

CIRDS: CTP Integrated Research Data System

Mark Gingrass¹, Vibha Kumar¹, Christine Wang¹, Wei Chen², Geetha Potluri², Alex Abramson², Nirmala uralikrishna², Wei Wu³
FDA/CTP/Division of Research Science Informatics¹, SAIC², MarkLogic³

CHALLENGES / BUSINESS NEEDS

Distributed information across both internal and external sources creates a significant challenge for the Center for Tobacco Products (CTP) personnel responsible for the Agency's tobacco regulation. Specifically, the CTP researchers and scientists need to

- search heterogeneous data sources for comprehensive studies (web URLs, different searching capabilities and configurations, import/export formats, etc.),
- access disparate systems use different technologies to express similar concepts (syntax and semantic level),
- maintain advanced search capabilities across heterogeneous applications, which leads to higher overall system costs, and
- find a User Interface (UI) and host environment for the iDAT (Industry Analysis Document Tool) data given that its Web UI will be decommissioned – this also posts technical challenges of ingestion of new data sources.

GENERAL INTRODUCTION

CTP Integrated Research and Data System (CIRDS) has been developed to address the above challenges. CIRDS employs MarkLogic's Enterprise NoSQL-based operational and transactional data integration platform to provide a Web-based UI for all CTP researchers and scientists to explore different data sources in a unified way.

CIRDS has gone through the following stages:

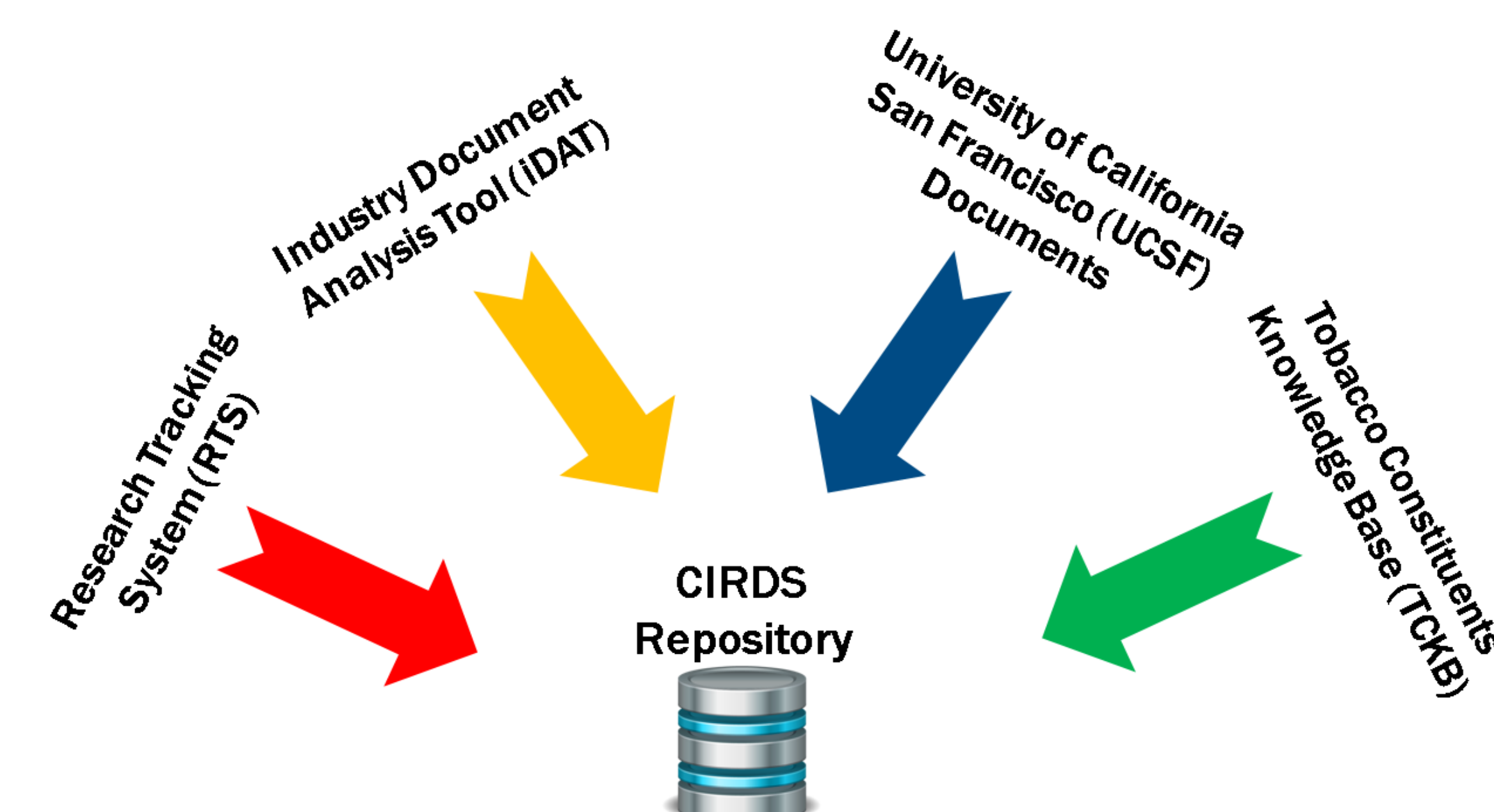
- From 2016 to 2019 CIRDS has been designed and developed as a Proof of Concept (POC) solution.
- In the past year CIRDS has been revised and improved to a state that paves a solid groundwork for advancing itself as a production system.
- Currently, CIRDS is in a transition process from POC to production.

OBJECTIVE

The objective is to introduce CTP's Integrated Research Data System (CIRDS) platform that allows users to search, explore, discover and disseminate information from multiple data sources.

So far, four disparate data sources have been ingested:

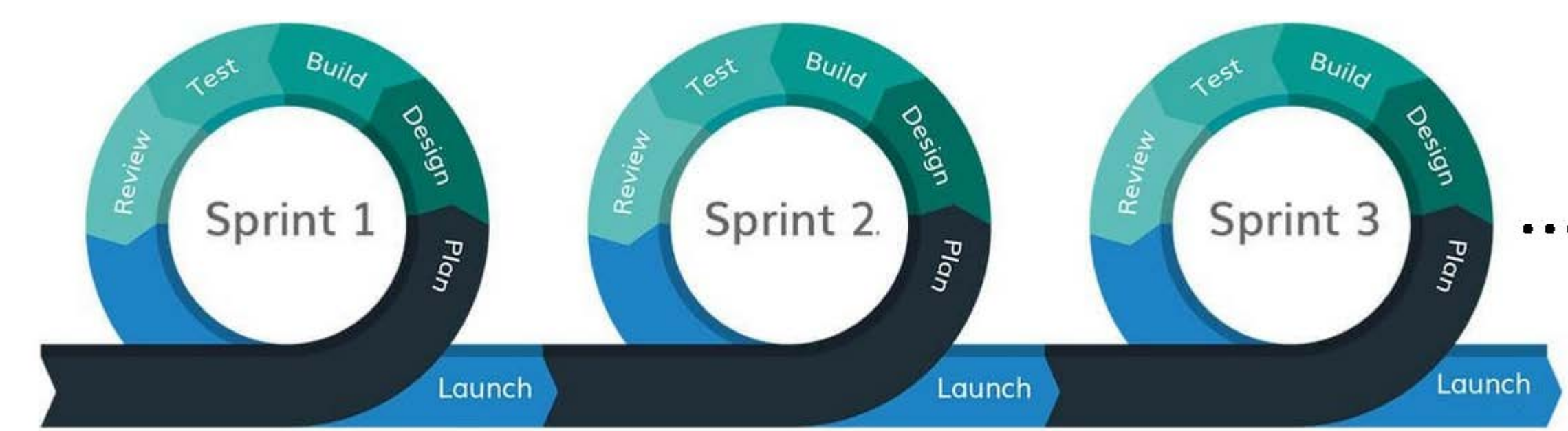
- CTP's Internal Research Tracking System (RTS),
- CTP's internal Tobacco Constituent Knowledgebase (TCKB),
- CTP's internal Industry Document Analysis Tool (iDAT),
- and the publicly available Truth Tobacco Industry Documents from UCSF.



SOLUTION APPROACH

Agile Methodology

CIRDS development and requirements gathering follow an Agile approach with small steps starting from November 2016 until September 2019. The development effort between September 2019 to September 2020 has been extremely rigorous following an Agile Scrum methodology with close communications with the COR and the CIRDS and iDAT users, streamlined tasks, adequate and dedicated QA/testing, ambitious and successful software releases, and the adaption to changing requirements.



CIRDS High Level Architecture

The CIRDS integrates multiple data sources into one repository employing multiple workflow steps:

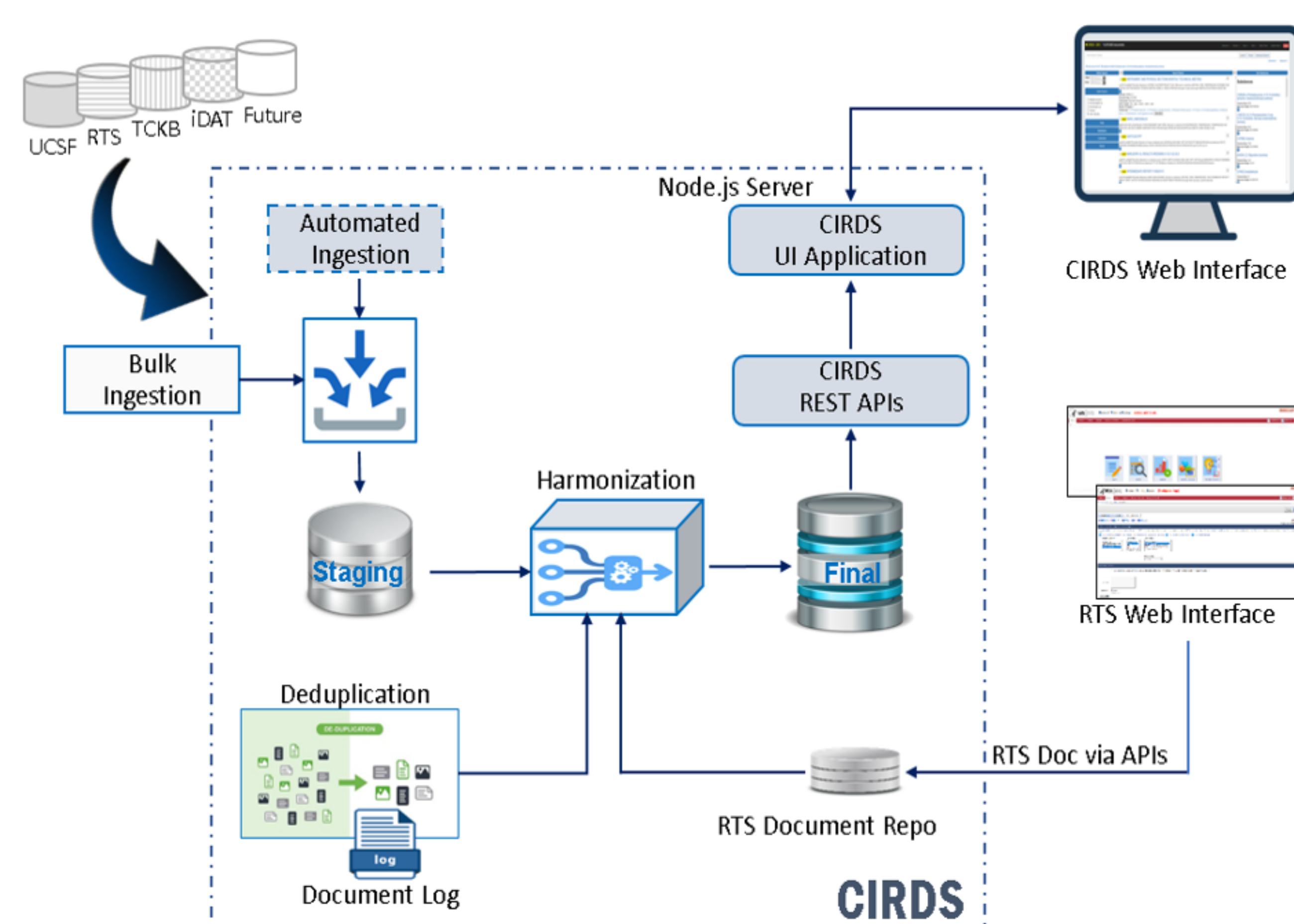
Data Harvesting- identifies the data source(s) and individual data sets within the source repository that is targeted for the system.

Data Ingestion- the data sets from the various sources are input as JSON formatted data records (a format for structuring data transmitted between a server and a web application). These records are stored in a "staging" database.

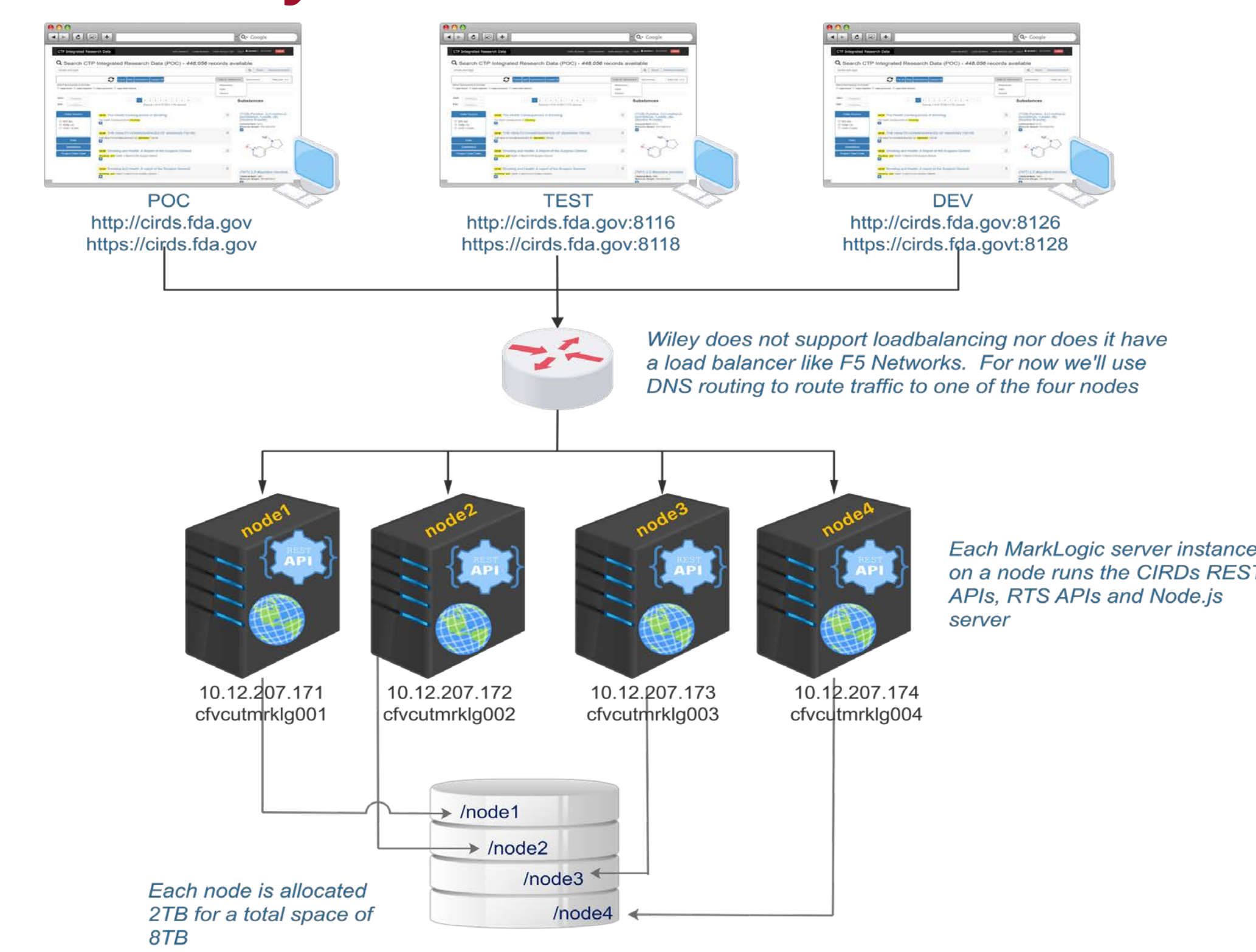
Harmonization- includes adding common fields for documents/records and identifying and indexing substances that are mentioned within each of the RTS projects and UCSF document records.

Cured Data Repository-the harmonization step produces the CIRDS Final Repository that is employed by the CIRDS Proof of Concept (POC) application. This repository is where records are stored for the query process and allow users to retrieve the records after completing a search.

CIRDS Components and Data Flow



CIRDS System Infrastructure

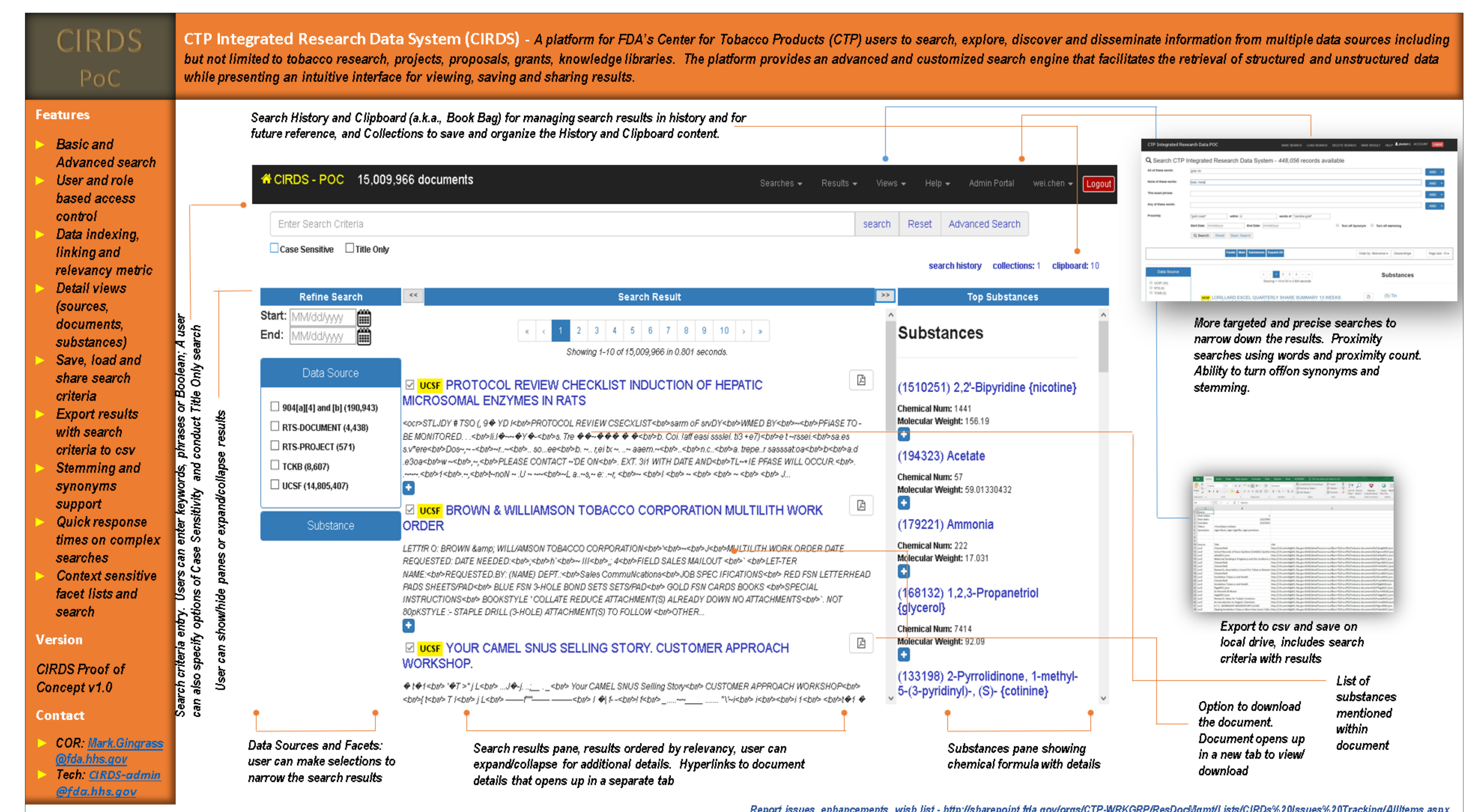


RESULTS

Features

- Unified UI displays all the data sources (pre/post search)
- Powerful search capability with customizable synonyms list and flexible options including Clipboard and Collections, etc.
- Effective system monitoring, management and security
- Highlighted search results with ability to Save and Export to MS Excel or Endnote.
- Developed a streamlined data ingestion process
- Employed MarkLogic as the integrated document search management system
- Adopted role-based user management

CIRDS Interface



CIRDS PoC
CTP Integrated Research Data System (CIRDS) - A platform for FDA's Center for Tobacco Products (CTP) users to search, explore, discover and disseminate information from multiple data sources including but not limited to tobacco research, projects, proposals, grants, knowledge libraries. The platform provides an advanced and customized search engine that facilitates the retrieval of structured and unstructured data while presenting an intuitive interface for viewing, saving and sharing results.

Features

- Basic and Advanced search
- User and role based access control
- Data indexing, linking and relevancy metric
- Detail views (sources, documents, substances)
- Save, load and share search criteria
- Export results with search criteria to csv
- Stemming and synonyms support
- Quick response times on complex searches
- Context sensitive facet lists and search

Version
CIRDS Proof of Concept v1.0

Contact
COR: Mark.Gingrass@fda.hhs.gov
Tech: CIRDS-admin@fda.hhs.gov

Search History and Clipboard (a.k.a., Book Bag) for managing search results in history and for future reference, and Collections to save and organize the History and Clipboard content.

Search Results: 15,009,966 documents

Substances list:

- (1510251) 2,2-Bipyridine (nicotine)
- (194323) Acetate
- (179221) Ammonia
- (168132) 1,2,3-Propanetriol (glycerol)
- (133198) 2-Pyrolidone, 1-methyl-5-(3-pyridinyl)-, (S)- (cotinine)

Export to csv and save on local drive, includes search criteria with results

List of substances mentioned within document download

Option to download the document. Document opens up in a new tab to view download

Search results pane, results ordered by relevancy, user can expand/collapse for additional details. Hyperlinks to document details that opens up in a separate tab

Substances pane showing chemical formula with details

More targeted and precise searches to narrow down the results. Proximity searches using words and proximity count. Ability to turn off synonyms and stemming.

Report issues, enhancements, wish list - <http://sharepoint.fda.gov/govcs/CTP-WRKGRP/ResDoc/Item/Lists/CIRDS%20Issues%20Tracking/AllItems.aspx>