# Agenda

1. *CCDI Data Federation Resource: Background*

2. *Data Harmonization: Aligning the Data to NCI Standards*

3. *Data Federation Resource API: Designing the API*

4. *Using the Federation Data Demo: Leveraging Jupyter Notebook*

5. *Future Applications*

6. *Q&A*

# Today's Speakers



**Geoff Lyle**
Technical Project
Manager, Treehouse
Childhood Cancer
Initiative

**Clay McLeod**
Director, Product
Development and
Engineering, St. Jude
Children's Research
Hospital

**Martin Ferguson**
External Consultant,
National Cancer Institute

**Allison Heath**
Director of Data
Technology and
Innovation, Children's
Hospital of
Philadelphia

# CCDI Data Federation Resource: Background

*Geoff Lyle*

# Why a Federated Childhood Cancer Data Ecosystem?

- Pediatric cancer data are currently siloed

  - Reduced and delayed access to data

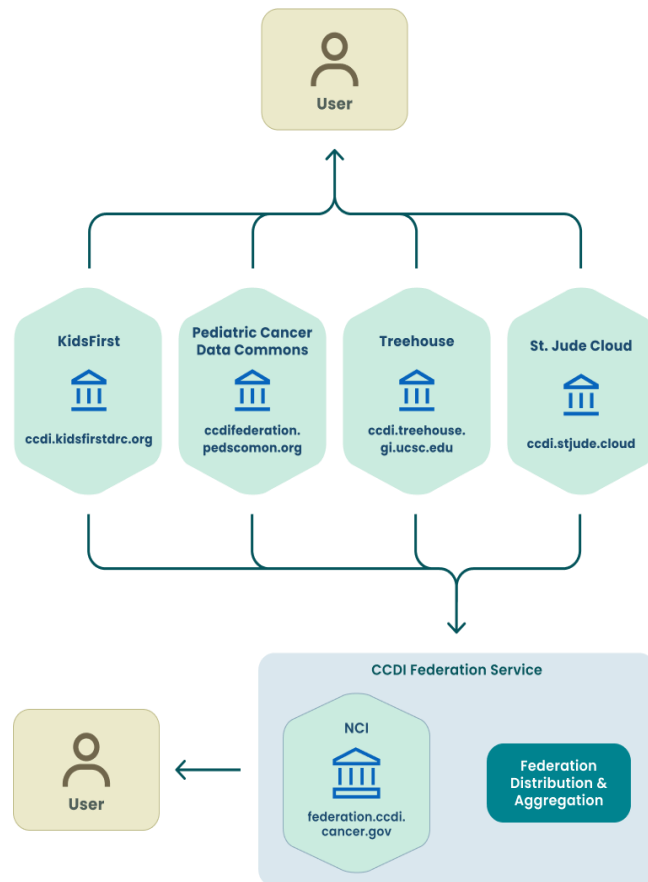  - Missed therapeutic opportunities

# Data Federation Objectives

**Facilitate large-scale biomedical research via a federated, real-time data search API**

- Develop and implement a common application programming interface (API) specification where deidentified, participant-level data from each member can be queried.

- Results from queries return responses, leveraging an ever-growing, harmonized set of metadata values.

- Rich, faceted search across the supported information.

- Data will not be moved or centrally warehoused; instead, users can access the data where it resides.

NIH NATIONAL CANCER INSTITUTE

# Data Federation Status

- Four current federation members:
  - Kids First Data Resource Center
  - Pediatric Cancer Data Commons
  - Treehouse Childhood Cancer Initiative
  - St. Jude Cloud
- Version 1.0 API implementation and demos
  - If you're interested in joining, please email NCIChildhoodCancerDataInitiative @mail.nih.gov

# Driver Scientific Use Cases

| Scientific Use Case | Description |
|---|---|
| **1. Disease and Genomic Variant Querying** | Search for diseases or genomic variants to gather data on alterations, uncertain variants, or mutations |
| **2. Participant Cross-reference System for Data Retrieval** | Retrieve comprehensive clinical and genomic information across institutions |
| **3. RNA Sequencing Map t-SNE** | Generate a global gene expression map using t-distributed stochastic neighbor embedding (t-SNE) analysis, providing visual insights into RNA expression patterns worldwide |
| **4. Flexible, Tiered User Query Handling** | Allow users to submit general queries and explore detailed matches to simplify data exploration and retrieval |
| **5. Harmonizing Multiple Batches of Samples** | Utilize internal patterns within datasets to align multiple batches of samples for a comprehensive analysis |
| **6. Biospecimen Metadata Integration** | Incorporate biospecimen information as an additional layer, expanding beyond clinical features and data types for users |

# Implementation Approach

Established two working groups:

1. **Data harmonization** – Ensures federation member data are harmonized to NCI standards

2. **API development** – Determines the best methods for delivering data that is accessible and useful to users querying information

   o Documents implementation guidelines for federation members to share data via the open API

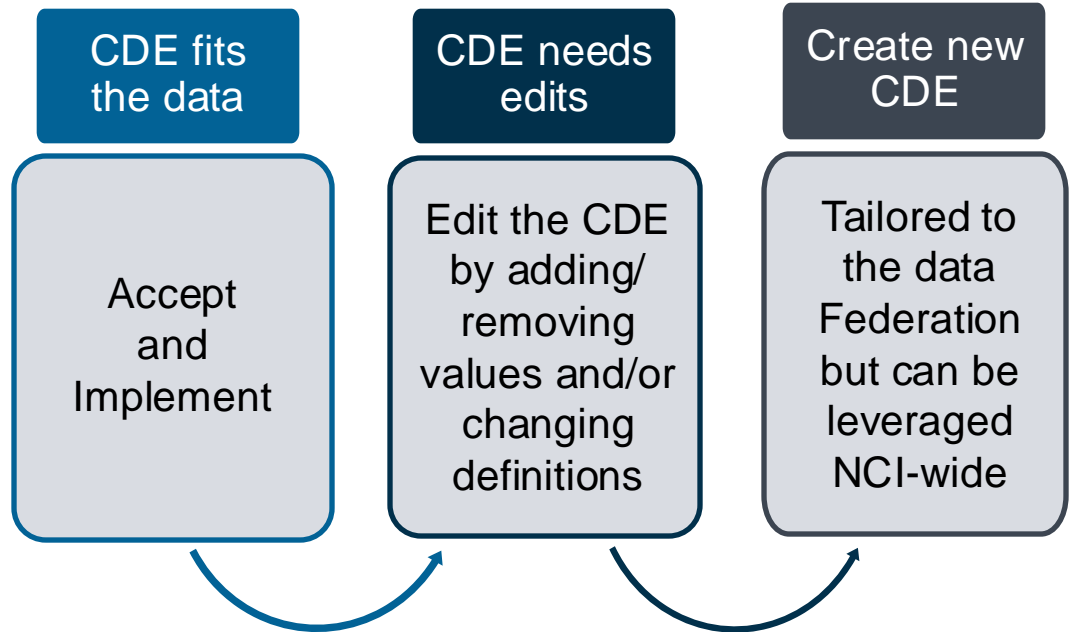# Data Harmonization: Aligning the Data to NCI Standards

*Geoff Lyle*

# V1 Federation Resource Data Summary

Through the API, the federation delivers **metadata** that will help users create a synthetic cohort across multiple institutions/data types.

| Source | Data Types | Data Level | Subjects (Participants) | Samples | Files |
|---|---|---|---|---|---|
| St. Jude Cloud | Genomic data, gene expression, imaging | Participant | 13,956 | 19,866 | 133,579 |
| UCSC - Treehouse | Genomic data, gene expression | Participant | 12,483 | 12,770 | 6 |
| Kids First - CHOP | Genomic data, gene expression, clinical data, imaging | Participant | 34,066 | 162,549 | 327,893 |
| PCDC - UChicago | Clinical data, imaging | Aggregate | 22,667 | | |

# Data Harmonization Approach

- Leverage existing CCDI standards

- Utilize caDSR Common Data Elements (CDEs) to map attributes and allowable values ([cadsr.cancer.gov](cadsr.cancer.gov))

- Develop harmonization guidelines when no NCI standard exists

- All discussions are tracked on [GitHub](GitHub)

| CDE fits the data | CDE needs edits | Create new CDE |
|---|---|---|
| Accept and Implement | Edit the CDE by adding/ removing values and/or changing definitions | Tailored to the data Federation but can be leveraged NCI-wide |

# V1 Federation Resource Scope – Harmonized Fields

Common Data Elements (CDEs) from Cancer Data Standards Registry and Repository (caDSR) (https://cadsr.cancer.gov)

| Subject (Participant) | Sample | Study and File |
|---|---|---|
| Sex (6343385) | Sample tumor status (5432687) | Study short title (11459812) |
| Race (2192199) | Tumor classification (12922545) | Study name (11459810) |
| Ethnicity (2192217) | Age at diagnosis (3225640) | dbGaP phs accession (11524544) |
| Vital status (2847330) | Age at collection (14473376) | Institution (12662779) |
| Age at vital status (14480965) | Library strategy (6273393) | File location (Link/Gateway) (11556141) |
| Subject ID (6867052) | Preservation method (8028962) | File description (11280338) |
| | Disease diagnosis (ICD-O; WHO CNS5) | File size (11479876) |
| | Disease phase (12217251) | md5sum (11556150) |
| | ICD-O morphology code & term (11326261) | File type (11416926) |

# Data Federation Resource API: Designing the API

*Clay McLeod*

# API Strategy

- **Select a standard that is purpose-built** for indexing the specific types of data we wanted to share (no enforced metadata standards).

- **Chose a scalable foundation** that will work for hundreds of thousands of samples and millions of files from day one.

- **Enable relaying information provided by source servers** with as little on-the-fly transformation as possible to enable high-performance aggregator services (e.g., the NCI aggregation server).

- **Ensure that joining the federation is as accessible as possible.**

  o The specification should not be more complicated than necessary.

  o The specification can be readily implemented using multiple open-source frameworks.

# API Strategy (cont.)

- **Considered three strategies/standards for creating an API:**

  o FHIR API

  o Beacon V2 API (GA4GH standard)

  o Bespoke API

- After careful review of the existing standards, discussion with the specification designers (for Beacon only), and internal discussion, we jointly decided that, today, **neither FHIR nor Beacon V2 met all our criteria**.

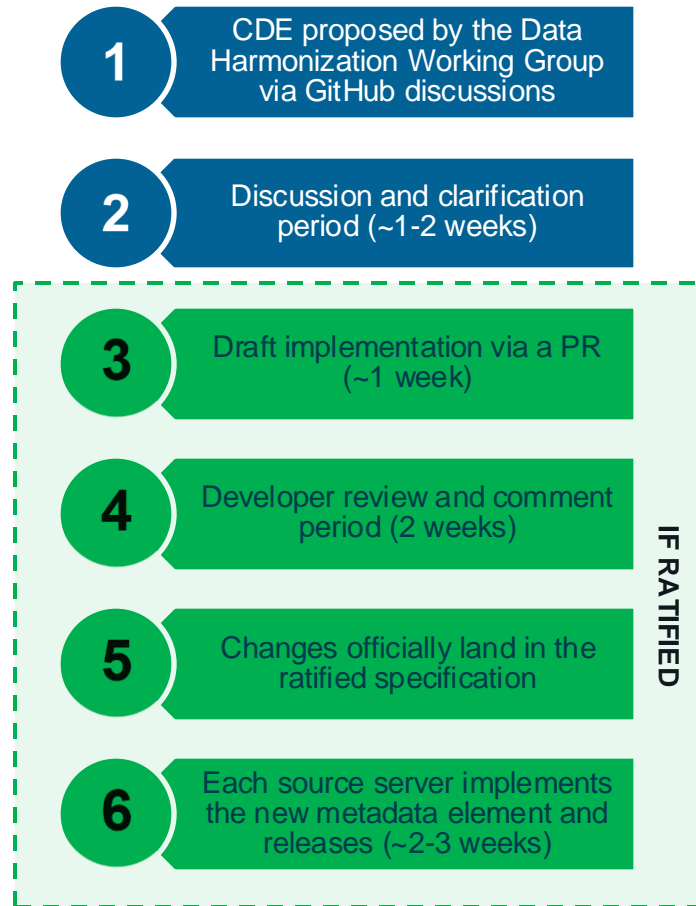- Given this, we decided to start by **creating a simple, bespoke API** that was purpose-built for indexing federation data and metadata.

# Development Approach

- The specification itself was designed to be robust and stand the test of time.

  - The OpenAPI specification itself, and all other source code, was written in Rust.

  - An example server with fake data that implementors can refer to.

  - A toolset for implementations to test compliance with the spec.

  - Testing to ensure the spec remains well-formed on every proposed change.

- All discussions regarding both the design of the API are open and searchable in the federation GitHub repository.

# Development Approach (cont.)

- GitHub was leveraged during metadata ratification, specification design, and development.

  - All CDEs followed a proposal to implementation lifecycle that was about ~1 month end-to-end.

- Slack was used for real-time, informal discussions amongst stakeholders as well as implementors.

- Each federation member, as well as the NCI aggregator, implemented the specification using their own framework/infrastructure.

**1** CDE proposed by the Data Harmonization Working Group via GitHub discussions

**2** Discussion and clarification period (~1-2 weeks)

**IF RATIFIED**

**3** Draft implementation via a PR (~1 week)

**4** Developer review and comment period (2 weeks)

**5** Changes officially land in the ratified specification

**6** Each source server implements the new metadata element and releases (~2-3 weeks)

# Information About the API

- Future blog post from API authors diving into the specification in-depth

- CCDI Hub (https://ccdi.cancer.gov/explore)
  - Information about the project
  - Link to GitHub

- CCDI Data Federation Resource GitHub (https://cbiit.github.io/ccdi-federation-api-aggregation/)
  - OpenAPI Specification
  - GitHub Wiki for metadata descriptions
  - GitHub issues/discussions for questions
  - Links to participating nodes API spec



https://ccdi.cancer.gov/data-federation-resource

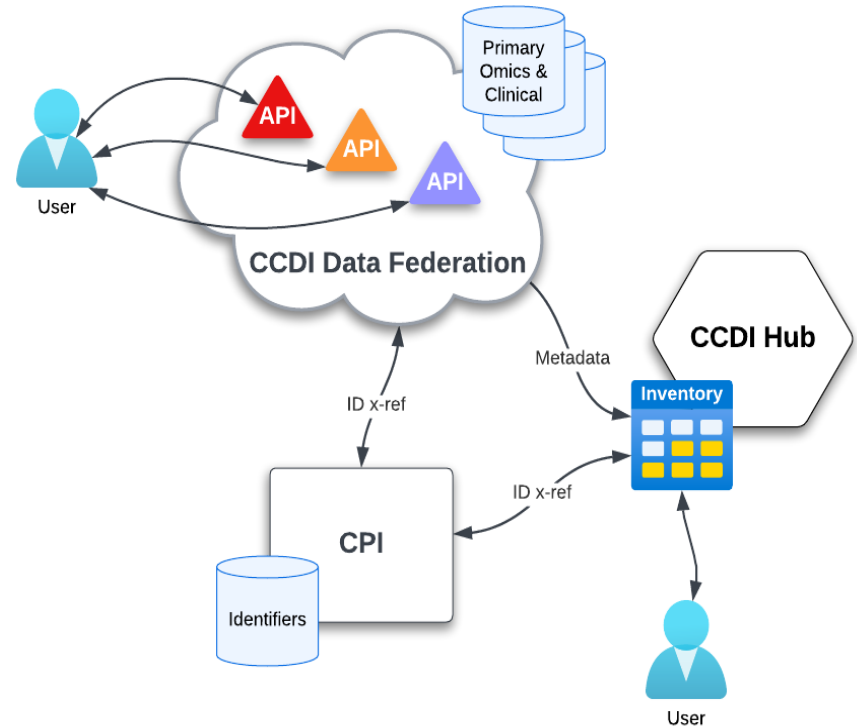# Using the Federation Data Demo: Leveraging Jupyter Notebook

*Martin Ferguson*

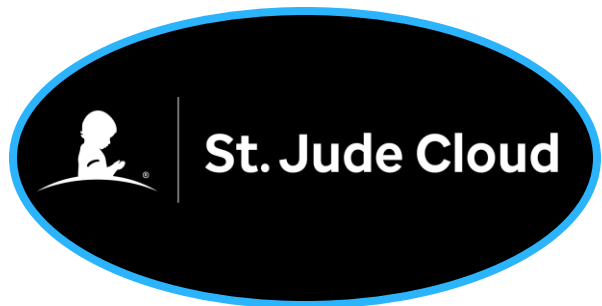# Future Applications

*Allison Heath*

# Further Enabling Scientific Use Cases: Childhood Cancer Data Initiative Participant Index (CPI)

- Link participants' data from the Federation and CCDI across:
  - Time (longitudinal)
  - Space (institutions, studies, trials)
  - Modalities (clinical, molecular, imaging)
- Minimize double counting
- Create a "Cohort of One"
  - Integrate data for a patient across federation
  - Live updates of new data for patients of interest

# API-based Integration for Discoverability

 = CCDI federated data of interest
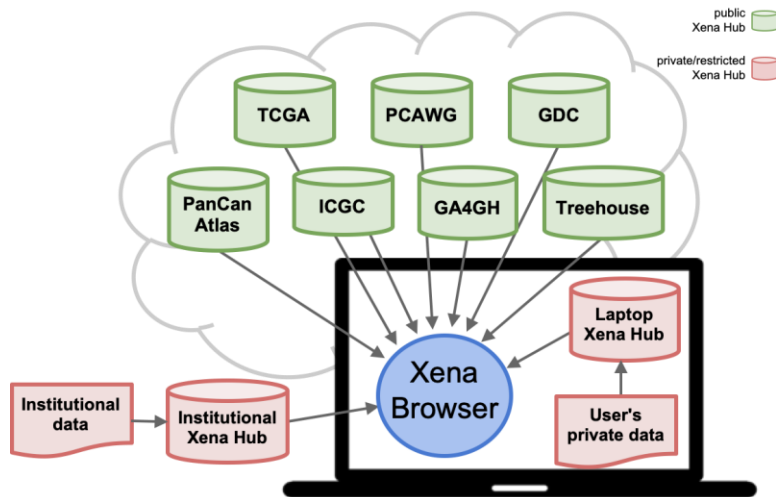

St. Jude Cloud


DATA FOR THE COMMON GOOD
Pediatric Cancer Data Commons

**Resources interested in CCDI data**


Gabriella Miller
Kids First
PEDIATRIC RESEARCH PROGRAM
Data Resource Center


Treehouse
CHILDHOOD CANCER INITIATIVE

# API-Based Integration with Local Data for Analysis

- Empower individuals to analyze their tumor data using the full scope of the CCDI federation data

- Data analysis tools utilizing CCDI API for rapid analysis by the research community

- Enable users to upload data and compare it to federation data (similar to matchmaker exchange

NIH NATIONAL CANCER INSTITUTE

# RNA-Seq "World Map": Foundational Layer for Scientific Use Cases

- Generate a gene expression "world map" using RNA-Seq data from all partners

  o Enables finding tumors that are "acting similarly" from an expression perspective

  o Individual institutions have already found these tools useful in diagnosing rare diseases as well as checking for high-quality data space (institutions, studies, trials)

- API-based mechanism would enable ongoing map updates and refinement as new data is generated

# Q&A

# Contacts

- Leveraging the federation resource

  - Email: [NCIChildhoodCancerDataInitiative@mail.nih.gov](mailto:NCIChildhoodCancerDataInitiative@mail.nih.gov) with questions related to CCDI federated data or accessing the CCDI Data Ecosystem

- Questions related to individual APIs

  - St. Jude Cloud: [support@stjude.cloud](mailto:support@stjude.cloud)

  - Treehouse: [treehousegenomics@ucsc.edu](mailto:treehousegenomics@ucsc.edu)

  - Kids First – CHOP: [nemarichc@chop.edu](mailto:nemarichc@chop.edu)

  - PCDC – UChicago: [lgraglia@bsd.uchicago.edu](mailto:lgraglia@bsd.uchicago.edu)

# How You Can Engage with CCDI

**Learn about CCDI and subscribe to our monthly newsletter:**
cancer.gov/CCDI

**Access CCDI data and resources:**
ccdi.cancer.gov

**Questions? Email us at:**
NCIChildhoodCancerDataInitiative@mail.nih.gov