

UNITED STATES OF AMERICA

FOOD AND DRUG ADMINISTRATION

+++

CENTER FOR DEVICES & RADIOLOGICAL HEALTH

+++

DIGITAL HEALTH ADVISORY COMMITTEE

DAY 1 CONFERENCE

+++

TOTAL PRODUCT LIFECYCLE CONSIDERATIONS FOR GENERATIVE AI-ENABLED

DEVICES

+++

November 20, 2024

9:00 a.m. EST

On-Site Conference

Transcript Produced By:



ACSI Translations

1025 Connecticut Avenue, NW, Suite 1000, Washington, DC 20036

<https://acsitranslations.com/>

**Participants**

Chairperson	Ami Bhatt, M.D.	Chief Innovation Officer (CIO) at the American College of Cardiology	Newton MA
Industry Representative	Diana Miller	Senior Director, Data Science Medtronic	Northridge, CA
Consumer Representative	Melissa Clarkson, Ph.D.	Assistant Professor of Biomedical Informatics, University of Kentucky College of Medicine	Lexington, KY
Voting Member	Jessica Jackson, Ph.D.	Founder and CEO Therapy Is For Everyone Psychological & Consultation Services, PLLC	Houston, TX
Voting Member	Thomas Maddox, M.D., MSc.	Vice President, Digital Products and Innovation, at BJC HealthCare, and Professor of Medicine (Cardiology), Washington University School of Medicine	St. Louis, MO
Voting Member	Chevon Rariy, M.D.	Chief Health Officer, and SVP Digital Health, Oncology Care Partners	Mequon, WI
Voting Member	Laura Stanley, Ph.D., CPE	Associate Professor, Gianforte School of Computing, and Director of Human Interaction Lab Montana State University, Bozeman, MT, and Clinical Associate Professor, Clemson University School of Health Research	Clemson, SC
Temporary Non-Voting Member	Taxiarchis Botsis, MSc., M.P.S., Ph.D.	Associate Professor of Oncology and Medicine The Sidney Kimmel Comprehensive Cancer Center Johns Hopkins University School of Medicine	Baltimore, MD
Temporary Non-Voting Member	Peter Elkin, M.D., MACP, FACMI, FNYAM, FAMIA, FIAHSI	UB Distinguished Professor and Chair Department of Biomedical Informatics Professor of Pathology and Anatomical Sciences University at Buffalo State University of New York	Buffalo, NY
Temporary Non-Voting Member	Rita Kukafka, DrPH, MA, FACMI	Professor of Biomedical Informatics and Sociomedical Sciences Department of Biomedical Informatics	New York, NY

		Chief Diversity Officer Columbia University	
Temporary Non-Voting Member	Steven Posnack, MIS, MS	Deputy National Coordinator for Health IT in the Office of the National Coordinator for Health Information Technology, US DHHS	Washington, DC
Temporary Non-Voting Member	Pratik Shah, Ph.D.	Assistant Professor, Pathology University of California	Irvine, CA
Temporary Non-Voting Member	Apurv Soni, M.D., Ph.D.	Assistant Professor of Medicine and Director of Program in Digital Medicine Health System Science Center for Digital Health Solutions, Department of Medicine UMass Chan Medical School	Worcester, MA
Temporary Non-Voting Member	Jagdish Khubchandani, MBBS, PhD, MPH	Professor of Public Health at New Mexico State University	Las Cruces, NM
Temporary Non-Voting Member	Thomas Radman, Ph.D.	Program Director Digital & Mobile Technologies Section, Division of Clinical Innovation National Center for Advancing Translational Sciences, National Institutes of Health, DHHS	Bethesda, MD
US Food and Drug Administration	Troy Tazbaz	Director of the Digital Health Center of Excellence, US Food and Drug Administration	Silver Spring, MD
US Food and Drug Administration	Sonja Fulmar, Ph.D.	Deputy Director, Digital Health Center of Excellence, US Food and Drug Administration	Silver Spring, MD
US Food and Drug Administration	Matthew Diamond, MD, Ph.D.	Chief Medical Officer, Digital Health Center of Excellence, US Food and Drug Administration	Silver Spring, MD
US Food and Drug Administration	James P. Swink	Designated Federal Officer US Food and Drug Administration	Silver Spring, MD

## Table of Contents

Digital Health Advisory Committee (DHAC) Call to Order.....	5
Conflict of Interest Statement.....	8
Welcome from FDA Commissioner.....	11
Opening Remarks .....	15
FDA Perspective - GenAI in Medical Devices .....	18
Sub-topic: Premarket Performance Evaluation .....	20
FDA Perspective – Regulatory Science Challenges for the Evaluation of Generative AI Applications in Medical Devices.....	20
Stakeholder Perspective – The Usage of Generative AI in Digital Pathology and Potential Challenges for Evaluation .....	23
Stakeholder Perspective – The Considerations for Multimodal Foundational Models & Generative AI Frameworks in Healthcare .....	37
Stakeholder Perspective – Measuring Performance of Generative AI – Methods and Lessons Learned .....	47
Dr. Bhatt’s Summary.....	57
Open Committee Discussion Q&A ( <i>Clarification questions</i> ).....	60
Open Public Hearing .....	81
Open Committee Discussion Q&A ( <i>Clarification questions</i> ).....	89
Committee Discussion of the FDA’s Questions ( <i>Deliberation and response to FDA</i> ) .	97
Sub-Topic: Risk Management.....	138
Stakeholder Perspective – Strategies and Controls to Mitigate Risks Associated with Gen AI Applications in Healthcare .....	138
Stakeholder Perspective - Narrow VS Generative AI: Risk Determination > Controls => Safe Innovation.....	149
Stakeholder Perspective – Safety from the Systems to Patient Levels: Risk Management for Large Language Models in Healthcare .....	160
Stakeholder Perspective - Risk Management for Generative AI-Enabled Medical Devices .....	169
Open Committee Discussion Q&A ( <i>Clarification questions</i> ).....	174
Committee Discussion of the FDA’s Questions ( <i>Deliberation and response to FDA</i> )	181
Day 1 Closing Remarks.....	199
Adjournment.....	200

1 Digital Health Advisory Committee (DHAC) Call to Order

2 Dr. Bhatt: Alright, I would like to call this meeting of the FDA's Digital Health  
3 Advisory Committee to order on November 20, 2024. I'm Dr. Ami Bhatt, I'm  
4 Chairperson of this Committee. It's an honor to be here. I'm a Medicine and Pediatrics  
5 Trained Cardiologist, practiced at the Mass General Hospital starting in 2009, ran  
6 outpatient cardiology, telecardiology, and most recently, three years ago, moved to be  
7 the Chief Innovation Officer at the American College of Cardiology, our global  
8 nonprofit. Again, it's an honor to be here. Important things about today. This is the first  
9 Advisory Committee to address digital health technologies and I think it's important to  
10 recognize that this Advisory Committee is the most formal and public way that the FDA  
11 can receive advice from the public on scientific matters.

12 Our Digital Health Advisory Committee has members with expertise in various  
13 fields of digital health, including machine learning, digital therapeutics, remote patient  
14 monitoring, software health data, user experience, real-world evidence, patient-  
15 generated health data, interoperability, personalized medicine and genetics,  
16 cybersecurity, implementation, and patient experience, along with a lot more. Just  
17 looking at the room, including our guests today, it's a really diverse group, and I want to  
18 congratulate the FDA on choosing such a diverse group to come together, because we  
19 recognize that we need to hear different opinions to make really good decisions.

20 I'll note for the record that the non-voting members constitute a quorum as  
21 required by 21 C.F.R., Part 14. I would also like to add that the Committee members  
22 participating in today's meeting have received training in FDA device law and  
23 regulations, and for today's agenda, the Committee will discuss and provide  
24 recommendations on how the use of generative AI may impact safety and effectiveness  
25 of medical devices enabled with this technology. The Committee will discuss premarket

1 performance evaluation, risk management, and postmarket performance monitoring for  
2 generative AI-enabled devices.

3 Before we begin, I'd like to set a few ground rules. We have a really busy  
4 agenda today, so we'll need to move through it efficiently and with purpose. I ask all  
5 presenters to stick to their allotted time. You will see a green, yellow, red button.  
6 Everybody else will see that too, and so that will help us make sure we stay on schedule.  
7 If any panelist would like to ask a question, please raise your hand or you can put up  
8 your card and I will address your question as we move through each session. To ensure  
9 smooth communication, please avoid speaking over one another, as this meeting is  
10 being transcribed for the official record. Finally, please remember to mute your  
11 microphone when you're finished speaking.

12 I would like to ask our distinguished Committee members and FDA experts to  
13 introduce themselves. I'll start by talking to Dr. Diamond here, and Dr. Diamond,  
14 maybe after you introduce yourself, we will just go ahead and go down the line and that  
15 might save us some time.

16 Dr. Diamond: Perfect. I'm Dr. Matthew Diamond. I serve as the Chief Medical Officer  
17 here at FDA's Digital Health Center of Excellence. Thanks.

18 Dr. Fulmer: Hi, I'm Sonja Fulmer. I'm the Deputy Director for the Digital Health  
19 Center of Excellence.

20 Mr. Tazbaz: Good morning, everyone. I'm Troy Tazbaz, the Director for Digital  
21 Health Center of Excellence at the FDA.

22 Dr. Radman: Good morning. Thomas Radman. I'm at the National Institute of Health  
23 and National Center for Advancing Translational Sciences.

24 Dr. Khubchandani: Good morning, everyone. Jagdish Khubchandani. Professor of  
25 Public Health, New Mexico State University.

- 1 Dr. Soni: Hello, everyone. Apurv Soni. Professor and Director of Program in  
2 Digital Medicine from UMass Chan Medical School.
- 3 Dr. Rariy: Hello, good morning. Chevon Rariy. I'm a Physician and Technology  
4 Executive with expertise in leading multidisciplinary teams really at the intersection of  
5 care delivery, data analytics, and digital health. I am excited to be here today and look  
6 forward to the discussion. Thank you.
- 7 Dr. Maddox: Good morning. Tom Maddox. I'm a Cardiologist and Professor at  
8 Washington University School of Medicine in St. Louis and I lead the Healthcare  
9 Innovation Lab at the school and its partner health system BJC HealthCare.
- 10 Mr. Swink: James Swink. I'm the Designated Federal Officer and team lead for the  
11 CDRH's Advisory Committee.
- 12 Dr. Jackson: Good morning. Jessica Jackson. I'm a Licensed Psychologist and I work  
13 with startups leveraging GenAI, digital therapeutics, and virtual reality for health  
14 innovation.
- 15 Dr. Shah: Morning. I'm Pratik Shah. I'm a Faculty Member and Professor at the  
16 University of California and my expertise is generative AI technologies and their  
17 deployment. Thank you.
- 18 Mr. Posnack: Good morning. Steve Posnack from the HHS's Assistant Secretary for  
19 Technology Policy. I serve as Principal Deputy.
- 20 Dr. Kukafka: Good morning. Rita Kukafka. I'm a Professor of Biomedical Informatics  
21 at Columbia University.
- 22 Dr. Elkin: Hello, I'm Peter Elkin. I'm a Distinguished Professor and Chair of the  
23 Department of Biomedical Informatics, the University of Buffalo, and also involved in  
24 doing research in AI in general, and I trained at MGH, so there you go. Thank you very  
25 much.

1 Dr. Botsis: Taxiarchis Botsis, Associate Professor of Oncology Medicine at Johns  
2 Hopkins University School of Medicine, also Director of the Biomedical Informatics  
3 Score at the Cancer Center.

4 Dr. Stanley: Hi, everyone. I'm Laura Stanley. I'm a Professor in Computer Science at  
5 Montana State University, Bozeman, and happy to be here.

6 Dr. Clarkson: I'm Melissa Clarkson. I'm the Consumer Representative. I'm also an  
7 Assistant Professor at the University of Kentucky in Biomedical Informatics.

8 Ms. Miller: Good morning. Diana Miller. Senior Director of Data Science at  
9 Medtronic. I'm here as an Industry Representative. Happy to be here and meet  
10 everybody.

11 Conflict of Interest Statement

12 Dr. Bhatt: Thank you all very much. Now James Swink, the Designated Federal  
13 Officer for the Digital Health Advisory Committee, will make some introductory  
14 remarks.

15 Mr. Swink: Thank you. Good morning. I will now read the FDA Conflict of Interest  
16 Disclosure Statement. So, The Food and Drug Administration is convening today's  
17 meeting of the Digital Health Advisory Committee under the authority of the Federal  
18 Advisory Committee Act of 1972. With the exception of the industry representative, all  
19 members and consultants of the Committee are special Government employees or  
20 regular federal employees from other agencies and are subject to federal conflict of  
21 interest laws and regulations. The following information of the status of this  
22 Committee's compliance with federal ethics and conflict of interest laws are covered by,  
23 but not limited to, those found at 18 U.S.C., Section 208, and are being provided to  
24 participants in today's meeting and to the public.



1 FDA has determined that members and consultants of this Committee are in  
2 compliance with federal ethics and conflict-of-interest laws. Under 18 U.S.C., Section  
3 208, Congress has authorized FDA to grant waivers to special Government employees  
4 and regular federal employees who have financial conflicts when it is determined that  
5 the Agency's need for a particular individual's services outweighs his or her potential  
6 financial conflict of interest.

7 Related to the discussions of today's meeting, members and consultants of this  
8 Committee who are special Government employees or regular federal employees have  
9 been screened for potential financial conflicts of interest of their own, as well as those  
10 imputed to them, including those of their spouses or minor children and, for purposes of  
11 18 U.S.C., Section 208, their employers. These interests may include investments;  
12 consulting; expert witness testimony; contracts, grants, CRADAs; teaching, speaking,  
13 writing; patents and royalties; and primary employment.

14 For today's agenda, the Committee will discuss how the use of generative  
15 artificial intelligence may impact safety and effectiveness of medical devices enabled  
16 with this technology. The Committee will discuss premarket performance evaluation,  
17 risk management, and postmarket performance monitoring for generative AI-enabled  
18 devices.

19 Based on the agenda for today's meeting and all financial interests reported by  
20 the Committee, the members and consultants have-- No conflict-of-interest waivers  
21 have been issued in accordance with 18 U.S.C., Section 208.

22 Ms. Diana Miller is serving as an Industry Representative for generative  
23 AI/large language models, and is acting on behalf of all related industry. Ms. Miller is  
24 employed by Medtronic plc. For the record, the Agency notes that Dr. Danielle  
25 Bitterman, who is an invited guest speaker with us today, has acknowledged possible

1 interests through her employer, Brigham and Women's Physician Organization, and in  
2 the form of federal research grants, and as a scientific advisor with an unaffected firm.  
3 Dr. Bitterman also acknowledged that she is an Associate Editor of the Radiation  
4 Oncology for HemOnc.org and she receives no financial compensation for that activity.  
5 Dr. Faisal Mahmood, who is also an invited guest speaker with us today, has  
6 acknowledged a possible interest through his employer, Massachusetts General  
7 Brigham, and has acknowledged interests with an unaffected firm in the form of a stock  
8 and as a scientific advisor. Dr. Pranav Rajpurkar, another invited guest speaker with us,  
9 has acknowledged interests with an unaffected firm as a Co-Founder of the firm. Drs.  
10 Keith Dreyer and Gabriella Waters are other invited guest speakers with us today and  
11 have reported no interests in relation to today's meeting.

12 We would like to remind members and consultants that if the discussions  
13 involve any other products or firms not already on the agenda for which the FDA  
14 participant has a personal or imputed financial interest, the participants need to exclude  
15 themselves from such involvement and their exclusion will be noted for the record.  
16 FDA encourages all other participants to advise the Panel of any financial relationships  
17 that they may have with any firms at issue. A copy of this statement will be available  
18 for review at the registration table.

19 Before I begin, I would like to make a few general announcements. We have  
20 transcripts of today's meeting; will be available from ACSI. Their phone number is  
21 202-599-8456. Information on purchasing videos of today's meeting and handouts for  
22 today's presentations are available at the registration table outside of the meeting room.  
23 The FDA press contact for today's meeting is James McKinney. I think he'll be-- Here  
24 he is, over here on this side. All written comments received were provided to the Panel  
25 and the FDA review team for their review prior to today's meeting. There is an active

1 docket where members of the public can post written comments. The link can be found  
2 on the FDA website and the registration table, and that will be open for another-- I think  
3 until January 30, so you have time to submit more comments. If you are presenting in  
4 the Open Public Hearing Session and have not previously provided an electronic copy  
5 of your slide presentations, please arrange to do so with Mr. Artair Mallett at the  
6 registration table. And in order to keep the transcriptionist-- Please identify yourself  
7 before you speak each and every time. Please, silence your cell phones and other  
8 electronic devices. Also, we made handouts for both days. If you took handouts, please  
9 bring them tomorrow as we have limited copies. Thank you.

10 Welcome from FDA Commissioner

11 Dr. Bhatt: Thanks so much, James, and thanks for all the work you did to get us all  
12 here today. So, we will proceed with welcoming remarks. We are honored to have the  
13 FDA's Commissioner, Dr. Robert Califf here with us. Dr. Califf.

14 Dr. Califf: Hey, it works. Alright. Good morning, everybody, and welcome to this  
15 inaugural Digital Health Advisory Committee meeting. I am pleased to speak with you  
16 today and I want to thank all the experts who are serving as members of the Advisory  
17 Committee; our CDRH team for their thoughtful planning of the meeting; and our  
18 speakers for sharing their perspectives. I also want to thank the many product  
19 developers, representatives of industry, academia, and consumer organizations, as well  
20 as the patient advocates and others who are participating today. Advisory Committees  
21 are a unique mechanism that help support the FDA's mission. They bring together  
22 leading experts who can advise the FDA about a particular field, in this case digital  
23 health technologies, giving us access to the latest thinking on broad issues and how to  
24 navigate them. These forums also encourage the public, including patients, healthcare  
25 providers, industry, and other interested people to share their views.

1           The Digital Health Advisory Committee will provide advice and  
2 recommendations on understanding the benefits, risks, and clinical outcomes associated  
3 with digital health technologies. We've established this Committee because we see great  
4 potential for digital health technologies to help address critical healthcare issues that we  
5 face today, and we need these technologies to be developed, deployed, and used  
6 responsibly in the best interest of patients and consumers. I've talked before about the  
7 current catastrophic decline in life expectancy in the United States, leaving us years  
8 behind our economic peers in life expectancy and ability of our people to function. We  
9 are even falling behind many middle-income countries like Thailand and Costa Rica,  
10 and just as worrisome and a critical issue for AI is that the negative trajectory is largely  
11 driven by disparities that are a function of race, ethnicity, education, and wealth, as well  
12 as where someone lives, particularly rural and underserved urban community status.

13           Artificial intelligence is changing how we think about health and healthcare, and  
14 it's one of the most exciting and promising areas of science because it's built to  
15 transcend boundaries. While we often describe the promise of AI in terms of helping us  
16 facilitate speedier delivery of new treatments, I believe the larger role for AI is to  
17 improve the efficient and coordinated delivery of care to patients in all facets of  
18 healthcare settings, including the operating room, the clinic, and even the home. Digital  
19 tools and AI have the potential to help improve quality of life and life expectancy,  
20 including prevention of illness in healthy people for all people living in the United  
21 States. To this end, I hope you'll consider CDRH's "Home as a Healthcare Hub" project  
22 during this meeting. The FDA regulates 20% of the U.S. economy and must  
23 continuously evolve to oversee the safe and effective use of AI across regulated  
24 industry, ensuring compliance while fostering innovation.

1 For us to be most effective and do our jobs as protectors of public health, it's  
2 essential that we embrace these groundbreaking technologies, not only to keep pace  
3 with the industries we regulate, but also to use regulatory channels and oversight to  
4 improve the chance that they will be applied effectively, consistently, and fairly. I can't  
5 get Ed Yong's phrase from his article in the Atlantic out of my mind. Technological  
6 advances drift into society's penthouses. Pandemics—we can substitute diseases for  
7 pandemics—seep into its cracks. Unless we have intentionality to reach technologically  
8 disadvantaged people, we'll continue to even-- Or even accelerate the current trajectory  
9 to further widen a huge gap in health and life expectancy between the wealthy and  
10 highly educated people, mostly living in urban areas, and people with less wealth in  
11 education increasingly dominated by the rural population. The FDA has been working  
12 for years to anticipate and prepare for challenges of AI and to harness the potential to  
13 enable major advances in the development of more effective and risky-- And less risky  
14 medical products.

15 The FDA's first approval of an AI-enabled medical device took place almost 30  
16 years ago in 1995. Since then, we've received approximately a thousand submissions  
17 for AI-enabled medical devices and more than 300 submissions for drugs and biological  
18 products with AI components. But I have to say, I'm also told on the drug and biologic  
19 side almost every application now has AI somewhere in the development process. These  
20 submissions have included drug discovery and repurposing, enhancing clinical trial  
21 design elements, dose optimization, endpoint and biomarker assessment, and  
22 postmarket surveillance. They also cover a growing diversity of medical devices that  
23 leverage AI to improve clinical workflows in patient experiences or outcomes in  
24 addition to sophisticated prediction algorithms. For example, the Agency recently  
25 reviewed and granted a new AI algorithm that supports earlier detection and

1 management of hypertrophic cardiomyopathy, something dear to my heart as a  
2 cardiologist. This granting created a new regulation for cardiovascular machine  
3 learning-based notification software and will help provide quicker identification of  
4 possible suspected HCM cases allowing for earlier patient treatment, hopefully before  
5 the first episode of sudden death.

6 We recently published a paper on how FDA's CBER, CDER, CDRH, and  
7 Office of Clinical Policy are working together on AI and medical products. The paper  
8 outlines and reaffirms our commitment to promoting the responsible and ethical  
9 development, deployment, use, and maintenance of safe and effective medical products  
10 that incorporate or are developed with AI.

11 Two final notes for your consideration. First, as I talk with clinicians across the  
12 country, I'm hearing increasing concern that the criteria for adopting an AI are almost  
13 purely financial. I hope this Committee will help us figure out how to balance patient  
14 outcome benefit with financial return. Financial return is relatively easy to measure  
15 these days, and AI is a supercharged financial tool. Clinical outcome measurement is  
16 hard. Unless you take this issue very seriously and form alliances of those concerned  
17 about improving health outcomes, this technology will improve profits at the cost of  
18 continued deterioration in our overall health status. At FDA, our statutes prevent us  
19 from considering cost in individual product decisions, so we're all in on this at the FDA,  
20 but we need your help on this issue because the tide is going the other direction.

21 Second, I'm convinced that we need a very different information ecosystem to  
22 ensure that local recurrent validation of AIs occurs as a continuous process in the midst  
23 of our ongoing healthcare system. This is the ultimate in real-world data and evidence,  
24 but I challenge you to envision how such a system would work and to focus on the  
25 specific examples you're discussing today. I have looked far and wide; I do not believe

1 there's a single health system in the United States that's capable of validating an AI  
2 algorithm that's put into place in a clinical care system. Beautiful operating  
3 characteristics of an AI in the premarket phase in no way should assure us of the  
4 performance in real life. Imagine that we developed apps for driving in a hypothetical  
5 space and didn't incorporate the evolving experience of people who are actually driving.

6 So, this is well beyond FDA and it should be. We don't have inspectors, for  
7 example, on every farm every day. But while not perfect, our farms adhere to the 10  
8 FSMA rules producing safe food for the country. We need something similar for AI and  
9 healthcare. You probably don't think of yourselves as farmers, but I believe that you  
10 really should in this regard.

11 I really do wish, and you could probably tell this, that I could stay for the two  
12 days to hear the robust conversations in this meeting. This Advisory Committee is a  
13 valuable resource and a forum for contributing to how we think about this technology to  
14 ensure that devices are safe and effective. We're seeking your input to help inform the  
15 FDA's thinking as we work to ensure that people living in the United States benefit, not  
16 just financially, but also with regard to their health outcomes, from AI's promise to  
17 improve healthcare.

18 Thanks for being here today. I hate to talk and run, but I will tune in. This is  
19 going to be on YouTube, right? Alright. Well, I take this in while I'm trying to sleep at  
20 night. Thank you.

#### 21 Opening Remarks

22 Dr. Bhatt: Dr. Califf, thank you so much for your remarks. Your insights are  
23 perfect, as we really think about this critically important topic. We'll proceed with  
24 opening remarks from Dr. Michelle Tarver, Director of CDRH at the FDA. Dr. Tarver,  
25 thank you for being here with us this morning.

1 Dr. Tarver: Well, thank you for having me. Good morning, everyone. I want to thank  
2 Dr. Califf for his remarks and I want to welcome you all; those of you on the Panel,  
3 those of you in the room, those of you tuning in online. Today's an exciting day. It's our  
4 first Digital Health Advisory Committee meeting. This is a milestone and it really does  
5 come in the context of how we're seeing healthcare, medical devices, and patients'  
6 expectations evolve over time. You've probably heard us say at the FDA "Patients are  
7 at the heart of what we do," and we really mean it. The Center for Devices and  
8 Radiological Health's vision statement begins with a commitment to the American  
9 people that patients in the U.S. have access to high quality, safe and effective medical  
10 devices of public health importance first in the world. And we don't say first in the  
11 world out of a spirit of competitiveness, but truly out of a pledge to timely access.

12 When we're talking about public health, every patient counts and every day  
13 matters. And it's by working collaboratively to foster innovation and to facilitate timely  
14 access that we can have true impacts on this nation's health. Now, you heard Dr. Califf  
15 talk a lot about the flood of medical devices that we have coming into our doors, and  
16 many of them are AI-enabled, and AI holds the possibility of transforming how we  
17 prevent, diagnose, and treat health conditions. And what's so exciting about it is that it  
18 continues to change and evolve each day, offering new opportunities for us to deploy it  
19 in the care of patients, but only if it's done responsibly and we can rely on it to be safe  
20 and effective. And it's going to be important not only in a traditional clinical setting, but  
21 also outside of the clinical setting, like the home and communities, closer, where they  
22 don't have clinics; where they don't have hospitals.

23 In fact, if you look at the CDC reports, currently there are 129 million people in  
24 this country that are living with chronic diseases, and 40% of them have two or more  
25 chronic diseases. That accounts for over 90% of the 4.1 trillion dollars we spend on



1 healthcare expenses. Huge population; huge cost. And that's in the setting where we  
2 have consolidation of health services, where we see contraction of providers and  
3 facilities. How are people going to get the care that they need? Well, AI-enabled  
4 medical devices offers a promise to extend into those communities and offer options for  
5 them to get care. When we're thinking about the burden of disease and who is most  
6 bearing that disease; that's small-town America. Those are rural communities. There are  
7 people who are of older age, racial and ethnic minorities, those with fewer resources and  
8 who live the furthest from healthcare facilities. So, technologies have to come meet  
9 them where they are. You heard Dr. Califf allude to our initiative on "Home as a  
10 Healthcare Hub", and it is really grounded on that principle of how can we democratize  
11 care in a way that everyone has access to it; that everyone can benefit from the  
12 healthcare-- Good healthcare outcomes that all of us would like to experience.

13 One of the things that I'm really excited about today's conversation is this is  
14 really a collective conversation on how do we move things forward in an equitable and  
15 ethical way. And that requires including everyone at the table. That means patients,  
16 providers, industry, regulators, developers to have a collective conversation. Well, today  
17 we're going to talk about generative AI and the total product lifecycle of these AI-  
18 enabled technologies, particularly looking at the benefits, the risks, the outcomes, and  
19 what are the considerations as we look at them both as regulators and ultimately how  
20 you all will be looking at them as care providers and patients.

21 One of the things that I want to definitely highlight is how important it is for us  
22 to be having the conversation now and not having it five years from now, when we're  
23 trying to figure out how to manage it in retrospect. So, being proactive, being deliberate,  
24 and being intentional in this conversation really will help us make informed decisions.

1 I'm really looking forward to the next two days. I'm really looking forward to  
2 the conversation, and I'm not going to keep you listening to me for very long, so I'm  
3 going to turn it over to Troy Tazbaz, who is the Director of the Digital Health Center of  
4 Excellence to give some remarks. So, thank you all very much for being here. Thank  
5 you all for your insights today and I'm looking forward to the conversation.

6 FDA Perspective - GenAI in Medical Devices

7 Dr. Bhatt: Thank you, Dr. Tarver. As Mr. Tazbaz starts, just for the Panel to know,  
8 we are now going to have the first five presentations of our agenda and we will then be  
9 able to ask questions after we have the next five Panel presentations. Mr. Tazbaz, thank  
10 you.

11 Mr. Tazbaz: Thank you. Good morning, everyone. I first want to extend my heartfelt  
12 gratitude for all of you, for being here and dedicating your time, your knowledge and  
13 expertise to these very, very critical discussions. I also want to thank the FDA team for  
14 working tirelessly to making this event happen. It was difficult, but it is going to be very  
15 well worth it, and we're very excited about the discussions that are going to be taking  
16 place today.

17 The Digital Health Advisory Committee has attracted significant interest from  
18 the entire healthcare community. We believe this is a reflection of the growing role  
19 digital health technologies will play in advancing patient care and overall health of our  
20 entire population. As the FDA's first Advisory Committee focused on crosscutting  
21 digital health issues, we hope DHAC, as we've been calling it, will guide us on matters  
22 that affect the wide spectrum of medical device products.

23 Today we are here to address one of the most transformative developments in  
24 healthcare: generative AI, or as we've been referring to as GenAI. We hope to explore  
25 and identify unique regulatory considerations necessary to ensure that safety and

1 effectiveness of these rapidly evolving technologies are maintained throughout their  
2 lifecycle. Generative AI brings incredible opportunities to enhance patient care but also  
3 introduces new regulatory challenges and oversight-- And operational oversight. As it's  
4 been highly covered, generative AI has the potential to reshape healthcare by generating  
5 contextually relevant outputs from vast amounts of data. However, unlike traditional AI  
6 models that predict or classify data, GenAI produces new outputs which can introduce  
7 additional layers of regulatory complexity for us. This is especially relevant when these  
8 outputs could directly impact the safety and effectiveness of medical devices.

9       To guide us throughout the discussions, we would like to take a total product  
10 lifecycle approach, as Michelle mentioned. The lifecycle perspective allows us to  
11 consider safety and performance from initial design and development all the way  
12 through the postmarket monitoring. Throughout the next two days, we will dive into  
13 critical questions surrounding generative AI, including premarket evaluation,  
14 performance metrics, and effective postmarket strategies. Additionally, we will discuss  
15 the complexities of monitoring generative AI models across multiple sites, as we need  
16 to ensure that they perform consistently while managing potential regional biases and  
17 data variations.

18       Our agenda is designed to give structure to these important discussions,  
19 beginning with an examination of information needed to evaluate generative AI-enabled  
20 devices, followed by in-depth sessions on risk management, and postmarket  
21 performance monitoring strategies, which is a very important part of the equation of  
22 when these things are deployed into real healthcare settings.

23       As we move forward, I'm confident that the insights and ideas shared here will  
24 shape the FDA's approach to balancing innovation with safety and ensuring that the  
25 patients and healthcare professionals can trust these powerful new tools.

1 I would like to say thank you once again for your participation and engagement.  
2 I look forward to the rich discussions ahead and the actionable steps we will identify  
3 here together. In a moment, you will also hear from Dr. Aldo Badano on some  
4 regulatory science challenges that FDA faces for evaluating generative AI applications  
5 and medical devices, and tomorrow FDA's Jessica Paulsen will provide an overview of  
6 approaches for managing changes for AI-enabled devices. Thank you, everyone.

7 *Sub-topic: Premarket Performance Evaluation*

8 FDA Perspective – Regulatory Science Challenges for the Evaluation of Generative AI  
9 Applications in Medical Devices

10 Dr. Bhatt: Thank you, Mr. Tazbaz. We will now proceed with Dr. Aldo Badano's  
11 presentation. Dr. Badano, as mentioned, is Director of the Division of Imaging,  
12 Diagnostics and Software Reliability in the Offices of Science and Engineering  
13 Laboratories at CDRH-FDA. Dr. Badano.

14 Dr. Badano: Good morning. Thank you for the introduction and it's an honor to speak  
15 to the distinguished members of the Panel in the audience.

16 Let me start by saying a few words about our office. OSEL in CDRH, Office of  
17 Science and Engineering Labs, is an ISO 9,001-certified research organization dedicated  
18 to promoting innovation for the development of new life science medical devices. Also,  
19 it is organized into about 20 program areas and AI/ML is one of the largest programs  
20 that we have. Also, outputs are regulatory science tools, which are innovative open-  
21 source tools for assessing the safety and effectiveness of emerging technology.  
22 Innovators can use these tools voluntarily into all stages of device development.  
23 Regulatory science tools have been cited in over 1,000 applications-- Premarket  
24 applications to the Agency, over a hundred product codes. If you would like to learn

1 more, there's some information about how to reach our regulatory science tools catalog  
2 with a website and the QR code.

3 GenAI-enabled devices in healthcare have a wide range of uses. Here's a table  
4 taken from the literature; an article by Meskó and Topol, which outlines some of the  
5 applications and how the community is thinking about GenAI in the medical space.  
6 Some of these products are clearly nonmedical devices, but there's a lot of excitement  
7 in terms of the potential performance and productivity gains. However, the same  
8 characteristics that make these devices and products have so much potential are the  
9 characteristics that might pose challenges for certain uses.

10 So, what are the regulatory science challenges that this technology brings about?  
11 Here's a preliminary list. I'm not going to go into detail in all of these; we don't have  
12 the time. But many of these challenges might exist in other technologies, but are  
13 exacerbated in AI applications, particularly in GenAI. So, difficulty in defining scope of  
14 the products intended use; some of the parts of the model like foundation models not  
15 being in the provenance of the manufacturers of the device; oversight of adaptive  
16 systems; and hallucinations, which are a new kind of error that has been introduced in  
17 recent years by deep learning and other technologies that bring about the new trade off  
18 between stability and accuracy due to concerns about reliability and trustworthiness.

19 There's also issues about data needs, including diversity; monitoring  
20 performance in the real world, including bias; transparency to users. So, let's narrow it  
21 down a little bit with one example use case. In this use case, a GenAI-enabled radiology  
22 device can produce a report when prompted with a medical image, as you see in the left,  
23 in the input side. The reference for this device might be a clinician-generated report. So,  
24 what strategies do we have at hand to evaluate these new technologies?

1 In recent years or months, I would say, three distinct classes of evaluation have  
2 emerged from the scientific literature. Benchmarking with the standardized reference  
3 datasets; expert evaluation, including a holistic evaluation of the AI-human interaction;  
4 and model-based evaluations, which include automated testing.

5 In the next slides, I'll describe a little bit the pros and cons of each one of these  
6 approaches. Benchmarking. Benchmarking is the evaluation of models on external  
7 testing datasets with predetermined evaluation metrics. This approach is practical,  
8 allows for head-to-head comparisons of models on same data and with same metrics,  
9 and can assess models at a large scale. However, benchmarks are limited in models and  
10 in tasks or applications, and there is this risk of that the models can train to the test,  
11 which is an overfitting technique for optimizing performance on a given dataset at the  
12 expense of performance for other data distributions, for instance, like the ones you see  
13 in a different clinical setting. On the right, you can see one of the many leaderboards  
14 available in open-source format.

15 Expert evaluation. Expert evaluation refers to using experts to establish the  
16 reference standard for a given task. For example, experts may determine if an error from  
17 GenAI will cause significant patient harm. Note that in the chart, experts can identify  
18 errors that might pass tests of a benchmark concern only with tech similarity metrics.  
19 For instance, right versus left, which is location of disease; fusion versus effusion.  
20 Expert evaluation is adaptable to a wide range of tasks with direct clinical relevancy.  
21 However, expert evaluation is resource-intensive and burdensome, and must consider  
22 the subjective nature of the readers and the highly variable effects.

23 And finally, model-based evaluation. This is about evaluating models with other  
24 models using a model-based approach that must have human oversight. Let me bring up  
25 the chart. The advantages are clear. This augments the human evaluation and is scalable.

1 However, the disadvantages include the fact that the validation needs to be robust and  
2 there is a potential risk of inter-model leakage, which means models in the device and in  
3 the evaluator can be correlated, therefore generating risks of the results. In our CDRH  
4 labs, we're actively investigating a model-based evaluator for factual accuracy of a  
5 GenAI-generated radiologic impression. A series of binary questions are used to  
6 determine the rates of different types of errors in the generated report. We're also  
7 working on performance metrics and associated statistical analysis plans. I hope that  
8 you'll stay tuned for updates in the coming months.

9 Let me summarize. For some GenAI, known evaluation strategies still may  
10 apply. However, new evaluation methodologies and new performance metrics may need  
11 to be developed for some other types of GenAI. Overall device performance evaluation  
12 requirements are governed by the intended use and the associated risk for the device.  
13 And finally, particularly in the case of GenAI-enabled devices, evidence may include  
14 pre and postmarket elements as the technology is meant to update often based on the  
15 availability of models that are trained on more and more data and the emergence of  
16 novel architectures and GenAI innovations.

17 In summary, rigorous evaluation of GenAI-enabled devices is necessary as a  
18 rapidly emerging technology developing a least burdensome approach to ensure safety  
19 and effectiveness will require partnerships like the one we're seeing here in this room  
20 and commitments to advancing regulatory science. Thank you for your attention.

21 Stakeholder Perspective – The Usage of Generative AI in Digital Pathology and  
22 Potential Challenges for Evaluation

23 Dr. Bhatt: Thank you, Dr. Badano. Next, we're going to hear a stakeholder  
24 perspective from Dr. Faisal Mahmood. Dr. Mahmood is Associate Professor at Harvard  
25 University Brigham and Women's Hospital. And just a reminder to the Panel that when

1 we are done with these, we'll take a short break and then we can ask questions of the  
2 panelists; that includes Mr. Tazbaz and Dr. Badano, so to save your questions for them  
3 as well as our next three speakers.

4 Dr. Mahmood: Okay, thank you for having me today. So, I'll be talking a little bit  
5 about some of the work that my group has done over the years in computational  
6 pathology and how it has led into generative AI and what we are currently focused on.  
7 Just a quick outline for the talk. I'll be talking about some of our older work and how it  
8 feeds into the generative AI-related work that we have done more recently.

9 So, the problem that we are trying to target is to analyze pathology images.  
10 Pathology images are hierarchical, they're large. They're very large gigapixel images  
11 that result from the digitization of glass slides. And the goal of our field is really to go  
12 from these images to everything in the red box here. This involves early diagnosis,  
13 prognosis, prediction of response to treatment, outcome prediction, patient stratification,  
14 and some form of biomarker discovery.

15 So, I'll start by talking about some of the "historical work," even though it's  
16 quite recent. So, this is some work that we did back in 2021. This allowed us to use  
17 these large whole-slide images and slide-level labels that existed in pathology reports  
18 while still being data efficient, so without requiring any pixel level and annotations.  
19 And this was published back in 2021 and we made the code and the models, everything,  
20 publicly available for the community to use. And these tools have been extensively used  
21 and have basically been applied to every major disease organ system. And we applied it  
22 to a number of different applications, including for cancers of unknown primary, where  
23 we tried to find the origin of the tumor. It's a very common and difficult problem in  
24 pathology, and we posed this as being an 18-class supervised classification problem  
25 using host-led images and labels that were taken from a pathology report.



1 Overall, we found that this approach works quite well. Internal and external test  
2 sets; we did a lot of evaluation assessing the variability. And then more recently we  
3 have tried to combine histology images with molecular data to see if we can improve  
4 outcomes and have shown that this is also possible. We're also interested in integrating  
5 histology with molecular information to improve outcomes in general, and have shown  
6 in this study from back in 2022 that you can improve patient outcome stratification into  
7 distinct risk groups, but then also go back and identify biomarkers that would be  
8 relevant to the diagnosis but also would correlate with response and resistance to  
9 specific therapeutic agents.

10 This is another example for endomyocardial biopsy assessments. Our approach  
11 has been that we can train these supervised models, but then targeting problems that are  
12 difficult to solve in pathology. So hence they would sort of encourage machine learning  
13 adoption because these problems are difficult to begin with, like cancer of unknown  
14 primary or endomyocardial biopsy assessment where there is large scale intra and intra  
15 observer variability.

16 But the purpose of showing all of these examples is that they all follow this very  
17 conventional approach that's commonly used in computational pathology, where we  
18 have whole-slide images that are digitized; we patch them because they're very large;  
19 we extract relevant features, or in this case, they're exhaustively extracted from a pre-  
20 trained encoder, leading to aggregation and eventually a prediction. This pipeline has  
21 been innovated on in various components over the years, both by the clinical  
22 community as well as in terms of how the data is used and by the machine learning  
23 community and how we feature, extract, and aggregate.

24 But over the years, we have found that the feature extraction component is the  
25 most important component, and that's what leads us into foundation models or self-

1 supervised models. In the examples I showed earlier, we were using a model that was  
2 trained just on a ResNet that was trained on real-world images from the ImageNet  
3 database and is used for feature extraction. But more recently we've seen that self-  
4 supervised models perform significantly better. So, motivating foundation models for  
5 pathology. So, foundation models are generic models that are capable of rich feature  
6 extraction. The particular use that we have here is that by having richer feature  
7 extraction, we can apply this to situations where we just don't have enough data; very,  
8 very common in pathology. For example, for rare diseases, clinical trials, and they're  
9 ideal for multi-task, multi-tissue kind of problems, and they're not necessarily meant to  
10 replace end-to-end trained supervised models. So, this thinking was sort of what led to  
11 these two studies. Those from March this year, where we have two foundation models.  
12 The first one, UNI, is a self-supervised model that's trained in a self-supervised manner  
13 on lots and lots of pathology data to extract rich feature representations from those  
14 pathology images. And CONCH is a visual-language model which further enhances this  
15 by contrasting with text. And this is in line with how the machine learning community  
16 is thinking about this, because they've shown that contrasting with additional modalities  
17 leads to better improved feature representation for both images and for text.

18         So, we collected a large dataset of about a hundred thousand whole-slide  
19 images. This came from a number of different hospitals including our own. And it's  
20 diverse; it was maximized for diversity. And this is also in line with how the machine  
21 learning community has found that diversity of data is much more important than the  
22 quantity of data in building these large self-supervised models. And we used the  
23 DINOv2 framework from Meta to train these models, and it was assessed against a  
24 number of different benchmarks, both conventional and pathology-related, and it was  
25 applied to 33 different downstream tasks. The radar plot is not the best way to show

1 this, but it does capture all the different tasks that this was assessed on, including for  
2 whole-slide level classification as well as for region-level or small-scale classification.  
3 And in particular for few-shot classification, where we attempt to train these models  
4 with very, very few data points, hence assessing its efficacy on situations where you  
5 would not have enough data available, for example, for rare diseases, clinical trials, and  
6 so forth. We also wanted to see how far-- I'm sorry, the colors on this slide did not  
7 show up as they originally were. But we also wanted to see how far can we push this.  
8 Can we have a single model that can cater to a really large difficult classification  
9 problem of 108 different cancer types in the entire OncoTree classification system? And  
10 the goal here is not to reach clinical utility—it would still take some time; more data,  
11 potentially, to reach to get there,—but to really assess that how rich the feature  
12 representations get if we use this form of self-supervised learning.

13 And in parallel, we were also working on contrasting images with text. And in  
14 this case we use 1.1 million image-caption pairs that come from the PubMed open  
15 access database. We could not use all the data in PubMed, only that that was publicly  
16 funded and is within the open access database with an appropriate license. But  
17 contrasting images with text is a very common approach used by the computer vision  
18 community to enhance feature representations for both images and for text. This is a  
19 distribution of all the data that was used in this particular case. And we assess this in a  
20 number of different scenarios, including for Zero-Shot classification, where the goal is  
21 to really assess the feature representation and not really achieve clinical utility, but also  
22 for Few-Shot classification, we're showing that fewer data points can lead to better  
23 performance and a variety of different disease models, both rare and common.

24 These models were made publicly available and we essentially did this to  
25 comply with NIH guidelines and what the requirements were from the journals where

1 we were publishing this. But the models were extensively downloaded and have been  
2 used all over the world for a host of different applications or for every major organ  
3 system all the way up to forensics. So, things that we never originally thought of. We  
4 didn't even know that this many people were working in computational pathology.

5         And they have also been assessed independently by a number of different  
6 groups. This is an evaluation from a group in Dresden, where they have shown how  
7 these models work against a number of other approaches, self-supervised and otherwise,  
8 for a host of different tasks including diagnostic tasks, as well as predicting what's  
9 commonly referred to as non-human identifiable features, like predicting molecular  
10 markers and other alterations directly from histology images and just prognosticating  
11 directly from histology images.

12         Here's another assessment that was done at Mount Sinai and at the Memorial  
13 Sloan Kettering Cancer Center using our foundation models and how well it fares for  
14 some of these very difficult pan-cancer predictive tasks. Since these models came out,  
15 we have further expanded them to a whole-slide level. So, the models I showed so far  
16 work at original level extracting rich feature representations from regions. We have  
17 since expanded this to extracting a single feature vector that corresponds to the whole  
18 slide, and we've done that on an image level. And then it's further contrasted with  
19 generated captions that come from our generative model that I'll talk about in a second,  
20 but also from the pathology report. So, contrasting both with a generative caption and a  
21 pathology report. And we find that this form of feature extraction, where we have a  
22 single feature vector for the whole slide, essentially enables a lot of different things that  
23 we can do downstream for incredibly difficult tasks and for common clinical tasks. For  
24 common clinical tasks, we can get clinical-grade performance and for more difficult  
25 tasks, we're further trying to improve the model to get there.

1 But I think what most people get very excited by is that-- Its ability to perform  
2 on Few-Shot examples in situations where diagnosis is incredibly difficult. Average  
3 diagnosis for rare disease is between six to eight months, and it really is a challenging  
4 problem where we can use these rich feature representations to find similar cases, but  
5 also train with very, very few cases. And then of course, because it has a text head,  
6 we're able to generate these reports. We compare it with the pathologist-given reports,  
7 and we've done quite some analysis on how good the model is for report synthesis, both  
8 at the level of the slide for the patient, what happens when we include electronic  
9 medical records, and so forth.

10 So, the story so far is that we started with lots of supervised models. We had  
11 self-supervised models or foundation models that relied on images themselves, and then  
12 we contrasted those images with text to improve the image feature representation. So,  
13 there are a lot of other modalities that we have in pathology and in healthcare in general  
14 that can be used to contrast data to further improve feature representation for each one  
15 of these individual modalities. So, we contrast-- In this example—this is an article from  
16 earlier this year,—we contrast pathology images with immunohistochemistry, and by  
17 contrasting we further improve the feature representations with just a single model that  
18 can extract rich features. And then it's applied downstream to a variety of different  
19 tasks, including IHC quantifications with the first time showing that you don't need to  
20 have pixel-level annotations for this very difficult task of quantifying IHCs. And it  
21 could be just this model trained on lots of IHC data that can cater to all of  
22 immunohistochemistry scoring as well as downstream applying it to survival and so  
23 forth.

24 And if this is true that we can contrast with IHC, an obvious next step would be  
25 that can we also contrast with transcriptomics? So, in this smaller example that was

1 published earlier this year in CVPR, we showed that we can do this and it can  
2 substantially improve Few-Shot performance for detecting rare diseases. But more  
3 recently we have expanded this to a very, very large-scale study where we are using all  
4 the molecular data from the Brigham and MGH that's been collected over the years and  
5 correlated with histology images. And by contrasting with our NGS data, with  
6 transcriptomic data, and with histology data, we were further able to improve feature  
7 representations. And these are a variety of different downstream tasks that it was  
8 applied to, but just to draw your attention, it was able to improve performance for these  
9 incredibly difficult tasks for predicting treatment response, so often for which we have  
10 very, very few samples available because at a tertiary medical center, like ours, it's  
11 often difficult to get pretreatment biopsies, because a lot of patients get referred rather  
12 than the hospital being their first instance. But given the data is so small and the  
13 downstream treatment efficacy is still low for a variety of different therapeutics,  
14 improving performance by using image-based biomarkers is of a lot of interest to drug  
15 companies and other stakeholders. So, this is quite exciting for us to see that for the first  
16 time we can just boost performance by using these contrasted foundation models.

17 So, what do we do after we have all these foundation models? We can of course  
18 train lots and lots of supervised models by using these self-supervised rich features. But  
19 what was most interesting to us is that we can also use them in a generative setting  
20 where we can have a singular model that can cater to all of human pathology. And this  
21 was the biggest undertaking that my group had done. So, our entire group was involved  
22 in collecting the data for this, because the data required is not collected essentially;  
23 essentially clinically. So, our goal is to build a singular, multimodal large language  
24 model that can cater to all of human pathology. And why we believe this is possible is  
25 because there are increasing number of large language models and multimodal large

1 language models that cater to very, very large amount of information. Like open AIs  
2 trying to build a single model that could cater to the world's information, and human  
3 pathology knowledge is much, much more constrained and specific. So, if the first is  
4 true, we should be able to build a model that caters to just all of human pathology. And  
5 what we need to do is that we need rich feature representation that can come from  
6 foundation models or self-supervised models; we need a large instruction dataset of  
7 images, text, and corresponding responses; and we also need robust evaluation.

8         The challenge we face is that we have these hierarchical large pathology images,  
9 but the annotations we have corresponding these images are at the level of the slide and  
10 not at region level, whereas an effective chatbot for pathology or a copilot for pathology  
11 or a singular model that can cater to all of pathology will have to operate at all of these  
12 levels. So, the single biggest challenge is that fine grained understanding of pathology  
13 regions and cellular level information leads to slide-level diagnosis. But that fine  
14 grained information does not reside in pathology reports. So, this required a large  
15 amount of manual data collection, which we had from multiple collaborators across six,  
16 seven different institutions. And the data collected also needs to be incredibly diverse  
17 and should cater to all known entities within, of course, neoplastic, infectious,  
18 inflammatory or otherwise, that would be added.

19         So, in this initial version that was published in June of this year, we had about  
20 999,000 question and answer turns, corresponding images, questions and answers that  
21 were used to train this model. So, this relied-- The image features came from the  
22 foundation models that we had developed earlier, and then it was fine-tuned on this  
23 large 1-million-instruction dataset. And then this led to PathChat, which is a generative  
24 AI tool that can cater to all of human pathology. It catered to a smaller subset back  
25 when the article was published, but now we have expanded this to about 13 million

1 question and answer turns and have really extended the data substantially with a lot of  
2 institutions volunteering data just by virtue of the fact that this made a bit of noise in the  
3 community. But you can ask basic questions like, “What am I looking at here?” So, it  
4 could be a very good training tool. But what we are most excited about is that it can  
5 basically work at a whole-slide level where it can generate the report and by the time a  
6 pathologist is looking at it, they would have the entire trajectory of the diagnosis  
7 worked out and that report already generated.

8 Another thing that we are very excited about is that it can ingest multiple images  
9 within the same context, which means that if-- by looking at an image, a pathologist  
10 decides that additional tests would need to be ordered. Those tests could automatically  
11 be ordered. And once those images from, for example, immunohistochemistry are  
12 received, they can be ingested into the same context and our report would already be  
13 generated. So, in this example here, the report is being generated. Of course, a big  
14 utility. We had a lot of interest from a variety of different foundations who would like to  
15 use this tool in lower source settings because there are very limited pathologists in  
16 resource-limited settings, and it’s the norm to send or ship out slides for a diagnosis to  
17 be made in resource-rich settings. However, that’s prone to a lot of delays and it’s been  
18 shown that such cases, depending on where they’re sent and how much data-- How  
19 much tissue material is available, also have a lower accurate diagnosis rate.

20 So, in this case, the user is just taking an image with a cell phone coupled to a  
21 microscope. It’s a very common approach used in lower-source settings, and people  
22 take these images and send them to their colleagues in resource-rich settings and trying  
23 to see how well this worked. We did a lot of domain adaptation analysis, trying to see  
24 that these models that are trained largely on whole-slide scanners from a variety of  
25 different vendors, how they would adapt to image generation in this form and how the



1 resolution would vary. And then went back to the drawing board and improved our  
2 training pipelines to cater to some of those variabilities that would improve domain  
3 adaptation.

4 Another example, and something that a lot of people are very excited by is our  
5 ability to use PathChat by just a microscope coupled to a camera. So only about 4% of  
6 pathology diagnosis in the U.S. is digital, or 90% of the pathologists still use a glass  
7 light under a microscope. Most microscopes in the U.S. have a camera attached to them.  
8 So, in this case, the user is asking that-- This is a lymph node mass, and what's the  
9 likely diagnosis? And the chatbot says that, well, it's likely a melanoma. And the user  
10 can ask that what IHCs or additional ancillary tests could be used to confirm this. And  
11 we have a number of different possibilities. And once those tests are in, you can take an  
12 image of the IHC and it would ingest it within the same context. And this capability is  
13 possible because during the training, we mapped out trajectory for over a hundred  
14 thousand cases. Of course, we had lots of one-shot training with lots of images,  
15 questions and answers. But also, we mapped out entire trajectory of cases where we  
16 start with an image, additional tests are ordered, and basically the entire diagnostic  
17 trajectory of the entire case, and use that as the training data. And we are further  
18 assessing this capability on a large dataset that's collected across multiple institutions.

19 The most important component in building this-- Of course, building this took  
20 about three years. But in building this, the most important component for us is the  
21 assessment and evaluation, and how the evaluation needs to evolve over time as this  
22 chatbot essentially gets better. So initially, this was largely based on-- The qualitative  
23 analysis is based on multiple-choice questions. We basically took inspiration from two  
24 different domains; inspiration from how generative machine learning is evaluated for  
25 conventional, like non-healthcare related machine learning, and how large models that

1 cater to multiple diseases have been evaluated on the healthcare side, like NGS essays  
2 and so forth, and brought the two together.

3 So, in this particular case, we're looking at lots of multiple-choice questions and  
4 how the model fares against other generative models that are not specific to healthcare  
5 or are specific to healthcare, but not specific to pathology. But I think because this is  
6 generative AI and it gives open-ended responses, it's very important for us to assess  
7 how well those responses are given and recorded. So, in this particular case, we had a  
8 Panel of seven Board-certified pathologists who assessed the entire trajectory of  
9 conversation and each response generated by the chatbot and compared it against other  
10 generative models. So, we collected both quantitative responses, whether the response  
11 was correct or not—the binary response,— but also a more qualitative response that  
12 how well does this-- Does in comparison to other models.

13 If we go and look at how other large language models are often assessed, it's  
14 often done in a competitive setting where multiple models are assessed against often a  
15 human. And we found that the models do quite well for their intended use, for  
16 diagnosis, and for describing morphologic descriptions, but don't do very well  
17 correlating back to clinical or just marginal decrease in performance, or for example,  
18 ChatGPT or GPT4V or GPT-4.0. And the reason could be-- We don't know, of course,  
19 because we don't know what GPT-4.0 is trained on, but likely OpenAI is training it on  
20 all of medical literature and correlating back to clinical might be much more easier in  
21 their case.

22 We're improving the capabilities in terms of the ancillary testing, and this is  
23 what was done for the article. But we have substantially expanded this to about 11,000  
24 cases versus 260 earlier and have also expanded the number of pathologists who are  
25 evaluating this as we continuously improve this model.

1           So, now that we have the generative AI model, a next step is obviously to look  
2 towards agents. So, we started with supervised models and then we built self-supervised  
3 models for literature extraction, and we used those self-supervised models for a  
4 generative AI tool. And the next step is to see if we can build an agent so AI agents can  
5 essentially do things for you. In our case, we do a lot of biomedical data analysis. We  
6 want to see that can we replicate what essentially we are doing every day and could be  
7 continuous discovery agents. So, I'd just like to show this one example that we have  
8 where the user is saying, "Here's my dataset and I have a cord of responders versus  
9 non-responders. Can you go train a model that would predict what the-- Can you go  
10 train a model to predict responders versus non-responders?" So, the model would  
11 essentially make a plan. So, this plan is essentially coming from the generative machine  
12 learning component where you would segment the tissue, extract the feature as train the  
13 models and the mechanism of evaluating this. But then it can make use of existing code  
14 and patch it together with some additional code that might need to be written. So, in this  
15 case uses PathChat for morphologic description, it uses GPT-4.0 for writing the  
16 additional code that might be needed, and a host of different libraries that are out there  
17 look at the-- To essentially conduct all of these independent tasks. So, once the model is  
18 trained, we can ask additional questions. This essentially enables discovery from people  
19 who would need a clinical or a computational collaborator to run some of these  
20 experiments. And then the other benefit of this is that it can be crawling through large  
21 amounts of data, essentially doing what we do with lots and lots of pathology data  
22 manually in pathology departments throughout the history of pathology. The way new  
23 morphologic diseases are identified-- Morphologic entities that correlate with disease or  
24 outcome are identified is that someone goes in, looks at lots and lots of different cases  
25 and says "This similar morphology across cases correlates with outcome." In this case,

1 it's doing this essentially on its own and is using PathChat to generate a report. So, this  
2 is quite exciting and we have this being tested at everything we are scanning at the  
3 Brigham and Women's Hospital.

4 Towards the end, I just want to touch upon bias and fairness in computational  
5 pathology. So, this is a topic that we are very interested in because we want to make  
6 sure that the models that we develop have the broadest reach and as we go from more  
7 specific models to more generalized models, we want to make sure that their  
8 adaptability across-- Is not just assessed across a variety of different hospitals and  
9 scanners, but is also assessed within protected subgroups, patients of different  
10 demographics. So, in this particular case, we're using very commonly used data, for  
11 example, the TCGA and other public datasets that have been used extensively in  
12 computational pathology literature to train these models. And then we're adapting those  
13 to data that came from MGH and the Brigham, but we stratified that data based on race,  
14 whether the patients had insurance or not, the income status that was extracted from the  
15 postcode, and a number of different other variables. This is a study from earlier this  
16 year.

17 And we basically want to assess that any of these shifts that are commonly  
18 known, whether it's acquisition shift or just demographic shift, leads to any bias and  
19 disparities. And we know that acquisition shift, for example-- Here we're showing three  
20 different examples of the same slide being scanned from a different scanner, or in the  
21 case of radiology slide, the same patient being scanned at two different hospitals. And  
22 these cohorts that are quite often used for training do have large disparities in where the  
23 data was collected and what the demographic makeup looks like.

24 So, our approach was that we would basically look at every component of the  
25 computational pathology pipeline where we train these models in a weekly-supervised

1 manner and assess for every modeling choice and the impact that it would have on these  
2 disparities. This includes from different pre-processing steps to the kind of features that  
3 are extracted to the aggregation profile and using commonly used techniques for bias  
4 mitigation, for example, using contrastive approaches or other modeling choices that are  
5 known for bias mitigation. And overall, we found that the feature extractor is the most  
6 important component in improving fairness for these models, and this was assessed for  
7 basically every common modeling choice in computational pathology.

8 So, I'll stop here in the interest of time and I want to leave you with this  
9 interesting poem that Judith Prewitt, who is a real pioneer in analyzing microscopy  
10 images using computational tools, wrote. She was doing some very exciting work in the  
11 1960s and in 1970s at the NIH and then at the University of Pennsylvania. She writes  
12 something that I think is very relevant today. So, "Optical illusions can deceive the  
13 subjective eye, but objective measurements and algorithms are assumed not to lie. It's  
14 often said that medicine could use such objectivity and thought that this justifies  
15 machine intelligence activity. Artificial intelligence is another craze that uses computers  
16 to cope with the diagnostic maze. Though the criteria for intelligence has never been  
17 resolved paper after paper claims that the problem has already been solved." So, I think  
18 it's very important to keep in mind that there's still several challenges that we need to  
19 go through, and they're mostly related to evaluation and assessment of these very large  
20 models. So, I'll stop here and thank you so much.

21 Stakeholder Perspective – The Considerations for Multimodal Foundational Models &

22 Generative AI Frameworks in Healthcare

23 Dr. Bhatt: Dr. Mahmood, thank you. Thank you for all the work that your lab does  
24 in this area and thank you for being so humble about the amount of work you're doing,

1 but the importance of getting it right for clinical care. So, next we have Parminder  
2 Bhatia, who is the Chief AI Officer at GE Healthcare.

3 Dr. Bhatia: Good morning, everyone. This is Parminder Bhatia. I'm Chief AI Officer  
4 at GE Healthcare. I'm super excited here to talk about generative AI and its potential to  
5 transform healthcare. But just as important, we'll also talk about some of the technical  
6 guidelines and methods to manage the risk associated with some of this technology.

7 Firstly, before we get started, I would like to thank FDA for bringing this  
8 Committee together to discuss this highly relevant topic that has the potential to  
9 improve the quality of care for the patients, as well as reduce the cognitive overload for  
10 the providers. Just like FDA led way path for medical devices powered by deep learning  
11 to be safe and effective, they'll have a critical component in ensuring GenAI is used in a  
12 safe and effective way as well.

13 So today at the outline we'll be talking about three major components. We'll  
14 start by diving deep into generative AI and foundation models, what they are, why they  
15 matter, how are they different from traditional AI models. Then we'll talk about their  
16 healthcare impact. We'll talk about how foundation models empower providers to create  
17 tailored adaptive solutions through smart data analysis. I'm convinced that cloud and AI  
18 together have the power to reshape healthcare more in the next decade than we have  
19 seen in the past century. But let's not forget we'll have to-- We are dealing with  
20 people's health. So, we need to implement these technologies that are safe and making  
21 sure we are thoughtful and respectful as well as responsible. So, I'll spend most of the  
22 time today talking about frameworks for managing risk associated with generative AI as  
23 well as these multimodal foundation models.

24 Let's start at the very beginning. A couple of years back with ChatGPT, which  
25 was one of the first technologies to showcase the incredible potential of generative AI.

1 To understand ChatGPT, we break it down into four major components. The first is the  
2 transformer architecture. We won't go into the details of the architecture, but this  
3 architecture lets capabilities or components like ChatGPT to understand and generate  
4 human-like language by recognizing complex patterns, relationships, and context as we  
5 saw in the previous session as well. Secondly, the advances in compute bar have  
6 allowed us to process huge amounts of data in parallel. Thirdly, alignment instruction or  
7 human reinforcement learning shapes generative AI by using human feedback to refine  
8 models' output for accuracy, relevance, as well as looking into the alignment with  
9 human values. As we'll see today, having humans in the loop is essential for rolling  
10 these technologies out responsibly. Finally, we'll talk about foundation models, which  
11 are the backbone of generative AI, and that's where we'll focus next on.

12 Foundation models at the core are trained on vast amounts of data, and that is  
13 often unlabeled, giving them a strong foundation to tackle a wide array of tasks without  
14 requiring significant datasets upfront. We are seeing these models make a huge impact  
15 where models like ChatGPT, Claude, for language processing as well as are also used in  
16 healthcare for a wide variety of medical imaging and multimodal tasks as well. These  
17 models started as language, but one of the key powers of these technologies is they were  
18 multimodal. They always started with multimodal. So, that's what we have seen over  
19 the last 18 months. It started with language. Now they're adding voice, they're adding  
20 images. We looked into the previous example how easy it was to even add pathology  
21 into some of these components as well. And because these technologies are built on  
22 diverse datasets, they are incredibly adaptive. They can be fine-tuned for specific tasks,  
23 whether it's answering complex questions, generating unique content, or detecting  
24 conditions in medical images.

1 To give an example, at GE Healthcare as well we have been doing research into  
2 foundation models in X-rays and MRIs, which could be game changers. For instance,  
3 these foundation models could help analyze images faster, more accurately, generate  
4 automated reports that support radiologists in their daily workload. Another major  
5 advantage of these technologies is that they can handle different kinds of datasets. It  
6 could be voice, it could be images, it could be text. Both at the input level as well as at  
7 the output level. This versatility to handle multimodal dataset at the input level as well  
8 as at the output level is kind of game changing and as well as invaluable, especially in  
9 complex fields like healthcare.

10 Next, talking about foundation models and how are they different than  
11 traditional AI models. Foundation models represent a quantum leap over traditional AI  
12 models in three key areas. Firstly, traditional AI models rely on limited and narrow  
13 datasets. Foundation models, however, are trained on vast datasets, capturing a wide  
14 range of knowledge. This results in actually more reliable and insightful predictions.  
15 Older models are often defined or designed for a single task. It could be text-only or  
16 image-only processing. However, foundation models being multimodal can actually  
17 handle multiple types of datasets simultaneously. This capability opens up new  
18 possibilities for richer and more contextualized AI applications. Finally, traditional  
19 models are usually built for specific tasks and lack flexibility to adapt to new contexts.  
20 This can lead to longer development cycle as well as massively increased cost for the  
21 healthcare cont. Foundation models by contrast, bring broad adaptability, where they  
22 can be fine-tuned for various applications, making them more cost effective and  
23 impactful across multiple domains.

24 I'll give an example over here of how foundation models can transform breast  
25 cancer care pathways. Today, patients move through a complex workflow. From



1 screening to diagnosis to treatment, as well as to the monitoring, where they're going  
2 through various devices and components such as ultrasounds, biopsies, MRIs, and PET  
3 scans. Because of the challenges associated with traditional AI that we just talked about,  
4 we see various challenges that can arise. For instance, that can lead to inconsistent  
5 screening adoption, complex biopsies, as well as various other challenges. It can also  
6 lead to low trial enrollment, which would slow innovation, while it can also lead to  
7 treatment selection, guideline adherence, and follow-ups, which remains inconsistent.  
8 Foundation models being flexible, multimodal, and adaptable can actually help  
9 streamline this fragmented process addressing these key pinpoints with data-driven,  
10 precise workflows that improve patient outcomes and care efficiently.

11 Now, let's bring all of this together and explore how these innovations can  
12 actually translate to tangible outcomes that can actually improve throughout the patient  
13 journey. Traditional machine learning crafts data from each clinical modality in separate  
14 silos, which creates significant inefficiencies. Foundation models, however, break those  
15 barriers by synthesizing complex multimodal healthcare, longitudinal healthcare data,  
16 from diverse sources into a unified view. These models are highly adaptive. So, over  
17 here we can see an example and we can actually think of an oncology-focused  
18 foundation model for breast cancer, for instance, can actually be adapted to other cancer  
19 related areas like prostate cancer. They can even go across multiple care areas, where a  
20 model from oncology can go into cardiology or neurology. Take Alzheimer's, for  
21 example. Imagine if a new drug receives regulatory approval. Foundation models could  
22 help develop the new care pathway solution in the matter of weeks versus months or  
23 years that it takes today, revolutionizing how quickly we can adapt advancements in the  
24 field of medicine.

1           Next, we'll talk about some of the considerations and measurements that we  
2 need to ensure to make sure we can actually bring some of these technology-- GenAI-  
3 enabled devices into the system. To ensure the responsible rollout of these technologies,  
4 I'll focus on four major pillars. Firstly, we'll discuss about the defined intended use,  
5 which is critical. By clearly outlining the product, what the product is designed to do,  
6 we can actually set the foundation for every subsequent step, from design to regulatory  
7 compliance. This clarity ensures that product aligns with the user needs and industry  
8 standards. Secondly, adopting robust risk management practices is essential. By  
9 identifying and addressing potential issues early, we can enhance safety, improve  
10 reliability, and build trust with users and stakeholders. Premarket evaluation comes  
11 next. This involves rigorous testing and validation before the product goes to the  
12 market, catching in on potential problems and ensuring it meets the quality and  
13 regulatory standards. And finally, establishing change control and lifecycle management  
14 is key. These systems allow us to manage updates and improvements efficiently,  
15 keeping the product compliant, high performing, and relevant throughout its lifecycle.  
16 With these measures, we can navigate the transformative potential of generative AI  
17 responsibly and safely.

18           Next, we'll go mainly discussing two major topics, which is-- First is the  
19 intended use and then around risk management. The concept of intended use is  
20 fundamental to the FDA regulatory system. For a medical device, its level of risk  
21 depends on its defined intended use. This principle also applies to GenAI-enabled  
22 products. To scope the intended use of a foundation model or a generative AI,  
23 developers can take several approaches. For example, they can control the prompts to  
24 shape the system's behavior. They can create boundaries on what model can generate.  
25 This allows developers to focus the product on specific use cases. Defining use cases or

1 intended use case initially is going to be critical for ensuring the trust and effectiveness;  
2 starting with a specific intended use case and expanding it gradually. From the  
3 regulatory standpoint, the process starts with first FDA authorization that corresponds to  
4 the scope of the initial use case, complemented by predetermined change control plan,  
5 or PCCP, that allows us to use these use cases to expand over time with the same  
6 intended use case.

7 In this slide, we can see an example of how we might want to start with, for  
8 instance, cardiac segmentation in ultrasound using GenAI or foundation model tool, and  
9 later expand to segment structures in kidney or other organs. We can do this as the  
10 intended purpose and workflow remains unchanged. However, detecting and  
11 categorizing lesions might represent a new intended use case, as it involves interpreting  
12 clinical significance and influencing diagnostic decisions as well. A similar process can  
13 actually be used even for GenAI or agentic AI systems designed to fulfill a specific task  
14 as we seek to expand the universe of possible use cases.

15 Medical device manufacturers use established risk management methods by  
16 assessing the impact of failure modes in the context of use, as well as designing  
17 mitigations to address the possible hazards. The approach equally applies for GenAI-  
18 enabled medical devices as well. However, GenAI and foundation models also  
19 introduce new failure modes, like hallucinations, where systems can generate fault or  
20 misleading information. The good news is by anticipating these use cases, we can  
21 proactively design effective mitigations, ensuring safe and more reliable use of these  
22 advanced technologies. This structured approach helps systematically to identify, as  
23 well as to address, potential failures to enhance device safety and reliability.

24 Just like we would do with medical devices, we have to be cognizant of the risk  
25 involved in the rollout of these technologies and think about strategies to address them.

1 There are some of these components which are common for even traditional AI models.  
2 For instance, misrepresentation of information can happen even with traditional AI  
3 systems, and we have to be aware of these limitations. We can apply techniques such as  
4 explainability and transparency to address them. For GenAI systems, which can also  
5 produce hallucinations, ontology-based methods or reasoning, something these models  
6 are really good at, can ensure accurate and reliable information generation. We can  
7 address the risk. One of the common questions that comes around is the inconsistent  
8 output that these generative AI solutions can build out. We can actually de-risk this risk  
9 associated with inconsistent output with mechanisms such as temperature control, which  
10 can help us to ensure that the model output remains consistent and repeatable over time.  
11 Finally, across the board we find that issues can be mitigated by keeping human in the  
12 loop, allowing for essential oversight and judgment and critical decisions. The level of  
13 human involvement is dependent on the intended use case that we look across these  
14 different areas as well.

15 Now let's look at an example of how we might mitigate risk. AI models trained  
16 mainly on common diseases may struggle to recognize rare conditions, leading to  
17 potential misdiagnosis. For instance, a patient with rare autoimmune disease might be  
18 misdiagnosed with a common condition like rheumatoid arthritis due to gaps in AIs  
19 training data. A machine reasoning algorithm can actually help by enabling the AI to  
20 use reasoning, which can be both deductive and inductive approaches, which can be  
21 enabled by chain-of-thought reasoning to better navigate complex medical cases. This  
22 can be extended with ontology-based disease models that can provide a structured  
23 framework for organizing and sharing medical knowledge, which can improve the  
24 diagnosis, treatment, and research efficiency. So, by bringing together machine

1 reasoning and ontology-based models, we can actually enhance AI's capability to  
2 manage both the common as well as the rare conditions in healthcare.

3 Now I want to give another example and another technique that can be used to  
4 mitigate risk. Visual grounding, which involves locating relevant objects in an image  
5 using natural language descriptions, enhancing image interpretation and diagnosis in  
6 medical context. This technology improves image understanding by allowing for more  
7 precise identification of abnormalities. In the example of this slide, a radiologist  
8 generates a report describing a condition observed in an X-ray, such as pneumothorax or  
9 lung consolidation. The report is processed by medical report grounding system, which  
10 highlights specific areas in X-ray that correspond to the descriptions, assisting the  
11 clinician by visually marking relevant regions and ensuring consistency in  
12 interpretation. One could imagine the same visual grounding system to also apply to the  
13 output of a foundation X-ray model generating a report that links directly to the content  
14 in the image, kind of creating a flywheel around that as well, which kind of eases the  
15 review by the radiologist. This is just one of the many examples and there are many  
16 more which we can't cover with a shortage of time, but we'll be covering more of these  
17 mitigation techniques in the docket that will be available in coming weeks as well.

18 One way to leverage the power of AI is to select the foundation model that best  
19 aligns with your development needs, and then fine-tune it using domain-specific data  
20 for a specific task. For example, with SonoSAM, we started with Meta's promptable  
21 image segmentation model, initially-- Which was initially trained on 1 billion  
22 segmentation mass from natural images. We then further tuned this model on 200 K  
23 more ultrasound image mass. In this case, we didn't know if SAM was trained on  
24 ultrasound images or not. So, our first approach was to understand its baseline  
25 performance. That's when it was determined that fine-tuning would be essential in this

1 use case. We eventually compared the performance of SAM and SonoSAM on an  
2 independent and representative dataset of ultrasound scans to understand the behavior of  
3 model on a variety of clinical images. It's important to note that the level of information  
4 about the foundation model or dataset used for pre-training does not significantly  
5 impact the process used for validation. We validated the fine-tuned model on an  
6 independent, representative and diverse dataset against the ground truth to ensure  
7 adequate performance as outlined in this paper as well.

8 As I near the end of my talk, I want to take you through three-step framework to  
9 think about change control and device management for GenAI-enabled medical devices.  
10 Firstly, change control, which is maintaining appropriate version control and end-to-end  
11 validation, ensures that GenAI-enabled devices produce consistent and reliable output.  
12 Adaptability and speed. We just saw that by leveraging PCCP, or pre-configured  
13 controlled plan, we can safely expand the model's use case within its intended scope.  
14 This setup can also be used for quick iterations, such as adopting new foundation model  
15 versions, to keep pace with advancements without compromising safety. And lastly with  
16 post model monitoring by applying best software practices are essential for capturing  
17 feedback during real-world use, and we have seen some of those today and we'll share  
18 more of those in future as well.

19 So, in conclusion, I hope you got the sense of how excited we are about  
20 foundation models, generative AI, and how they can have significant potential to  
21 transform healthcare. With the right controls, this can effectively be integrated into  
22 medical devices to enhance their functionality. Leveraging existing best technological  
23 practices within the current regulatory framework allows us to manage these  
24 integrations responsibly. What's more, combining foundation models with PCCP  
25 enables us to implement changes safely and efficiently, ensuring timely advancements

1 while maintaining high safety standards. Thank you for your attention today. Thanks for  
2 your time.

3 Stakeholder Perspective – Measuring Performance of Generative AI – Methods and  
4 Lessons Learned

5 Dr. Bhatt: Thank you very much, Dr. Bhatia. Our last presenter for the morning is  
6 Dr. Pranav Rajpurkar, who is the Assistant Professor at Harvard University. Dr.  
7 Rajpurkar, please, your presentation. Thank you.

8 Dr. Rajpurkar: Delighted to be here today. Thank you so much for having me.  
9 My goal today is to share some of the lessons that we have learned in trying to evaluate  
10 generative AI applications. Here's my disclosure.

11 There are many applications where we can have generative AI make big impact.  
12 One of them is the ability to generate text from images, be able to summarize clinical  
13 notes, be able to record conversations between patients and doctors, and to be able to  
14 summarize them. It's not just about text outputs. We can enhance images, we can  
15 generate visualizations, which can all be helpful to different stakeholders. A few months  
16 ago, I was shadowing a GI colleague of mine at one of the Boston hospitals and they  
17 use these narrow AIs to be able to highlight colon polyps on the screen. There was a  
18 time at which one of these detections came and the GI doc didn't know what it was  
19 about. Wouldn't it be cool if we would be able to ask an AI model a question, explain  
20 the object appearing on the screen, and the AI model could understand the intent and be  
21 able to answer that query in real time? These are the kind of technologies that I'll be  
22 talking about today.

23 The reason that we are seeing a lot of advances in generative AI is because of  
24 technology advancement on three fronts. One of them is the ability to leverage on-label  
25 data; another, the ability to leverage different kinds of data or multimodal data; and then

1 the ability to communicate that in natural language. So, the main question I want to talk  
2 about today is if we have a generalist AI system that can generate these text outputs for  
3 a variety of inputs, then how do we think about evaluation? And I want to break it down  
4 into two fronts.

5 One is how do we think about automatic metrics which allow us to do things at  
6 scale. And the second is what is the role of humans in the loop, of experts in the loop, as  
7 part of this evaluation. So, I want to start with the question of metrics. What metrics  
8 allow us to determine, in this case, whether a report that's written by an AI matches a  
9 report that's been written by a human?

10 Here are the examples that we can see on the left and we're going to look at  
11 different metrics on the right. So, we started working on this problem in 2020, and one  
12 of the literature fields which has looked at this question a lot, has been the field of  
13 machine translation, where we have translation produced by a model and then a  
14 translation written by a human. And there one of the most common metrics are metrics  
15 like BLEU. And what BLEU does is it looks at word overlap to determine similarity.  
16 For example, here there's a perfect match on certain terms like air bronchograms, and  
17 this is reliable when identical terms are used. However, this often fails where different  
18 words are used to convey the same meaning, so "left lower" and "left base,"  
19 "consolidation"- "opacity," "effusion" - "fluid collection". This is not something these  
20 can recognize. So, how do we solve this? So, the time we turn to thinking about  
21 embedding-based metrics. Embedding-based metrics decide how to embed every word  
22 such that we capture the meaning of that word. Here, "consolidation" and "opacity"  
23 would be considered similar; and we were using this metric called BERT Score.  
24 "Effusion" - "Fluid collection" now becomes similar. However, one pitfall of these  
25 metrics that often happens is that opposite meanings can be represented very similarly.



1 So “with” and “without” mean two very different things. We want them to be far apart,  
2 yet these metrics often don’t distinguish between the negation and the presence of  
3 something.

4 So, in 2020 we started to think how do we design a metric that captures what  
5 matters in clinical medicine. And so, we came up with this clinical accuracy metric  
6 called CheXbert to evaluate the medical context within the text rather than the text  
7 itself. So here understanding things like “with and without mean opposites” correctly,  
8 mapping synonyms to standard terms, preserving the negation, are all the positive  
9 things. However, this took us several months to build and was limited to a set of limited  
10 conditions. We also could not link these findings to the locations in which these  
11 findings were occurring. So, if you have a left pleural effusion versus a right one, this is  
12 not something this metric would be able to help us disentangle. So, how do we solve it?

13 So, we started to think how can we extract medical knowledge from these  
14 generations. And we developed a metric called RadGraph-F1, where the idea is: each  
15 one of these reports contains medical context; can we explicitly extract the medical  
16 knowledge here? And so, we generate these graphs off of these two reports like you see,  
17 which captures the different entities and the relationships of those entities to their  
18 absence or presence, but also to the anatomical location. And that’s very helpful. Now  
19 we can compare these two graphs with respect to each other. The advantage of this is  
20 that we can capture a bunch of findings, we can preserve the relationship between the  
21 finding and the anatomy and we can handle negation, but this is quite hard to extend  
22 across modalities. RadGraph-F1 was developed for chest X-rays. We collected a bunch  
23 of human labels; we created a model that would be able to automate part of that. And  
24 then the second challenge is it’s very hard to normalize across entities to be able to

1 understand that one medical term is actually a subset of another medical term. It's a  
2 challenge that's very hard to solve.

3 So fast forward three years later, we have been working very actively on this  
4 front. A metric that we developed, called HeadCT-One, is able to compare two reports  
5 using knowledge ontologies. Here, the original report and an AI-generated report have  
6 both been tagged with the anatomies, the descriptors, and whether or not an observation  
7 is present or absent. So, what you can see on the right is we can now understand that the  
8 extracted entity of "brain" and "cerebral" mean the same thing. "Atrophic changes" and  
9 "atrophy" mean the same thing. And the "small vessel disease changes" means the same  
10 thing as "microvascular ischemic changes". And we understand that the content of these  
11 two reports, where it matters, is the same.

12 But I think one question that we should start asking when we think about  
13 automated metrics is, "Does this align with what we care about?" In particular, if we  
14 asked experts to score this report to say, "This is the AI report, this is the generation. Do  
15 you think this is correct?" And they say there's one significant error here, there's one  
16 insignificant error here. Do any of these metrics allow us to really get at that salient  
17 question? So, how do we design this experiment? We thought in 2022, "Let's go ahead  
18 and understand the correlation between each one of these metrics and the expert rating."  
19 And so, we have the radiologist score on the x-axis and the metric score on the y-axis  
20 here, and we're going to evaluate these four metrics. And here's what we found. We  
21 found that the expert scoring revealed the highest alignment with RadGraph-F1  
22 compared to BERT Score, CheXbert, and BLEU. But what we also found is that these  
23 different metrics are really capturing different ideas, and maybe it makes sense to think  
24 about how can we get the best of all worlds in creating a really good metric.

1           And so, we thought, could we compose these metrics? Could we take the  
2 weighted combination to be able to create the alignment with experts? And that's what  
3 we did with the metric that we proposed called RadCliQ, where we showed it has the  
4 highest correlation to experts and actually combines these different metrics together.

5           But often—and certainly is going to be the case clinically—a single number that  
6 tells you how well a report matches another one is often not fine-grained enough to tell  
7 us what exactly is wrong here. And we believe there is a strong need to understand the  
8 source and the type of the errors in the generations are made. And we break it down into  
9 two questions. One is “What parts of the generations actually have errors?” And two,  
10 “Is this a significant error or is this a stylistic one that won't make a difference to patient  
11 management?” And so, we developed this metric called FineRadScore.

12           But before I talk about the metric, I want to talk about the idea here. The idea is  
13 “Let's start with an AI-generated report and let's ask experts to fix it line by line.” Now,  
14 they might decide to remove certain lines, they might decide to edit certain lines, and  
15 they might decide to add certain lines. Each one of these edits has a cost associated with  
16 it or a level of significance associated with it. For example, the first two errors might be  
17 non-actionable or actionable non-urgent errors, while we might have an emergent error.  
18 And we really care about catching those emergent errors; we don't want to make worst-  
19 case mistakes.

20           Here's a comparison of what happens when you look at the average error made  
21 by an AI generation and a human. If you look at the different error severities on  
22 average, you can see that if we take a look at something that's a clinically significant  
23 but non-emergent error, which is at two, 68% of AI reports and 69% of human reports  
24 have errors of two or less. And this seems great; it's almost equivalent. But it doesn't  
25 tell you, and doesn't tell us, the full story here, which is if you don't look at the average

1 error, you look at the maximum error, what is the worst error that is produced by the  
2 AI? You see that there's actually a significant difference between the AI generation and  
3 the human generation. And this goes to show and say that we should not be thinking  
4 about averages. We should also be thinking about worst-case analysis.

5 One of the things that our lab is actively working on is "Can we automate this  
6 process of fixing AI-generated reports given access to the expert reports?" such that it's  
7 not something we have to go out and seek experts for in order to be able to do and  
8 FineRadScore performs this task pretty well.

9 So, I want to switch gears a little bit and talk about current state of the arts in  
10 terms of building these generalist models. And I want to share one work from our group  
11 called MedVersa. The idea of MedVersa is to think about building a model that goes  
12 across different modalities. Typically, we have different models that work on chest X-  
13 rays, on dermoscopy images, on CT images. Would it be that crazy to have a single  
14 model that could look at all these kinds of images and be able to perform a host of  
15 different tasks? These might be tasks like the ones we've talked about today, where  
16 we're generating a report, but it might also be traditional tasks like being able to classify  
17 a disease off of the image or even being able to outline the area of interest, draw a  
18 bounding box.

19 What we found was that when we trained a single model to be able to do these  
20 tasks, we found that MedVersa was able to outperform the previous state of the art on  
21 medical report generation—we're here using RadCliQ;—on being able to classify  
22 images—we're using the F1 score;—on being able to draw bounding boxes—we're  
23 using the Intersection over Union. And as we've seen earlier today, lots of images are  
24 used to be able to actually train such a model. And this idea that we are working  
25 towards generalist models that can actually perform better than the specialist models is

1 fairly new, because it's been long assumed that if you want something that works really,  
2 really well for this specific task, we should only go ahead and train for that specific  
3 task. That no longer seems to be the case.

4 So, I want to hopefully show a demo here of this model. So, I want to show here  
5 a demo of MedVersa as it goes through and generates a report on a chest X-ray image.  
6 You can see it outputs the output here for this image that should hopefully be visible  
7 now. One of the things we expect these models to do is look at the context. So, here it's  
8 a 70 to 80-year-old female, and we're going to generate a report here. One of the  
9 findings is moderate pulmonary edema. We should be able to ask a question particularly  
10 on that finding, so to say "How severe is the edema?" And it should say an answer  
11 that's consistent with the report that it's generated. Now we might want to do a task like  
12 being able to outline the heart. It's not quite right here, but I like to say its heart's in the  
13 right place.

14 We'll have to think about how to evaluate such systems, but it can look at  
15 dermoscopy images, be able to make a diagnosis, be able to perform tasks like  
16 segmentation on 3D modalities—here being used to segment out the liver,— and having  
17 a new task be solved, be as simple as simply changing the prompt used in this particular  
18 case.

19 So, that's MedVersa. And I wanted to share this example as an example of the  
20 capabilities that we're thinking about when we're talking about the next generation of  
21 multimodal generative AI models and what it takes for us to evaluate them. So, I want  
22 to now talk about human-centered evaluation, because that's such a key area. And one  
23 thing I'll mention before we do that is if we look at automated metrics, what I just  
24 showed you today, MedVersa is currently the best benchmark model in terms of  
25 generating reports on chest X-rays. If over the course of the next year or two, I expect

1 we'll see significant progress here and benchmarks are going to be very important. And  
2 we have the ReXrank benchmark that we have generated for this now.

3 So, onto human evaluation, how do you leverage expert capabilities to be able to  
4 evaluate such a generalist system? Well, one of the ways that we thought about it was to  
5 say, "We're going to ask experts to choose which report is better." One of them is going  
6 to be generated by another expert, one of them by the AI model. Now, the expert might  
7 say, "Both of them are equivalent clinically," or might say, "Report A is better" or  
8 "Report B is better." What we found was that overall expert-written reports are  
9 preferred; AI-written reports are preferred some of the time; and they're equivalent  
10 more than 50% of the time. We have a breakdown of how things perform on abnormal  
11 cases as well as normal cases. And on abnormal you see the human preference being  
12 greater, and for normal you see they're pretty much equivalent most of the time, and an  
13 equal split of the remainder between whether the AI was preferred or the human was  
14 preferred.

15 One of the other ways in which we think about evaluation is to think about what  
16 is the setup in which they're going to be used in practice. One of the exciting  
17 opportunities here for radiologists is to be able to spend less time on cases. The idea of  
18 working off of a draft and being able to edit it to sign off is not new. That's how  
19 academic radiology often works. And so, could we have a system in which these AI  
20 drafts are given to experts for them to modify and then sign off?

21 The way we set this up is we say, "We're going to measure the time spent on a  
22 study, the mental effort and the confidence in each report as radiologists are going  
23 through X-rays." The control here is going to be starting off with a normal negative  
24 template and the AI model as well as the radiologist is going to be supplied with the  
25 demographic, the indication, and also the comparison.

1           And what we found here in a small pilot study right now is that there is, in fact,  
2 significant time savings and mental effort reduction, as well as increased confidence in  
3 each study. But the caveat here is that it's going to be very important to think about how  
4 this technology can work for everyone. Even in small-scale trials, we're seeing some  
5 people experience big benefits from this technology while others show no improvement  
6 in their performance or their time taken as a result of these technologies.

7           I want to switch gears a little bit and talk about one of the things that we've all  
8 been very excited about. It's the ability of large language models to be able to have  
9 conversations. Now, there's a big disconnect in the way that these systems are evaluated  
10 and the ways in which these systems are expected to be used in practice. And I want to  
11 share one example of that. The primary evaluation for LLMs right now is based on  
12 answering medical exam questions, and we've seen incredible numbers in terms of the  
13 ability of LLMs to outperform humans on this task. Several studies, including the ones  
14 this week, that even show that the AI alone performed better than the clinician being  
15 assisted by the AI.

16           But these exam questions and the setup is broken in fundamental ways. One is  
17 that no patient ever presents with "Here's a summary of my symptoms, here's my  
18 vignette;" no patient talks about their medical terminology in such precise terms; and no  
19 patient gives options for a doctor to choose between. And these are ways in which we  
20 need to change the way we're thinking about the evaluation of these clinical large  
21 language models. One of the ways in which we're thinking about this is "How can we  
22 have benchmarks for conversational assessment of LLMs?"

23           Now, we can get these LLMs to have conversations with lots of patients, but that  
24 would be unethical in many ways. "Could we have LLMs simulate patients to have  
25 these conversations with clinical LLMs and see in that setting whether or not these

1 clinical LLMs do well?” is how we’re thinking about this. So, I want to present our  
2 results here on evaluating clinical LLMs on 12 medical specialties, both text-only and  
3 multimodal LLMs, and talk a little bit about continuous monitoring care.

4 The setups we have is, we have these case vignette evaluations, which we know  
5 are unrealistic, and what we do is we take these case vignette evaluations and ask a  
6 patient to share information with the clinical LLM. The clinical LLM asks questions to  
7 the simulated patient. The patient answers using the available information in the  
8 vignette and we evaluate this in a multi-turn, in a single-turn, in a summarized  
9 conversation setting. The main point was this, that when we actually get these clinical  
10 LLMs to have conversations, to ask the right questions, they actually do much, much  
11 more poorly than when they’re given the vignette. A large part of this has to do with  
12 incomplete gathering of relevant medical history for these clinical LLMs to actually  
13 have all the information to arrive at the right diagnosis. Similarly, multimodal LLMs are  
14 severely limited in their image interpretation capabilities and this becomes worse when  
15 we try to evaluate them in this conversational setup.

16 I want to end with talking about an application of generative AI that has several  
17 different components, but goes to show the potential value in terms of helping patients  
18 understand their own medical information. A lot of us, or our family members might  
19 have had experiences in which we get scans and we get radiology reports associated  
20 with those scans, and the scans come on CDs and the reports come over email, and it’s  
21 very hard to understand what any line of the report means, especially for those who  
22 might not be medically trained. So, we thought about, “Could we solve this problem by  
23 creating a system which allows us to understand different lines in the salient  
24 information in the radiology report and in particular be able to link it to regions of the  
25 image and also regions of the anatomy?” This would allow a patient or their family



1 member to understand what’s going on in their image, what was found, what does it  
2 look like, be able to compare to what a normal scan looks like, be able to see the organ  
3 and also have this virtual avatar do some explaining. Now, this is a system that can be  
4 built with a combination of existing components and components which we developed  
5 to segment out abnormalities. I now want to play you a demo of what this system looks  
6 like.

7 Okay. So, in this particular setting, there’s no audio, so I’m going to be the  
8 virtual avatar in this case. But, the first finding that’s being explained is that of aortic  
9 calcifications and as the system scrolls through the aorta, you can go ahead and have the  
10 virtual avatar talk about this. We are not going to go through the whole video, but if we  
11 fast forward it, we’d also have a comparison to the reference normal image, right there,  
12 and then also the organ rendering on the right. And this is our vision for what might be a  
13 world in which we have the patient be empowered to make their own medical decisions.

14 We can go back to the slides. So, in conclusion today I wanted to share some of  
15 our lessons of evaluating LLMs and generative AI. And what I hope to have shared with  
16 you is that they’re challenging, that there are automated metrics that take us a long way,  
17 that there are human-centered evaluations that take us a long way as well. The space of  
18 opportunities is immense and it’s a tremendously potentially helpful technology, and we  
19 need to think about the best way and the most safe way of bringing it to practice. I’d  
20 like to thank my lab members who made all this possible. Thank you.

21 Dr. Bhatt’s Summary

22 Dr. Bhatt: Thank you, Dr. Rajpurkar. Surprisingly, and I want this formally noted  
23 for the record, we are ahead of time. Thank you very much to our speakers in the room.  
24 I’ll take a minute to just go through the whirlwind we just heard one time in summary  
25 and then we’ll go to a break. I think it’s quite a feat we have ahead of us, but we can

1 already see that we have really talented people we can rely on, and so it feels good to  
2 have had this morning happen.

3 We started with Dr. Califf, who reminded us it's imperative to bring technology  
4 to those who need healthcare the most and have limited access. But then, he specifically  
5 said the criteria for adopting AI are largely financial and we need to think about clinical  
6 outcome measurement. I think that our lectures today helped us in that. He also  
7 challenged us that we don't have an information ecosystem. Local ongoing validation of  
8 AI is something we need to focus on.

9 We then had Dr. Michelle Tarver and Troy Tazbaz speak with us. And they  
10 reminded us that we can have a collective conversation to move forward equitably, that  
11 we need to think about AI for access. It's interesting. They've mentioned home health  
12 quite a bit from the CDRH, and when I was training, we used to talk about the last mile  
13 of healthcare, and that was the patient in their home. And I think what I'm hearing is we  
14 really need to focus on the first mile of healthcare and how we're going to diagnose  
15 patients earlier and help them earlier. So that's part of the conversation we're here for  
16 today.

17 We then had Dr. Badano speak to us, and he talked about regulatory science  
18 challenges. He talked about the challenge of defining scope, the fact that foundation  
19 models actually don't live necessarily with a device manufacturer, how oversight is  
20 important, hallucinations must be controlled and diverse data needs to exist and are  
21 sometimes challenging to find. And then reminded us that performance assessment has  
22 been established. We know what benchmarking looks like, we know what expert  
23 evaluation is, and now we are adding model-based evaluation, where there may be some  
24 challenges like inter-model leakage, but also a lot of potential benefit.

1 Dr. Mahmood next gave us a real impressive look into pathology and how far  
2 ahead they are. He reminded us that input patching, feature extraction, feature  
3 aggregation and prediction is a model we've always used, but that feature extraction is  
4 actually where AI can make a significant difference, and showed us what self-  
5 supervised models look like. Very impressive to see OncoTree go from 43 to 108, but  
6 really humble in saying, not necessarily ready for clinical use when we expand that  
7 much, that we still have work to do. So, thank you for that.

8 Interesting to see clinical reports versus generated reports in both his  
9 presentation and the presentations following and start thinking about what that looks  
10 like when we not only talk about the models and how they're going to work, but how  
11 we're going to understand their output. Because it's not only whether that model works,  
12 is whether or not we can do something clinically with it that is also correct. We then had  
13 the opportunity to hear about PathChat. Thank you for thinking about equity and global  
14 equity, not just United States equity.

15 Moving forward. Important to think about why foundation models. And  
16 Parminder Bhatia talked to us about training. How we can go from narrow to large  
17 training sets, modality, the ability to do multimodal analysis that didn't exist before  
18 other than in the human brain, and adaptability. Important to think about intended use,  
19 robust risk management. And thank you for bringing up the FDA's pre-control change  
20 plan. It is an excellent way that people have been able to look at devices. It may be  
21 applicable as we think about generative AI in those devices.

22 Dr. Rajpurkar gave us an excellent look at what we are doing now for  
23 performance measurement in generative AI. I think it's so important to recognize that  
24 we can leverage unlabeled data, that we can leverage multimodal data as we saw. And  
25 again, coming back to communication, how do you communicate that in a natural

1 language way that those who are going to use these tools will use them as correctly as  
2 the tools are able to work? I think as I look at all of this, one thing I will ask people to  
3 do is to think, importantly, but not talk about during the 10-minute break, “How does  
4 this help us think about premarket evaluation, ongoing risk management, and  
5 postmarket surveillance?” Because we’ve heard some very different fields, but we  
6 recognize that each of them applies to what are we doing premarket to assure a right?  
7 What are we doing in terms of risk management? And what are we doing when we are  
8 thinking about postmarket surveillance? So, just while it’s fresh in your head,  
9 sometimes thinking about it by yourself.

10 We will take a 10-minute break and maybe we can take more time for lunch  
11 later. We’ll take a 10-minute break now. Panel members, again, joking aside, please  
12 don’t discuss the meeting topic during the break amongst yourselves or with any of our  
13 panelists and guests yet. We will have the chance to talk about this when we come back  
14 from the 10-minute break. Lastly, yes. I’m going to be writing with pen and paper all  
15 day. I realize this is a Digital Health Advisory Committee, but this is how my brain  
16 works. So, as you see me scribbling, some things maybe you’re never going to change.  
17 Thank you and enjoy your break.

18 Open Committee Discussion Q&A (*Clarification questions*)

19 Dr. Bhatt: Well, I would like to thank all of the speakers for presenting this  
20 morning. We now have an opportunity for Open Committee Discussion and clarifying  
21 questions from the Committee. As a reminder, although this portion is open to public  
22 observers, public attendees may not participate, except at the specific request of the  
23 Committee’s Chair. That’s me.

24 Additionally, we request that all persons who are asked to speak, please identify  
25 yourself each time. This helps with transcription. So, for the panelists, please just

1 quickly say your name before you have your clarifying question, and then when you ask  
2 it to one of our five speakers this morning, if the speaker before you answers that  
3 clarifying question, could again say your name and then answer the question? That will  
4 help our transcription process.

5 With that-- Again, a reminder, if you have a question, you can either raise your  
6 hand or sometimes it's just easier to put your card up vertically like this. And James and  
7 I will keep track of the order of questions. And we have about 30 minutes for this  
8 process. Panelists, all yours. Okay. Dr. Apurv Soni.

9 Dr. Soni: Hi. Apurv Soni from UMass. One of the points of consideration that I've  
10 been thinking about as we looked at different approaches for evaluation of generative  
11 AI is, "Where does AI assurance labs fit into this ecosystem?" and "What role do we  
12 imagine AI assurance labs having?" And I'll just kind of add my perspective on it. I  
13 think there has been a lot of momentum behind Coalition for Health AI as an  
14 organization that has grown. And there is also a lot of emerging momentum behind  
15 development of AI assurance labs that may allow for either human-in-the-loop  
16 evaluation or model-based evaluation. And how does that interplay with some of the  
17 metrics that need to be evaluated for FDA's determination.

18 Dr. Bhatt: You didn't want to address that to anybody specific. So, I will now  
19 maybe say, do any of the Panel members, or more importantly, do any of the five  
20 speakers who perhaps have been thinking about human-in-the-loop, have a thought  
21 about that? Otherwise, we can take this as a comment as well for later. Okay. Such an  
22 important area to think about with AI assurance labs, and where they work. Let us  
23 continue with questions for now, and I think we can bring that conversation in a little bit  
24 later. Thank you so much, Dr. Soni. Anybody else? Yes, Dr. Shah.

1 Dr. Shah: Hi. This is Pratik Shah, University of California. So, my question is  
2 specifically directed at Dr. Aldo and Dr. Mahmood, or Dr. Rajpurkar, either of the three.

3 So one is that when we have GenAI models, we know that they will generate  
4 information that they necessarily were not trained on for clinical decision-making. And  
5 at the FDA, we have three different pathways. We have the De Novo pathway,  
6 premarket approval pathway, and then we have the 510(K) pathways. Right? And for  
7 the top two, we have those pathways activated when there is no other predicated device  
8 in the market for those pathways. So, my question is, when you have these kinds of  
9 models where there are no predicated devices in the market and they're making  
10 decisions or determinations based on training data that they haven't necessarily been  
11 completely trained on, for example, if you have a weekly-supervised learning or  
12 multiple-instrument learning or statistical learning, it hasn't seen all the data and it  
13 makes a determination. What actual regulatory science pathways should be considered  
14 these models fit into? That's one sub-part.

15 And the second part is, "What is the best way to do a randomized controlled trial  
16 on some of these models?" If that's one of the ways we want to go. Maybe a human-in-  
17 the-loop model where the AI model doesn't do clinical decision-making and the human  
18 continues doing clinical decision-making. Or a second pathway where it's an assisted  
19 care model. Right? So, these are some of the bigger concentrations to Dr. Aldo or any of  
20 the speakers, Dr. Mahmood or Dr. Rajpurkar, I can take. Thank you.

21 Dr. Bhatt: Maybe let's take that second question first, right? When we talk about  
22 level of testing necessary for us to trust something. Randomized controlled trials come  
23 up in clinical medicine a lot, but how are we proposing to test things? Do any of the  
24 speakers have kind of a thought or two? Yes, please. Thank you.

1 Dr. Rajpurkar: I think of this in two ways. One of them is, “If we have a  
2 generalist device, could we scope it to tasks that we can evaluate really well?” As an  
3 example of this, you can have a multi-disease detection model that can detect a bunch of  
4 different diseases, and we can apply existing frameworks to be able to do evaluation of  
5 those narrow paths. So yes, there are a lot of them, but we can evaluate each one of  
6 those individually.

7 The second aspect is if we leave it in its unconstrained format where there’s a  
8 bunch of inputs, there’s a bunch of possible outputs, one of the ways in which we need  
9 to transition our thinking is from thinking about sensitivity, specificity for individual  
10 models to thinking about preferences. And this change to preferences has already been  
11 something that outside of medicine, people have thought deeply about.

12 And the idea behind preference-testing is the following. We have, in a lot of  
13 cases, reports generated in the clinical settings by people, and we have reports generated  
14 by the AI model. Let’s ask a bunch of experts, “Which one of these do you prefer?” and  
15 “If you are preferring one over the other, is it because one of them would lead to  
16 clinically significant differences in terms of patient management?” Thinking about both  
17 the preference and the reason for the preference. We can set a non-inferiority margin to  
18 say that if we fulfill this non-inferiority margin, then this is at least as safe. And then  
19 downstream, we can think about it as at least as effective where the increased accuracy  
20 might come into play or the efficiency might come into play.

21 But I think this transition to thinking about preferences and evaluation on a  
22 dataset, which is representative of the kinds of settings in which we want to deploy, as  
23 well as with the kind of people that we want to deploy, too. For instance, if something is  
24 targeted at radiologists, then those should be the users. If that is particularly for  
25 residents, then those should be the users of that.

1           And I think we don't necessarily need to be in a setup in which we need a  
2 prospective double-blinded randomized control trial. In a lot of the settings, the  
3 blindness of the randomization won't even work because it'll be easy to see what's  
4 actually a normal negative template versus what's an AI generated report.

5           And I think we can actually go a long way, even with the retrospective  
6 evaluations to tell us, based on preferences, how well certain models perform. That's  
7 my two cents.

8 Dr. Shah:     Right. Quick follow up, just to clarify. At the FDA, we have this intended  
9 use for medical devices where we want to label the product by the people for what the  
10 product was intended to do. For some of these things, which are so broad-- I think you  
11 answered, Dr. Rajpurkar, some of those questions. For some of these products, which  
12 are so broad, it may be hard to come up with an intended use, a priority. Because once  
13 you have an intended use, later on you can have tracking and then software product  
14 lifecycle where you can track the model's performance. So, I think one of the things that  
15 you mentioned is to actually have a very narrow definition of who the end-users are,  
16 what the product label should be, and how the performance should be evaluated in the  
17 software product lifecycle. Correct?

18 Dr. Rajpurkar:     Yeah, and just to talk about that intended use piece-- I think if  
19 you're looking at, let's say, a modality and you're looking at the natural distribution of  
20 whatever scans are to be read, then that would be the scope of the usage. So yes, it's  
21 broad, but if in the intent that we're seeing retrospectively or with human usage, it's  
22 representative of the distribution that we want to deploy in, then that to me makes sense  
23 as a way to evaluate.

24 Dr. Bhatt:     Excellent. Thank you. For the transcriptionist, that was Dr. Rajpurkar  
25 who was giving the answer as the speaker and the additional comment was coming from



1 Dr. Shah. What I'm hearing, just to clarify, is that when we're thinking about device  
2 scope, we are also thinking about who is using the device and the setting in which it's  
3 used and defining all of those are equally important and sometimes that could be limited  
4 in a randomized blinded mechanism. It does bring us forward to the question for later,  
5 which is "What does postmarket surveillance look like then if we are putting these into  
6 the market and then requiring an understanding of what is happening once they're being  
7 used in the right setting with the right people, and how can that be expanded?" Let's go  
8 to Dr. Peter Elkin next, please.

9 Dr. Elkin: Thank you so much, Dr. Bhatt. So, first of all, I'd like to say that for  
10 every-- To help answer Dr. Shah's query. For every model that would be even  
11 considered for approval, we need a model card that has on it the exact dataset in which  
12 it was trained, what the demographics of the people whose data was used are, and so  
13 that we have a sense of the breadth of the training. If there's an underlying foundational  
14 model, we'd need to have the properties of that model that's defined in the card about  
15 what is in the potential model. We need the results of evaluations that were done prior  
16 to this, including the sensitivity and specificity of the model for particular tasks. And we  
17 need to understand what those evaluations have been, how they've been done, what was  
18 the methodology that was done in the trials that led to the submission.

19 I personally believe, and this is my opinion, that in between a clinician in some  
20 specialty or a patient and a model should be a clinical informatician. Someone who  
21 understands how the model was trained with the inclusion/exclusion criteria, for the  
22 studies, what the sensitivity and specificity of the model for different tasks are. So that  
23 they can give good recommendations, interpretations of how they're used. Very similar  
24 to an echocardiogram; it's used by cardiologists, it's not used by every doctor, it's used

1 by specialists in the field. We do have a Board-certified specialty in clinical informatics.  
2 These people can be a link to help us to make sure that it's safe.

3 And then to answer the second part, because I have strong feelings about that  
4 too. There's two kinds of evaluations that we need to put in place. One kind of  
5 evaluation are for rare but significant errors that happen, so that we do biosurveillance  
6 of problems that come up with the model. And then, second is for drift. Over time, as  
7 medicine changes, as the population that's used on it changes, the accuracy and  
8 effectiveness of the model can drift. So, we need particular agents that are looking at  
9 this on a regular way.

10 So, we need to create machine learning algorithms, probably, not just human  
11 research capabilities, that can actually be constantly monitoring this kind of information  
12 and sending feedback to some kind of a central registry where we can have researchers,  
13 like the people who presented today so eloquently, look at the data and understand when  
14 we have a problem. And there should be Data Safety Monitoring Board for these, like  
15 any other large study that would, in a phase four trial, be able to bring it back to this  
16 group, or whatever the appropriate body is at the FDA, in order for them to gain some  
17 kind of headway of when they should do something to keep the public safe. In other  
18 words, have thresholds of action where, if the model drifted beyond a certain point, it  
19 would be recalled. Or if there were a specific subpopulation that was at pretty high risk  
20 of some secondary bad effect, that those would be notified and that the use of that  
21 model would be restricted not to involve people who are at high-risk for harm. And if  
22 anybody wants to respond to that, I'd be happy to hear. Thank you.

23 Dr. Bhatt: Thank you so much. That was very helpful. Dr. Kukafka, did you want to  
24 comment?

1 Dr. Kukafka: Yeah, thank you for all these wonderful presentations. And, it got me  
2 thinking, which I think was the purpose of these talks. And I just want to-- It might be--  
3 We'll have more discussion in the afternoon about this, but as I'm thinking-- Because I  
4 do a lot of RCT type work in this space, and I just wanted to differentiate the way I'm  
5 thinking about some of the evaluation that we're discussing. A lot of what-- I think I've  
6 heard two things. A lot of what you just brought up or what I would term as process  
7 measures. And then I heard one reference to more endpoints, and the endpoint that I  
8 heard was around preference; the clinician preference, patient preference, I'm not sure,  
9 but preference.

10 I'm just curious at this point, and I know we'll have more discussion about this.  
11 What would be-- I think the process measures might be more well defined, I think. And  
12 I'm curious about what-- In evaluation, in a clinical setting, what would be the expected  
13 endpoints? Since it is very atypical to the very traditional kind of RCTs that we might  
14 have done in the past. So maybe just start to think about what would be the endpoints in  
15 this space? Not the process, but the actual clinical endpoints. And they don't have to be--  
16 - Preference would be one, because we do use patient preference, we use some decision-  
17 making metrics, which might be another clinical endpoint, like improvements in some  
18 clinical status. I'm just curious if there's any thinking about this area.

19 Dr. Bhatt: I love those comments. We have a deliberation section that will happen  
20 after lunch, and one of the ways that we were hoping to organize it, is to really think  
21 about structure, process and outcomes and how we're going to think about generative  
22 AI in that setting.

23 Just a quick reminder that I just got from our FDA colleagues, which is that we  
24 don't have this morning speakers beyond this morning necessarily. So, if we could use

1 our time to direct our questions to them, we will have more time for deliberation. And  
2 I've taken notes on all of those really good points.

3 If you could say your name before you start speaking, that would be great. And  
4 Diana Miller, you are up next. I did see that.

5 Dr. Radman: Thomas Radman. Thank you for asking your question on endpoints  
6 because I also wanted to ask some clarifications on this idea of preference-testing. And  
7 I'll give you my interpretation and then I'll invite the speaker to clarify. But I guess I  
8 wanted to say that-- So, we're here to discuss generative AI. And AI—let's call it  
9 maybe for the sake of this Panel, pre-generative AI—has been in existence for a while  
10 and there's some really great guidances and regulations that FDA has used to regulate  
11 that thus far. So, we can learn a little bit from the past to inform the future of generative  
12 AI.

13 So, when I think of preference-testing, and please correct me if I'm  
14 misinterpreting you and I have a question for you directly actually, but-- The FDA  
15 currently, essentially has two types of diagnostic devices. You have devices that directly  
16 provide a diagnosis, and I think you guys are allowed to correct me as well if I'm  
17 misspeaking, right? As the FDA. But, there's diagnostic and there's an adjunct to  
18 diagnosis, which is where I interpret this phrase of preference-testing because  
19 adjunctive diagnosis is-- The clinician basically is-- It's just helping him to-- And he  
20 could reject or accept the recommendation by the machine learning algorithm, pre-  
21 generative AI algorithm typically thus far. And then when we talk about preference-  
22 testing, you can imagine a trial where you have half of, let's say radiologists, if that's  
23 the example we want to use, use this device to adjunctively help them provide the  
24 diagnosis. Does it augment their ability to diagnose? And then you have half doing  
25 standard of care. And, so if I'm interpreting you right, what is preferred would be

1 whichever of those two randomized groups in a clinical trial performing the diagnosis  
2 better by some significant margin.

3 The question I had for you before you start speaking is, because we have this  
4 track record, if you will, of adjunctive diagnostic devices being cleared for marketing  
5 authorization by the FDA, what would be your recommendation for the number of users  
6 in such a study in, let's say it's a two-group study, with and without a generative AI  
7 device? And does generative AI, which should be our focus here, provide new questions  
8 above just a pre-generative AI machine learning device where you might want more  
9 users than the FDA is used to?

10 Dr. Rajpurkar: This is Pranav Rajpurkar. Let me try to be very concrete about the  
11 trial setup, just so we can have a working point. The trial, I'm imagining-- Let's start  
12 with the intended use case of drafting. So, there is a draft that's generated. A clinician  
13 will be responsible for editing that draft and signing off on it rather than writing  
14 something from scratch. Let's use that. So, the idea that we've had and performed at  
15 this point of a randomized controlled trial is one in which you have a certain set of  
16 clinicians that's starting from the current clinical standard, that might be from scratch.  
17 So, they have those and then we have another set that's randomly assigned to start with  
18 an AI draft and they edit and we get two final sets of reports that have been generated.  
19 Now, a panel that's independent of this reviews each of these reports and decides  
20 whether they preferred one that was generated with the AI human combination or we  
21 can do the classic human in their standard, or we can even have the AI by itself in the  
22 case that we're thinking about having autonomous generations. And then the preference  
23 is essentially what fraction of the time was one preferred over the other. And this takes  
24 into account, a lot of the time, clinically significant errors being made. If those are  
25 made, that shouldn't be preferred compared to the clinical standard. And so, the idea

1 behind a lot of this is to say we're not aiming for perfection, we're aiming for better  
2 than existing clinical standard and this is one way that we can think about setting this  
3 up.

4 The second part of your question had to do with how many users. And I think  
5 it's a very challenging question. Primarily because it matters what the diversity in the  
6 intended users looks like. And what we have shown in a natural medicine publication  
7 earlier this year is that things like years of experience, things like specialty, things like  
8 previous experience with AI tools were actually not good predictors of whether  
9 someone was going to experience benefit with the use of AI. So, it's unclear right now  
10 in literature what is the sample size needed from the number of users. One of the ways  
11 that we can set primary outcome is to think about the average. It's to say, "On average  
12 as a group, the clinicians that are representative of the intended users should experience  
13 improvement" and recognize that there is going to be a distribution where some  
14 clinicians are not going to benefit and some clinicians are going to benefit a lot by the  
15 use of that technology. How we make that variance small, we don't know, but I think  
16 that's going to be the reality of the evaluations today.

17 Dr. Bhatt: Thank you so much. We'll go to Diana Miller next and then Dr. Stanley  
18 after that.

19 Ms. Miller: Thank you very much for this session. I'd like to commend the FDA for  
20 organizing this and getting ahead of technology. As far as I recall, I think FDA  
21 approved over 950 AI-enabled devices as of past August. So, we did a good leap in that  
22 technology in the past and I'm sure we're going to be doing a good leap in GenAI  
23 technology in the future. So, this is a great start.

24 I have very three concrete questions to some of the presenters that are related  
25 more on how we design these models. So, first is Dr. Mahmood and computational

1 pathology. My question is, and maybe I missed that when you mentioned-- Sounded  
2 like you used a pre-trained model; what pre-trained model did you use at the beginning  
3 and how did you-- Did you have any criteria of selections of a specific pre-trained  
4 model? So that's question number one. How did you pick that?

5 And the second question is for Dr. Bhatia, and I'm unsure if I mispronounced  
6 the name. You mentioned that we use foundational models and we can transfer them  
7 from one domain to another domain, and we talk about the intended use a little bit and  
8 how do you define intended use in that space if you move from one domain to another.  
9 Maybe you can clarify that a little bit.

10 And then the last question is for Dr. Rajpurkar about metrics, and I know he's a  
11 popular one today, this morning. He gets all the questions, but one more from me. Did  
12 you look at our metrics outside radiology space and if not, how would you advise  
13 designing a metric for GenAI in other space?

14 Dr. Bhatt: Wonderful. Let's take those in order and start with the selection of the  
15 pre-trained model. Thanks, Faisal.

16 Dr. Mahmood: Hi, this is Faisal Mahmood. So, we developed our own pre-  
17 trained model. It was trained on a lot of internal data from Brigham and MGH, but also  
18 data that was sent from clinical collaborators at other hospitals. And we attempted to  
19 maximize for diversity just based on what we knew from conventional computer vision-  
20 based applications, that diversity of data is much more important than the quantity of  
21 data in training these large self-supervised models. And then we use that model for our  
22 generative AI tool. Thanks.

23 Dr. Bhatt: Thank you. Dr. Bhatia. Your question was to define intended use for  
24 foundation models if we think that there may be multiple potential uses.

1 Dr. Bhatia: Yeah. Hey everyone, this is Parminder Bhatia. Thanks for the question.  
2 So, I'll first actually try to just answer the first question that you just asked as well on  
3 how do you decide the foundation model. It depends on the use case that the developer  
4 or anyone who's working on it would look into. For instance, one of the examples I  
5 gave was ultrasound segmentation model SonoSAM that was built using SAM, which is  
6 built by Meta. And there you start evaluating it at the baseline, how is it performing on  
7 those benchmark datasets, which are more universal, and then you start to see "Is there a  
8 value to fine-tune that or use that model off the shelf?" If you go into multimodal  
9 models, you could even adapt it to models like ChatGPT or Anthropic and tune them for  
10 specific use cases as well. In some of the cases it might be beneficial, as I think Dr.  
11 Mahmood also talked about, training the model from scratch as well, the foundation  
12 model from scratch. So, I think it depends on the use case and the accuracy bar that you  
13 can get from these off-the-shelf models as well as trying to adapt from those use cases.

14 Now to your second question on how do you go from one domain to another.  
15 One of the examples which I would like to give is on summarization. You have data for  
16 cancer patients that could be thousands of clinical notes that can be created over the  
17 years. Now you could start with an intended use case on prostate cancer for  
18 summarizing that data. And now you can even extend as you're starting to think, "Can I  
19 add more domains into that?" Now, "Can I add from prostate cancer to breast cancer  
20 into those same use cases?", which is summarization. And same thing goes across on  
21 the report generation, as we looked into as well. So, you start with that. You're going  
22 from image, you're going into report generation, you start with a specific disease-based  
23 care area, and then you start adding more components and that's where you can actually  
24 leverage capabilities from PCCP and others to add those as part of additional  
25 components, which are again aligned with the same intended use case where you're



1 either trying to do summarization or report generation or even segmentation where you  
2 start with few organs, add more organs later as well. So, the intended use case is still the  
3 same. You're just increasing the horizon of how either you can adapt more organs or  
4 you can adapt for more disease care areas as well. Thanks.

5 Dr. Bhatt: Great. Thank you. Dr. Rajpurkar. And while you're coming up, so from  
6 Dr. Bhatia, just thinking a little bit about defining intended use, still very specifically in  
7 the area in which that was trained, the idea of the pre-control change plan maybe now  
8 applied to generative AI, which is a little different than the way it's used now, but might  
9 be something we want to pin for our discussion later. Dr. Rajpurkar.

10 Dr. Rajpurkar: This is Pranav Rajpurkar. The question was around the usage of  
11 metrics outside of the field of radiology. So, we've done an evaluation of generative AI  
12 models across 12 medical specialties, and one of the ways that we've set up metrics is  
13 diagnostic accuracy. And on diagnostic accuracy, in one of the setups, we actually have  
14 an LLM that's responsible for determining whether the diagnosis according to the  
15 model was the same as the one made by an expert. I think that's where the subtleties of,  
16 and the brokenness in some ways, of using medical exam questions as a way of  
17 evaluating these models comes in. Because our main finding was that, "Yes, you can  
18 have these LLMs, have these wonderful answers when they're given options, when  
19 they're given these very clear questions," but things break down once you think about  
20 them in a conversational setting. So, we have this benchmark called Kraft M.D., which  
21 essentially has patient agents that are simulated, that have conversations with these  
22 clinical LLMs. And the evaluation ultimately is across the cases; what was the  
23 diagnostic accuracy? And then we perform subgroup analysis to think about "Are there  
24 specific subgroups?", which are either age, sex, we could do ethnicity or even specialty  
25 where there are greater performances than others. We've also investigated this idea of

1 preference testing once again, where if you have a report that's generated, even on  
2 cases, by models and by humans, then how much do we prefer one versus the other?

3 Dr. Bhatt: Thank you. We'll go next to Dr. Stanley, then Dr. Rariy. I'll ask you to  
4 put your placards down once you ask your question to help me. Thank you.

5 Dr. Stanley: Hi, this is Laura Stanley. So, this is for the presenters. I'm loving this  
6 mashup of the academics with the industry folks and the regulators. This is fantastic. So,  
7 thank you. So, my question-- I might've missed this and I apologize if I did. So, I kind  
8 of sit at this human tech integration, human AI space, and what I wasn't hearing, and  
9 perhaps I may have missed it, I apologize if so, is this notion and there's a huge body of  
10 science around trust and automation.

11 Have you collected-- What are some of the metrics? What are some of the  
12 things? What are your thoughts about trust and automation? I spent 10 years working  
13 with autonomous vehicles. We have levels of automation which we might want to  
14 consider, if it hasn't already been considered, in terms of-- I think of this as like human-  
15 out-of-the-loop, human-in-the-loop, human-on-the-loop. That's a level of automation  
16 we might want to consider. And now, what about the trust? And when I think about  
17 trust, we have to think about "Are we over-relying?" I heard that come up. Over-  
18 reliance. "Are we vigilant? Are we situationally aware?" So, the need for human trust is  
19 critical in these machines as we're interacting. So, I'm curious from the presenters, if  
20 you collected that, what are some of your results? What are some of your metrics?  
21 There are validated instruments that do exist that people have modified from aviation,  
22 surface transportation, nuclear plants, social robots, etc.

23 Dr. Bhatt: Okay. Validated metrics for trust and automation when it comes to  
24 generative AI giving us answers of results. Yes. Dr. Bhatia. Thank you.

1 Dr. Bhatia: Hi, this is Parminder Bhatia. So, I think one of the areas where generative  
2 AI or foundation models is different than models we have seen in the past is  
3 explainability or what's called a chain-of-thought reasoning. So, these models are not  
4 just giving out a number which is positive or negative or confidence score, but they can  
5 also break down the problem and explain as to why they came out with the output as  
6 well. And I think that adds more to the explainability as well as trust as well.

7 So, you're getting the output which could be the summary and you can actually  
8 tweak it in a way to provide reasoning as to why it generated that output as well. So, I  
9 think that's where some of these technologies being multimodal as well as coming,  
10 thinking more from how a human would think, both from cognitive side of things and  
11 inductive side of reasoning as well, breaking down a bigger problem into smaller  
12 problems and coming out with rationales for that. I think that's one of the areas where  
13 generative AI can actually help decipher some of these things and make explainability  
14 more as a norm than we have seen in the past as well. It can actually be used for  
15 traditional AI systems as well. In some ways, we are starting to see research where  
16 generative AI or LLMs are being applied for evaluation or reasoning even if the output  
17 came from some other system as well. So, I think that's one of the areas where it'll help  
18 to bring in more trust into these systems as well. Because at the end you still have  
19 human-in-the-loop, but you're providing not just the output which is positive/negative,  
20 but reasoning behind that as well.

21 Dr. Bhatt: Wonderful. Thank you. We have five questions remaining, 10 minutes  
22 remaining. So, my panelists, if we could make our questions directed directly at the  
23 specialist who spoke. Dr. Rariy, then Dr. Maddox.

24 Dr. Rariy: Hi. Dr. Rariy. Thank you so much. As I trained within the Harvard  
25 system, in the Brigham system specifically, so it's great to see so many colleagues here.

1 One question I had was for Dr. Mahmood, specifically as a physician. One of the  
2 questions that I have, and I'm very curious to hear your response to, is around the end-  
3 user information or end-user labeling. So as an example with the PathChat, the  
4 information that's provided back to the doctor I think is as important in terms of the  
5 references as an example, or in clinical practice. I'm an endocrine oncologist, so I focus  
6 a lot on precision medicine and the information obtained back from that includes  
7 references, includes specificity, sensitivity, accuracy, and that very much impacts how  
8 we're presenting this information to patients.

9 And so, while as an example, ChatGPT, when we're asking questions, it just  
10 provides the answer back. There's not really that rigor or that layer underneath around  
11 "Where did this information come from?" or "What specific information from?" as an  
12 example of that pathology impacts the sensitivity or the specificity or how the user  
13 uploaded that information, whether it's magnification changes or things of that nature.  
14 So, I'm curious to hear from your perspective and recognize you with the understanding  
15 that it's still in research phase and not necessarily ready for end-user prime time. But  
16 again, with that physician perspective, that's going to be pretty important. Dr. Elkin  
17 already stated that the necessary labeling for FDA approval is required, but I would say  
18 and articulate that from the provider standpoint, that's going to be extremely important  
19 as well.

20 Dr. Mahmood: Yeah. This is Faisal Mahmood. So that's a very important  
21 question. So, there are a couple of components. The first one being that how do we  
22 quantify-- "How well the model is essentially doing?" The other one is, "How do we get  
23 to the underlying references where the information might be originating?" In terms of  
24 quantifying how well the model is doing or the model generating, how confident it is in  
25 a said prediction that can then lead to as a sort of a trust-score for the physician. It's

1 currently slightly difficult to do that directly from generative models. It's just based on  
2 where the machine learning research is. But what we think would happen down the line  
3 for these very specific generative models, for example, one that's focused on pathology,  
4 we would have a generative model that focuses more on explainability and morphologic  
5 description and we would have a supervised model that does crude diagnosis and the  
6 two would work essentially hand in hand. And for the supervised model, we would be  
7 able to generate a confidence score as to what the diagnosis is. Whereas for the  
8 morphologic description, we would not have that kind of a score.

9         The other aspect around the references is that there's quite a lot of research  
10 showing that you can somehow trace back what the model used in getting to a particular  
11 determination. Most of it is *post hoc* and is not based on what the model actually  
12 generated. So, there are a number of commercial tools, and I think ChatGPT, the latest  
13 version also now gives some references in instances where the information is retrieval-  
14 augmented generation, where it's directly pulled from the Internet. It's easy to do.  
15 Whereas with crude synthesis, it might be more difficult to generate what the origin  
16 might be. That's just based on how the feature space is organized. It might not be  
17 possible to map those abstract feature representations directly back to where the origin  
18 was. So, in my opinion, the current approaches for how you get to what the origin of  
19 those references would be are not holistic to be used in the healthcare space. We can get  
20 ideas or hints, but on the other hand, this is also true for a lot of the machine learning  
21 models, including those that have received FDA approval where interpretability  
22 techniques are giving hints and not crude information for what the abstract feature  
23 representation essentially is looking at.

24         So, I think in the short term, a combination of supervised models working with  
25 generative models for specific use cases like for pathology could be the solution in the

1 long run. We should keep an eye open for what's happening in the conventional  
2 machine learning space for enhancing interpretability and so forth. Thank you.

3 Dr. Bhatt: Great. Thank you so much. With the remaining few minutes, Dr. Maddox  
4 and then Dr. Jackson.

5 Dr. Maddox: Thank you. This is Thomas Maddox. Dr. Bhatia in particular. I'd just be  
6 interested in your thoughts specifically around operationalizing some of the user-output  
7 of generative model use cases. I think you make the very important point that humans-  
8 in-the-loop is a critical point of auditing. I think we've heard from you and others that  
9 we do need to think about making visible things like confidence scores, the cognitive,  
10 the cognition the machine is employing when it provides a particular diagnosis or  
11 recommendation. But I think we all know, and I'm certainly guilty of automation bias  
12 and I think it's probably unrealistic to expect every frontline clinician to pay attention to  
13 each of those metrics in the course of them going through their busy clinical day. So, to  
14 that end, what is the proper role for auditing? Are we envisioning that we want a  
15 periodic charge to some sort of auditor—informaticians, clinicians, some  
16 combination—that uses that information and ensures that we are dealing with  
17 information of sufficient confidence, sufficient cognition behind it. What are thoughts  
18 about that auditing process that allows us to think through that?

19 Dr. Bhatia: Thanks for your question. This is Parminder Bhatia. So, I think  
20 operationalizing these models, this is going to be key as we look across the spectrum of  
21 things. You touched upon areas of bringing human-in-the-loop. Initially, as we go into  
22 some of these use cases, it will be a lot of offline metrics. Some of the things which we  
23 talked about from preference side of things, getting extensive post-model monitoring of  
24 some of these things, so making sure it's validated at multiple sites to make sure it's  
25 kind of meeting the bar as well. So, I think that's going to be key to make sure they

1 meet to the standards as well. So it's going to require more benchmarking across  
2 multiple sites, kind of making sure it meets and actually meets the bar in a lot of those  
3 areas.

4         And then subsequently, I think in post-model deployment monitoring is going to  
5 become key. In some form, how do we capture that. And then have more periodic  
6 reviews of those. And I think that's where we need to have more machine learning and  
7 science behind that as well, like having mathematicians looking into that, providing  
8 inputs into that, where it'll help us to know not just-- Yes, there could be multiple ways  
9 in which we can get the output. It could be confidence score, it could be reasoning  
10 behind that, as well as multiple metrics as well, but not all of them might be relevant  
11 when the clinician is actually looking into the output, where the whole idea of these  
12 technologies is to reduce the overall cognitive overload by adding these additional  
13 components. We are in a way increasing that. So, I think a lot of this needs to move to  
14 offline and once we start meeting the accuracy as it is there with I think current models  
15 as well, we go with sensitivity/specificity, we are meeting 80, 90 and then we start to  
16 deploy those components as well. I think in similar form you'll look over here as well  
17 where initially you would want to do a lot of validation of these things, want to look  
18 into post-model monitoring and some of that will be discussed later today as well. And  
19 then kind of creating that flywheel going across, as well. And that's where in some form  
20 or capacity, it has to be a combination of AI and cloud capabilities as well, because you  
21 would need to create that flywheel that can go across multiple areas as well. Thanks.

22 Dr. Bhatt:     Great. Thank you. For our panelists, I want the Committee to think--  
23 Sorry, not panelists, the Committee to think a little bit about this as we move forward.  
24 The who, what, when, where, how of auditing, right? Just keep that in your head as

1 we're having more conversations. Let us end with Dr. Jackson as our final comment for  
2 this section.

3 Dr. Jackson: Thank you. This question is for any of the presenters today. One of the  
4 things we're tasked with thinking about risk management is controls for training. And  
5 I'm curious if you all have thought about "Is this going to be-- "Do people have to  
6 specialize in being able to use these models?" What are some of the training controls  
7 that will be needed for people who want to deploy these models in their hospital  
8 systems?

9 Dr. Rajpurkar: Pranav Rajpurkar. I think training's going to be important. We  
10 recently performed a study with 227 radiologists to find out what is the impact of AI  
11 assistance and we found that there was just incredible variability in the help that was  
12 produced. And what that goes to show is that there's something missing in the  
13 collaboration; recent studies now even showing that human plus AI performed worse  
14 than AI. All of this gets to the question of how do we enable that collaboration to  
15 succeed? And I think training's going to be a very important part of that. I think part of  
16 the problem, though, shouldn't just be on the users. The problem is often with the user  
17 interface and the developers and I think the onus in a lot of cases is if the impact of the  
18 tool isn't as great as we want it to be, then what do we need about the way in which  
19 we're interfacing that tool?

20 And regulatory wise, I think one of the things to think about is if we're thinking  
21 about separate regulations for different functionalities, for instance, if we want to have  
22 the ability to look at regions of the image which are being highlighted and have that  
23 same system be able to prioritize, are there existing regulations for that? And what do  
24 we need to do to be able to support the translation of systems that can create the sort of



1 trust and enable the sort of quick training to happen to really meet the user's needs and  
2 be able to improve their performance? So that's one way in which we think about it.

3 Dr. Bhatt: Thank you so much. It's now 12:02. We will now break for lunch. Thank  
4 you so much to our speakers this morning. For our audience who is here taking  
5 aggressive notes and to our Committee. Panel members, please do not discuss the  
6 meeting topic during lunch amongst yourselves or with any member of the audience.  
7 Panel members, please look at your meeting folder. It'll include information regarding  
8 your lunch arrangements. We will reconvene here at 1:00 p.m. to resume the Panel  
9 meeting. When we come back, it will be the Open Public Hearing portion of this  
10 Advisory Panel Meeting, so if you are registered to speak at the Open Public Hearing,  
11 please make sure that you're checked-in with the registration table just outside this  
12 meeting room. This meeting is now in recess for lunch. Thank you.

13 Open Public Hearing

14 Dr. Bhatt: Alright, 1:01 p.m., we can all start to trickle back in and get to our seats.  
15 Okay, welcome back everyone. This meeting is now reconvened. Before we proceed  
16 with the Open Public Hearing of this Advisory Panel, I just want to address our Digital  
17 Health Advisory Committee. You'll see in front of you three little post-it notes in order  
18 to, when we get to the question answer part, help us a little bit. If you have questions  
19 that arise as someone is speaking, feel free to write it on the note and pass it over to me  
20 after the speaker is done. That way if the questions are similar, we can compile them  
21 into one question and ask them to the speaker. And when those who speak come up  
22 afterwards, we'll just actually have them each come up one at a time and direct  
23 questions to them and hopefully that'll help us with our flow. So, thank you for  
24 cooperating with that.

1 Now, to proceed with the Open Public Hearing of this Advisory Panel Meeting.  
2 For the record, all Panel members have been provided written comments received prior  
3 to this meeting for their consideration during the Open Public Hearing, public attendees  
4 are given an opportunity to address the Panel to present data, information, or views  
5 relevant to the meeting agenda. Mr. Swink will now read the Open Public Hearing  
6 Disclosure Process Statement.

7 Mr. Swink: Thank you. Both the Food and Drug Administration and the public  
8 believe in a transparent process for information-gathering and decision-making. To  
9 ensure such transparency at the Open Public Hearing Session of the Public Advisory  
10 Panel, FDA believes that it is important to understand the context of an individual's  
11 presentation. For this reason, FDA encourages you, the Open Public Hearing speaker at  
12 the beginning of your written or oral statement, to advise the Committee of any  
13 financial relationships that you may have with any company or group that may be  
14 affected by the topic of this meeting. For example, this financial information may  
15 include a company or group's payment of your travel, lodging, or other expenses in  
16 connection with your attendance at this meeting. Likewise, FDA encourages you at the  
17 beginning of your statement to advise the Committee if you do not have any financial  
18 relationships. If you choose not to address this issue of financial relationships at the  
19 beginning of your statement, it will not preclude you from speaking.

20 Dr. Bhatt: Thank you, James. The FDA and this Panel place great importance in the  
21 Open Public Hearing process. The insights and comments provided can help the agency  
22 and this Panel and their consideration of the issues before them. We ask that each  
23 presenter speak clearly to allow the transcriptionist to provide an accurate transcription  
24 of the proceedings of this meeting. The Panel appreciates that each speaker remains  
25 cognizant of their speaking time.

1 We have four requests to speak. The first one was to be a virtual presentation by  
2 Sophia Phillips. I don't know if we have that video. We do not. We will move on then.  
3 Our next presenter would be Annie Soh from the American Institute for Minimally  
4 Invasive Surgery. Is Annie Soh present today?

5 Okay. The next request is from Dr. Zimmerman, [Corrected 04:34:29 - VP of  
6 Translational Science] at Tempus AI. Dr. Zimmerman, are you here? Dr. Zimmerman,  
7 you have five minutes. Thank you so much for speaking with us.

8 Dr. Zimmerman: Is this for advancing? Alright. Good afternoon, everyone. And my  
9 thanks to the Committee and for the FDA for the opportunity to discuss this important  
10 topic. My name is Noah Zimmerman. I am a biomedical data scientist by training and  
11 the Vice-president of Translational Science at Tempus. At Tempus, our mission is to  
12 develop technology that allows every patient to benefit from the experiences of the  
13 patient that came before. To this end, we ingest and structure huge amounts of clinical  
14 data, integrate this data with data from our own clinical laboratories, and use it to create  
15 and deploy artificial intelligence-based tools that clinicians and researchers use to  
16 advance healthcare treatment and discovery. We see tremendous potential for the  
17 application of generative AI to the vast amount of data that exists within the US  
18 healthcare system. And I'm here today to talk about how we believe regulatory  
19 frameworks can accommodate innovation in generative AI and best serve clinicians and  
20 patients.

21 In a recently published survey of physicians, one out of five reported that they  
22 are already using ChatGPT for clinical tasks including documentation, differential  
23 diagnosis and even therapy selection. This rapid uptake of a new technology, especially  
24 for an industry that still loves its fax machines, is a testament to the tremendous benefits  
25 that frontline clinicians see in leveraging AI for patient care. And the potential benefits

1 are immense, including use cases currently under development such as medical note-  
2 taking, consultations and virtual second opinions and results interpretations for patients.  
3 I think the message here is that the cat is out of the bag. The train has left the station.  
4 The alternative is not a world where patients and providers lack access to generative AI.  
5 The reality is that they will turn to general purpose tools often without the appropriate  
6 guardrails in place. That's why it's so important for this Committee to convene today to  
7 catalyze a critical discussion about how to establish those guardrails, ensuring safety  
8 and effectiveness while fostering continued innovation. FDA has deep experience in  
9 premarket review of medical devices with different technologies and intended uses.  
10 What's clear is there's no one-size-fits-all approach to determine the evidence needed to  
11 support premarket review. Like so many advances that have come before, generative AI  
12 is just a technology that medical devices will use to achieve specific purposes. Focusing  
13 too much on the technology, likely misses the forest for the trees.

14       It's also important to note that not all generative AI used in healthcare is likely  
15 to be a medical device. For example, at Tempus we are exploring the use of large  
16 language models for the task of information extraction from a longitudinal patient  
17 record. We may ask the model to review the case and extract the stage of cancer at  
18 diagnosis or the results of a particular diagnostic test. Here the task is discreet; take  
19 information from a clinical document and the results are specific; a yes or a no, or a  
20 stage at diagnosis. Used in this context, the output is easily verifiable and evaluation can  
21 be conducted just like any other machine learning algorithm.

22       However, many other applications will likely be considered devices and when it  
23 is, the FDA should apply the same approach to authorize almost a thousand safe and  
24 effective AI enabled medical devices to date. FDA's current framework considers both  
25 the technology and its intended use. And FDA should continue to apply sound

1 regulatory and scientific judgment to evaluate these factors and consider specific  
2 benefits and risks when determining the appropriate information to support review.

3 All of this is not to say that there is nothing unique about generative AI. One  
4 particular challenge facing the medical AI research community is testing and evaluation  
5 of healthcare applications of generative AI. Unlike predictive AI, which works with  
6 defined boundaries, generative AI produces original content (text, sound, or images)  
7 that can appear highly human-like. It doesn't have the degree of obviousness that it is a  
8 computer generated, which may lead to an increased risk of automation bias. And so,  
9 we think that this introduces two risks that are particularly important with respect to  
10 generative AI; the risk that the output is inaccurate and the risk that the output is  
11 mistaken for human judgment.

12 Early evaluation efforts for GenAI have relied on well-defined tasks such as  
13 answering multiple choice questions from a medical board exam. These tasks have the  
14 benefit of being clearly defined with inputs and validated answers, making them easy to  
15 assess. However, they fall short in capturing the complexities of real-world use cases. In  
16 practice, healthcare applications of generative AI often operate in unbounded  
17 environments where the inputs are diverse, unpredictable, and context-dependent. For  
18 example, a model used in clinical decision support might need to interpret incomplete or  
19 ambiguous patient data or a chat bot for mental health support may need to respond  
20 empathetically while avoiding harmful advice. These real-world scenarios present  
21 unique challenges. How do we evaluate performance when the correct output isn't  
22 always clear? How do we ensure reliability when the same input might generate  
23 different outputs depending on subtle contextual factors?

24 These challenges underscore the need for more dynamic and real-world  
25 approaches to monitoring generative AI systems, especially as their behavior can evolve

1 over time and vary depending on context. This is where robust postmarket surveillance  
2 mechanisms, like a modernized MAUDE database can play a pivotal role in ensuring  
3 transparency, safety and continuous learning across the total product lifecycle. The  
4 medical device reporting database is a critical tool for postmarket surveillance, allowing  
5 manufacturers and the public to report adverse events associated with medical devices.  
6 However, as we think about the unique challenges of generative AI, we believe there's  
7 an opportunity to modernize MAUDE to better meet the needs of this evolving  
8 landscape.

9 Generative AI introduces risks that are more dynamic and context-dependent  
10 than traditional devices. For example, the nature of errors such as hallucinations or  
11 inaccuracies at output can vary depending on the input prompt, the user's interpretation  
12 or the context in which the tool has been used. Modernizing MAUDE could serve  
13 several important functions including increased transparency, crowdsourcing insights,  
14 and continuous learning.

15 Further, a robust postmarket database can be used as a resource for research and  
16 development, allowing manufacturers and academic institutions to analyze patterns and  
17 design safer, more effective iterations of their devices. In short, MAUDE has the  
18 potential to evolve into an active data-driven platform that supports transparency,  
19 postmarket surveillance and ongoing innovation. By embracing this vision, we can  
20 ensure that the entire ecosystem from regulators to developers to end users remain  
21 vigilant and adaptable in the face of rapidly changing technology landscape. Thank you.

22 Dr. Bhatt: Thanks very much. Our next speaker will be Dr. Bernardo Bizzo,  
23 Associate Chief Science Officer, the American College of Radiology, Data Science  
24 Institute.

1 Dr. Bizzo: Thank you. Good afternoon, everyone. Thank you very much for having  
2 me here today. It's a pleasure to be speaking on behalf of the American College of  
3 Radiology about our Assess-AI program. I'm a diagnostic radiologist by training,  
4 Associate Chief Science Officer at the college and Senior Director of the Mass General  
5 Brigham AI business. Those are my disclosures.

6 And just for a brief context, the ACR exchanges images and data with over  
7 40,000 imaging entities across the country through its technology, the ACR Connect  
8 infrastructure for accreditation, registries, research, education, and very importantly  
9 today, AI programs. The Recognized Center for Healthcare-AI or ARCH-AI was  
10 launched recently and is the first national AI quality assurance program for radiology  
11 facilities with, at the moment, 18 recognized centers and almost 50 centers in the  
12 pipeline for engagement including international sites.

13 And as part of this Arch-AI umbrella, Assess-AI is a registry that provides real-  
14 world monitoring of imaging-based AI tools deployed in the clinical workflow. And  
15 Assess-AI allows for comparison of local performance metrics of these models to  
16 national benchmarks and also to facilities with similar characteristics such as region,  
17 facility type, trauma level, if it's urban or rural.

18 And it's currently being piloted with 15 sites with four AI vendors and platforms  
19 engaged for two FDA cleared clinical use cases—intracranial hemorrhage, ICH, or  
20 pulmonary embolism, PE on CT cases. And just at a high level, the current state, how  
21 it's being applied for AI monitoring. Each one of these sites or facilities are part of the  
22 program with existing ACR infrastructure can connect PACS data, including images  
23 and associated information, as well as medical records information and the AI  
24 manufacturer and platforms to ACR Connect and reach the ACR Assess environment  
25 where currently introduced two narrow AI use cases.

1           We collect site information about the final radiology report for each case that ran  
2 through the model, the AI model information, the manufacturer, the version of the  
3 software used and the output. This case-- Those are triage devices, just the binary  
4 classification of positive/negative for ICH or PE. And information about the case; the  
5 hospital, the facility, the patient demographics—age, sex, ethnicity,—the DICOM  
6 header information with technical details about the case that the model ran on.

7           And then at the ACR Data Science Institute, we are using large language models  
8 prompts to bias subject-matter experts to extract information from those recharge  
9 reports, in this case, if they are positive or negative for those two findings, if they are  
10 adequate or suboptimal in their quality as reported, and if the regionals had higher or  
11 low certainty when reporting those findings. And that data helps to generate analytics  
12 and reporting. So we can compare the report information with the model outputs and  
13 create dashboards with subgroup analysis. For example, assessing how AI impacts  
14 based on demographics, technical parameters, and also generate data regarding, for  
15 example, national incidence of findings.

16           And then, how this data can be distributed to different stakeholders at different  
17 levels. For example, for the site's local performance dashboards, benchmark reports for  
18 other sites that have similar characteristics I mentioned, and, of course, a national AI  
19 registry that can serve for much more information generation about these tools across  
20 the country. And that's the current state using narrow AI tools.

21           But of course, the future state is using GenAI use cases and draft reporting is the  
22 first one that we are thinking about going next. Instead of using the CNNs, of course,  
23 using visual language models to draft a radiology report. And then that information can  
24 be assessed using the same framework that is already established to assess draft report  
25 and final report differences.



1 And we're doing that through the Healthcare AI Challenge powered by the AI  
2 Arena. And you're going to hear more about that later today by Keith Dreyer. But in  
3 summary, it's a collaborative community for healthcare AI validation monitoring that  
4 includes the American College of Radiologists stakeholders, as well as Emory  
5 Healthcare, Mass General Brigham, University of Washington and Wisconsin.

6 So this data can be aggregated, healthcare professionals can be exposed to these  
7 AI tools and you can assess the local validation of these models as well as monitor their  
8 performance. With that, I want to be respectful of your time. Thank you very much and  
9 happy to answer any questions later.

10 Open Committee Discussion Q&A (*Clarification questions*)

11 Dr. Bhatt: Thank you so much, Dr. Bizzo. Will you stay there if you don't mind?  
12 That might make it easier. And for our Digital Health Advisory Committee, any  
13 questions for Dr. Bizzo? I'll start with one. First of all, excellent work and thank you for  
14 sharing how you think about it. As we're thinking about generative AI rather than  
15 machine learning per se, the way we've been using it with radiology, are there specific  
16 elements that you think about monitoring with regard to safety or efficacy that you  
17 would include in the kind of work that you're doing?

18 Dr. Bizzo: Yeah, sure. As you all know, this new technology acts in a probabilistic  
19 way. So, if you run the same data through the same model a hundred times, you may  
20 have different answers, hallucinations-- So we need to really focus on the content of the  
21 more output and how that, as was discussed previously by other presenters, can impact a  
22 patient safety and the healthcare providers not being biased if they're not experts in the  
23 area that they're focusing on. So, basically the content in the use case of drug reporting,  
24 I think there's a lot that we need to focus on with regards to the risks of the changing  
25 care management in a way that can affect patient safety and is an ongoing effort to think

1 about the specific variables that need to be accounted for. And that's the work we're  
2 doing at the college.

3 Dr. Bhatt: Are you able to share anything about those specific variables that you  
4 include to measure whether or not that is happening?

5 Dr. Bizzo: Sure.

6 Dr. Bhatt: Or if you're able to.

7 Dr. Bizzo: This is of course ongoing work. I think the number of changes compared  
8 to a final report when you look at the draft report in how clinically meaningful those  
9 changes are, because there are stylistic changes, there are content changes. So, we are  
10 really focusing on how these tools can help with efficiency gains of the healthcare  
11 profession while maintaining patients' safety. So, if you look at the number of changes  
12 that could impact patient safety and patient care, if those are 1, 3, 5-- How much that  
13 would impact the management of that patient, how that impacts the efficiency of the  
14 healthcare professional using that tool. So, we are really starting to scratch the surface  
15 about the specific metrics as was well discussed in prior speakers' talks. But I think the  
16 number and the clinical relevance of those changes are the initial points.

17 Dr. Bhatt: Thank you so much. Dr. Botsis.

18 Dr. Botsis: Taxiarchis Botsis. I think you more or less responded to my question that  
19 I'm about to ask. Are you also talking about specific metrics that should be used for  
20 specific use cases? And I think this is very much linked to Dr. Elkin's point before  
21 about having a card for each of those models with particular characteristics, but I want  
22 to really focus on the evaluation piece and that card.

23 Dr. Bizzo: Yeah. For sure, there's a lot to think about the specific metrics, and I  
24 think they're going to be use case, they are going to be model task-specific to a certain  
25 degree. The way we are envisioning—and Keith is going to talk more about the

1 Healthcare AI Challenge,—how different events, how this generative AI technology is  
2 applied to different types of data, generates different types of challenges. So, in a  
3 nutshell, I think we're going to need to expose those AI outputs to healthcare  
4 professionals at scale, to a lot of local validation, a lot of monitoring based on risk  
5 determination of how well that tool can impact the clinical care activity by those  
6 healthcare professionals.

7 Dr. Botsis: Taxiarchis Botsis. A follow up question on that. Do you specify certain  
8 thresholds for your use cases, thresholds that those models would exceed?

9 Dr. Bizzo: Yeah. So, we are thinking about threshold specifically about how well  
10 those models are performing. So, the initial thinking is having healthcare professionals  
11 rate the clinical skill level of, in this case, a draft report. How it was created by one or  
12 many of these generative AI tools. And then let's take the radiology example. Let's take  
13 chest X-ray. As a radiologist, I can look at a foundation model generated chest X-ray  
14 draft report and say, "This is at the level of a radiologist attending or at the level of a  
15 fellow." Good, but I need to make changes that change the clinical significance of this  
16 report, not just stylistic changes or at the level of a resident. It's okay, acceptable, but  
17 there is a lot to be changed to make sure that it's safe and effective. Or at the level of a  
18 student, where it's, with all respect, poor quality at that specific task—not at the level of  
19 an attending radiologist,—or it's unacceptable. So based on those rates of how well an  
20 expert professional judges the draft report of that use case, I can set thresholds where  
21 specific models that go below that threshold will not be exposed in subsequent steps of  
22 this challenge, of this event until we get to a point where a clinically safe threshold is  
23 defined for actual clinical use. Of course, bear in mind the whole regulatory aspect that  
24 needs to be taken into account before reaching that point.

25 Dr. Bhatt: Great. Thank you, Dr. Bizzo.

1 Dr. Bizzo: Sure. Thank you.

2 Dr. Bhatt: Does the Panel have any questions for Dr. Zimmerman? Oh, still for Dr.  
3 Bizzo. I'm sorry, I missed that. Dr. Bizzo, will you please come back? Dr. Elkin has a  
4 question.

5 Dr. Elkin: And you thought you were off the hook.

6 Dr. Bizzo: That's okay.

7 Dr. Elkin: So, one of the things with these that I always worried about, and I  
8 wondered how you're handling it, is that, first of all, there's trigger text that's used by  
9 many radiologic systems. So, you could say normal chest X-ray in a whole, your normal  
10 pops out and that's biasing in some way in the training data. But the thing that has been,  
11 I think, very difficult is that as people age, as they-- Different races and ethnicities may  
12 have different what looks like normal on, for example, a chest X-ray image. And how  
13 are you taking into account the breadth of normal? Because if you were to call  
14 pathology in all of the things that were normal, because it's very subtle in the way that it  
15 acts, it could be a burden to the practice. And I'm sure you're trying to improve  
16 efficiency. So, I'm wondering, how do you handle the routine variation in normal across  
17 the population?

18 Dr. Bizzo: Yes. Great question. And the same way that radiologists do it today. So,  
19 the actual images are being presented to the radiologists while they look at a draft report  
20 in this simulated virtual environment. And they can look at if those are generative  
21 changes, normal variations of anatomy, and how the generative AI draft report  
22 represents those changes. If called them pathologist disease findings that are worth  
23 mentioning on the report or as a radiologist that we say those are changes expected for  
24 the age of the patient. So of course, we need the expertise of the end-user to drive the  
25 significance of those findings and make sure that that's not biased in the results.

1 Dr. Bhatt: Thank you so much, Dr. Bizzo. Dr. Zimmerman, if you would rejoin us.

2 Dr. Radman has a question.

3 Dr. Radman: Yes. Thomas Radman. Thank you. In your presentation you had these  
4 differences between generative AI and traditional machine learning, if you will, or  
5 algorithms. And I'm really interested to learn from your expertise in developing  
6 generative AI algorithms. So, I guess I just want to get a feel from your developer's  
7 perspective on the necessity for generative AI to have unbounded outputs, as you said.  
8 Or is it feasible to have more singular narrow scope or intended uses of a generative AI  
9 product? And I guess in particular, what is the cost-benefit trade-off? Because I could  
10 imagine an unbounded output could maybe-- If it's just a one-shot regulatory approval,  
11 you've got something that could impact a really broad range of conditions potentially.  
12 But perhaps it's not validated as thoroughly in some respect, and certainly in terms of  
13 development costs, in terms of a business case for going condition by condition versus  
14 unbounded product. And what are the trade-offs there?

15 Dr. Zimmerman: Yeah. It's a great question. So, I think to start, I think you hit on a  
16 really important point, which is, I think, sort of the canonical examples that most people  
17 in the room might think of with generative AI because of the way that it was rolled out  
18 in the form of a chatbot, have that unbounded nature. So, it's a natural language  
19 interaction. And so, as such, the unbounded part is you can write any English sentence  
20 in there and then it can produce any sentence as output. Of course, there are more  
21 narrowly defined cases as you sort of alluded to, one of which I think I brought up in  
22 my remarks, which was using it for very specific kinds of tasks. So, taking the context  
23 of a patient and asking about the stage of a disease or the presence or absence of a  
24 particular finding. And so, I do think that there are very clear ways that you can bound  
25 those environments. And I think really it's about the intended use of the device and how

1 you contextualize the GenAI. What is the component that the GenAI is doing and how  
2 does that relate to the intended use of the device?

3 Dr. Bhatt: Great, thank you so much. Dr. Soni.

4 Dr. Soni: Apurv Soni from UMass. The two things that you mentioned, cats out of  
5 the bag and trains leaving the platform, whichever one we want to stick with. Really  
6 curious to your opinion and considerations of how can we study that cat and train. And  
7 this is in the context of Commissioner Califf advocating for more pragmatic clinical  
8 trials. Director Tarver commented on real-world evidence. So, what would be some of  
9 the considerations where we can observe digital breadcrumbs or other things? Because  
10 if one in five providers, which may be an underestimate due to social desirability bias,  
11 are using some of these GenAI tools, what can we learn from it in vivo?

12 Dr. Zimmerman: Yeah, I mean, as I think that this is the opportunity for the FDA to  
13 operate as a central clearinghouse. So, I'll give an analogy, is the Adverse Event  
14 Reporting System for drugs. So, that system has existed for a long period of time. One  
15 of the things-- One of the many things that such a system enables is third parties to  
16 come in and study these things. So, there are entire academic labs and PhD students  
17 who have devoted their lives to mining that Adverse Event Reporting System. You  
18 understand things like toxicity profiles, drug-drug interactions, having a similar kind of  
19 resource for generative AI in practice would provide an enormous opportunity for the  
20 community to collectively study what it looks like in the real world. I think in the  
21 absence of that, and without the FDA's gravitas behind such an effort, what you'll find  
22 is local places that have this information and can get iteratively better. But in terms of  
23 pushing forward the entire community to create more safe and effective devices, I think  
24 it would be a real contribution.

1 Dr. Bhatt: When you were talking earlier and you mentioned “level the playing  
2 field” you were referring to also being able to report what doesn’t go right, what are the  
3 errors, what happens. Can you talk a little bit more about error reporting and not in the  
4 academic sense, but in the industry and business-minded sense. Is that something that  
5 we can realistically ask of companies? What are your thoughts on that coming from the  
6 industry side and error reporting?

7 Dr. Zimmerman: Yeah. I think that postmarket surveillance is going to be a really  
8 critical part of this puzzle because as I mentioned in my opening-remarks, I think that  
9 the premarket presents unique challenges in evaluation. And so, I do think that industry  
10 will benefit from clear pathways towards how to get devices into market. And if that  
11 means additional guardrails in postmarketing surveillance, I think that that could be  
12 accommodated. I think there’s ways we could do that, that the burden could be not  
13 super high on reporting. For example, you could imagine an API-like interface that a  
14 manufacturer could decide to build into such a device that would allow people to report  
15 unusual outcomes and sort of have that be part of the postmarket process.

16 Dr. Bhatt: Thank you Dr. Zimmerman. And just quickly, Dr. Bizzo, will you come  
17 back for one more question from Dr. Jackson? And that will conclude our  
18 question/answer session.

19 Dr. Jackson: Thank you. Not to put you on the hot seat, but you were talking about the  
20 data from over 40,000 entities that you’re working with. And I’m curious because  
21 humans are fallible and clinicians are not excluded from that. Are there any controls in  
22 place to make sure that the images are real images, that they’re not AI generated  
23 images? Because a lot of this is based on what is the input if it’s an image. So, what are  
24 controls in place to make sure people are not creating images but using what was  
25 actually given to them?

1 Dr. Bizzo: Yes, of course. And just to clarify, the 40,000 facilities across the country  
2 are part of the ACR network. We have a subset of them engaged in the Assess-AI  
3 program that's currently being used for everyday cleared devices. So, only data from  
4 real clinical workflows where radiologists are looking at patient images, creating the  
5 report, and that information is being used for clinical care at those local facilities, are  
6 being transmitted to the ACR identified with all the proper guardrails regarding access  
7 to that data. And that is the data that's being used at the current state for the assessment  
8 with the AI results. So, you can see if those patients had or didn't have, based on the  
9 radiologist report, a finding of interest—intracranial hemorrhage or pulmonary  
10 embolism.—So there is no new data being created by those institutions, by other  
11 technologies, meaning AI generated data. But there could be, and I think it's a very  
12 valid question. In the current guardrails, since we're only using clinical data that's really  
13 being provided for patient care, that really simplifies that point. But we definitely need  
14 to keep that in mind for future endeavors with terms of AI and so forth.

15 Dr. Jackson: Can I ask a follow-up? So then I'm curious of your thoughts, thinking  
16 about scaling this. What would we need to think about? What are controls to think about  
17 so that people don't, right? Sometimes people misplace something, maybe they made a  
18 mistake because that happens and so they need to go back and correct even if they're  
19 using the image but they make a change to it. What are controls or guardrails to make  
20 sure that people are not manipulating the images that are the input that go in?

21 Dr. Bizzo: Yeah. We should leverage the existing controls in place in all those  
22 accredited facilities that are providing patient care and they are audited by the college,  
23 by other societies or other organizations, making sure that they are compliant with all  
24 the regulations. And we can add, on top of that, any additional guardrails if needed. But  
25 I do believe that those are very strict guidelines that already are in place so those



1 facilities can provide care for their patients. And we use secure connections to make  
2 sure that no data is interjected between that facility and the ACR. So yeah. I think  
3 relying on existing practice and make sure that they are enough for this new wave of  
4 technology of AI being used is certainly something that we keep in mind very, very  
5 critically.

6 Dr. Jackson: Thank you.

7 Dr. Bizzo: Sure.

8 Dr. Bhatt: Wonderful. Thank you both so much for your participation in this open  
9 session. So now, I now pronounce this session of the Open Public Hearing closed and it  
10 is now time for Panel deliberations and the discussion of FDA questions. The FDA has  
11 generated a series of questions for the Panel to consider. Ms. Aubrey Shick will present  
12 the questions and the Panel will deliberate amongst ourselves. Ms. Shick, please present  
13 the first question. After the Panel has discussed the question, I will summarize the  
14 recommendations for FDA before moving on to the next subpart of the question. Team,  
15 just in a less formal way, the goal is each of our questions that we received in advance  
16 have a subpart. What I'd like to do is achieve for the FDA three clear points of  
17 recommendation at a minimum for each subpart before we move on to the next one. So  
18 that's kind of our goal within this discussion. Oh, and they're in your folders as well.  
19 They might be in your email if you have your laptops open. I'm going to give you a  
20 minute to get those out. Okay. And Ms. Aubrey, you may present the first question.

21 Committee Discussion of the FDA's Questions (*Deliberation and response to FDA*)

22 Ms. Shick: Thank you, Dr. Bhatt. Our first question area is in regard to premarket  
23 performance evaluation. We would like the Committee to please discuss what specific  
24 information related to generative AI should be available to FDA to evaluate the safety  
25 and effectiveness of GenAI-enabled devices considering, for example, that foundation

1 models leveraged by the GenAI-enabled device will change over time and that there  
2 may be limited information available on the training data utilized for these pre-trained  
3 generative models.

4 So, there are four sub-questions. I will read the first one now. So, question 1.a in  
5 this premarket performance evaluation topic. What information should be included as  
6 part of a device's description or characterization in the premarket submission when the  
7 device is enabled by generative AI? For example, when a human is/is not intended to be  
8 in the loop, or if a device is intended only to recall information versus to generate new  
9 recommendations. What information is particularly valuable to evaluate the safety and  
10 effectiveness for devices enabled with generative AI in comparison to non-generative  
11 AI?

12 Dr. Bhatt: So maybe let's start with the first part of that question, which is, "What  
13 do we think should be included in the description or characterization in the device  
14 enabled by generative AI?" And if you have a present answer for one of the other parts,  
15 feel free as well. But just to help us organize.

16 Dr. Radman: Thomas Radman. I just had a clarifying question. So, I think the first part  
17 says, "Considering that GenAI devices can change over time," and I think we should  
18 distinguish that there's two different cases. There is a possibility as far as I'm  
19 concerned, that as far as I understand, GenAI devices can be locked. For example, we  
20 have GPT 3.5, GPT 4, etc. And then there's another separate question of continuous  
21 learning adaptive algorithms, which is a problem we need to figure out still for just pre-  
22 generative AI, right? And there's guidances on change control plans. So, I think it might  
23 help our deliberations if we separate those two questions like GenAI generally and then  
24 GenAI devices that continuously learn. Am I correct in that? GenAI does not have to  
25 continuously learn or is that a necessity?

1 Dr. Bhatt: I think that is very fair for us to break down as part of our response,  
2 which is lock GenAI versus allowing continuous learning. Is that okay? Alright. Dr.  
3 Maddox, you had a response?

4 Dr. Maddox: Yeah. So, I'll build on that. So just the three things that occurred to me  
5 and then I'd love people to edit or pile on is what's the use case or use cases envisioned  
6 by the particular device? Is it dynamic or not? And then, is it autonomous or not? Which  
7 speaks of the human-in-the-loop. And to me those are some of the things that I've heard  
8 to be important characteristics when it's time for the FDA to evaluate, its need to  
9 evaluate and the risk kind of category we put in place. But above reactions to that or  
10 building on that.

11 Dr. Bhatt: Tom, will you tell me the third of those intended use or use case dynamic  
12 or not, and-

13 Dr. Maddox: -Autonomous.-

14 Dr. Bhatt: -Autonomous. Thank you.

15 Dr. Stanley: If I can-- Oh. Can I build on that?

16 Dr. Elkin: Did you put me or her? Okay, so--

17 Dr. Bhatt: Why don't we go with-- Yes. You and then her. That's okay. Thank you.

18 Dr. Elkin: I'm sorry. So, I would like to build on Thomas's comments. I agree with  
19 them completely. I do think we also need to know the population on which the machine  
20 was trained and we need to know any results from prior testing, the methodology for the  
21 testing and the accuracy and the breadth of the testing. So, what populations was it  
22 tested on and what was the accuracy in those tests? Both and the safety of it, meaning  
23 that what were the error rates and qualitatively, although this is rarely reported, the  
24 things that-- Because they always answer, so the false positive rate and the  
25 confabulation rate are equivalent. So, the idea is that from those erroneous answers, how

1 much harm could come from those deviations from truth and that kind of an assessment  
2 would be lovely to have as you're making a decision about the safety profile of the AI  
3 and device. By the way, my large language model, when I put the question in, doesn't  
4 believe it's a device, it believes it's a software and it made a point of it to me. So, I  
5 thought I would give you that the LLM disagreed with me. I do believe it's a device.  
6 Thank you.

7 Dr. Bhatt: Dr. Stanley.

8 Dr. Stanley: Thank you. This is Laura Stanley. I was just going to quickly add on to  
9 that. I think that level of automation, that level of autonomy, I think perhaps needs a  
10 breakdown. We think about human-in-the-loop, human-out-of-the-loop, human-on-the-  
11 loop and how we classify that—if it's those three or if it's not,—I do think it's  
12 important that we know kind of where that human sits in that process and is it  
13 anticipated or is it not anticipated? If it is, how much? So, we can kind of think about  
14 those things of overreliance and vigilance, etc.

15 Dr. Bhatt: Dr. Botsis.

16 Dr. Botsis: Taxiarchis Botsis. So, to build a little bit on that about the use cases and  
17 how autonomous GenAI technology might be, I would think that it might be important  
18 really discuss the possibility of having a list of the particular functionalities or particular  
19 components that this particular GenAI supports for a certain use case because there may  
20 be a series of steps that one has to take there.

21 Dr. Bhatt: I have Dr. Soni next. Before that, I'm just going to, every couple  
22 comments, summarize what I'm writing so we're all on the same page. I have intended  
23 use or use case of the device with generative AI as well as the intended population for  
24 use is important to this Panel. I have the concept of "Is it dynamic or not?",  
25 "autonomous or not?", and importantly thinking about autonomy as human-in-the-loop,

1 human-out-of-the-loop, or not necessary, or on the loop. We'll come back to that in a  
2 little bit. And the anticipated human interaction that we expect. I have the definition of  
3 which model was trained, properties of the model and demographics of the training  
4 population from Dr. Elkin earlier in our conversation and understanding the  
5 functionalities and specifically the sensitivity or specificity of the device for those  
6 particular tests or functionalities. Combining comments from before with now. Dr. Soni.

7 Dr. Soni: In addition to who and what, I think we also need to consider where or  
8 the setting in which the medical device is being used, especially as a lot of our work and  
9 growing momentum is towards bringing healthcare to home. If a medical device that  
10 uses GenAI is going to be deployed in a decentralized setting, that needs to be pre-  
11 specified.

12 Dr. Bhatt: Dr. Shah next.

13 Dr. Shah: Right. Pratik Shah. So, I'm going to respond to the specific part A of  
14 what should be included. So, one thing that works for us is datasheets for datasheets, but  
15 before we accept a product in our lab to do any research, we have a standard pro forma  
16 datasheet where we have to fill out all the metrics of the data—the size of the data, the  
17 type of the data, the number of examples, where it came from, how long back it was  
18 collected.—So one recommendation could be that at the FDA we standardize this  
19 datasheet pro forma of accepting something to review, like, “Fill this datasheet out so  
20 we can track where all your data came from” and that could be just standardized  
21 completely. So there is no ambiguity when somebody submits something, what do they  
22 need to include? What do they not need to include? So that would really standardize the  
23 premarket. Just a characterization of what came in. That's one.

24 The second part is generative AI. I think a lot of people talked about it. When  
25 we do research on generative AI, there are different types of generative AI. So, if one

1 model converts one image to complete different image domain, that that model has  
2 never seen, that's a very unbound generative AI model. But if it converts an image that  
3 it has seen and generates some sort of notification or segmentation or some text report,  
4 that's kind of a different lower-bound of generative AI. So, it's very important that we  
5 have a clear distinguishing of what are the upper-bounds and lower-bounds of the  
6 models' generative capabilities on the tasks it is being asked to perform. Because you  
7 can have a lot of ambiguity in what the model is generating. Is it generating something  
8 it has already seen or it's generating something it has never seen? That's a second.

9 And the third is uncertainty estimations. So we use that all the time. We also  
10 need to know where the model is most uncertain or where your data was most uncertain  
11 that went into training the model. And there's Bayes by Backprop, Monte Carlo  
12 dropout, the technical ways to do it, but also the statistical and the common sense way  
13 of doing it is that asks counterfactual questions. Like what has this model not seen?  
14 What would it see that would break it, right? So that the uncertainty estimations are also  
15 very important deterministic things that we should really consider. So, we ask the tough  
16 questions upfront. What do you not know?

17 Dr. Bhatt: Dr. Kukafka?

18 Dr. Kukafka: Yeah. I think this might have been mentioned, but if a model is trained in  
19 a single institution and has a specific performance and it moves to another institution,  
20 usually the performance would certainly not be the same and most likely decline. So,  
21 whether that's part of the expectation that it would be trained in-- Depending on how  
22 it's being proposed to be deployed, I think that would be important to know.

23 Dr. Bhatt: Dr. Jackson next.

24 Dr. Jackson: I just wanted to add on more specifically to some things that were already  
25 shared. I think it'd be important to specify if they see it as a primary or secondary use.

1 Is it meant to support or is it meant to generate something brand new in the application?

2 I also think it could be useful to indicate what is the percent of hallucinations or error-  
3 rate that they have found when they're submitting so that there's a general  
4 understanding of that when it is submitted.

5 Dr. Bhatt: Dr. Rariy, please.

6 Dr. Rariy: Yes, thank you. I think some of the other things that I'd like to consider  
7 is-- And some of it was touched on, but if we're looking at risk-mitigation-- So if the  
8 individuals who are submitting for FDA clearance or approval of the medical device  
9 already have a sense of what the limitations are of the model, then that needs to be  
10 specified or if there are specific requirements that are needed for the intended use to be  
11 executed. As an example, are professionals supposed to have specific training? Or end-  
12 users, do they need specific training to leverage the device in a particular way?

13 I also would recommend that we consider the cybersecurity and the data privacy  
14 aspect as well as part of that general standardized framework. So, what has gone into  
15 ensuring that there's not an ability to tamper with the product in current state or that  
16 there was that ability already in place when the pre-training was conducted. Similarly,  
17 the privacy standards that they applied for that data needs to be included. And then  
18 lastly, to include in this framework, perhaps the FDA can maybe come up with  
19 examples of what type of stress testing would be recommended so we can clearly  
20 determine, or at least articulate, the importance of stress testing to some extent that  
21 would need to happen for the device for submission.

22 Dr. Bhatt: We'll do Dr. Soni, then Dr. Elkin. Sorry, Dr. Radman first, Soni, Elkin.  
23 And we'll pause for a second to listen to the second part of this question because we're  
24 drifting a little into it, which is great. Dr. Radman, then Dr. Soni, then Dr. Elkin.

1 Dr. Radman: Thank you. Thomas Radman. And I apologize that pretty much all of our  
2 comments seem to be in multiple parts because it's such a broad topic. But I did want to  
3 reiterate the point Dr. Zimmerman actually brought up on the ability to build in a  
4 method to report for users to report any errors. And I'll add to his comment that ideally  
5 the particular report would go to a centralized independent database analogous to the  
6 FDA's Adverse Event database. Just any bias detection or hallucination, I think is a  
7 large part of this because of the challenges of doing an equally rigorous premarket  
8 evaluation as we currently are able to conduct with pre-generative AI devices.

9 And then I had this comment about-- We're talking about several Panel  
10 members correctly pointing out the need to describe the composition of the databases  
11 used, but we should recognize that the truly largest language model such as an open AI  
12 product, they currently do not disclose what's in their databases. And I'm sure there's  
13 industry representatives here that are interested in just directly using that foundation  
14 model in their products and perhaps that's using those products as a way of achieving  
15 equity where a smaller company can use these costly models. So, I think we should  
16 consider if we're talking about database description, are we talking about like an open  
17 AI or other very, very large language model and even are they able to describe their  
18 databases?

19 Dr. Bhatt: We'll do Dr. Soni, Dr. Elkin. Then to keep me honest, Diana Miller, I'd  
20 like you to come and just respond a little bit thinking about revealing foundation models  
21 in industry and error reporting in industry. And then we'll pick up with the other two.

22 Dr. Soni.

23 Dr. Soni: I'll keep my comment brief because Dr. Radman covered a major part of  
24 it, which is I'm all for having nutrition labels for these models to know what's inside it  
25 and what we should expect from it. But the ingredients, if those ingredients are coming



1 from repurposing of foundation models, we also need to be practical about what is an  
2 onerous task that's unachievable versus what's for the safety. Recognizing that a lot of  
3 these tools are already starting to be used.

4 Dr. Bhatt: Dr. Elkin.

5 Dr. Elkin: So, you guys really made me think. And one of the things is that  
6 obviously these models can do a lot more than they're being tested to do. That would  
7 come before the FDA and I worry that there's variability in the accuracy of those  
8 ancillary tasks that have not been formally tested. And I think that this body, at some  
9 point, should discuss whether they would ban off-label use of the AI or as an  
10 alternative, some kind of postmarketing surveillance of off-label uses that would  
11 continue to monitor closely what they're doing because there is the chance of doing  
12 harm in situations like that. So, you guys are really bright. You're making me think a lot  
13 here and I think those are things that this Panel should at some point consider. The other  
14 piece of information that you'd have to gather from people premarket is any ways that  
15 they've used their technology that have failed miserably, crashed and burned, that  
16 definitely are unsafe. And I would like to see us have them report those upfront before  
17 approval so the indication can be accurate in terms of what it's used for.

18 Dr. Bhatt: Thank you. I'm going to ask everybody to look at the current question  
19 one more time. What information should be included as part of a device's description or  
20 characterization? For example, human in or not in the loop, device intention, whether  
21 it's a recall or new recommendations and safety and effectiveness. If you wouldn't  
22 mind, Ms. Shick, would you read part 1.b? Because I think we might be traveling  
23 between the two.

24 Ms. Shick: Thank you, Dr. Bhatt. So, we have 1.b up on the screen. And that  
25 question is what evidence specific to generative AI-enabled devices should the FDA

1 consider during premarket evaluation regarding performance evaluation and  
2 characteristics of the training data during the total product lifecycle to understand if a  
3 device is safe and effective?

4 Dr. Bhatt: Thank you. So, I think we've already begun answering this as well. Let  
5 me tell you where we are and then we will go on to more comments. I have broken us  
6 down into intended use, use case of the device with generative AI as well as indented  
7 population. We've also added the setting in which the AI is going to be used, whether  
8 that is a decentralized model or central, and primary versus secondary use. After use, we  
9 have the data itself—size, type, number of examples, standardized a pro forma datasheet  
10 perhaps to help us track our data.—We then have the model itself, the definition of  
11 which model was trained, the properties of said model, the demographics of the training  
12 population, where was it trained in terms of location in which it was trained, and does it  
13 include both a diverse dataset but also diverse location of care from where the data  
14 came? Does the model create bounded or unbounded results? And what are those upper  
15 and lower bounds? Requirements to leverage the device, logistic or training, and then  
16 cybersecurity and privacy precautions. I have those written out in longer sentences. I'm  
17 saving us time. We then talked about uncertainty estimations. Where your data or model  
18 is most uncertain, hallucination rates, sensitivity and specificity of the device for  
19 particular tasks or functionalities. We then addressed “dynamic or not,” “autonomous or  
20 not,” previous failures of the technology. That's where we are to date. Dr. Miller. I'm  
21 sorry, Ms. Miller.

22 Ms. Miller: Yes. Hello. You can hear me right? So, Diana Miller from the industry  
23 representative here. So, listening to all the comments from all of you, something comes  
24 to mind that the existing framework for approval devices from the FDA think covers a  
25 lot of these things. And I think a risk-based approach is feasible for this. And risk versus

1 benefit analysis should be developed when we submit this kind of medical devices and  
2 the risk mitigations apply to the intended use and the context of use are very important.  
3 So, it's very-- I wouldn't call it difficult, but I don't want to speculate about open-ended  
4 use cases. I think each use case has its value and each of these devices that deploy this  
5 technology, they should be reviewed in the context of use and in intended use. And with  
6 the risk/benefit in mind. And incorporating harm framework, like estimating harm and  
7 quantifying and special controls to address a harm is a framework that we already use  
8 today and we can leverage and apply that.

9         And going back to the risk/benefit, the comment that somebody made earlier  
10 that physicians are using this today, I just want to point out that we need to look forward  
11 and to be flexible in this new framework we developed to allow adoption because  
12 otherwise people will use it anyway and then the harm is worse. So, I'm not sure if I am  
13 clear enough, but since we don't have this kind of technology in medical devices today,  
14 and it's a technology that's coming to us, we should be a little flexible in deploying this  
15 framework at the beginning so we can all learn and climb in the same time. And the  
16 current framework with the risk assessments, it can work and we have cybersecurity  
17 assessments and all of those things that you talked about in place already.

18 Dr. Bhatt:     Okay, thank you so much. So, two quick notes. Our question number two  
19 that we'll discuss later today will be on risk management. So, we'll come back to some  
20 of those thoughts there. As we go through the next couple comments, I want us to focus  
21 again on specifically generative AI as we give our answers for the sake of transcription  
22 in this discussion. And so a lot of our answers are becoming slightly more broad and  
23 moving towards AI in general. And we really need to help the FDA think specifically  
24 about generative AI. So, if you could please orient your comments towards that for the  
25 sake of transcription, that would be very helpful. Thank you.

1 Dr. Kukafka: Dr. Kukafka. So, I just wanted to push a little bit more on hallucinations  
2 because that was mentioned. And in the context of adverse events, which might-- I think  
3 it more traditionally are fairly transparent. I'm wondering if there's a need for defining  
4 how hallucinations are detected because they can be much less transparent.

5 Dr. Bhatt: Okay, thank you very much. Back to Dr. Stanley.

6 Dr. Stanley: Yes. This is Laura Stanley in hopes of helping to narrow in on this need.  
7 I don't think we've spoken a whole lot about the user interface design. I think we need  
8 to think about what that is and bringing in some fundamental human factor needs, and  
9 that surrounds ourselves around back to the user-trust safety. So, what does that user  
10 interface look like and what can it include to help that explainability need, that  
11 transparency need? Does the interface show regions that help explain the output? Giving  
12 that feedback to that user throughout will help gain trust? And ultimately, hopefully will  
13 gain trust and lead to safer outcomes. So just thinking about really what the interface  
14 needs to include. There's a lot of research in, more recently, in explainable AI, human  
15 AI that we could perhaps tap into.

16 Dr. Bhatt: Okay. Yes, exactly. So, Dr. Botsis, do you have comments on B?  
17 Because I feel like we've started to drift towards C. Okay.

18 Dr. Botsis: Yeah. I want to specifically refer to this sub-question here and maybe  
19 take a step back because we are talking about training here. And traditionally, training  
20 really requires some annotated and labeled datasets. And for GenAI, this process may  
21 be a little bit different. So, for example, when we try to train an LLM by doing some  
22 prompt engineering, there are certain steps and we have to understand what doesn't  
23 really go well in all those steps. We need more information about that and why the  
24 developers made the selection a specific selection to improve the performance through

1 prompt engineering. And I think this is an important part of the process. And what is the  
2 type of the training data that is used through that process?

3 Dr. Bhatt: Can I push on that just a little bit more? For the sake of this specific  
4 question then, what kind of evidence would get at the comment that you just made?  
5 What could we ask device manufacturers with GenAI to give us to help better  
6 understand that? And if you don't know now, we can come back to you.

7 Dr. Botsis: Like a suggestion that comes in mind right away is about the correction  
8 measures they took to improve the performance of the GenAI.

9 Dr. Bhatt: Thank you so much. Next we have--

10 Dr. Botsis: Oh, there's something else.

11 Dr. Bhatt: Dr. Soni.

12 Dr. Soni: Specifically commenting on the effectiveness part of GenAI. I think a lot  
13 of the predominating use cases for GenAI is to help improve the day-to-day workflow  
14 for providers in the setting where there is a mass exodus of healthcare providers. There  
15 is a lot of momentum and opportunity to use generative AI tools to improve their day-  
16 to-day workflow. So, one of the metrics need to be provider facing and how does that  
17 impact providers satisfaction or providers work to help make the decisions and make the  
18 clinical care more efficient and traditionally that has been seen as a preference or as a  
19 soft outcome, but maybe for generative AI, that may be one of the more opportunities to  
20 centralize on provider centered outcomes.

21 Dr. Bhatt: Okay. I'm going to take that comment. I'm just going to say something  
22 earlier today we also talked about the importance of having individuals who do a lot of  
23 kind of information science work and we do have the tap program through the FDA that  
24 puts people in touch with the right people to understand what you're doing with your  
25 information. I wonder if this is one of those upstream type programs that may be less of

1 a requirement and more of an opportunity. So, I'm just saying that to myself for the  
2 record. So, we have it for the transcript later.

3 Let's do Dr. Radman and then Dr. Rariy, and then let's look at C and see what  
4 we've covered for C.

5 Dr. Radman: Thomas Radman. I have a few comments and I'm happy to get cut off if I  
6 go too long or off topic, but someone mentioned restricting off-label use. The one thing  
7 that brings to mind is something I've heard frequently is that FDA does not regulate the  
8 practice of medicine. So, with that understanding, we should recognize that GenAI  
9 products have different risks than pre-GenAI products. For example, much more  
10 versatility to be used off-label than just a yes-no diagnostic that has a fixed input that  
11 we're used to how it's regulated. And then we've talked about breach of privacy,  
12 regurgitating someone's personal information, bias, misuse is part of the off-label  
13 hallucinations. So, if we agree that the risks of GenAI is higher than pre-generative AI  
14 products, now we're talking about FDA regulation where there's a substantial  
15 equivalence pathway. So, I think we maybe as a Panel should determine should the  
16 benefit of increase of accuracy have to be better than a pre-generative AI product if we  
17 are accepting that the risk is higher. It's not just the accuracy, if someone comes in and  
18 the accuracy is the same, but they're using GenAI, are we going to call that equivalent  
19 now accepting these risks. So that's one point.

20 And also, what I'm understanding from our panelists from the research  
21 community, that it's still kind of early days in the evaluation of these products and in  
22 terms of many aspects, determining bias, detecting hallucinations and just evaluation in  
23 general. I don't think that being early days should prevent someone from attempting to  
24 come on to impede progress, but we should recognize that that's also an increased risk  
25 that may justify at least a more rigorous evaluation process than a pre-generative AI

1 product would have. Then in terms of particularly speaking to evaluation, so because  
2 these devices have these nuances such as the way you engineer your prompt, how you  
3 ask generative AI, what you want could change your output and that could require  
4 different levels of training and also the outputs can vary as well because it can be  
5 conversational, if we're talking about a chatbot type of function that leads to broad  
6 interpretability by different users. So, I think the evaluation should be at an end-user  
7 level and there should be a stated diversity of the user base in terms of experience,  
8 demographics, like you name it, it should be very diverse and a sufficient number.

9 Dr. Bhatt: Great. Can we put up C? Because I think we've covered A and B well,  
10 and maybe Ms. Shick once you read C, I will let the group know what I think might  
11 answer, see what they've already said, we can add anything else to that segment. One  
12 thing I do want to mention is we're coming up against this quite a bit, and this is what is  
13 so challenging about generative AI. There are specifically the data, the model, the  
14 outputs that we are talking about, and those are all very kind of logistic parts that we  
15 can offer advice to the FDA. There is then this other side of the usability, who the user  
16 is, the flexibility in its use, the indications for use, which may not necessarily fall under  
17 our ability to enforce or state and yet is important to all of us. So, I am thinking about  
18 how we best convey the parts that may be easier to ask for require and the parts that are  
19 perhaps more suggestive about the practice of medicine, but very important to our  
20 group. And so maybe when you are telling me what you're thinking and I'm typing, feel  
21 free to say, here are some concrete logistic things that we have to have included and  
22 here are some concepts that are worrying me about the implementation of this in  
23 practice and that will just help me in my typing because as you know, I prefer to hand  
24 write. So, this typing thing is it's going okay. Ms. Shick, can we talk about C please?

1 Ms. Shick: Thank you, Dr. Bhatt. What new and unique risks related to usability  
2 may be introduced by generative AI compared to non-generative AI? What, if any,  
3 specific information relevant to healthcare professionals, patients and caregivers is  
4 needed to be conveyed to help improve transparency and/or to control these risks?

5 Dr. Bhatt: Wonderful. You see why I went here now? I think we've gone there. So,  
6 we have heard earlier from Dr. Stanley that user interface needs to help explain the  
7 output and give feedback throughout that product used to gain trust. We have also heard  
8 that it is important for individuals to recognize what the on-label indications are,  
9 therefore that there are off-label reasons that are not actually approved and be very clear  
10 on what those are. We've also spent time talking about training and the importance of  
11 understanding how things like changing your system prompt may actually change how  
12 you get an answer and that goes back to training recommendations that we may make to  
13 help improve both transparency, but help control risks. Are there other things that are  
14 new and unique risks related to usability that we have not yet talked about with regards  
15 specifically to generative AI that we should include in our response to this section? Dr.  
16 Rariy, I believe you were up next.

17 Dr. Rariy: Thank you. And I think this falls more in the prior around evidence, but  
18 one thing that I wanted us to consider is an opportunity for the device company to  
19 benchmark against existing AI devices or AI-enabled models that already exist and to  
20 look at if there's value added or if they're at parity or just some sort of looking around  
21 at the general device atmosphere. I think that would be helpful for information to come  
22 back with to the FDA on. The other thought is as we're looking at what is required in  
23 premarket-performance evaluation, I think it's always helpful to have a forward-  
24 thinking plan and so to include what could be your postmarketing strategy or your  
25 postmarketing plan, I think would be important specifically for generative AI where we



1 know that there could be subject to drift, there could be subject to other areas in  
2 information that is gathered from the real world. So I think that should be included in  
3 the premarketing evaluation.

4 Dr. Bhatt: Can I ask a quick clarifying question? An opportunity for the device to be  
5 benchmarked against other models, other devices with the same intended outcome and  
6 use that may include or may not include AI, that includes machine learning or that  
7 specifically includes GenAI, like theirs or all of the above.

8 Dr. Rariy: Probably all of the above, including if they're first to market, they may  
9 not have a specific comparative device, but perhaps they have the AI-generated models  
10 which they're currently leveraging and so there's that comparison, but at least an  
11 opportunity to think through what that comparison looks like.

12 Dr. Bhatt: Great. Dr. Jackson, and just let me know if you're on C or if you're  
13 thinking about A and B still that is okay for something additive.

14 Dr. Jackson: Thank you. I did have one additive for B and then I do have something  
15 for C. For B, and I may have missed this. I think it would be helpful to have information  
16 on the performance evaluation specific to different demographics and settings so they  
17 will have the overall performance, but how does that differ for the intended use, right?  
18 There's a difference in gender if there's a difference in if it's a rural clinic or something  
19 else. For example, foresee I think a new and unique risk, not just that the output cannot  
20 be reproduced, but as what was shared earlier that the reasoning will not, so they might  
21 still get the same answer, but the reasoning will not. The reasoning will not be  
22 reproduced and that should be clear that how they got to that answer may not be able to  
23 be regenerated. For what needs to be conveyed to help improve transparency, who is the  
24 intended user of the model? I think this goes to the training. Will it be techs? Will it be  
25 physicians, nurses, so that it's very clear who will need any training.

1 Dr. Bhatt: Dr. Shah.

2 Dr. Shah: Pratik Shah. I'm kind of B and C additive together kind of segueing. So,  
3 one is like it's very hard to have one regulatory framework for all generative models, so  
4 one recommendation for risk and usability would be for generative AI models that don't  
5 operate on images versus generative AI models that operate only on text, right? Because  
6 generative AI models that generate images do a very different risk and usability  
7 assessment than something that's just a chatbot that generates new text and then the  
8 hybrid models that kind of go sit in between both of them that use an image that's  
9 generated and then generate a text. So that's one usability aspect that's very important  
10 to point out. I'll give you a quick example, you could have a locked or a generative AI  
11 model that's not using chat GPT, but it's again a generative adversarial neural network  
12 that generates new images at the point of care. It's a generative model, but not a large  
13 language model or chat GPT that's still generated AI right? People still use that. So not  
14 everything is chat GPT or large language model. So that's one very important  
15 distinguishing factor to make for the regulatory framework to understand what we are  
16 reviewing and where it comes from and what the intended outcome is.

17 The C part which is related to it is one suggestion would be to actually include  
18 evidence to patients that this result was generated by an AI model eventually. I know  
19 that's a complex topic, but I feel like in the sake of transparency, if there is downstream  
20 clinical decision making happening from an AI-generated prompt, it may be good to  
21 label the initial AI-generated result that was generated. So, if it goes into the healthcare  
22 system and we lose track of it, we know where the AI-decision making was made. So,  
23 we can reverse engineer the path up and find out the node where the bifurcation  
24 happened. And that goes into the postmarket surveillance part and the trust and

1 transparency because if we don't tell people where AI intersected in their healthcare, we  
2 have no mechanism to track it back and fix it.

3 Dr. Bhatt: Dr. Elkin.

4 Dr. Elkin: Thank you. Peter Elkin, University of Buffalo. I remembered her  
5 instructions. So, I do have one for B, but Dr Botsis has made me think of this that when  
6 you're looking at getting a metric of whether the data was overtrained for a particular  
7 purpose, the number of epochs that it was trained on can be helpful. Normally,  
8 [Indiscernible - 05:48:26 (experts)] recommend you do about three to five epochs. I've  
9 started to see people who use 20 or a 100 epochs in order to train and in that case there's  
10 going to be high specificity for a particular task and a great fall off in functionality when  
11 the tasks differ from the one that's been overtrained. Then for the usability thing, I got  
12 to tell you a little story in order for my comment to make sense. In the years around the  
13 2000, a bunch of us who were human factors engineering people went to NIST and we  
14 got them to recommend to you guys that every device be usability tested. And sure  
15 enough, the guidance came out from FDA that all devices need to be usability tested  
16 and they stopped there, no implementation guidance on what that meant.

17 So, most of the things were just, we showed it to some customers. There was no  
18 formal think-aloud method, Jakob Nielsen based a usability evaluation with a standard  
19 protocol in an unbiased way that took a number of participants and were watched  
20 unbiasedly for this. And there is a real protocol around usability testing that can be  
21 employed. I really urge the Committee to think about the last mile of this in putting in  
22 place a requirement for real usability testing before approval. Thank you.

23 Dr. Bhatt: Dr. Botsis.

24 Dr. Botsis: I'm not exactly sure whether this comment is so much related to one C. I  
25 think it is really important to appreciate the types of our output. So, for example, we

1 have a binary classification problem and we want to classify a case as A or B. In a non-  
2 generative AI approach, we would probably have a misclassification for a case A as B,  
3 but in the case of GenAI, you may have a third label unexpectedly. So, you really need  
4 to know that. The second example, when you want to summarize some long text, it is  
5 not a matter of whether you have a summary that may or may not fully represent the  
6 original text. It is, as we all know, like an output hallucination. You may have a  
7 summary that is completely different from what is described in the original text. I think  
8 that those matters are critical.

9 Dr. Maddox: Tom Maddox, WashU in St. Louis, I'm going to push back a little bit on  
10 Elkin because I'm thinking about this really important distinction that the FDA has  
11 about regulating devices but not regulating the practice of medicine. When we think  
12 about drugs and typical devices, the FDA appropriately stops at that point of how they  
13 approve it for a use case, they ensure the device your drug is safe and effective and then  
14 they give it to our clinician community to then use it as they see fit, recognizing that  
15 some things are off-label, some things are on-label. I'm trying to think of a similar  
16 framework here. I think what's important maybe is to stop at dictating use, I think that's  
17 sort of a fool's errand anyway, just given how heterogeneous this will be. But I actually  
18 think it's actually not appropriate for the FDA to be in that space just given how the  
19 statutes are set up.

20 But what does occur to me is one unique thing about generative AI is it very  
21 much interacts with the user in a way that's very different and changes the output in a  
22 way different than a pharmaceutical or a traditional device ever has been. So, it's made  
23 me think that interface, I'll call it the user interface between the GenAI device and its  
24 ultimate use might need some clarification by the FDA. The two things that I've heard,  
25 but just to outline them, are that we need to think about the user information that is

1 provided in a helpful way, I think Dr. Stanley's been very thoughtful on this,  
2 explainable AI, confidence intervals, etcetera. There may be some specifications there.  
3 And then, also user training, which Chevon and others have said, and we would need to  
4 think probably about that. I'll just leave it with this.

5 I think there's an analogy here from my world and NAMI's world of cardiology,  
6 and there is an antiarrhythmic medicine called Dofetilide, which has a lot of risks and  
7 because of that, the FDA has said, I must as a clinician prescribing, go through a  
8 prescribed training on how to use this med and I can't prescribe it without it. The  
9 industry I forget who makes the medicine, but they actually run the training so I had to  
10 go through a training before I could prescribe that med. So maybe there's a similar  
11 certification pathway to train me to interpret the user information to then equip me to  
12 use this in practice. So anyway, just a thought.

13 Dr. Bhatt: Perfect. I have Radman, then Clarkson, and then Shah. However, I want  
14 to give great credit to the FDA people who put these questions together. They appear to  
15 know what we're thinking about next before we think about it next. So, let's go ahead  
16 and show what question D or part D is, and then feel free to just tell me A, B, C, or D as  
17 you give your comments so that I can categorize. Ms. Schick.

18 Ms. Shick: Are there prospective performance metrics that are particularly suited or  
19 most informative for these technologies, given their complexity? What kinds of  
20 performance metrics are needed for multimodal systems for example, text image models  
21 where either inputs, outputs, or both could be multimodal? Performance metrics will  
22 typically vary with device intended use.

23 Some examples of known metrics to support discussion may be modality-  
24 specific such as generative text (perplexity, quantitative comparison to reference text)  
25 for generative images (Fréchet Inception Distance (FID), Structural Similarity Index,

1 Measure (SSIM) or for generative audio (Log-Spectral Distance, Perceptual Evaluation  
2 of Speech Quality) or may be functionally-based such as frequency and types of errors  
3 made by the generative AI-enabled device.

4 Dr. Bhatt: Okay, so we will go to-- You're still up, right? Yeah, Dr. Radman. And  
5 then just a reminder, A and B were particularly about evidence specific to GenAI  
6 devices premarket, what they should have shown us C is new and unique risks related to  
7 usability with these devices. And now D, really thinking about performance metrics that  
8 we would like to see.

9 Dr. Radman: Thank you. And giving credit to the FDA, I'll have to give them credit  
10 for taking this very robust and kind of varied conversation and disseminating into  
11 guidances and final regulations. So thank you. Thank you, Tom, you set up my points  
12 perfectly on labeling. I think that the labeling should be very clear that it should be in  
13 your face, this output, this chat advice. Especially when you're considering patients or  
14 lay users, has been generated by a chatbot, artificial intelligence chatbot. In terms of  
15 some of the training points that Tom was getting at. So the repeatability or the  
16 temperature setting of the device and its propensity to give different outputs for the  
17 same input is important part of the training for a user, but it should also be clear the  
18 input, the premarket information to the FDA should also have some-- And this is  
19 already included for traditional algorithm review, like a repeatability and reproducibility  
20 that really hasn't been spoken about yet. But how consistent is it with the same inputs  
21 on a large scale.

22 We should recognize that some of the latest literature shows that repeating the  
23 same inputs is now a mechanism to recognize what part of a response is hallucination.  
24 Usually the varied part if you're repeating, and then the part that stays constant is the  
25 truth of the response. There's also literature that shows that repeating prompts actually

1 increases accuracy so if you take a generative AI for a task and you have it answer  
2 something once versus answering something three times, it does better in the three-time  
3 case. So that should be part of the premarket information, how it's going to be used,  
4 education about repeatability, what's the temperature, what are the guardrails, what's  
5 the complexity of the model that's directly linked to its propensity to hallucinate. So the  
6 larger the model, the more a premarket review should be looking for concerns on  
7 hallucinations.

8           So that's maybe D and getting it to-- Sorry, that was C and maybe getting it to  
9 B. I also wanted to mention, I think this is either B or D, we talked about broad-use case  
10 generative AI or unbounded, and Dr. Mahmood gave the example earlier that he had  
11 this ability to classify now 108 cancers, that sounds so amazing, but we should  
12 recognize that some of those cancers are very rare cancers. So if you consider a device  
13 like that and you just have a broad indication of it can detect these 108 cancers, but one  
14 of those cancers occurred 1 in a 1000 just making up a number, we didn't really validate  
15 that one cancer. So we should, even though it is a broad use, we should be breaking it  
16 down and particularly paying attention to the rarest cases, those are probably the hardest  
17 to classify.

18           That same point goes for users when you have broad demographics using a  
19 device, whether that be clinicians or patients. I'll give the example of language, the  
20 rarest language used in the population has been documented to perform the worst in  
21 generative AI use. So I think a focus on the rarest class is sort of like this idea of  
22 asymptotically arguments. If you're thinking a philosophical construct to prove  
23 something, take the hardest, most extreme case and that would be in this general of AI  
24 analogy, the rarest classification or user group. Then I just wanted to mention there's a  
25 publication many of us might be aware of it, co-authored by actually someone who used

1 to be at the FDA Digital Health Center for Excellence, Bakul Patel, where he talks  
2 about regulating generative AI as doctors. It's unbounded, doctors are also unbounded.  
3 So how do we regulate them? They have to take a test every year and there's challenges  
4 with giving something that's ingesting everything in the web, constantly on testing it.  
5 But I think that's just a really interesting idea for this Panel to consider, or at least I'm  
6 just stating that idea for the record.

7 Dr. Bhatt: Dr. Clarkson, and then Dr. Shah.

8 Dr. Clarkson: Hi Melissa Clarkson, consumer representative. I want to follow up on Dr.  
9 Elkins comment about, and I'm back on C, I believe it was talking about risks. It  
10 mentioned patients and caregivers, and I think that was a reminder that we've mostly  
11 been talking about in-clinical use of generative AI, but this is going to end up in homes  
12 or it might be assigned in a clinic and then you bring it home. I mean, this is going to be  
13 used by people who do not have clinical training and I think it is at the fuzzy boundary  
14 of-- Is this really within the FDA's boundaries? But studies to understand how non-  
15 clinical people seeing that generative-AI output understand it? How they react to it?  
16 How do they interpret its authority? Because it could be like if I go to a very prestigious  
17 medical center, I might say, oh, this physician has much more authority than my little  
18 hometown doc. How do they perceive the authority of generative AI? Especially, if it's  
19 using my name and it seems very personalized and it knows all about me. Do I  
20 understand that it's even more authoritative as my doctor I haven't seen in three years  
21 because it's too expensive. I really have no other access to healthcare. So I think some  
22 sort of study going in that direction of how do non-clinician viewing this output  
23 understand the output.

24 Dr. Bhatt: We'll go to Dr. Shah and Dr. Soni. I just want to make one comment.

25 Thank you Dr. Clarkson for that. There was a day about a year ago where I called a



1 chatbot. Really nice. I was like, wow, that's a really nice chatbot. And then I just  
2 thought about what I did and so I love the comments you just made. Okay, over to Dr.  
3 Shah and then Dr. Soni.

4 Dr. Shah: Okay, Pratik Shah. I'm answering specifically 1D and great technical  
5 questions. I really like it. Thank you for the FDA team for putting it up. So I think this  
6 is a specific example where a generative model has created a new image that potentially  
7 did not exist and you're trying to characterize the quality of that image, and you're  
8 using Pearson correlation coefficient, SSIM, etcetera. And we have dealt with this  
9 problem. And I think my recommendation in general for the regulatory science  
10 perspective of this is that we need to do all of it. We need to have a really good  
11 computer vision validation metrics, which has no human involved. So it could be pixel  
12 by pixel quality assessment of the generated images by Pearson correlation coefficient,  
13 SSIM, PSNR and once the image is signed off by those computational metrics, the other  
14 aspect to consider would be would a clinician sign off on that image?

15 And we have done that work because we have to do a common-sense approach  
16 that once we have the generative AI model image validated, we then have to go to the  
17 clinician and say, okay, will you sign off on this image? And what would be the  
18 standard of care you would give using this image? Right? So that's the second sub-part,  
19 which really needs to be considered. The errors that will compound could come from  
20 the computer vision problem where the Pearson correlation coefficient or SSIM wasn't  
21 good enough, and that compounded in the clinical decision making, or it could be the  
22 clinical decision making was independently errored in its own errors, which had nothing  
23 to do with the image, right? Because clinicians sometimes don't agree on rare subtypes.  
24 So that's the two part, and then the very important part is how do you communicate it to  
25 people and understanding the multimodal part.

1           So I think for the multimodal part, the communication and the interpretation of  
2 this model is very important. And I'll give you a very quick example: People confuse  
3 interpretability with explainability in machine learning models and AI models.  
4 Explainability means saliency maps or some sort of grad-CAM where you can show  
5 what the model was looking at when it got activated, but it doesn't have much to do  
6 with why it chose that activation path to make the decision. So I think it's very  
7 important to kind of communicate that explainability of the model may be unlinked to  
8 interpretability and saliency maps or activation maps that generative models use may  
9 not have much to do with why the decision. So the solution to that problem is  
10 potentially looking at the saliency activation maps and then teasing out specific features  
11 from those saliency activation maps, which actually led to the deterministic decision  
12 made by the model. Because if we don't have that link to the output, and just an  
13 activation, we are confounded in the technical side.

14 Dr. Bhatt:     Okay. We are slowly coming up on time, so we have Dr. Soni and Dr.  
15 Maddox. To the FDA team, I will ask a question which is we are going in order of the  
16 agenda given today. I'm curious as to whether the break can be given before my  
17 summary to Mr. Tazbaz. Is that allowed or possible? Great, thank you very much. I will  
18 say probably for the sake of the FDA, nobody can talk to me during the break. I will just  
19 sit here by myself with my computer so that we don't have anything that was not  
20 recorded for public record.

21 Mr. Tazbaz:    So, the question, I guess, is are we going to have an opportunity to ask  
22 clarifying questions?

23 Dr. Bhatt:     Yes, clarifying questions first would be great. So, I will do a version of a  
24 summary, clarifying questions and then a final summary.

25 Mr. Tazbaz:    Okay, wonderful. Thank you.

1 Dr. Bhatt: To Dr. Soni please

2 Dr. Soni: A comment on 1D. Those metrics are really important and I consider  
3 them as safety metrics because that helps understand and continue to prospectively  
4 evaluate whether or not the models continue to operate as intended. Very challenging  
5 with generative AI because there is variable input and potentially variable output or by  
6 design variable output. But I do also want to highlight that, and this may vary based on  
7 intended-use case, but there may be a need to not lose track of some of the traditional  
8 metrics which we have learned to understand a lot of devices and decision-making tools  
9 which is sensitivity, specificity, positive predictive value, negative predictive value,  
10 precision and some of those metrics still need to be part of the discussion because if  
11 metrics are only understandable by very highly-specialized workforce, those become  
12 challenging. Again, this may become a concept part, but I think the safety and  
13 effectiveness we need to maybe think about them as separate metrics.

14 Dr. Bhatt: Thank you. We'll do Dr. Maddox then Dr. Elkin.

15 Dr. Maddox: So to continue my role as a contrarian, Dr. Shah, you talked about  
16 evaluating the use in this case of an image example. Why is that not getting into the  
17 practice of medicine?

18 Dr. Shah: Sorry, repeat your question again.

19 Dr. Maddox: So we can't regulate the use of the practice of medicine.

20 Dr. Shah: Correct.

21 Dr. Maddox: And if I understood you right, maybe I didn't, was that you said one  
22 thing we need to think about is once an image is generated that we then need to evaluate  
23 how the clinician is using that and that struck me as a little bit out of bounds here, but  
24 maybe I misunderstood.

1 Dr. Shah: Right, so I think different settings require different oversight, I assume.  
2 So my specific comment was that if you have an image that was generated to train a  
3 deep-learning model and the image was then used by clinicians, validated by clinicians  
4 to help the deep-learning model perform better or performance metrics, that's a valid  
5 oversight of the model. That was my point.

6 Dr. Maddox: Thank you.

7 Dr. Bhatt: Dr. Elkin, next.

8 Dr. Elkin: Thank you. I've been thinking a lot about all the comments and there's  
9 some really wonderful things to build off from Dr. Maddox's comment. There's a  
10 difference between drug companies and the companies that are marketing the AI in  
11 terms of their usual practices. I think that in order to above all else do no harm, we  
12 really want them to have good guidance from the FDA about how they can market these  
13 things. I think it falls short of regulating practice, but it does give guidance to the people  
14 so that it's not a slippery slope to this model's approved for one thing, but it can do  
15 everything and because these are very broad and interesting, it made me think of an  
16 example that might help your case, Dr. Shah, which is let's say you have an MRI and  
17 you take from the MRI and generate an image, an augmented reality for the  
18 neurosurgeon taking out a brain tumor. You could see that as a generative image that  
19 would come out of one of these things that would've real effect on patients. Just  
20 because it was proven in that circumstance doesn't mean it's going to help a cardiac  
21 surgeon find all the calcifications or a colorectal surgeon find the extent of the tumor,  
22 right? So we really want to give good guidance to how these things can be marketed.

23 Dr. Shah: So I think to quickly, so that's why I suggested labeling where it came  
24 from and what it was intended to do and if you do something else with it, then you are

1 taking the risk, not the FDA, not the people who made the model because you decided  
2 to use a drug for an off-label recommendation.

3 Dr. Bhatt: Okay, with that lovely spirited discussion, we're going to move for a  
4 minute now to give a short summary and then ask the FDA if they have clarifying  
5 questions for the group. So those who didn't get to speak yet another time, I think  
6 there's still opportunity for discussion. So I'm going to go back to the beginning with A.  
7 When we are looking at this, we are looking at-- I'm going to do the short form, tell me  
8 if they want longer. Intended use that includes use of the device use cases intended  
9 population, where it's going to be used, primary versus secondary use. We're talking  
10 about data, we would like to make sure that the FDA recognizes we need the size, the  
11 type, the number of examples. We need to think about standardized proforma data  
12 sheets that include all of this information and allow for tracking.

13 We do need uncertainty estimations where your data model is most uncertain  
14 and how it was trained and whether it was overtrained. For the model itself, we have the  
15 majority of information that we are providing to the FDA about our opinions. This  
16 includes a definition of training of the model properties, demographics of the  
17 population, and does it include a diverse dataset but also location of care. Importantly,  
18 bounded or unbounded is unique to generative AI and important part of what we would  
19 like people to be able to report in the premarket state. Cybersecurity and privacy  
20 standards as with other AI. However, recognizing that these are models that may change  
21 over time, so will have unique needs that may not otherwise exist when these algorithms  
22 are changed over time. Then for B specifically to generative AI, we do think that there  
23 might be an opportunity for the device to state how it benchmarks against other models  
24 that may have existed previously without generative AI or other generative-AI models  
25 similar to it. We do think that it would be nice for them to include a little bit about

1 postmarketing plan in their premarketing conversation, but we will get to postmarket  
2 surveillance tomorrow. We would like information on the performance in different  
3 populations and settings as well as information on repeatability and reproducibility that  
4 is important and unique to generative AI as if they have information on how different  
5 temperatures offer different responses that would be interesting to those who are going  
6 to either use the device or build off of it.

7 Also important to recognize how narrow the indication for which it was created  
8 is how many epochs were actually used or evaluated. Hallucination rate, we want to  
9 understand sensitivity and specificity of the device for particular tasks and hallucination  
10 rates for those tests and functionalities. And then lastly, we have not had a thorough  
11 conversation yet, but we will about the breakdown of autonomy, which also goes to the  
12 degree to which it affects direct patient care, bringing in what Dr. Elkin just ended with  
13 talking about neurosurgery versus cardiology.

14 Lastly, for C specific to usability, we are thinking about the user interface and  
15 it's explaining the output and giving feedback throughout to gain trust. We are thinking  
16 about the fact that the output and reasoning may not be reproduced as with prior  
17 technologies and therefore that should be something that users are aware of. We  
18 recognize that generative AI models may operate on images versus text versus  
19 multimodal, and we likely have specific metrics that are coming in D depending on  
20 what they are actually operating on.

21 We do want to talk about requirements to leverage the device logistics and  
22 training, and we may recognize that outputs of generative AI may not be what clinicians  
23 are used to of A or B, but may create a novel response, again, necessitating section D,  
24 the metrics for those novel responses. Lastly, we commented on C how we need studies  
25 to understand how non-clinical people understand these outputs because the location of

1 use will also likely influence who is using this for direct healthcare, and that is not  
2 always a super-trained clinician.

3 Lastly, for D four metrics, we agree with the known metrics for safety that you  
4 have suggested here. We do think we need to continue to evaluate models and models  
5 remain accurate without drift over time and we need a metric to be able to do that. Then  
6 finally, we do think that it's important that these companies demonstrate good computer  
7 vision metrics without human input first. Then what does the human input of clinical  
8 evaluation say? And then lastly, how do they clearly communicate those outputs. FDA,  
9 any clarifying questions of what our group has discussed?

10 Mr. Tazbaz: Well, thank you Dr. Bhatt. Give it up. Thank you for such a rich  
11 discussion. I know this is a very complex topic and it is very engaging and we're  
12 actually having a hard time not to engage as we're instructed to do. But we do have a  
13 few quick hit clarifying questions and maybe three additional questions that probably  
14 are not going to be answered here, but it's something that we should definitely take  
15 back and address at a later time. I'll maybe just go through the quick hit questions and  
16 one is we've heard the term dynamic AI models being used a few times and we're  
17 trying to understand what that is mean because we're trying to align on terminology as  
18 well. Can't remember who used that, but there's been the term dynamic AI, and I want  
19 to make an assumption, which I don't really want to do, but is that adaptive AI? Maybe?  
20 Okay.

21 Dr. Maddox: I used it and you're right, that's a more precise term, but I used it initially  
22 and that's a more precise term, but I also would look for validation on that. Is that  
23 consistent with what others are thinking? Okay.

24 Dr. Radman: Yeah, I think one of my first comments was about continuously learning  
25 there and then that brought us into the dynamic.

1 Mr. Tazbaz: Okay, thank you. Thank you. Yeah. Okay, so adaptive artificial  
2 intelligence. Okay, the next one, hopefully Dr. Rariy, you mentioned about  
3 benchmarking against other models. Where should those benchmarks and performance  
4 metrics should be stored? I guess who's responsible for maintaining those benchmarks?

5 Dr. Rariy: I think it's a good question. The only thing I can lean on is where  
6 historically medical device benchmarks have been stored. So perhaps looking at the  
7 literature as an example. That's where I would start, but it might need to adapt over time  
8 because while we've seen quite a bit of literature in the generative AI space that may or  
9 may not be representative of the devices that we would see through the FDA.

10 Mr. Tazbaz: Wonderful, thank you. So here are some questions that are very  
11 generative AI specific that we captured from the discussions and these are really  
12 important topics that we do need to address. And the first one is around hallucination  
13 and drift detection. Dr. Rariy you mentioned about stress testing models and coming up  
14 with methodologies and perhaps tools to be able to stress test. From a regulation  
15 perspective, should the requirement for breaking these models and anyone who's done  
16 software development understands ethical hacking, right? You have ethical hackers in  
17 your staff that tries to break these applications, should that be a requirement for  
18 generative AI where there is a stress testing requirement for regulatory approvals that  
19 allows you to understand the breaking points when something is going to hallucinate  
20 versus drift? I would love to hear some comments on that.

21 Dr. Rariy: I have an additional comment there I suggest. Yes, and perhaps we could  
22 base it based on risk. So one of the comments that I was going to propose is specifically  
23 the metric would be an error-rate analysis. So not only describing the type of error, the  
24 frequency of the error, but the severity of the error as well. So this comes back to the  
25 stress testing requirements can be based on risk of the medical device, risk of harm if



1 you will. Is the risk high enough to cause death or is the risk high enough where there's  
2 misinformation and there's different stress testing recommendations dependent on the  
3 risk.

4 Mr. Tazbaz: Okay, so the requirement for stress testing or the level of stress test needs  
5 to be aligned to the risk associated to the model.

6 Dr. Rariy: That would be my recommendation, but I would certainly want to hear  
7 others.

8 Mr. Tazbaz: Thank you.

9 Dr. Jackson: I would say that the ownness of responsibility should be on the company.  
10 I do think that it does need to happen, the stress testing and I think about, I know a lot of  
11 our presentations today we're talking about radiology and different things, but coming  
12 from behavioral medicine and mental health perspective, we've seen that with chatbots.  
13 I think sometimes we conflate generative AI and we forget generative AI in healthcare  
14 needs to be vastly different than just thinking about how we produce GenAI in general.  
15 So I think companies do need to show along the same line of responsibility for  
16 clinicians of do no harm, how far it would take or how long it would take for  
17 hallucination so that people can mitigate the risk of harm to patients that they're  
18 referring to or that they're using the GenAI for.

19 Mr. Tazbaz: Okay. And which is going back to Dr Bhatt's question around the  
20 industry's responsibility and some of these performance management.

21 Dr. Elkin: So I think it's absolutely necessary and I'll give you a couple of good  
22 examples of why and then I've got another recommendation aligned with that. So as  
23 you're debugging these, as you put them up and debug them, there's some  
24 unpredictability in the results that you get and the wider the audience that you can stress  
25 test it on, the better you'll be able to adjust the parameterization of the models because

1 in addition to the training of the models, you get to parameterize them if you really  
2 want, I can do a teaching session on it, but the idea is that there's a good lot of choices  
3 that you can make. They have real impact in the confabulation rates and so forth, but the  
4 more people you get to try it and you tell them to try to break it, the more likely you  
5 find a good share of these upfront before you would release it in any kind of way.

6 We do this routinely, and I tell you in my department, people have real fun  
7 trying to break the models I build because they love to tell me if they did it because  
8 they're real smart and it shows they're smart and so they have fun. It's actually a fun  
9 activity for the companies. I'm sure they would enjoy it and it would help them to hone  
10 the parameters that they put in. The other thing you didn't say, but I think might be  
11 aligned with that is maybe people should report Shapley curves for the models. These  
12 are deep-learning models and they're not explainable completely, but we can back  
13 propagate through it with Shapley curves and try to get a sense of what the most  
14 important features are. I know a lot of people care a lot about explainable AI, I would  
15 love to see that. Of course I was just in San Francisco and the Waymo were there, self-  
16 driving cars and nobody who got into one of those cabs actually was able to ask ahead  
17 of time how the model was trained in the car that allowed them to drive through San  
18 Francisco safely.

19 They were happy that they had a good safety profile and they gave their credit  
20 card and they took them where they wanted to go and they got in and out and even  
21 people I know are very sensible and careful people. Matter of fact, you guys know  
22 David Bates, he was one of the people who did this. He got in one of these things and  
23 he's a safety person and he couldn't know how these things were trained, but he had a  
24 good assurance that safety profile was good and he was willing to trust that. That's sort

1 of interesting for someone who's so well-known in patient safety, right. Anyway,  
2 thought I'd share.

3 Mr. Tazbaz: And by the way, it is a good example of using different industries who  
4 are applying autonomous use to different industry use cases and can we learn from the  
5 guardrails that they've established and something that we have been asking. Just to be  
6 respectful of time, this one may not necessarily be something that we're going to answer  
7 today, but there was a discussion on third party foundational model usage in  
8 applications and from our point of view, we do wonder about what information should  
9 be actually available to us during submissions if a third-party foundational model was  
10 used as the underlying model for applications to be applied to a specific kind of  
11 healthcare use case. We've been calling it either multi-tiered or multi-layered  
12 architecture, so that's in an area that we are very interested in kind of pursuing  
13 additional information. That said, I think the last few comments kind of goes back into  
14 or ties directly into the question that we have around usability.

15 So it's about the interaction with a device, I think is what we're talking about,  
16 the usability. The question for us is how should we thinking about performance  
17 management and the guardrails that should be applied to usability considerations around  
18 how someone is interacting and whether the performance changes based on the  
19 interaction with a specific device.

20 I don't think we talk about this often, but in the early days there were a lot of  
21 actually psychologists looking at what people were doing and part of it was that they  
22 wanted to break the model and so they were trying to see what that break point was and  
23 that's a behavioral pattern in engaging in. So the question now for all of us to consider  
24 here is that in medical devices, how should we be thinking about this thing from an

1 engineering point of view of how the user is interacting and whether performance  
2 changes based on that interaction? Would love to hear any comments on that.

3 Dr. Bhatt: So the group at NYU is doing some really good work in this and entering  
4 healthcare, which is, it's almost like the behavioral, I don't want to call it economics,  
5 but the behavioral manipulation with AI that happens. Right? Which is they're in fact  
6 studying if you are using the device, does it make you better at it? We talked about this  
7 earlier, right? Sometimes AI plus human is worse, but also what is it telling you about  
8 yourself, about your knowledge level, about your confidence in decision making? How  
9 does it affect your ability to care for that patient if something tells you something  
10 different than what you had? I think it's a growing field, but it is a young field and so I  
11 too would be interested in what the group says, but I would also be careful. I think  
12 there'll be a lot more data about this and this field is about to grow and luckily we'll be  
13 together for three years, but I think this field is about to grow.

14 Mr. Tazbaz: And I guess a more concrete question on this one is should that be part of  
15 the evaluation criteria.

16 Dr. Bhatt: As far as performance metrics that may influence outcomes for patients?  
17 It goes back to what Dr. Rariy was saying. I think a lot of our conversation the group  
18 can agree or disagree entails what exactly are you doing for the patient, how much harm  
19 could you possibly cause by doing this thing using generative AI? And I think when we  
20 base it on that, are we triaging? Are we giving an exact diagnosis, are we affecting the  
21 operating room? I think there are different categories and based on that there is probably  
22 a different level of stress testing. I would say how a human interacts with it and  
23 changes their clinical care probably falls in that realm. If you're going to change  
24 something that's not going to really affect long-term care or short-term acute care, that  
25 might be different than acute difference in what you would do based on interaction, I'd

1 almost maybe think about including that in stress testing. Dr. Maddox, you can be  
2 contrary if you want.

3 Dr. Maddox: No, I agree with that. Shoot, I wish I could. I think that's right and it also  
4 makes me think maybe this is a role for a calibrating temperature to risk and so the more  
5 risky it is, maybe it's more deterministic, which in turn will lead it less dependent on the  
6 user interaction with it about the output. So this may be something that you evaluate on  
7 a spectrum along with how you proposed it on risk.

8 Dr. Radman: This is a good time to bring up a point that we haven't discussed is that  
9 the support provided by the provider of the device, the GenAI, is part of the usability.  
10 So I'm just thinking of a study where they were using online smoking cessation, which  
11 is not GenAI but it's a digital health product and they had one group just using it. One  
12 group was given broadband access and one group was given broadband access and  
13 access to a digital health literacy coach. The group with the coaching had almost double  
14 the treatment response. So just like simple customer support provided by a company  
15 could really increase the efficacy of a device and so that's just a concept we might want  
16 to consider too.

17 Ms. Miller: Thank you. I think I have comments through a lot of things that I didn't  
18 get to talk to. So there might be a little bit in the random order, but I'm going to start  
19 from-- Because how you summarize, I think, structured the way I wanted to talk about  
20 it. So, going back to Dr. Thomas. When manufacturing, we are very careful and mindful  
21 to serve the patient when deploying this kind of technology and our doer is best  
22 approach. So to your point, we don't use the technology just because the technology is  
23 available. So the fact we look at the problem and I try to find the best solution to the  
24 problem, and I think that still applies for GenAI, they probably be different way to solve

1 the same problem and not using GenAI. With that being said, the risk that the device  
2 entails also targets in a way how we pick the best solution to it.

3 So in the same fashion, the stress-based testing that we are talking about should  
4 be targeted based on risk. I'm worried of enforcing a stress-based test for all devices and  
5 then they're going to be companies trying to pop up, "Okay, give it to me, I'll test it for  
6 you," and that kind of stuff. So, I want to be mindful of that and when we do need to  
7 enforce a stress requirement in a device. Now going back and what the benchmarking  
8 would be, and this is a new technology and that kind of stuff, so it might just hinder  
9 deploying our devices if we go just blindly. On a matter of usability and transparency,  
10 when we submit devices right now and we talk about this in PCCP and in our regulatory  
11 frameworks, we also include the targeted-intended population and to what you  
12 mentioned earlier, the intended population for which a device is intended to be used as  
13 part of the device package in a way.

14 So therefore there are different usability considerations and different usability  
15 studies that the manufacturer does when they are targeting a population. And we will  
16 need to be mindful of it. I don't know how to explain this, but I'm going to maybe use a  
17 little different words or too much confession. For example, for the sake of  
18 transferability and usability, you might disclose too much information when it becomes  
19 harmful because it's very confusing. So I call this a confession disease because you and  
20 I apologize, I want to make my point clear. It has to be targeted to the population you're  
21 intended to target for because you have a different kind of transparency for a patient and  
22 different kind of transparency to a physician and different, and you mentioned very  
23 greatly that practice of medicine varies and in cities versus its undeserved places, the  
24 practice of medicine varies. So depending on what your target population is you need to  
25 design it differently. And with that I think also the benefit of how do you design this

1 device to allow adoption to healthcare and underserved population, may outweigh some  
2 of these risks that we don't exactly know today how to design force. And I had one  
3 more comment that I'm trying to remember. I gather all of them when you guys are  
4 talking. No, I think the usability was the one. Thank you.

5 Dr. Soni: I'll just say in relation to usability testing, we do have some opportunities  
6 with use of simulation labs that have been more widely used for learning, but in many  
7 ways there is a great need for learning and a very steeper learning curve on how we are  
8 using and how the end-user is using GenAI. So usability testing data that comes from  
9 simulation environments where it's more contrived-use cases and we're able to get at  
10 specific components of GenAI and its implications on interactions and decision-making  
11 may be helpful and dovetailing that with where some of the benchmarks should sit. I  
12 think the idea that was brought up by Tempest spokesperson of having clearinghouses  
13 where we can continue to be able to look at prospective validation as well as ongoing  
14 monitoring is helpful to be able to continue to learn from ongoing practice and use of it.  
15 I'm not sure how much clearinghouse in that concept fits into the roadmap for the data  
16 that gets generated for its evaluation of GenAI, but that may be an important  
17 consideration.

18 Dr. Jackson: To your question about usability influencing performance metrics. I  
19 believe that it does, especially, I'm just thinking from mental health examples, there's  
20 research that shows if it's a person that end-user is a patient who is using the generative  
21 AI, the relationship feels different and they see it more as an authority to Dr. Clarkson's  
22 point compared to if it's a clinician who's using it and how they're thinking about it. So  
23 I think then the perceived performance could be different even how a patient may say  
24 this is being more effective than it actually is because of how they're feeling about it.  
25 The relationship compared to a clinician and the outcomes that they may get. I also

1 think when you're saying should it be considered, I think one of the things we haven't  
2 thought about is consideration to what weight. There's lots of things to be considered,  
3 but it is something that is a 90% weight in the consideration. Is it a 10% weight? So I  
4 think that is something also to be considered when evaluating how much weight do  
5 some of the considerations should have or do have.

6 Dr. Bhatt: I think Radman, then Miller and then Shah.

7 Dr. Radman: Thomas, I just want to make a comment on the question on ethical or  
8 white hat hacking. I think that it's actually very important for the reason that there's this  
9 concern of the ability to generate AI models to ingest any tests or benchmark that we  
10 might be able to put forth. So once a benchmark becomes a target, it's no longer  
11 effective. It was mentioned, actually the first comment today, this idea of these quality  
12 assurance labs, or CHAI, and I think Dr. Bizzo's organization might be related to this  
13 assurance function, if I'm understanding correctly-- ARCH-AI. And FDA has some  
14 familiarity in they already do a third-party review for some products, so passing off  
15 some review functions to an independent source and I wouldn't want a company  
16 themselves doing their own quality assurance for this white hat hacking function, but  
17 hopefully their profitability would be able to fund such a function in an independent or  
18 not-for-profit framework. Thank you.

19 Dr. Bhatt: Dr. Shah.

20 Dr. Shah: Right. I want to specifically answer Troy's question about data drift  
21 hallucinations and also what to do with it, right? So I think this is like a no unknown,  
22 question that you're asking, right? If people who are submitting a new weather data drift  
23 was then they would've fixed it. So that's the big problem that everybody's going to  
24 face and it's going to be really hard to put that problem back on the submitter of the  
25 evidence because then they will have to do all this extra work to Dr. Miller's point



1 about not burdening the submitters with all this extra work. And neither do we want  
2 internally at the FDA people have legions of reviewers reviewing all this. It's just not  
3 pragmatic. So the problem is that if it should all go back to the label of the product, if  
4 you're labeling your product to do something very specific, let's say you're  
5 characterizing segmenting prostate cancer, Gleason grade 5 tumors for example, that's  
6 the label, that's what they're sending.

7 Then the data questions, there should just be a questionnaire to start with to say  
8 how many examples of each type do you think your model has seen for which your  
9 product is labeled? And if they say, we have seen 10,000 examples of this tumor, 5,000  
10 examples of this tumor, 3,000 examples of this tumor, then at least we know the upper  
11 and lower bound of their data and their performance and we can at least anticipate  
12 where it'll go. Then also ask them the confession thing that Dr. Miller was saying,  
13 where do you think the model might drift in your evaluation of training it? And they  
14 might say, we have not seen too many Gleason grade 5 tumors before. So that's one  
15 thing. I think your life is significantly more harder when you have third-party  
16 foundation models because you really don't know what the underlying training data was  
17 on which the medical AI product was built.

18 And that's going to be a big problem to sort out. And my recommendation  
19 would be that we ask them that, does your underlying foundation model in your  
20 evaluation, has it ever seen the types of examples the product has been labeled for? If  
21 the answer is negative, then the significance of data drift is high. And if they say yes,  
22 then it's like okay, at least they have a notion. And then you ask them about the fine  
23 tuning of the foundation model on the specific product label and say how much fine  
24 tuning has been done and go from there. The second part was the hallucination part you  
25 asked, and that's a specific problem for text-based models as well, not just vision-based

1 models. So for vision transformers or their tokenization problems or for other models,  
2 we need to know what hallucination looks like for it to colon hallucination. That also  
3 goes back to asking a simple questionnaire versus all this evidence, where do you think  
4 this will drift and what do you think are the obvious stress points that we should be  
5 careful of?

6 Dr. Bhatt: Wonderful. Thank you so much. We are a little bit over time. So Mr.  
7 Tazbaz, is this adequate for now and are we ready for a break?

8 Mr. Tazbaz: I mean with that transition? No. Yes, it's adequate. Thank you.

9 Dr. Bhatt: Excellent. We'll take a quick 10-minute break. When we come back  
10 we'll hear from a couple more experts and the next question only has one part, so that  
11 should be good. Thank you all.

12 *Sub-Topic: Risk Management*

13 Stakeholder Perspective – Strategies and Controls to Mitigate Risks Associated with  
14 Gen AI Applications in Healthcare

15 Dr. Bhatt: Thank you all so much. That was a wonderful discussion by the Panel  
16 and we'll all slowly move to our seats to resume the meeting. I see people eating apples.  
17 I'm very proud of your healthfulness. I personally have a cookie, but it made for a good  
18 break. Thank you so much for taking your seats. If there are conversations that need to  
19 happen, please feel free to step out of the room for just a moment and then come back in  
20 if you are not actively participating.

21 Okay, first of all, I want to thank the Panel for their thoughtful recommendations  
22 in the previous session, there was a lot to go through and they really not only hit the  
23 highlights but went into great depth. We'll now resume the meeting and we will hear  
24 four more presentations after which the Panel may ask questions to each of the  
25 presenters. Our first presenter is Dr. Michael Schlosser, who is Senior Vice President of

1 Care Transformation and Innovation at HCA Healthcare. Dr. Schlosser, you may  
2 approach the podium. Oh, you're already there. Thank you.

3 Dr. Schlosser: Excellent. Thank you. And thank you for the opportunity to  
4 address the Panel. We at HCA Healthcare are excited for the possibilities that AI can  
5 bring to our healthcare teams and our patients. The opportunity to make healthcare more  
6 efficient and effective and a better experience is really a generational one. So we  
7 appreciate the FDA and the federal government's thoughtful approach to ensuring that  
8 common-sense regulation is advanced that supports responsible use of AI while not  
9 stifling innovation. I also want to clarify that I'll use some examples from HCA  
10 Healthcare in my presentation and while these are not medical devices, I believe the  
11 framework we'll use to evaluate them or we are using to evaluate them is broadly  
12 relevant. Just a little bit of context, which I think is important as I go into my  
13 presentation, HCA Healthcare is one of the largest provider systems in the United  
14 States. We have 188 hospitals across 20 states, around 2,500 other sites of care and we  
15 have the privilege of providing about 43 million patient encounters a year.

16 Why that context is relevant is when you start evaluating generative AI use  
17 cases in particular and you start thinking about things like hallucination rates and error  
18 rates and you start talking about 1% and 2%, the mass starts to add up really quickly.  
19 I'll give one example in my talk about nurse handoff, the handoff between one shift to  
20 the next in a hospital. We do that about 400,000 times a week, about 22 million times a  
21 year. So those 1% error rates start to look very different at that kind of scale.

22 So three real topics: two that I'll speak to in a little bit of depth and then one that  
23 I'll just make some quick comments on. The first is really around governance, which I  
24 think speaks directly to the topic at hand of risk mitigation. And I'll speak to the  
25 governance structure that we have stood up and are using and studying and learning

1 from across all of HCA Healthcare. And then the second topic is around a specific risk-  
2 mitigation strategy that we have used with some frequency, which is the human-in-the-  
3 loop and there's been a lot of conversation around that today and so I think a little bit of  
4 a deeper dive into that is really useful and warranted. Then lastly, I want to make just a  
5 quick comment on data and data use. All AI strategies are data strategies, and so I think  
6 that's an important part of the conversation.

7         So this is HCA's responsible AI framework: We developed this and put it into  
8 place about a year ago. We have a policy which is available on our public website that  
9 describes this in some detail. That policy went into effect earlier this year and really  
10 requires all of our 300,000 colleagues, including our practicing physicians and nurses  
11 and others, to engage in the responsible AI program if they're going to use AI use cases  
12 within the bounds of our healthcare system. So whether they're using it on their  
13 personal devices like we heard earlier today, people just pulling up Chat GTP on their  
14 phone to ask questions or whether it's an AI system that we have deployed or are  
15 planning to deploy, the policy really directs that the first step is to engage with the  
16 responsible AI program and go through this risk mitigation strategy before the AI is  
17 really used at scale.

18         The backbone of the strategy really is a risk-based approach. So as we've talked  
19 about a lot today, the types of AI vary dramatically from simple algorithms to these  
20 very large complex transformer based generative AI models. So a nuanced approach  
21 that really looks at risk and benefit and then understands whether we want to take on  
22 that risk and the benefit outweighs it. And then what are the mitigating strategies that  
23 we can use to decrease the risk is really the approach that we saw as the most successful  
24 one. I'm not going to go through everything that's in that wheel, that wheel describes  
25 our program, but I want to highlight a few of the key areas. You can see the seven

1 domains. I will talk about those on the next slide in a little more detail. Safe and secure,  
2 private, transparent, and again, I'll go into those just in a little bit more detail.

3 Those are the domains or areas of risk and it's our framework. There are many  
4 other similar frameworks out there that we use to identify what are the unique risks that  
5 any specific AI-application if deployed could create. Underneath those the individual  
6 wedges are sort of sub-domains so more detailed questions in each of those domains  
7 that help us really get to the specific nuanced risks that may be created based on the AI  
8 solution we're considering deploying. Around the outside of the model are various  
9 mitigating strategies or levers that we have to pull in order to try to mitigate that risk, so  
10 we talked earlier, the Panel talked earlier about re-skilling and education and how  
11 important that would be. And so you can see that's a key part of our strategy, machine  
12 learning operations, which really are the systems for what we would call here  
13 postmarket surveillance, but how do we continue to observe and manage models that  
14 are deployed at scale and ensure that they continue to perform at the level we want them  
15 to is another critical lever when we think about deploying these kinds of models.

16 And then the other areas are our ethics, our cultural norms. These are ways that  
17 we expect our team members and our colleagues to show up and engage in these  
18 processes and do things in an ethical way, in ways that are consistent with how HCA  
19 Healthcare practices in general. My last comment on this slide is that the team that  
20 represents this program that administers this program and it's a team of leaders and data  
21 scientists, really understand that their job is to enable the use of artificial intelligence,  
22 it's not to limit it. It would be very easy to manage the risk by just saying no, that's a  
23 pretty easy lever to pull, but that would again miss that generational opportunity we  
24 have to transform healthcare. So they understand their job is to identify the risk to

1 mitigate it and then only to say no if we feel like we're in a zone where we're not ready  
2 to take on a level of risk that a product potentially creates.

3           So this is the seven domains of risk that we study with each of our AI  
4 applications before we would pilot it or deploy it. Just to comment on that, many times  
5 the risk assessment really goes in multiple phases where we may only understand a  
6 piece of the picture of the risk it creates when we're just studying it on paper or  
7 studying other folks' use case or use of that product. You actually have to deploy it in a  
8 pilot and study it to really fully understand how these risks may show up, that gives you  
9 within a second layer of coming back to this process and then restudying and  
10 determining do we have the right mitigating strategies for a scale deployment versus a  
11 pilot deployment? So there really are no surprises here, in fact, I think this is very  
12 similar to the FAVES construct that the ONC has used though these don't organize  
13 themselves into a nice acronym, so the government couldn't make use of this  
14 framework.

15           But otherwise, I think very similar concepts. So safe and secure, we will  
16 establish careful and intentional guidelines to protect patient data against risk and ensure  
17 data is trusted and interpreted accurately. I'm not going to read you the entire slide,  
18 obviously you have it in the record, but each of these questions or statements describes  
19 to the business user, the end-user or the person responsible for this AI application, what  
20 it is that we expect of them as we think about how the application or the model is  
21 deployed and the mitigating strategies that we put into place. I would say that if you  
22 think about safety and security, privacy, really fair and impartial, these tend to be pretty  
23 technical. You're in the weeds there in terms of how does the model connect to our data  
24 systems and what are the ways in which those connections are governed and how do we  
25 ensure that the data is being transmitted in a safe way and that the model can't be used

1 to access the data, that it has access to very, very technical conversations. When you get  
2 into transparent and explainable and robust and reliable I think these are actually the  
3 more interesting domains of risk and it's where I'll spend the rest of my time discussing  
4 here today. These are areas where often we don't have the complete answer and we're  
5 still moving forward. I'm going to show you examples of how we're using AI, but I  
6 think the right answers are not completely determined yet in those areas.

7         Just one more point on this risk model that we have. So underneath those seven  
8 domains are actually 43 individual questions that make up a Risk Register. So those  
9 questions each map to one of those domains. Again, you have this, so I won't read it to  
10 you, but this is an example on privacy and so as we're considering deploying a model as  
11 a business unit or a user is considering using a model, they're asked to go through these  
12 43 questions and determine how does their model map to this Risk Register. Is this a  
13 risk that that model particularly has? Is it a high likelihood that this risk shows up when  
14 the model's deployed? Is it low-likelihood? And then there's some data already in the  
15 register that suggests what the likely impact of this risk would be if it does show up.  
16 And then that all can be used, all 43 questions in conjunction to come up with a score so  
17 that we understand is this a low-risk use case, is it a medium-risk, is it a high-risk or a  
18 critical-risk use case?

19         And then we have a governance structure that oversees how we respond to those  
20 four levels of risks. So we have a very senior-level governance committee, which  
21 involves some of the senior leaders in the company like our chief lawyer, chief of ethics  
22 and compliance, and our chief clinical officer. The highest-risk use cases have to be  
23 reviewed all the way by those levels of folks before they would be used. If it's a  
24 medium or low-risk use case, there's faster paths to at least piloting and testing those  
25 use cases. But the rigor of actually going through this process and asking these 43

1 individual questions, understanding those seven domains of risk, mapping it to a level of  
2 risk, a bucket of high, medium, low, we believe is a framework with the governance  
3 structure to review that, that will help us deploy AI safely and effectively while not  
4 putting our patients or our care team members at risk in a way that we wouldn't want to.

5 A big piece of this strategy also comes very much down to use case selection. So  
6 a lot of what we end up wrestling with is do we want to take on this level of risk yet or  
7 not? And I'll tell you HCA Healthcare's approach to this is that we're very much  
8 tiptoeing into the use of generative AI and for a few reasons. One is I think generative  
9 AI, we still have an enormous amount to learn about these models and how they  
10 perform and why they do the things that they do. So putting a patient's treatment or  
11 diagnosis or other aspects of their care sort of directly at the end-result of a generative-  
12 AI model, I think is something at least our organization is not ready for. I think there's  
13 risk in taking on that too soon. Also in that if we take on those kinds of risks and we  
14 make mistakes and patients get hurt, we could end up setting ourselves back years  
15 before we ever get to the point of being comfortable doing it again. So this incremental  
16 acceptance of risk over time I think is a way that we can build trust with models, we can  
17 get people to use them in a way that's low-risk, but it builds comfort, it helps them  
18 understand how to use them, what they're good at, what they're not good at, which will  
19 make them better prepared to take the next step. And then the next step.

20 So this is just an example of what the risk register output of a generative AI  
21 model could look like. I'm not going to go through this all in detail here. I'm going to  
22 focus just on one bar or maybe two, which is the hallucination. So, of course generative  
23 AI models and in this case we were thinking about a use case where the model is  
24 producing an output like a medical note or a note summary. There's the risk of  
25 hallucination, which when you have the human-in-the-loop, we would rate that as a



1 moderate risk because the human is there to correct the errors. Then there is that risk of  
2 over-reliance on AI systems, which we have not really demonstrated yet, but certainly is  
3 something that we would be concerned about with these types of systems. So I want to  
4 focus just in my last few minutes here on this concept of human-in-the-loop, which I  
5 think is an important aspect of generative-AI models.

6 In fact, I'm going to make the argument that I think it's a required aspect of  
7 generative-AI models at this point in time. We often talk about human-in-the-loop as if  
8 it's just sort of one thing, which is the model produces an output and then we ask the  
9 human to review it, they make their edits and then they accept the output and then it  
10 goes into the medical record or whatever the next step of that model is. I think to do this  
11 really well is actually far more nuanced than that because at scale, again, if I go back to  
12 the 22 million handoffs, if the end users are not really engaging in this human-in-the-  
13 loop process in a robust way, you could end up with a huge number of errors making it  
14 through those humans and into the medical record. And then it's very hard to calculate  
15 what the downstream impact of that could be over time, especially as it accumulates  
16 over time.

17 So, here are six concepts that I think are important when designing systems that  
18 include human-in-the-loop as its safety measure. Promoting transparency in whatever  
19 way we can. We've had a good discussion about where you can be and can't be  
20 transparent in showing how a model made its decisions when you talk about generative  
21 AI, but there are lots of good transparency levers you can put into place, like asking the  
22 model to provide citations that show where it found the data that it displays. Displaying  
23 levels of uncertainty where it doesn't have a citation or can't show you exactly where  
24 that information comes from. Make sure that it's in the interface, that the user can see  
25 that level of uncertainty. Encourage active participation. So, designing the UI / UX in a

1 way that draws the end user into that activity. This is where your UI / UX designers are  
2 really worth their weight in gold. Provide training. We have a requirement in our  
3 responsible AI program that anyone who uses AI has to go through at least a minimum  
4 level of training. We have a series of micro-learnings they're exposed to which teaches  
5 them at least what generative AI is and what responsible use of AI looks like. So, they  
6 have some baseline understanding of what these models are and what they can do.  
7 Design for trust and calibration. And then last, implement feedback mechanisms. And  
8 I'm going to double click on "feedback mechanisms" here for just a second. I know I'm  
9 running short on time. I'm going to use this handoff example. So, we, partnering with  
10 Google, designed a nurse handoff tool to assist in that process. It's a very important  
11 time in care delivery, as all the clinicians know. We used a standard framework, the I-  
12 PASS framework. This is something that's been used for a long time. And we ventured  
13 to teach an AI model how to do this process, how to use an I-PASS framework to read  
14 the medical record and come up with a consolidated summary, a timeline of the  
15 progression of care and a prioritized list of tasks and key considerations for the  
16 oncoming shift.

17 We built into the user interface of that handoff tool the ability for the end users to  
18 provide direct feedback on five what we thought were the most critical variables in their  
19 assessment of the quality of that model: factuality, which speaks to hallucinations,  
20 coverage, organization, conciseness and helpfulness. And they would evaluate all five  
21 of these on this red, yellow, green scale. And then, there would also be a composite  
22 score where if they rated all five as green, then that would be a green note. It's  
23 important to note that these five features are actually, sometimes, in conflict with each  
24 other. So, you can have really good coverage, and you can have good organization if  
25 you allow the model a lot of creativity, but that usually leads to hallucinations. If you

1 make the model highly factual, it will often omit a lot of things, and its organization  
2 sometimes can be poor. So, balancing these things against each other is critically  
3 important to the usability of the model and the safety and helpfulness.

4 So, in this last slide, I'm just showing you one of the variables. This is factuality.  
5 We have this data for all five of them, but in the interest of brevity, I'm just showing  
6 you one example. What I think is interesting about this is if you look at model version  
7 zero, which was still at the time we did this work, which is this past summer, we used  
8 some of the best foundation models that were available. We tried multiple different  
9 ones, including Google, Gemini and OpenAI 4.0. And you can see that the model  
10 actually performed pretty poorly. And this was despite the fact that at the time the  
11 literature would suggest these models were really good at summarization tools, but  
12 when you asked it to do something just slightly different than summarize, and you  
13 applied a really critical eye to it, the nurse, and said, "No, I don't want just a summary. I  
14 want the critical elements, the most contextually relevant elements from the data to be  
15 presented," the model immediately did poor. And then, we had to go through a series of  
16 steps. And we didn't change the foundation model actually through this entire  
17 progression, we changed the system that sits around that foundation model. And I'll just  
18 editorialize. That system includes the person that's using the model; it's part of that  
19 system. And we created feedback loops where they would score on that red, yellow,  
20 green. That information was passed back to the system, sometimes directly to  
21 developers, to the prompt engineers, sometimes in the form of recording the output as  
22 green, the highest quality output, and using that as examples that we fed the model. And  
23 using all of those techniques in an iterative way, we got to the right-hand side of the  
24 slide where we were over 90% factual, and the nurses told us at that point the model  
25 was good enough. It was also over 90% helpful to be able to use in at least a production

1 pilot. This was a very controlled environment. And so, we've now moved this into the  
2 next step.

3 So, a couple of final thoughts and then I do want to mention my comment on data.  
4 So, first is-- I think when you evaluate these kinds of systems, it's very hard to detangle  
5 the model itself-- In this case, the model itself is really just a foundation model. It's  
6 Gemini 1.5 Pro. From the system that you put around the model, which includes the  
7 user interface, the feedback mechanisms, and in the end the people that are using the  
8 model. And so, my opinion is if we're going to evaluate the performance of a generative  
9 AI system, you have to look at the performance of the system, not just the technical  
10 pieces of the model.

11 I also believe that in order for AI to be successful, it actually has to be used. And so,  
12 there's a certain threshold that we may want to see before we were deployed something  
13 broadly. As you can see, we did that ourselves. But I think that we're going to have to  
14 come to an acceptance that the answers to these questions are never going to be 100%.  
15 And so, we're going to have to study how do you deploy a model that's 98%  
16 successful? How do you ensure that it's designed such that humans are monitoring that  
17 2% in a successful way? And then, how are we returning that data that the humans are  
18 providing us back to the system and forming a continuous loop so that the model is  
19 constantly correcting errors and getting better? And that should protect against some  
20 forms of drift and other changes in the environment that the model may encounter. But  
21 just understanding the premarket performance and assuming that it will somehow  
22 provide safety in the real world I think is going to be incredibly challenging for  
23 generative AI.

24 My last comment on data is that these foundation models are incredibly powerful,  
25 but it's the context-relevant data that makes them useful in workflow. And so, these are

1 always, in my opinion, going to be a combination of technology companies who have  
2 the ability to generate these foundation models and folks like ourselves who have rich  
3 context-sensitive data. And it's the bringing those two together that's really going to  
4 create successful outcomes. And so, while I know FDA doesn't necessarily regulate or  
5 doesn't regulate patient data, I think making sure that we don't create dogmatic rules  
6 when it comes to how data is used that makes it impossible for us to do that ongoing  
7 context-sensitive training, I think it's also going to be important for success. So, I'll end  
8 there. Apologize for running over time and thank you.

9 Stakeholder Perspective - Narrow VS Generative AI: Risk Determination > Controls =>

10 Safe Innovation

11 Dr. Bhatt: Great. Thank you, Dr. Schlosser. Next on the agenda is Dr. Keith Dreyer,  
12 Chief Data Science Officer and Chief Imaging Information Officer at Mass General  
13 Brigham. Dr. Dreyer.

14 Dr. Dreyer: Thank you and thank you all for the opportunity to present. I'm a  
15 diagnostic radiologist, but I have spent the last 10 years of my career dedicated almost  
16 100% to the creation and deployment of AI in the pre-Gen state and in the current  
17 GenAI state.

18 I'm going to talk about narrow AI versus generative AI with respect to risk  
19 determination, controls and safe innovation. I think it's important for this Committee to  
20 be able to get a good purview of what the FDA has done successfully and some of the  
21 challenges that have been in place through the narrow AI era, if you will, in thinking  
22 forward, as to some of these pluses and minuses might progress into this generative AI  
23 state. I want to look at technology, regulations, usage, and monitoring as kind of this life  
24 cycle. When the technology CNNs came out in the early 2010s, it was seen as a  
25 transformational technology that computers for the first time could do things that they

1 couldn't do before at a percentage that was acceptable. And so, that entered the foray of  
2 the concept of putting this into healthcare. I would say that this challenges the  
3 regulatory system in this country and several others around the notion of should we  
4 create a software as a medical device and treat these as devices? They're obviously not  
5 food or not drugs. So, if you have to have something that already exists, device  
6 probably makes sense. But I still think that this Committee should even question today  
7 if that rubric of a device is the right way to take a look at GenAI, like what was  
8 mentioned with Bakul Patel's recent article talking about this is more human-like. And  
9 so, I think some of the challenges that we're going to see historically looking back is  
10 that we maybe made this decision too prematurely. The adoption of AI in this narrow  
11 AI state was also somewhat limited, and I'll talk about that, for a variety of different  
12 reasons. And, as you probably all know, monitoring is almost exactly zero throughout  
13 the United States in this by device manufacturers as well as those that have deployed the  
14 technology.

15         So, I just want to talk about what happened through this stage of good  
16 technology that ran into some regulatory hurdles. First, I know the decision is: is it AI?  
17 Is it used in healthcare? If the answer is yes, it falls into HHS. If it's a medical device  
18 that the answer is yes by 21st Century Cures or whatever mechanism you want to use to  
19 be able to determine that, then it falls into FDA, this group here, and makes the decision  
20 falls into SaMD or Software in a Medical Device, and then gets weighted for risk. The  
21 reality is Software as a Medical Device is essentially all one risk's number two. So,  
22 there is no risk stratification, as Dr. Schlosser described what they're doing. And then,  
23 there's different ways to get them approved. Probably most of you are familiar with  
24 this, but almost all have come through the 510(k) pathway, being that these devices are  
25 predominantly imaging-based, where 75% are radiology, but upwards of high 80% are

1 imaging of any kind, optical, cardiac, etcetera. And they fall into these known  
2 categories. And these categories also kind of determine risk, but now at the level of  
3 specificity of risk, category I, II, III. That has resulted in 900 AI medical devices for  
4 imaging of any type, a thousand total, another hundred that are Software as a Medical  
5 Device for non-imaging. And I want to focus on these. This is really because that's what  
6 I do and also this is the predominant area that has been used for the last 10 years and  
7 approved and deployed, because I think it's a good lesson to look at it retrospectively to  
8 see some of the things that could happen going forward.

9         So, basically there are two different categories, despite the fact that there are  
10 multiple risks and it's really determined by-- These six categories are defined, but yet  
11 when it comes time to be able to impose a control to limit that risk, there's really only  
12 two choices. There's the concept of standalone performance testing, or comparative  
13 effectiveness or reader study. And that has created some challenges in the market today--  
14 - I'm going to put my glasses on. Challenges in the market today in that the cleared  
15 devices of those 900 look like this. About 13% are those high-risk, but high-value  
16 solutions. And the bulk of the products that are available today are these CADt, CADa/o  
17 and MIMPS-type solutions that had a lower bar to get through. So, it was much cheaper  
18 as opposed to millions of dollars to go through a CRO-approval process. It was  
19 thousands-- Hundreds of thousands and took months, instead of years. So, the result is  
20 there's a predominance of triage software and non-diagnostic solutions, and there's a  
21 very small number of diagnostic solutions available today. The result in that is-- And  
22 through the American College of Radiology survey of radiologists asking about  
23 adoption rates of AI over the last five years, it's about 2% today, and there's reasons for  
24 this. The FDA requires that triage not to show imaging results. So, because triage is  
25 seen as a lower solution and standalone performance test is okay, I believe the belief is,

1 or could be discussed, that if you do show images, now you're making a prediction and  
2 as opposed to just triaging the patient, but that also has an unintended consequence of  
3 making it much more complex for that human to interpret that. So, if I know there's a  
4 problem in the image, but I can't see where, and I haven't been told where, I'm  
5 naturally going to take longer and be more concerned that I'm going to miss it, but I  
6 won't know what it is that I'm missing. And since all AI products require human  
7 oversight with the exception of one-- So, that 900 require that. That triage AI makes the  
8 workforce burnout worse for radiologists. A very important point here is that if I see  
9 that there's-- if I don't see what it is that the AI is saying, but the AI is saying there's  
10 something there. It's going to slow down those that are doing the interpretation.  
11 Therefore, many experts won't use this type of AI, and that represents 87%. CMS has  
12 created fewer than five CPT codes for payment of AI. And there's no commercial  
13 coverage for FDA cleared AI products today. So, there's no incentive to deploy. That  
14 means that it's off the back of practice performance that needs to be purchased. And  
15 these same products in other countries are considered diagnostic solutions, providing  
16 improved patient care, but the AI coverage of diagnostic tasks in the United States is  
17 about 1%.

18 And how would I come to that number? If you look at what a radiologist does, a  
19 fully capable general radiologist, of which there are a few, because everyone's  
20 subspecialized, there's about 3,500 findings that you would have to find to be able to  
21 read all modalities, all organ systems. If you look at the landscape of AI in the United  
22 States today approved for CADe or CADx-- It's this limited number that are there  
23 inside of those boxes, as you could see predominantly in mammography, chest X-ray,  
24 etcetera. But there's about 42 of those 3,500 that are necessary, which is less than about  
25 1%. Also, of note though, if you look at the chest X-ray as an example, seven findings



1 that are CADe or CADx for diagnosis are-- The same AI products in the UK and Asia  
2 provide full diagnostic coverage. Those CADe, CADt solutions here are full diagnostic  
3 solutions elsewhere in other sovereign nations.

4 So at the end of that, because all AI requires human oversight, but not by an  
5 expert, arguably if a non-qualified individual is using that AI to make an interpretation,  
6 it's arguably coming out sometimes with a worse diagnosis than if a human was able to  
7 do it completely independently, separate from the other lesser qualified individual that  
8 would be using AI to make the diagnosis. And we've had these problems in the  
9 organization-- At Mass General Brigham, in our organization. Limitations in SaMD risk  
10 classification in my mind are preventing AI from providing patient benefit at this point.

11 So, if that's the state today, what happens as this technological innovation  
12 transforms into generative AI technology? Clearly, again, an absolutely incredible  
13 technology that appeared on the scene just in the last few years. There is no regulatory  
14 pathway to have these medical devices approved. Dr. Schlosser approved or rather  
15 discussed different ways for non-medical device AI to be used, but there is a  
16 tremendous challenge right now to try and use GenAI inside of a regulatory or medical  
17 device, even if it's created inside of an organization. So, yeah, usage is somewhat high,  
18 and monitoring is still bleak. And what do I mean by that? Well, there's pathways for  
19 people to be able to go straight to the solution, like GenAI, drop an image in or do  
20 various things, that is happening today. And so, it's not this life cycle. It's like a direct  
21 line; a lifeline approach to this is what's happening today.

22 So, let's talk about this regulatory challenge now in light of narrow AI and also  
23 generative AI. If you look at the differences between the two, at least from my  
24 standpoint of how I see the challenges of regulation, which have clearly been discussed  
25 here, traditional neural networks are deterministic, they retrieve, they don't generate,

1 and they have narrow capabilities, whereas these are probabilistic, they're generative or  
2 can be generative, they can also retrieve and they have emergent behavior, unforeseen  
3 activities that they can perform. So, the question that I think is reasonable to ask is this  
4 notion of risk and that process that has been capable of controlling these deterministic,  
5 retrieval, narrow-capable devices with standalone performance tests and comparative  
6 effectiveness tests. Could that same kind of thing be applied to probabilistic generative  
7 behavior? I think the answer is yes. I think that you could just extend the risk curve and  
8 be able to use the same controls that were considered before with some caveats. If you  
9 look at this as all medical devices independent of the technology, but just looking at  
10 them as their intended use, their users and their risk, then I think you get a kind of good  
11 analogy. But the thing I would say, and I'll talk about this more, is I would not use  
12 comparative effectiveness. Running a CRO and watching the use of reader studies  
13 compared to standalone performance tests, I don't see the value. I don't understand the  
14 value in comparative effectiveness testing. I don't see that it does anything stronger than  
15 a standalone performance test, as compared to continuous monitoring and things you  
16 can do far beyond premarket testing. So, what I would do is look at standalone  
17 performance tests required for everything, as well as clinical validation. And I'll talk  
18 more about that in a second, what I mean by clinical validation, but then I would also  
19 look at the notion of ongoing monitoring. And those things should be based on and  
20 controlled by the amount of risk that the device imposes.

21 So again, if I look at this flow chart and look at these CAD solutions that are out  
22 there today, standalone performance tests handle these and comparative effectiveness  
23 tests handle these, but I would argue again that standalone performance tests could  
24 probably do adequate at that. These are all premarket testing with no monitoring. So, if  
25 you add GenAI onto this list and it falls through as a SaMD or SiMD, then probably, I

1 hope, the knee-jerk reaction is to just do comparative effectiveness testing because I  
2 think it's going to have the same challenge that it's had previously. And what I would  
3 do instead is start to impose standalone performance tests plus validation, plus the  
4 concept of monitoring. And now you have all these new devices that have to go through  
5 that rigor to various levels and degrees, but there's also-- As Dr. Schlosser mentioned,  
6 there are non-device GenAI applications that also need to be considered. I'll talk about  
7 those in a second, but those are not medical devices. So, somewhere up in the level of  
8 HHS, maybe not FDA, maybe that also needs to be managed and considered from a risk  
9 standpoint and for mechanisms of controls.

10 So, now we have devices that are not able to be used inside the United States,  
11 and we have devices that are able to be used inside the United States. And so, it's a very  
12 confusing time for providers like us to decide where the line is drawn, what these things  
13 are. And then how exactly do you enforce these internally and how do you deploy  
14 them?

15 Let me just use one example of voice-to-text transcription as well as diagnostic  
16 draft reporting. And it was interesting to see this used so many times as an example  
17 today, but let's just do that. So, text-- Just like we do speech recognition, as opposed to  
18 image recognition or image ability to create a text. So, the process is, as a radiologist--  
19 And we have been doing this for decades. I look at an image, I dictate the interpretation  
20 into AI of some sort. Previously, they were RNNs, now they're more and more LLMs  
21 or foundation models. And it does voice-to-text and I get a report. That's a draft, I sign  
22 it off and I'm ready to go. What I'm talking now about is the notion of taking an image,  
23 putting that into essentially the same AI, but now I get a draft report out and now, as the  
24 radiologist, I look at the draft and I look at the image and I make an interpretation. And  
25 now that's how I get the report. So, transcription is considered to be not a device, but

1 *transvision* or whatever you want to call it, it just made up the word, is a medical device  
2 if you look at the rules, 21st Century Cures. But it's the same AI solution.

3 So, the question I start to ask is if transcription is okay to use and *transvision* is  
4 not okay right now to use unless it's-- As it's considered a medical device, and it has to  
5 get approved through a regulatory process to be determined, what about summarization?  
6 Today, there are applications that say, "You dictate half the report, and you let AI  
7 complete that report for you into the impressions". That has at least been deployed and  
8 used extensively throughout the country and has not had any kind of FDA restriction on  
9 it. So, I would say it's probably not a medical device, but unknown. What about quality  
10 checking? Others are talking about the ability to be able to say, "You dictate whatever  
11 you want and then I'll take a look at what you dictated. I'll take a look at the image, and  
12 I'll produce the difference between that to help to try and decide if there was anything  
13 that was missed." Again, unknown. So, these are things I think that this group has to  
14 decide as to what are the bounds and confines of generative AI that falls under the  
15 rubric of a medical device? And does that even matter anymore today? And then finally,  
16 consumer usage, where there's no claims and no regs, seems like it's okay to use. So,  
17 my biggest concern and what I see happening are-- Healthcare, yes. We welcome AI.  
18 Come on in, generative AI. Please, come in. And then, "Oh. It's a round peg." And all  
19 of the AI that's already there is going through a round hole, and this says, "Wait... I can  
20 use the consumer entrance." No claims, no regs, no problems. That's what's happening  
21 unfortunately. I know that no one wants this to happen, but this is what is happening  
22 until we find a better way to solve this problem.

23 I want to use diagnostic draft reporting as a GenAI case study of what I think  
24 was happening. This is how we're studying this up at our organization. So, if you think  
25 of this as the same model, it's doing both, then you can think of it like this. So, I,

1 basically as the radiologist, am looking at the image and I'm doing dictation, just like I  
2 have been for decades. That's being handled by AI, and it generates a draft report. What  
3 I'm talking about now-- Sorry. Then, I read that draft report, say, "Yes. It says what it's  
4 supposed to say," and I have a final report. What I'm talking now about is also adding  
5 the image into that AI so that it now can add additional information to that draft report  
6 so that I can do the same thing I've always done, but now that's a medical device. Now,  
7 if this was accurate, and this has been described in other presentations, it would add up  
8 to about 150% efficiency. We've seen no efficiency from the AI that's available today  
9 through the 900 products that are available to radiologists. This would clearly just  
10 improve efficiency. So, the question is how do we test this component for accuracy? By  
11 clinical validation and ongoing monitoring of this device.

12 So, Bernardo, Dr. Bizzo, talked about the healthcare arena, but basically this is a  
13 collaborative community in the form of the FDA's collaborative community concept to  
14 be able to create a solution that allows people to be able to look at these large language  
15 models that are for medical devices or not for medical devices and see if they're  
16 applicable to the environment that we have, to the data that we have, to the subject  
17 matter experts that we have, and the involving foundation models that are being used.  
18 So, what it does is it takes de-identified data and AI models and experts and makes them  
19 available here. In this example, it's a draft report. The important thing to note here is  
20 that, when I talk about a model and the instructions, the instructions are the prompt and  
21 the prompt is actually baked in here, right? So, it's not an open prompt. I can't ask  
22 whatever I want. This has been specifically trained and designed to be able to create a  
23 draft report. And so, I take de-identified data, I run it through that generator, that draft  
24 report generator, and I get an AI result. I get a draft report. As the expert, I look at that,  
25 and most important, I look at the input data to that model as well. So, now I see the text-

1 - I'm sorry, the image that the model saw, and I see the output that the model generated,  
2 and I rate that. I rate that from a range of acceptable all the way to unacceptable. And  
3 through that process of multiple experts with multiple data being able to run this, I can  
4 now start to extract discordance, performance, generalizability and bias, because we  
5 have ways to pull back to the data to say, "This subpopulation is biased," etcetera,  
6 etcetera. And so, the arena gives you this geographically and demographically diverse,  
7 embargoed data; unlimited clinical use cases (this one is the draft report generator);  
8 geographically and institutionally diverse experts, and expert evaluation of AI agents  
9 with input data. So, I can look at the input data and look at the agent's output data, and  
10 you wind up with these analytic evaluations of agents in real-world environments that  
11 help us to make a decision of whether we're going to deploy or not. Independent of  
12 FDA approval, we have to make that decision, and as well create performance ratings in  
13 a public forum.

14 This is what it looks like. So, in this particular example, I say, "Create a draft  
15 report." I can select the areas that I want to look at, emergency pulmonary, and I can  
16 now see there's clearly a fracture here. I can also see the output of the model that  
17 created that report that I'm now rating and evaluating. And I say, "I feel that this is at  
18 the level of-- At a Likert scale, resident three." I've got another example here. I read  
19 through the report. Another example. I rate it, etcetera. And now we have the ability to  
20 create these structured outputs where these stand and all of those statistics about these  
21 models that we need to decide to determine whether we want to use or are able to use.

22 The one thing that we've seen early in the deployment of this is the ability to  
23 reduce risk of these models. So, one thing I would petition the Committee to think about  
24 is if there were pathways to be able to solve this problem while you're still trying to  
25 solve the entire problem, but it allows models to be able to come out that are safe or low

1 risk, it'll help us all learn much more quickly how to go to the next ones that are higher  
2 and higher risk. I wouldn't try and solve everything at once. This is what I guess I'm  
3 trying to say. So, maybe if an open query isn't the best way to go or there's no clarity on  
4 open query, at least a closed query using LLMs might be something that's acceptable.  
5 So, in this case, the ability to create this information with the user and interpretation  
6 with an open query is pretty high risk. But what we've seen that the risk is much lower  
7 if you have fixed instruction prompts, if you show the input data, if the user is expert in  
8 that input data. So, if it's a radiologist looking at MRIs, if it's a pathologist looking at  
9 pathology images, as well as the AI. And the interpretation, the usage is monitored, then  
10 that risk drops considerably. So, fixed instructions, input data shown, user expert data  
11 and usage is monitored, is the big key. On this graph, it looks like this. You want to  
12 monitor that usage.

13         So, we see in this process of validation, of testing the model before it goes into  
14 deployment, it informs the manufacturer of potential defects that determines the site  
15 accuracy and generalizability; it lowers the risk of failure on site; it promotes early and  
16 safe innovation, and it provides a mechanism for ongoing monitoring as well. You can  
17 use the same solution for that.

18         I'll finally close with this thought. Is there a thoughtful way to regulate  
19 healthcare AI beyond a medical device definition? So, if you think of this complex  
20 drawing that's got all these different pathways is the way to maybe think of this to say  
21 performance tests, site validation and site monitoring are three controls that can be used  
22 for medical devices or should they really be things that are be able-- Accessible for  
23 devices, whether they're medical devices or not, just the fact that they do use AI?

24         Thank you very much.

1 Stakeholder Perspective – Safety from the Systems to Patient Levels: Risk Management  
2 for Large Language Models in Healthcare

3 Dr. Bhatt: Thank you so much, Keith. Our next presenter is Dr. Danielle Bitterman  
4 who is the Assistant Professor at Harvard Medical School. And then, we'll have one  
5 more presentation directly after that. Thanks, Dr. Bitterman.

6 Dr. Bitterman: Thank you so much to the Panel for inviting me to speak today. I'm  
7 Danielle Bitterman and I am at Mass General Brigham and Dana-Farber Cancer  
8 Institute, Harvard Medical School. And today I'm going to be speaking from my dual  
9 roles both as a practicing clinician, I'm a radiation oncologist, as well as a researcher in  
10 natural language processing. My lab is focusing on developing natural language  
11 processing, including large language model methods, to support and enhance safe  
12 clinical care delivery. And we're very interested in developing methods to improve the  
13 automated evaluation and robust evaluation of these models. Here are my disclosures.

14 So, a lot has been talked today already about how large language models,  
15 foundation models, are developed, but it is really difficult, as others have mentioned, to  
16 disentangle the quality of the foundation model from the safety and risks of a device. I  
17 really do think it's essential to understand and continue to monitor the underlying large  
18 language model that will be integrated into any device. So how do we get to the helpful  
19 and useful large language models that have really catalyzed this explosion of  
20 innovations and medical AI? So, today's advanced large language models are generally  
21 trained by taking large amounts of text curated from the Internet and having the model  
22 predict the next word in a sentence. So, most of this text is not medical; it's standard  
23 text. And from that we get a model that is able to transform and handle language. It can  
24 answer nonmedical questions pretty well.



1           There also is a lot of information about biomedicine on the Internet. And so,  
2 though language models happen to be trained on that in the process of their general  
3 training, and as we've already heard, this has allowed us to develop-- This has led to  
4 models that are large enough and turn on enough data that can now pass a variety of  
5 different medical benchmarks, including the US Medical Licensing Exam. So, it really  
6 is quite impressive, but there is still this outstanding question of how does performance  
7 on these benchmark evaluations in the general biomedical realm translate into safe and  
8 effective applications that truly help a clinician in practice and are safe for our patients?  
9 So, while large language models have a lot of potential, they also bring unique risks that  
10 require proactive controls. And the goal here is to balance innovation and safety so we  
11 can realize the benefits for patients and providers as quickly as possible, but also as  
12 safely as possible.

13           And again, the training processes through which we arrive at a large language  
14 model is the foundation for any medical device that uses that model's performance,  
15 behavior and risks. So, we already discussed language model pre-training. Most of these  
16 large language models also go through additional tuning processes, known as  
17 instruction tuning and our preference tuning. So, instruction tuning is when a model is  
18 given an additional fine-tuning dataset with instructions and the output of the desired  
19 instruction. They're tuned on that. That makes the model able to, of course, follow  
20 instructions. Preference tuning is a process by which a model is shown two different  
21 examples of a model output that have been human labeled for a preference. And while  
22 these make our models helpful, they also introduce new risks, such as leading to models  
23 that might be overly compliant, easy to jailbreak, as these models will often try and  
24 always give an answer, even if they may not "know the right answer" and that is risky  
25 for medical misinformation.

1           In our group, we're really interested in understanding whether we can start to  
2 understand the risks and behaviors of a large language model from its pre-training data.  
3 We carried out a study on a smaller publicly available large language model known as  
4 Pythia, for which the pre-training corpus is completely publicly available. And we were  
5 able to show that the pre-training corpora do reflect realities of current medical  
6 knowledge. Here you can see. These are the number of mentions of various diseases. I  
7 highlighted COVID-19 and infection in the pre-training corpora that these models are  
8 trained on. And you can see, as expected, COVID-19 and infections spike around 2019.  
9 So, that's encouraging. For this model, it is, at least at the time of 2019, up to date with  
10 general medical knowledge. However, we also showed that biases that were present in  
11 the statistics of the pre-training data led to biased diagnostic decision-making of the  
12 ultimate model that used those data to train on.

13           So, that's actually very valuable to know. If you're able to-- If you have access  
14 to a language model's pre-training data, you might be able to anticipate risks and its  
15 behavior in a more generalizable fashion than if you're relying on the inherently limited  
16 sets of benchmark datasets and red teaming prompting. However, this requires having  
17 transparency into the fine-tuning datasets-- The pre-training datasets and fine-tuning  
18 datasets used to develop a model. It also requires transparency into how the model is  
19 trained. And I do want to really emphasize, from a clinician and researcher's  
20 perspective, it is a lot of risk for the clinicians and patients to take on-- To onboard  
21 models for which we do not have enough information into how they were trained. I  
22 know this is likely going to be a challenge, especially from the larger companies where  
23 this information is currently proprietary, but when we're talking about clinicians'  
24 liability and most importantly patient's healthcare, it really is incumbent on us and we

1 have, I think, a really unique opportunity right now to require more transparency so that  
2 we can ensure that we develop these safe models.

3         Alright. Moving a little bit more into the translational side, there are risks at  
4 various levels of clinical application once we start to think about moving these models  
5 into clinic. These risks range from risks to systems privacy and security, risks arising  
6 from the device and also risk to the device, things like cybersecurity concerns; risks to  
7 patient safety and clinical effectiveness related to device performance, safety and impact  
8 and outcomes; new risks related to workflow integration for which we need better  
9 understanding of human factors, better ways to collect user feedback in real or near-real  
10 time and better methods for scalable monitoring; and finally, a lot of ethical and legal  
11 concerns regarding the transparency, equity, accountability and responsibility when one  
12 works with these models. And again, the behaviors and knowledge learned during  
13 model training and tuning modulate, and I would argue, really form the backbone for all  
14 of these risks.

15         In terms of systems-level risk and controls, some of these are unique to  
16 generative AI and some are more broadly applicable to software integrated into a  
17 healthcare system. But large language models might interface with many different  
18 points of patient data, clinician data within the Electronic Health Record system, as well  
19 as on the Internet and patients interacting with these models via their personal devices.  
20 So, data governance and clear definitions for data governance for the data input into the  
21 model and output from the model, as well as appropriate data masking, is important for  
22 understanding risk privacy. Security protocols, such as encryption, audit trails and  
23 cybersecurity protections, will be important from a systems-level. And, as we've heard  
24 a lot during today's talks, deployment controls for on-label use are really important.  
25 While pre-training and tuning form the backbone for a model's behavior, the prompting

1 methods are also very-- Are of course the kind of final step to determining what that  
2 model will output. So, role-based access controls might be one approach, if clearly  
3 defined, to help ensure on-label use, beyond simply role-based access controls, being  
4 clear about how a device is going to standardize, clean and ensure that the data input  
5 into the model, because those become the healthcare data which become the prompts,  
6 are consistent and will be important, and including issues about-- And ways to monitor  
7 for the type of data modalities, languages, and tasks that the model is being used for.

8         Now, underlying all of this is ensuring the adequate clinical evaluation of these  
9 models. Again, current benchmark datasets that are used to assess large language  
10 models, general biomedical knowledge, have clear gold standards and have reliable and  
11 automated evaluation methods, which is good. However, they do not reflect most of the  
12 really useful real-world applications that we care about and where there's a lot of  
13 excitement for advancing healthcare. So, I give one example of a real-world application,  
14 which is using large language models to help assist clinicians in responding to patient  
15 portal messages. This is a big source of clinician burnout, and many institutions are  
16 piloting this application. This is just one example of a question a patient might ask in  
17 their patient portal: "I've been experiencing hot flashes and night sweats for the past  
18 week. How likely is this a side effect of my prostate cancer treatment? What should I  
19 do?" Now, if I asked three different physicians or oncologists to respond to this  
20 message, you would get three very different answers. If I showed those three  
21 oncologists a draft response to that message, you would likely get three very different  
22 assessments of the quality of the response. And this really makes it hard to assess  
23 language model performance. Currently, there are no or very few gold standards or  
24 reference sets to automate the evaluation of such generative-type outputs and no way to  
25 reliably automate these evaluations. Human evaluation still really is the backbone, but

1 as I just mentioned, it has its own complications and requires clear definitions and  
2 transparency into how those evaluators were instructed to evaluate the quality of a  
3 model's output.

4       Okay. So, because of the limitations of existing evaluation benchmarks and the  
5 practical challenges of evaluating a small set of exemplars, a small set of inputs for your  
6 given task, I really do think that a multi-level holistic evaluation is needed to start to  
7 understand a language model's risk profile. So, one initial approach is to start-- Maybe  
8 understanding the general safety of that model. In the general domain, there are  
9 benchmarks to evaluate how likely a model is to be truthful. How likely is it to provide--  
10 - To produce toxic language? Evaluating these general harmful behaviors might start to  
11 give us a baseline understanding of the more explicit riskiness of a model that's being  
12 integrated into a medical device, but it's not adequate on its own of course. Our lab was  
13 interested in understanding just how these benchmarks might be-- The extent to which  
14 they might truly be reflecting a model's "knowledge" versus just kind of stochasticity or  
15 whether the model just was pre-trained on data. So, we did just about the simplest thing  
16 we could do. We took the USMLE exam, the benchmark most commonly used to  
17 evaluate a large language model's clinical reasoning capacity, and we switched all  
18 mentions of a brand name to a generic name of the same drug and vice versa. Now,  
19 when we asked-- So, this was an example. A mechanism of action of leuprolide is  
20 GnRH agonism. That's the correct answer. When we switch those two terms referring to  
21 the same drug, the models often responded differently. And that's despite the fact that  
22 those models in most cases were able to accurately match the brand and generic drug as  
23 being the same. So, benchmarks themselves need to be audited. This might be because  
24 the models were trained and overfit at this point to the USMLE Medical Licensing  
25 Exam. But it also speaks to the fact that language models "reasoning" or "knowledge",

1 and I'm sorry I'm anthropomorphizing a little bit, is very different to how humans learn  
2 and reason. Language models are probably much better than humans at memorizing  
3 huge amounts of data, but they have flaws in knowledge manipulation that humans may  
4 not have. And this complicates using benchmark dataset or testing credential-- Existing  
5 credentialing methods to monitor language models' knowledge because you might  
6 really just be testing memorization and not a true understanding of the concept.

7 We, in conjunction with several other NLP researchers internationally, have  
8 developed the TRIPOD-LLM reporting checklist to provide standards for what should  
9 be reported in biomedical research according to large language models. We lay out  
10 information that should be reported regarding the data used to develop the model or  
11 whether that researcher did not have access to information, to that data, how the model  
12 was developed, how the evaluation was carried out, who carried out the evaluation, and  
13 what were the guidelines by which the evaluation was conducted, and finally what were  
14 the results and endpoints for that study? Now, speaking to the kind of open-ended  
15 nature of large language models, we ended up using a modular format that adapts to  
16 different use cases, but this is one approach that we've arrived at to try and develop a  
17 more consistent method to review-- To evaluate these models.

18 Now, finally, as we've seen, many AI models work incredibly differently *in*  
19 *silico* than they do when integrated into a busy, complicated, messy clinic. When large  
20 language models started to be rolled out into EHRs for patient portal messaging, we  
21 were really interested in understanding how that may impact clinician reasoning and  
22 risks when used with a human-in-the-loop. So, we carried out a preclinical study where  
23 we asked oncologists to respond to patient messages asking about new symptoms  
24 related to cancer, things that we often respond to in a patient portal. First, we had the  
25 oncologist's manually written responses as usual. Then, we had GPT4 draft response,

1 and asked the oncologist to revise it until they thought it was appropriate and safe to  
2 send to a human. And what did we see? We saw that, in most cases, a language model  
3 response was safe and in about 60% of cases, oncologists thought that the draft could be  
4 sent to a patient without additional editing. However, there was a small, but clinically  
5 significant risk of severe harm or death without editing. So, that's concerning. However,  
6 if humans are affected at overseeing these models, we could reduce that risk and  
7 potentially it might be a safer implementation method. That said, we did see evidence of  
8 automation bias and over-reliance that could expose patients to those rare or more  
9 harmful risks of a model. So, we saw that clinicians-- The content of responses was  
10 significantly different when clinicians manually wrote the response, compared to when  
11 they use the LLM draft. The LLM just tends to reflect quite-- Tended to closely reflect  
12 the LLM draft, suggesting that physicians might be relying on the LLM's clinical  
13 reasoning more than just using it as an efficiency assist. So, even these more  
14 administrative approaches might impact clinical reasoning in unexpected ways.

15 So, when we're thinking about workflow integration controls, it will be really  
16 important to understand and to see evidence of how medical devices that use large  
17 language models impact effectiveness and safety when used as a part of the human-  
18 machine team. Automation bias and over-reliance are common and are likely to  
19 increase. And we can't rely only on human oversight to control these methods, to  
20 control these risks.

21 There are various different approaches to post-deployment and real-world  
22 performance monitoring ranging from regular quality assurance, prompt versioning and  
23 controls, ongoing audits, monitoring for off-target use and evaluation for shifts and how  
24 a user interacts with a model. Are they spending more time editing a response? Are  
25 more words being added or deleted from a report that's generated by an LLM? Those

1 may give us early hints of faulty performance and devices that integrate these controls  
2 are likely to be safer.

3 I'm running out of time so I'll skip over this and just touch on some of the ethics  
4 and legal risks, which is really important for clinicians and patients to feel they can trust  
5 a model. LLMs come with a host of equity concerns. They learn biases from the model  
6 they're trained on. Humans and machines may introduce new biases together, and they  
7 perform less well on lower resource languages. Transparency will be essential for  
8 ensuring trust from both the patient and clinician standpoint. Patients and clinicians  
9 really do deserve to know when they are interacting with a large language model that  
10 might not be always clear because some of their output is so human-like, and when their  
11 data is being used to further refine a model. And finally, clearly defining who was  
12 accountable and responsible for any risks and errors of a model is an important risk  
13 control. If a clinician doesn't know that they're going to be accountable for the output,  
14 they may not use it in as safe a way as is necessary.

15 We've defined some levels of human-in-the-loop previously, drawing from the  
16 Department of Transportation, categorizing models from assistive all the way to fully  
17 autonomous, which might be a helpful guide for further defining this.

18 So, in conclusion, language models have potential to advance health but have  
19 risks at multiple levels of healthcare. Challenges persist in scoping, monitoring and  
20 mitigating risks. And human-computer interaction modulates the benefits of risks and  
21 need to be investigated and risk-mitigated. There are several emerging approaches that  
22 should be considered when refining and developing oversight methods, including  
23 automated risk and performance assessments, interpretability methods and innovations  
24 in usability design. And the goal of this overall more measured approach will be to  
25 facilitate a durable and sustainable implementation of AI that advances human health in



1 the long-term in a way that patients and clinicians feel trust and comfort with our  
2 changing healthcare system. Thank you very much.

3 Dr. Bhatt: Thank you, Dr. Bitterman. We are going to forego our next break.  
4 Everybody gets 30 seconds to stand up for a second and not get sciatica. Do not leave  
5 your seat. Please just stand up and stretch. And while you're doing that, Gabriella  
6 Waters please come on up. Artificial Intelligence Evaluation Research Associate at  
7 NIST. You will have 15 minutes for your presentation.

8 Stakeholder Perspective - Risk Management for Generative AI-Enabled Medical  
9 Devices

10 Dr. Bhatt: Alright. Sciatica now being prevented. Everybody, come on and have a  
11 seat. Yeah. Okay. We now welcome Gabriella Waters, AI Evaluation Research  
12 Associate at NIST and the Director of the Cognitive and Neurodiversity AI Lab. Dr.  
13 Waters, you may begin, and you may notice me using two arms when things turn red up  
14 at the podium. So, my apologies if I start waving at you.

15 Dr. Waters: I will just wave back. Hi. My name is Gabriella Waters. Thank you so  
16 much for this opportunity to speak. I was furiously deleting notes as the presenters  
17 before me were going, so this will be quicker because we were kind of walking in the  
18 same lane on what we wanted to talk about.

19 I don't have to tell you what generative AI is, but as a theoretical AI researcher,  
20 as someone who develops these models, what I want to say about them in response to  
21 some of the comments that were given earlier is that LLMs are a kind of technology that  
22 you should be more amazed when they get it right versus when they get it wrong. The  
23 probabilistic stochastic nature of them lends them to furiously confabulating. So, the  
24 example given by Dr. Bitterman where the generic drug had the reversal of the names is  
25 very logical when you consider that it is rolling a die on what the next word should be.

1 It is only used to seeing the words in one order. So, when it's not in that order, it's  
2 rolling a die on what the next word should be. And so, when you have a model that  
3 behaves in this way, it's very challenging to deploy it in a high-risk situation, like in a  
4 clinical setting. It's also sort of risky to do any of this testing post-appointment; that is  
5 too late. All of the testing needs to happen beforehand. So, I want to put a good order of  
6 operations in place before we continue on through.

7 So again, I'm a theoretical researcher, so I'm going to always champion  
8 innovation and how we can move the technology forward and what that's going to look  
9 like. So, generative AI is really good at helping us to do personalized treatment  
10 planning, helping to enhance that medical imaging capabilities. We've seen a lot of this  
11 work in diabetic retinal neuropathy in terms of how models can actually perceive  
12 images and make diagnoses and look at decision support. We've got lots of examples of  
13 that. I'm not going to go through it. I told you I promised I wasn't going to go over the  
14 same ground, but there are definitely risks that we need to consider: data quality and  
15 bias, harmful bias. As a developer, we like bias in AI. It's the harmful bias that we want  
16 to get rid of. Good example: if I am looking at a dataset and I only want women in the  
17 dataset, the men in the audience might feel that's really biased. But if I tell you it's for a  
18 women's clinic, then that model needs to be biased. So, bias can enhance performance.  
19 It's when the bias leads to harmful outputs that we have to be on alert. So, just making  
20 sure we understand how the model is biased and under what situations that bias  
21 emerges, and then looking at the risks associated with that. We want to look at model  
22 transparency and explainability, and I keep those two separate. People who hold IP do  
23 not like transparency. So, we want to do a tradeoff between transparency and  
24 explainability. Full transparency means I know the entire suite of your ensemble. I  
25 know every single algorithm. I know how its architecture looks. Real-world

1 transparency means I know when a model is being used. Explainability means that any  
2 person in this room can tell you how the model arrived at its decision. That's important  
3 for building trust; it's important for deploying things in the real world. And we're also  
4 looking at performance variability in the real world. LLMs and generative AI are non-  
5 deterministic. So, you heard from Dr. Dreyer about deterministic narrow networks. We  
6 know that if we pull a piece of information from that network, it's going to look the  
7 same every single time. In generative AI, that's not the case. You can ask the exact  
8 same question on 40 different machines and get a slightly different answer every time.  
9 And if that tiny bit of variance leads to a misdiagnosis even 2% of the time, that's  
10 dangerous. And so real-world performance has to be something that's evaluated and  
11 tested for. At NIST, we have the AI Risk Management Framework. So, Dr. Dreyer also  
12 mentioned validation. We also need reliability. It's got to be reliable, and it's  
13 challenging to do that if you haven't tested for it beforehand. It's difficult to go back  
14 once it's out in the wild and then try to redress some of the challenges that you're seeing  
15 within the model.

16 Dr. Bitterman mentioned the governance controls. We definitely need some  
17 form of compliance as these models are developed and deployed and implemented and  
18 what that looks like. We need some sort of Review Board. As a researcher, I have to  
19 submit to the IRB; as an AI developer, I do what I want. That's not right in terms of  
20 clinical applications of this technology. So, these are things that we need to consider.

21 We talked about training, validation and controls. Robust training protocols:  
22 what does that mean? It is literally a free-for-all. You can develop your model in any  
23 way that you like. Do you remember there used to be a slogan "have it your way"?  
24 That's literally what it is when we're developing models. We just do whatever it is and  
25 then present it. We need protocols for this. We need those rigorous validation

1 procedures. We need to have TEVV as a standard part of the AI implementation  
2 process. The design process has to be there. Testing, evaluation, validation and  
3 verification.

4 Feedback mechanisms. It can't just be-- I have the thumbs up there for a reason.  
5 It can't just be helpful and harmful. Binaries do not provide context. Why is it helpful?  
6 Why is it harmful? What happens when it's good? What makes it good? The model  
7 can't learn outside of the context, and it can't update with just a thumbs up or a thumbs  
8 down. We love Likert scales. We really do. What else can we use to leverage training  
9 these models for more robust performance? We need user feedback integration. How do  
10 we get all of the developers who are developing models for people to actually like  
11 people? Because at the end of the day, people are interfacing with these agents. And the  
12 number one complaint I get when I bring this up is, "People are messy." Like, "Yes.  
13 That's why you need to test with the actual people." This shouldn't be a radical idea.  
14 We need to get that continuous feedback integration going so that we can understand the  
15 surfaces for which a risk can emerge, what the impact is of that kind of risk, and then  
16 what's the mitigation strategy for it? If we are only observing these things after  
17 something has gone wrong, we are again too late.

18 Postmarket surveillance and performance benchmarking. This is where I usually  
19 take a step back from the microphone and take a deep inhale. Benchmarking does not  
20 tell us about risk. Benchmarking tells us about capabilities and performance, and those  
21 things are not concurrent with risk. That just tells me that the model is doing what you  
22 said it was going to do. And in the process of doing that, it's making misdiagnoses. It's  
23 making judgements that are inappropriate. It's culturally irrelevant. It's doing all these  
24 things, but its F1 score is great. Those things don't tell us enough about the context of  
25 the use of the model. And so, we want to make sure that we're not just hung up on

1 KPIs, that we're also looking at how robustly we can measure the risks and not just the  
2 performance.

3 So, I'm trying to go fast. I'm going fast. So, trends. I am going to fly through  
4 these slides because you don't really need all of the details, but you need to know that  
5 people are thinking about how they can integrate generative AI into everything. As an  
6 AI researcher, I get two questions. What is AI? Is it chat GPT or is it the Terminator?  
7 Those are the two. So, everyone thinks that generative AI is AI, and businesses are  
8 trying to capitalize on integrating these systems into whatever use case they have. And  
9 so, you're going to see through these slides that a lot of businesses are looking at an  
10 adoption curve to really drive home how they can get the best ROI on GenAI. It's not  
11 great right now, but the promise of a robust return is there. And so, we're seeing a lot of  
12 interest in these areas. So, I'm just going as fast as I can. Looking at how we can scale  
13 up, that's another thing.

14 Theoretical researcher hat is on. What is on the horizon? Lots of potential use  
15 cases for what the patient will experience in terms of how GenAI will touch their  
16 healthcare continuum as they journey through the system. From pharmaceuticals to their  
17 providers, the papers, the medtech, the services, public health agencies, all of these  
18 areas are touched by GenAI. We're going to start seeing mobile devices as medical  
19 devices. We're already seeing apps that are already on board. Most of these devices that  
20 can track things like your heart rate, like your respiration, all kinds of-- You have things  
21 that will alert you if any change happens from your baseline. So, we're already starting  
22 to see this. And we're seeing a lot of online doctor interactions, online medical  
23 professional interactions. And the risks are-- Again, we heard about over-reliance, but  
24 we're also thinking about that as a two-way street because the more distance we place  
25 between the healthcare provider and the patient, the more the patient may distrust the

1 healthcare provider's advice, the more they may have trouble adhering. So, we have to  
2 look at how these systems can augment, instead of supplant. So, there are ways where  
3 you can have a LLM that is parsing and summarizing information that's being pulled in  
4 from the online portal, or if it's in a mental health situation and someone's submitted a  
5 report overnight and all of these kinds of information. But you still need to do that as  
6 part of a team. So, it's you get this information, it's digestible, and then you're able to  
7 follow up as a human. So, we need the human-in-the-loop. We need the human  
8 feedback. We need to very critically test and evaluate these systems. We need to  
9 understand their mechanics. If I rolled a die in front of you each time I said a word, you  
10 would not trust that what I was going to say to you was going to be coherent. But that is  
11 the nature of generative AI. It's a dice roll, every single word. And as long as we have  
12 these systems in this way, and as long as the nature of them is so opaque, it's really  
13 challenging. And I'm not saying Black Box because I love Black Box models. I love  
14 them. They're really good at finding patterns that humans don't understand. Opaque  
15 models are not really good at helping us understand why something is going wrong. So,  
16 it comes back to the explainability, it comes back to the testing so that we can  
17 understand. I have successfully finished early.

18 Open Committee Discussion Q&A (*Clarification questions*)

19 Dr. Bhatt: That was absolutely fantastic. All right. We are going to take a couple of  
20 minutes to ask some clarifying questions of the speakers that we just had. I'm going to  
21 take a minute, even though we're over time, to just remind us of where we've come  
22 from. So, the things we've heard are the following. We must create a safe infrastructure  
23 to employ the equitable use of generative AI devices in healthcare with guardrails and a  
24 robust evaluation mechanism. We are not trying to inhibit but to promote safety. There  
25 has to be strategies and controls to mitigate risk associated with GenAI applications,

1 governance frameworks, designing model systems for safe deployment and evaluation,  
2 and mitigating strategies that range all the way from ethics to reusability to cultural  
3 norms and re-skilling. We then learned that evaluation of the performance of an AI  
4 model involves evaluating the performance of the system and all its parts, including its  
5 deployment, its users, its interpretation, its iteration. And critical variables within a  
6 framework should include factuality, coverage or admission, organization, conciseness  
7 and helpfulness. But we very recently learned that benchmarking tells us about  
8 capabilities and performance, not necessarily risk. So, for risk mitigation, we do want to  
9 think about user feedback for GenAI model improvement and educate users about AI  
10 limitations and potential errors. Finally, we need to train clinicians and patients, when  
11 applicable, to use these devices because human-in-the-loop is essential, since at this  
12 stage we're looking largely at qualitative outputs, and we have few automated measures  
13 that assess closeness to truth in the absence of a human expertise. Equally important,  
14 however, is role-based access to AI models. If this AI model is not in your division, in  
15 your area, should you have access to that? Automation bias and over-reliance  
16 complicate the good part of human oversight. However, therefore, human training is  
17 still, again, essential for the success of the system. Lastly, we learned that premarket  
18 testing to date includes standalone performance testing or comparative effectiveness.  
19 However, previously, it was deterministic, involved retrieval and narrow capabilities.  
20 GenAI, we've heard quite a bit now, is probabilistic, generative and it has emergent  
21 behavior. So, maybe we look at standalone performance, clinical validation, and  
22 ongoing monitoring particularly. Lastly, we learned that generative AI, I love this, is  
23 furiously confabulating, a die roll for the next best word. So, we need to focus on real-  
24 world transparency, explainability. And real-world performance variability must be

1 evaluated prior to deployment in the premarket stage and not be relied on only in  
2 postmarket surveillance.

3 So, I think that is what we heard today. With that reminder in mind, do we have  
4 questions for the panelists? I will ask our team here to please keep your question limited  
5 to only a question because after this we have a chance to discuss question number 2  
6 where you can give your opinion on the matter. So, specific questions. We'll take five  
7 minutes for questions of the Panel that just spoke. Dr. Soni, who are you asking your  
8 question to so that they can come up?

9 Dr. Soni: The question is for Dr. Dreyer, but really thought-provoking commentary  
10 from everyone. Dr. Dreyer, you really made us recognize some of the complexities of a  
11 lot of different pathways. You highlighted something that was really important, which is  
12 one of the outcomes that potentially should be evaluated is efficiency. I'm curious if  
13 you can expand on that a little bit more, and how do you measure it? How do you  
14 compare it to non-users of GenAI? And how do you continue to monitor it on an  
15 ongoing basis?

16 Dr. Dreyer: Yeah. So, this is a great topic and it's going to take longer than a couple  
17 minutes. So, I'll just paraphrase. I think I question if it's the FDA's job to measure  
18 efficiency because I think the market should do that. With a thousand devices today,  
19 people look at these solutions and say, "Is this making me faster? Can I afford to  
20 purchase this? Is there some return on investment here?" And that's usually through  
21 efficiency. It's assumed that FDA measures the safety and efficacy of this, makes sure  
22 that this is a safe device. But I think where reader studies get challenging is to try and  
23 demonstrate that someone is more proficient or efficient in the process of with versus  
24 without. I just don't-- I see that as something that we have to do every time we choose a



1 device. So, I don't know-- It seems very redundant and very expensive to try and do that  
2 *a priori*, as opposed to at a particular facility.

3 Dr. Bhatt: Dr. Dreyer, will you stay there? Do we have other questions for Dr.  
4 Dreyer specifically? Yes.

5 Dr. Radman: So, you talked a lot about your, I guess, advocacy of postmarket  
6 monitoring.

7 Dr. Dreyer: Yeah.

8 Dr. Radman: And I think we've talked a lot today about model drifting. I think we all  
9 agree that they're going to drift. So, it has me thinking about that fact and the FDA's  
10 framework of, as you mentioned, class II devices being the most common, which use  
11 substantial equivalents.

12 Dr. Dreyer: Right.

13 Dr. Radman: So, if you take a new device that has validation data, like Timeway this  
14 year.

15 Dr. Dreyer: Yeah.

16 Dr. Radman: And then you're comparing it to a previous device, the ecosystem which  
17 is tested is different. And that's the whole-- It's one of the components causing data  
18 drift.

19 Dr. Dreyer: Yeah.

20 Dr. Radman: Right?

21 Dr. Dreyer: Yeah.

22 Dr. Radman: Model drift. So, do you-- Considering your position on the importance of  
23 data drift and continuing monitoring, do you also agree with the substantial equivalence  
24 framework that it's equivalent to compare to a device at a previous time? Or should you  
25 do a head-to-head bake-off at the same time?

1 Dr. Dreyer: No. That's a great question. I feel like everyone goes through-- Everyone  
2 takes the easiest way to market. Right? And so, an equivalency test is what everyone  
3 tries to achieve. I've looked at some equivalency tests and looked at the two devices that  
4 are considered equivalent. I don't consider them equivalent, but I'm not in a position to  
5 make that decision. So, I think even beyond what you're saying, there's other reasons  
6 why I don't know that that really works well. And I think in GenAI time, it's going to  
7 be much more complicated to try and figure out that equivalence. The thing that I would  
8 say is, in my space-- And I try to remind folks that these devices are not acting on  
9 humans; they're acting on data that's coming out of a device that came from a human.  
10 So, for example, the 900 devices that are being used today in the old-fashioned narrow  
11 AI don't drift. These are deterministic, but the data that they're looking at are coming  
12 off of CT scanners, MRI scanners, EKG machines. Those things are changing all the  
13 time. So, the data is changing all the time. We have 2,000 scanners at our facility and  
14 none of them are the same. So, the concept of trying to have one device work on every  
15 single scanner, another device, that then acts on humans is just not realistic. And that's  
16 why I think that validation has to take place, as opposed to comparative effectiveness or  
17 premarket of anything, because it's a moment in time with fixed experts and fixed data.  
18 And I just don't think that scales, especially in the generative AI space.

19 Dr. Bhatt: Great. Thank you so much. Ms. Miller, who is your question for? And a  
20 reminder that we are focused on risk mitigation and management with our current  
21 panelists.

22 Ms. Miller: My question is for Dr. Schlosser. I think I said the name right. This is a  
23 very quick question. What is the intended use for your model? I'm not sure you  
24 described that well or I missed that or-- Exactly, how that's going to be used, how it's  
25 being used.

1 Dr. Schlosser: You're talking about the nurse handoff model that we talked about. Well,  
2 so, let me first clarify that I do not believe that model is a medical device.

3 Ms. Miller: Okay.

4 Dr. Schlosser: And so, we've not specifically defined an intended use the way that the  
5 FDA would, but if I had to come up with an intended use, it would be to facilitate the  
6 handoff of a patient from one shift to the next.

7 Dr. Bhatt: Great. Dr. Kukafka, who is your question for?

8 Dr. Kukafka: I think it applies to all three speakers, but whoever wants to answer it.

9 Dr. Bhatt: Since you're there.

10 Dr. Schlosser: Okay.

11 Dr. Kukafka: Okay. So, I think all three speakers talked about the need to train users  
12 and patients to mitigate risk. And as someone who's worked in more generally  
13 biomedical informatics and more recently AI, training clinicians is next to impossible  
14 because they're busy. I mean, we just can't do it. And certainly, training patients, all  
15 patients, seems less than feasible. So, I just wanted some response to that. And what's  
16 the strategy to do this training to mitigate risk? How is that actually going to happen?

17 Dr. Schlosser: Well, so, I would say that those two populations I would view very  
18 differently because I agree. I think if we are thinking about direct-to-patient AI models  
19 that don't have a clinician-in-the-loop in between the AI and the patient, then you are  
20 going to have to deal with the patients, which is a very large diverse population with  
21 different understandings of what that model is. And I think that's one of the challenges  
22 that we have right now about thinking about direct-to-patient generative AI solutions.  
23 It's that we can't expect to train them and educate them on how to use those models,  
24 which is one of the reasons I think that's high risk at this point. On the provider side, I  
25 see this maybe a different way, which is I think providers who leverage AI and

1 understand how to leverage AI are going to replace those that don't. And so, I see this  
2 as an opportunity more than a challenge. I think that it should be a barrier. In our system  
3 it is. You can't use AI applications, at least in the ways that we can control the use of  
4 AI, if you haven't completed at least minimum training. And I think those who choose  
5 not to go through that training and don't use the AI will over time become obsolete.

6 Dr. Bhatt: Great. Thank you. Dr. Elkin. Last comment and then we'll move to  
7 discussion. Our last question rather, pointed question.

8 Dr. Elkin: Yeah, so question. Two questions. Sorry, I apologize. Two. One is  
9 several of you mentioned that you had a risk assessment survey that you use. Can you  
10 share them? And then the second question is, and something I've been wrestling with--  
11 People here were talking about different models for different modalities and that makes  
12 sense of course, but there's also the issue that these models may not work equally well  
13 for all populations and I wonder if we can find a highly accurate model that works  
14 equally well for all populations or we need to train many models and use the right  
15 model on the right population person. I don't know if any of you can respond to that  
16 one.

17 Dr. Schlosser: I'll grab the first one. I'm going to let someone smarter than me do the  
18 second one. We have shared our risk assessment survey and strategy and in fact I didn't  
19 mention in my talk, but we're actually collaborating with the FDA on studying that risk  
20 framework to determine if it actually achieves the outcomes that we think it will  
21 achieve. We think it's a good framework. It's been used in other industries, but we need  
22 to really study it. And so, we're going to take a number of AI use cases through the  
23 framework, study what the impact was on our system and our outcomes, and then,  
24 hopefully collaboratively with the FDA, publish those results. But we're happy to share  
25 it in its current form as well.

1 Dr. Dreyer: I'll go next unless someone's smarter and answers for me as well. This is  
2 a fundamental question I just think that you've raised here. It's that when we create AI,  
3 sometimes we create it as a product for the United States and sometimes we create it for  
4 ourselves. It's much easier to create a solution for a smaller population and you can get  
5 much higher accuracy by doing that. So, the question to this group is: is the concept of  
6 fine-tuning at the site level something that should be considered? Because if you think--  
7 The goal is to try and have high accuracy devices. You're probably going to get them  
8 higher accuracy if you're focusing that onto a population that you're going to be using it  
9 for.

10 Committee Discussion of the FDA's Questions (*Deliberation and response to FDA*)

11 Dr. Bhatt: Okay. It is now time for Panel deliberations and the discussion of the  
12 FDA question number 2 about risk management and mitigation. Ms. Shick, if you'll  
13 please present the next question. Fortunately, it only has one part. We will have 10  
14 minutes for focused discussion amongst the Digital Health Advisory Committee. We  
15 will then have 10 minutes for the FDA to ask us clarifying questions. I left a little  
16 wiggle room that can extend to about 15 minutes. Ms. Shick.

17 Ms. Shick: So, the question here is: "Risk management: What new opportunities,  
18 such as new intended uses or new applications in existing uses, have been enabled by  
19 generative AI for medical devices, and what new controls may be needed to mitigate  
20 risks associated with the generative AI technologies that enable those opportunities?"  
21 For example, controls related to governance, training, feedback mechanisms, and real-  
22 world performance evaluation.

23 Dr. Bhatt: Great. As you are putting up your placards, I will start with comments  
24 that were made earlier that probably belong in this section. To begin, the Panel  
25 generally believes that the output of generative AI-enabled devices should be

1 categorized according to risk, as we heard in the earlier discussion. Whereas this may be  
2 challenging, the basics of direct acute care provision, clinical diagnosis or triage of risk  
3 level may allow for differential stress testing suggestions by the FDA. I thought that  
4 comment belonged here. A centralized database was mentioned earlier today. It would  
5 be essential, in generative AI, to build in a method for users to report errors for risk  
6 management, similar to the Adverse Event database with particular attention to bias  
7 detection and hallucinations. Earlier, we also mentioned, despite FDA labeling and  
8 marketing requirements for device companies, even the best guardrails may not prevent  
9 natural human drift when it comes to generative AI use. Therefore, it's essential for us  
10 to think about postmarketing surveillance when we talk about off-label use reporting for  
11 risk mitigation. And lastly, human-in-the-loop feedback was essential. It allows  
12 promotion of transparency using explainable AI to understand AI decisions and to  
13 demonstrate confidence levels to inform users of the reliability of outputs. Those were  
14 the four comments I pulled from earlier from our team that I thought perhaps belong in  
15 risk management. Dr. Elkin.

16 Dr. Elkin: Thank you. So, I think that first we have to realize that there's at least  
17 two kinds of risks. One are errors of omission and the other are errors of commission  
18 that happen with these models. I think-- In terms of the new intended uses, I think the  
19 biggest one that I think most of us have talked about is that these present a new way of  
20 presenting information that we accumulate within a model that seems more humanistic.  
21 And because it seems more humanistic, it gives the impression to the user of real  
22 intelligence, whether there is real intelligence or not. I think there's lots of work being  
23 done toward making these more intelligent systems, but the truth of the matter is that  
24 the impression to the user that these are intelligent-based on their conversational  
25 capabilities makes them seem more authoritative in the way that they deliver

1 information. So, even if that information is available from other sources, the impact in  
2 on-prem and toured the information may be actually appreciated differently. So, I think  
3 that governance needs to be done at many levels. Obviously, local governance is  
4 probably the most important thing, rather than national governance of these things  
5 because value sets are often local. And I think we heard that by some of the speakers  
6 this afternoon. I agree with Dr. Kukafka that, although we'd love to be able to train  
7 people to use these, doctors are not-- They're very busy and it's not likely to be able to  
8 enhance things because in-- We know from EMR implementations that if you plan on  
9 your EMR requiring physician training, that it just won't work, and the best approach  
10 has always been at-the-elbow training. So, having clinician champions and clinical  
11 informaticians around that people can go to when they're not sure about things may be  
12 quite helpful.

13 Feedback mechanisms are expensive, and you don't know who to take feedback  
14 from in terms of the model. So, there's two ways to think about that. How you improve  
15 the model for the feedback, but then of course you just want to mirror best practice.  
16 Long time ago, when I was at Mayo, one of the top radiologists said, "Peter, what I  
17 want you to do is come up with a model that would tell junior clinicians exactly what  
18 the five smartest people in our department would've done with this case and that would  
19 help them to learn more rapidly." And I've grown to appreciate his comments over the  
20 years.

21 And then, how we use this as an educational tool to increase the quality of our  
22 care across the board. So, we often talk about best practice, but we still have a problem  
23 with the bottom. Right? We want to raise up the bottom to minimal standards of care,  
24 but then we don't want to stop there. Clinical guidelines do that. We want to get beyond  
25 that and move closer and closer asymptotically to best practice. So, I would look

1 forward to people's thoughts about how we can accomplish that with these new  
2 technologies.

3 Dr. Bhatt: Dr. Clarkson, I'll have you go next. I just want to make a differentiation  
4 again. There is generative AI that is inside devices, which is what is our purview here.  
5 However, appropriately we have discussed after Keith Dreyer's presentation that there's  
6 a lot of other generative AI that is affecting clinicians and patients. I understand that.  
7 But if you could maybe start with the GenAI device part and then if you have thoughts  
8 on the other, we'll record that as well. That would be helpful for anybody who goes  
9 next. Dr. Clarkson.

10 Dr. Clarkson: Hi. I'm Melissa Clarkson, the Consumer Representative. I want to pause  
11 a moment and ask: risk of what? So, our first presenter I believe had seven categories  
12 and one of them was data security and such. And those are all absolute valid risks to pay  
13 attention to. Risk of patient harm I hope is up there. As a patient advocate, we have a  
14 tense relationship with the word "risk management" because we know how this plays  
15 out in the real world, which is first try not to harm patients, but if you do, it becomes  
16 risk management of financial liability and institution reputation. So, with that said, I see  
17 an increased risk for lack of accountability in healthcare. It's already really difficult for  
18 harmed patients to get accountability, but if there's now generative AI, which seems  
19 very trustworthy and authoritative and your hospital has implemented it and it came  
20 from an expensive company and you're a hurried doctor and you used it and now a  
21 patient is harmed, tracing that accountability and getting accountability for harmed  
22 patients, I see this as an increasing problem.

23 Dr. Bhatt: Thank you so much. Dr. Soni.

24 Dr. Soni: So, in thinking about the feedback mechanisms in real-world  
25 performance evaluations and controls over that, perhaps one of the phrases that has been



1 unspoken, but on many people's mind, has been the learning health systems as enabling  
2 a lot of this generative AI when implemented as a medical device. How is that going to  
3 continue to inform evidence-based practice and practice-based evidence as more  
4 intended uses develop? And I think that may be an opportunity where-- Again, different  
5 agencies within HHS have championed learning health systems within FDA for ongoing  
6 monitoring. Where does that come in? An ongoing monitoring of real-world  
7 performance evaluation is easier said than done. But if every implementation is going to  
8 be better informed by local validation, then perhaps that ongoing local validation, and  
9 some sort of an audit report may be helpful.

10 Dr. Bhatt: Dr. Miller. Ms. Miller.

11 Ms. Miller: Yes. Thank you. Diana Miller from Industry. So, I like the way the  
12 question is phrased and I'm going to start from there and piggyback on what you said.  
13 How we tie the risk with the existing intended uses or new applications for existing  
14 uses. And I want to describe a little bit how we do in the industry, how we judge the risk  
15 and the importance of intended use for the medical device. So, even highlighted by the  
16 presentations you heard, think about the controls and existing controls you already have  
17 in place for an existing use and think maybe, you start from there, a more traditional--  
18 An existing use, an existing device, and maybe you incorporate a new technology into  
19 it. What I want to point out is that a lot of times, I think the majority of times in the  
20 industry, when we think of risk and we think of controls, we think of the patient's risk  
21 and intended use and not target specific to the technology. So, you start-- And I'm going  
22 to describe very short the process. You start from assessing the failure modes for that  
23 particular device. So, I have this device, this is a use, what's a failure mode and what  
24 harm can you do, related to that intended use? So, that's agnostic of the technology and  
25 a lot of times what we heard today was close to that, but a lot of times we can apply this

1 kind of framework, and I think the first presentation was very close with what we do in  
2 the industry with red, yellow, and green. We do the same thing and assess and score the  
3 risk and the harm. That can be applied to this technology as well. So, after you evaluate  
4 the clinical implication of that harm and potential failures, then what do you do? Do you  
5 add special controls to mitigate them? Okay. So, you say for this particular harm, for  
6 this particular failure mode, I'm going to do this control, and this is going to mitigate  
7 that harm. And then after that, after you add all those specific controls, at the end of it  
8 all, you have a lady we call residual risk. After all of these controls and all the  
9 mitigations I put in place, what are my residual risks that I still need to mitigate? So,  
10 this is a very simple framework that-- I mean, it's not simple. You saw it's very close to  
11 what we heard today, but this is an existing framework we use, and we focus on the  
12 clinical harm first and it transcends the technology. I don't know if that's a word. I  
13 might have made up the word, but I think it does exist. But it transcends the technology.  
14 Thank you.

15 Dr. Bhatt: Great. Thank you. Dr. Jackson. And if you're done speaking, if you put  
16 your placard down, that would help me. Thank you.

17 Dr. Jackson: Thank you. One of the things that I was thinking about is-- I was trying  
18 to synthesize; it's so much data listening to everyone and everything that we read today.  
19 It relates to controls related to governance. And I'm starting to wonder if we need to  
20 think about having different governance controls because, to Dr. Shah's points, there's  
21 different GenAI models that will be used in medical devices. And even thinking about--  
22 In the documents the FDA gave us, all the uses are not always going to be an image that  
23 goes in. Sometimes, if we have user interface, there's going to be a different  
24 interpretation that somebody is putting in as a part of the human-in-the-loop in different  
25 things. And I think having governance that is the same across both, I don't know that

1 we're going to have a solution, but if there's an image that's just going in, that's very  
2 easy to say, "This is a fact," but if there's a part of it that has human-in-the-loop, but  
3 there's still some interpretation and we're saying that has the same governance, I think  
4 we get into how do we determine which one is valid and reliable and all those other  
5 things that we were talking about? And so, I think that something that needs to be  
6 considered is having different governance controls depending on intended use and  
7 maybe even the model that is being used.

8 Dr. Bhatt: Dr. Radman.

9 Dr. Radman: So, I'm reading this question, and it has this point about real-world  
10 performance evaluation, and it leads me to want to just say that the risk of misdiagnosis  
11 is not new to generative AI. That's something we live with already with traditional AI,  
12 and that risk's getting more likely as the model drifts is also not new with generative AI.  
13 But I think-- As I'm building my framework of the collective thinking here, I'm  
14 realizing that the new challenge with generative AI is that the-- Is the exponential  
15 expansion of the input and output space of the data and the responses. And so, I'm not  
16 sure it's actually possible to adequately test in a satisfactory way in a premarket setting  
17 everything that could occur in that exponential input and output space. So, that really  
18 increases the importance in my estimation of real-world performance evaluation. And I  
19 think it is-- A new kind of regulatory framework is going to be needed. And that brings  
20 us back. We already talked about quality assurance labs and the day started with Dr.  
21 Califf talking about learning health systems, and I think he was right on it when he said--  
22 - He brought up learning health systems and he said he's not sure in his searches if there  
23 is a learning health system that can adequately validate an algorithm in the real-world  
24 setting. But I think that's where we need to go. And anything in the FDA's-- As  
25 someone put it, their gravitas of requiring, for example, closely related to the

1 predetermined change control plans, a plan to continue to monitor. And then, what  
2 happens when you detect you've drifted too much or there's new things that you didn't  
3 expect coming up, like we're discovering with other large language models giving us  
4 crazy responses. Every day we hear about new things. So, I'm just 'plus-one-ing' the  
5 real-world performance evaluation. I hope we could go further with it.

6 Dr. Bhatt: Dr. Shah, Dr. Maddox, and then questions-- Clarifying questions from  
7 the FDA, please.

8 Dr. Shah: So, I'm going to respond specifically to the part about risk and  
9 governance, which is on the screen Aubrey put. I think at the FDA, we have a very  
10 good-- Very good recall, market watch and safety alert system for drugs. Right? And  
11 that seems to be working pretty well. So, one would be for the FDA internally to think  
12 what could be a recall mechanism, what could be a safety alert mechanism and what  
13 could be a market watch, like notifications for people who are deploying these  
14 algorithms, to protect patients and avoid risk? And then, the FDA also launched a  
15 sentinel initiative, which I'm sure many of you are familiar with during COVID, which  
16 actually allowed reverse data collection from the field from patients. So, integrating  
17 some of those existing mechanisms, maybe not as regulatory heavy for some of these  
18 things, but having those safety alerts, etcetera, market watch and sentinel, should  
19 probably be in the framework of risk mitigation if things go wrong eventually,  
20 unfortunately.

21 And one point I want to bring, which is in my opinion not related to this  
22 question, is that Dr. Dreyer mentioned about CPT codes that there are only five CMS  
23 providers underwriting insurance for patients for these algorithms. I think that's a very  
24 big point that eventually patient-centered care is not going to work if you don't have  
25 providers-- If you do have providers being compensated by insurance companies. So,

1 then really what is the delta about helping patients get better care if their insurance  
2 providers are not going to pay for these tests? So, that can be tabled for now, but  
3 something to be considered. Thanks.

4 Dr. Maddox: Tom Maddox from WashU. I don't know if this made it to the list, but  
5 just to make sure, output variability to me is such a unique aspect of these tools,  
6 particularly when it's having a conversation with either a clinician or patient. And I just  
7 think we're going to have to think about how we both assess and control for that. I'm  
8 looking at that final phrase of real-world performance evaluation and I think there's  
9 probably a role in postmarketing evaluation as well, which we'll talk about tomorrow.  
10 Thank you.

11 Dr. Bhatt: Okay, really quick here, really quick there. And then FDA. No more  
12 cards up, please. Thank you.

13 Dr. Radman: I just wanted to add-- I think Tom is getting-- Is riffing on the point I was  
14 making about real-world performance evaluation as well. And I think the output  
15 variability increases the cost and burden of premarket evaluation. That's the issue here.  
16 There's just too much to evaluate. So, it really does necessitate the real-world  
17 evaluation to eventually do it thoroughly if you continuously monitor once it's  
18 deployed. And we could only do the best we can do in the premarket space, but just  
19 really resource constraint now that there's just so many outputs.

20 Dr. Botsis: Taxiarchis Botsis. A very brief comment inspired a little bit by Ms.  
21 Miller's comment about the failure mode. Obviously, we need to know the failures of  
22 the GenAI methodology. Absolutely. I wonder though whether it might be easier to  
23 specify and define a safe mode for those technologies. Just to give you an example, and  
24 imaging is not my expertise honestly, but if we have specific characteristics, image  
25 characteristics that may specify the generation of a text report for that particular image,

1 a CT scan for example, whether that set of features really take us to a safe mode for that  
2 particular technology.

3 Dr. Bhatt: To our FDA colleagues for clarifying questions. And then I'll  
4 summarize.

5 Mr. Tazbaz: Thank you so much. So, we have one clarifying question around when  
6 the Committee is talking about clinical validation, if we could just maybe describe what  
7 you mean by clinical validation, see if we're actually aligning from a regulatory  
8 language perspective. And then, we have a couple more questions around. Do you see  
9 any difference between a generative AI-- From a risk perspective or risk management  
10 perspective, a generative AI tool that is patient-facing versus clinician-facing? That's an  
11 important part of the equation of risk management. And then, there were some  
12 conversations about accountability. So, when something does go wrong, who's  
13 accountable for it? Right? Is it the developer? Is it the hospital system? Is it the  
14 individual clinician? And so, risk management kind of takes a different consideration  
15 there. And this is really not maybe pertaining to this conversation immediately, but I  
16 also would like you all to think about this concept of risk management and the  
17 complexity that we just saw, and even Dr. Jackson talked about multiple risk  
18 management frameworks or governance frameworks. And so, one of the things that  
19 we've been talking about is health equity and this is incredibly difficult to go  
20 implement. And so, how should we be thinking about this more collectively as an  
21 industry around this concept of risk management to ensure that these technology and  
22 tools can be used by everyone, not just the ones that have the infrastructure to be able to  
23 deploy and operate them? Thank you.

24 Dr. Bhatt: Chevron, may I ask you to comment on some of these?

1 Dr. Rariy: Yeah, absolutely. I was just trying to put my thoughts together, but I  
2 think obviously they're excellent questions around the risk management. I remember  
3 there were a few speakers who brought up the IRB process in your traditional clinical  
4 setting and it's more locally based. And I think we would benefit from considering  
5 something of that nature as well, not only a governance process, but some sort of  
6 regulatory oversight compliance, etcetera, that is locally based, and it will lend itself  
7 well to risk management in a similar way as the current state of risk management on the  
8 local level. So, I think there's benefit in having multiple layers, if you will, of risk  
9 mitigation strategies embedded within the-- In a systematic way. And part of that  
10 governance and risk mitigation strategy, perhaps the FDA can lend itself to creating at  
11 least an overarching scaffold of a framework. And that framework could include high  
12 level-- We saw the seven areas including equity as one, but at least providing various  
13 localities with the opportunity to adopt this framework and expound on it I think will  
14 help with accountability as well at the local level and perhaps provide more rigor in  
15 clinical care.

16 Secondly, I think it goes to your question, Mr. Tazbaz, around the patient-facing  
17 or direct to patient-facing generative AI versus direct to clinician. I think the key  
18 consideration here, as we heard it from some others, is that if we are reliant on expert  
19 interpretation for whatever medical specific use case of the generative AI output, then  
20 we need to take that into consideration that as a lay patient, they will not have that  
21 expert medical consideration. And so, they are two very distinct use cases that likely  
22 require two very distinct considerations.

23 Dr. Bhatt: Thank you. Going back to your patient-facing versus clinician-facing  
24 question, I just want to echo what Dr. Clarkson said just a little while ago. At the end of  
25 the day, does it lead to patient harm? I think sometimes we carry the assumption that

1 something that is clinician-facing may not directly harm the patient, and that's one of  
2 the risks we run in that terminology. And so, I think I would just make sure that we  
3 extend that to how does it affect the patient, whether the device is facing a patient or a  
4 clinician? Let's do Dr. Radman. Dr. Khubchandani, it would be great to hear from you  
5 after Dr. Radman, and then Dr. Jackson.

6 Dr. Radman: I just have a simple comment about the patient-facing that-- And I guess  
7 we were talking about-- I think Peter mentioned maybe we need a model for this group  
8 or this case and so forth. And so, we should also consider, while we want to protect  
9 patients and everything, there's also a case to be made for the benefits that patients may  
10 get, like consider a rural population that lives two hours away from the closest  
11 specialist, and whereas a generative AI can be a surrogate for that. So, there's a balance  
12 between the risks and the benefits. And the benefits can be far-- Increase for just  
13 thinking this particular rural populations, but I'm sure we could think of other use cases  
14 like that, whereas if you compare a benefit where a hospital is rolling out AI to cut costs  
15 and become more efficient or something like that, that's great, but it may be a lower  
16 benefit level there and therefore the risks should be accounted for accordingly.

17 Dr. Bhatt: That's great. Thank you. I'll piggyback off that for a minute and then go  
18 to Dr. Khubchandani. When thinking about the health equity part of your question,  
19 Troy, one challenge we face, because generative AI is oftentimes related to clinical  
20 guidelines or information or clinical decision support, is what the gold standard is. And  
21 when we think about whether or not we're delivering gold-standard clinical guideline-  
22 derived treatment throughout the United States throughout all of the different  
23 specialties, the answer is generally no. And so, it's how close does generative AI bring  
24 us to the gold standard versus how much better will the care get versus what it is now?



1 And so, just based on what Dr. Radman said, I think that's an important differential as  
2 well for us to think about.

3 Dr. Khubchandani and then, Dr. Jackson. And Troy, we haven't forgotten that  
4 you asked about clinical validity and accountability.

5 Dr. Khubchandani: Yeah. I think from the public health perspective, risk  
6 management-- When we submit IRB applications to do a study, we take the liability, the  
7 university takes the liability. Same for the healthcare system. They had a choice to adopt  
8 an AI technology or not to adopt it, so they would be in a way liable. My other  
9 comment is regarding how I see generative AI. It's like-- As Dr. Clarkson said, how do  
10 you define risk? So, I bought a wearable device. I did not lose weight. Risk. I lost some  
11 money. I could have invested it in a gym. But it kept misinforming me about my heart  
12 rate and I died. Now, that's a different type of risk. And so, risk has to be defined. Risk  
13 assessment has to be defined and tied up with risk management. That's my comment  
14 here. Thank you.

15 Dr. Bhatt: Dr. Jackson.

16 Dr. Jackson: I wanted to speak to the accountability question. I think it's a really  
17 interesting question, and I'll be honest, I'm not sure that it's necessarily fully under the  
18 FDA's purview, right? There's an article in the New England Journal of Medicine from,  
19 I think, January, February of this year, and they were looking at generative AI and  
20 accountability and liability for providers, and they laid out multiple legal cases of when  
21 software's medical devices come up and who ended up being accountable and even  
22 liable. And I think risk management and liability are two different things, but that's  
23 what they were talking about. How do we differentiate between the two? And it showed  
24 that every time there was a different answer, depending on what was happening. And  
25 so, I don't know that we can have a definitive accountability in this case, but thinking

1 about it from a health equity lens, which for me has been very important in thinking  
2 about all these discussions. I wonder if we could borrow from some of the value-based  
3 care work that has been done in thinking about shared risk management and what that  
4 looks like when a healthcare system or a practice or a clinician decides to deploy GenAI  
5 in their work, and how do we look at shared risk management as a part of that that  
6 incentivizes all parties to make sure that we're looking at safety and quality making  
7 decisions? And I think that is a definite pathway to also improving health equity if  
8 everyone is ensuring that. Because I think one of the things that could happen is people  
9 will decide, if the risk management is too high, not to deploy some of these tools, which  
10 would increase disparities because they feel like they have too much accountability.

11 Dr. Bhatt: Thank you. Dr. Kukafka.

12 Dr. Kukafka: I just wanted to talk a little bit more about the patient-facing that you  
13 brought up in the context of risk mitigation. I don't know if the FDA considers  
14 exasperating health disparities as a risk or is that out of context, out of your purview? I  
15 mean, my point is if we do patient-facing, we have to design for different levels of  
16 literacy and numeracy to be-- So that we don't exasperate disparities. Interesting. I just  
17 completed a study where patients were able to download their own EHR data, review  
18 their data, and we found that the patients with lower literacy levels, lower education  
19 levels were more likely to do that than patients more educated. Now, we have  
20 hypotheses that they might be more-- Higher level of information-seeking because they  
21 do have issues with obtaining and interpreting health information. But again, is that a  
22 risk for the FDA to consider exasperation of health disparities? That's a question. And  
23 then, if the answer is yes, then the usability has to address that.

24 Dr. Bhatt: I think that's an incredibly helpful comment. And I would say I think it's  
25 essential for us as a Digital Health Advisory Committee, not just today, for our other

1 sessions as well, to think about digital and health literacy for the patient-facing parts of  
2 whatever we are going to look at, whether it's generative AI or something else in the  
3 future. So, thank you for bringing that up, and I think it dovetails nicely with what Dr.  
4 Jackson was saying as well about exacerbating the discrepancy in access or usability,  
5 and we can probably break those down as we think about this further. Dr. Elkin.

6 Dr. Elkin: Yeah. I'd just love to build on that. It touched something in my heart as a  
7 clinician. Many years ago, when people started getting information on the Internet,  
8 they'd come into your office, and some of my colleagues really didn't like that, but I  
9 loved it because it was getting involved with their own healthcare. And the more they  
10 were involved in their healthcare, the more we could have a great conversation and  
11 draw them in and try to help solve the problems that they're coming in with. And it  
12 turned out that, at least from my practice, this was very helpful, even if the information  
13 that they were worried about didn't really apply to them. I could explain why and  
14 alleviate fears and tell them when they needed to come in and encourage them to  
15 continue to be involved in their own healthcare. And I see the ubiquitous nature of these  
16 devices where they can be on your phone. Right? It doesn't have to be a very expensive  
17 computer that you have. It's going over the Internet and the main hardware that's  
18 running is in the cloud. So, it could give people that don't have access to high quality  
19 medical information much more medical information for them to bring to their  
20 clinicians to get involved in their care and proactively, I think that will help to increase  
21 the quality of care and decrease disparities.

22 Dr. Bhatt: Okay. Mr. Radman, maybe last comment. I'm going to look over at Dr.  
23 Clarkson and Diana Miller for any final comments they may have. And then I will read  
24 to you our concerns, but not our summary in the interest of time, Mr. Tazbaz.

1 Dr. Radman: Thank you. Regarding accountability, I think there should be strong  
2 requirements for labeling studies because, for example, chat GPT already has a  
3 disclaimer that it's not a medical device, but we've already seen data that shows its  
4 adoption and use by clinicians already as a second opinion. And we all know how  
5 ineffective the listing of the side effects on drug commercials can be at the end of those  
6 commercials. So, I think if-- We discussed already the importance of having some kind  
7 of label saying, "This output is generated by a generative AI or large language model,"  
8 and maybe even saying-- I've seen people recommending-- "Which has a propensity to  
9 hallucinate." Perhaps even we could have a premarket evaluation to determine  
10 hallucinations occur X percentage of the time or something like that. But do people  
11 understand that label? Whatever the label says. So, that's really important to verify with  
12 an adequate study with a diverse user group that reflects who's going to be using the  
13 device.

14 Dr. Bhatt: Great, thank you. Let's do a rapid fire. Just a quick round. Does anybody  
15 have a 30-second answer to clinical validation? Who used that phrase? And do you want  
16 to claim that you didn't explain to Troy what you meant? No?

17 Dr. Maddox: I actually think it came up in Dr. Dreyer's presentation. So, first of all,  
18 can you just expound on your question real quick to make sure we're answering it?

19 Dr. Diamond: Dr. Dreyer, I think you were saying that you didn't see a lot of value in  
20 some types of comparator studies, and instead we're advocating for a different type of  
21 clinical validation. Could you explain?

22 Dr. Dreyer: Yeah. I'll talk about this in the open part tomorrow to talk about  
23 monitoring as well where I'll get to this. But if you think of software testing, just  
24 straight software testing, the next step would be you would do a performance test on  
25 that software. That means that you have data, but you don't have experts. The next thing

1 you would want to try and do is consider doing a comparative test with experts and the  
2 data to say with and without the AI. Next, as we deploy every time we deploy AI, we do  
3 clinical validation. What I mean by that is think of it as a one-time monitoring. So, we  
4 deploy this in a mode where we can see the exact output of what's coming at the facility  
5 using our data and our experts with the algorithm. So that's the distinction between a  
6 comparative effectiveness of premarket testing with a pre-deployment testing at a  
7 facility. And then, as that continues on, either intermittent or continuous, that's the  
8 monitoring process. So, what I mean by clinical validation or validation of the test is  
9 that it's validated on that site to prove that it is at a level of acceptability at that site.  
10 And I make that distinction also because what we've seen continuously with the  
11 previous models, and I think it's going to be the same here, is that because this works at  
12 five sites, it doesn't mean it'll work at the next five. And then the question that comes,  
13 as we've discussed, is if it's working at 50 sites but it's not working at 10 sites, is it  
14 recalled or is it not? Is that a good thing or a bad thing? And so, I don't think you can  
15 premarket decide that. I think you have to validate at each facility to be able to make  
16 that distinction.

17 Dr. Diamond: So, you're talking about site-specific validation in the context of  
18 real-world use.

19 Dr. Dreyer: Or even pre-real-world use, but essentially real-world evidence by not  
20 maybe the clinical deployment of it, but the preclinical deployment of it to test it with  
21 site users, site patient data and the algorithm.

22 Dr. Bhatt: Any last 30-second comments from the Panel? Ms. Miller.

23 Ms. Miller: I want to address a little bit about the patient-facing and the access to  
24 healthcare. And I think there are two important aspects to think about when we design a  
25 medical device targeting patients. You obviously have to think of different personas and

1 who do you design the device for? And I think it's a real benefit, and that's what we  
2 mean by intended use population target. But a real good benefit that we still need to be  
3 mindful of is the access to healthcare. It's increased ever since phone's been around.  
4 There are medical devices, applications on the phone that contribute to a lot of benefits  
5 to the patient's access. And even though maybe we don't get to try it first-- We keep  
6 iterating onto it. And I think the FDA acknowledges this with the PCCP framework and  
7 all of that, that allows faster iterations on these devices because something we need to  
8 think about when we design for patients is that we need to have a consumer mindset.  
9 And in a medical device industry, we're a little handicapped in a way that we can't  
10 really act that fast like tech companies and we cannot deploy that fast our medical  
11 device apps to those devices and technologies. So, we really appreciate the PCCP, and  
12 we need to realize that we need to give the same experience like the Apples of the  
13 world-- I don't want to call names, but the big tech companies. We need to be on the  
14 same bar and that's what it's hard to strive for.

15 Dr. Bhatt: Okay. Thank you so much. So, in the interest of time, I have the  
16 documentation here of the summary that I will hand in for the record. And I will leave  
17 us just with some concerns that the Panel has according to this question that include the  
18 following. This is a new way of presenting information that seems more humanistic and  
19 gives the impression of real intelligence. That is concern number one. Number two is,  
20 whereas we are thinking about national guidance, local governance may be most  
21 important and value sets appear to be local. The third is clinicians are busy and may not  
22 engage in training the way we are expecting. So, our reliance on training to ensure that  
23 generative AI systems work may not work as intended and it may be challenged. The  
24 risk of patient harm must be central to this discussion. There's an increased risk of lack  
25 of accountability in healthcare with the addition of generative AI to care delivery.

1 Governance may need to differ across different types of generative AI. This may be  
2 differentiated across type of input, type of model and potential clinical outcomes. And  
3 there's a suggestion that we could perhaps use infrastructure for drugs and other devices  
4 that already exist at the FDA and expand them to generative AI-related devices. Digital  
5 health and literacy is essential to include in risk evaluation, monitoring and mitigation if  
6 we are to keep health equity central to our efforts. And although accountability may not  
7 be under FDA purview, responsibility may lie with health systems that deploy GenAI  
8 devices or the companies that create them and others at this time, because there is not a  
9 large body of clear precedents. Shared risk management may likely be the best concept  
10 for now as it will incentivize all parties to look at risk mitigation as their responsibility.

11 Mr. Tazbaz, is this an adequate summary for now? Thank you.

12 Day 1 Closing Remarks

13 Dr. Bhatt: With that, I would like to thank the Committee and the FDA for their  
14 contributions. I would also like to thank again the open public hearing speakers, patient,  
15 industry, healthcare provider, academia and FDA for their remarks.

16 And before we adjourn, I want to ask the FDA staff seated at the table if they  
17 have any final comments.

18 Mr. Tazbaz: Thank you, Dr. Bhatt. I would like to echo the same gratitude towards the  
19 presenters, the public commentary and of course the Committee itself, but also to you,  
20 Dr. Bhatt, for running a very tight and productive session. So, we're looking forward to  
21 seeing-- We've been able to capture several action items that could come out of it, but  
22 we're looking to collaborate further to drive more of these action items that we might be  
23 able to tackle collectively. So, thank you.

1 *Adjournment*

2 Dr. Bhatt: Great. Thank you. The November 20th session of the Digital Health

3 Advisory Committee is now adjourned. We'll continue this agenda tomorrow at 9:00

4 a.m. sharp. Have a wonderful night and thank you.