# DISCLAIIMER

The findings and conclusions in this presentation are those of the author(s) and do not necessarily represent the views of the Centers for Disease Control and Prevention.
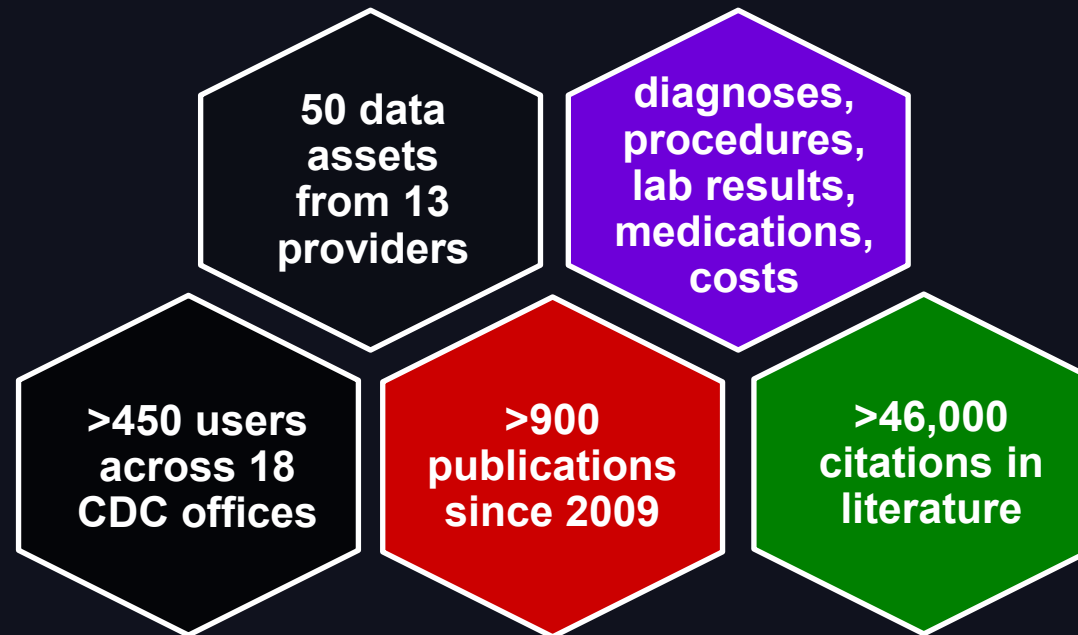
The information shared here is based on implementation experience of a CDC use case and is no way, shape, or form a direct endorsement of the product.

# CDC DATA HUB HEALTHCARE DATA ECOSYSTEM

Program acquires, evaluates, manages, and supports the use of healthcare and related data sources for public health research and action—*A "one stop shop" for the agency.* Leverages Databricks within CDC's cloud environment for data management and intelligence capabilities

- 50 data assets from 13 providers
- diagnoses, procedures, lab results, medications, costs
- >450 users across 18 CDC offices
- >900 publications since 2009
- >46,000 citations in literature

# PRESENTATION OBJECTIVE AND AGENDA

- Explore an accelerator methodology to build faster and more reliable visualization data products

Agenda

1. **Our Journey in Big Data Visualization in Public Health**

   - **Process** Modernization: Key changes and impacts

   - **Technology** Modernization: Advancements and implementations

2. **Example Dashboard:** Respiratory Virus - Performance Insights

3. **Conclusion:** Achievements and future directions

# Hi, I'm John

| | | | |
|---|---|---|---|
| ● | 1988-1992 | 🎓 Graduated UGA | BA, Fortran Punch Cards |
| | 1993-1995 | 🏭 Southern Company | Mainframe, DB2 |
| ● | 1995-1997 | 🌲 Georgia Pacific | Web, SQL |
| | 1997-1998 | 📶 MCI | Web Portals, SQL |
| ● | 1998-2001 | ⚡ Clarus Corporation | SQL, .NET, .COM Startup |
| ● | 2001-2007 | 💊 Pfizer + Pharma (Ctr) | SQL, XML, BizTalk |
| ● | 2007-2018 | 🤝 Microsoft (Ctr) | Data Architect |
| ● | 2018-Present | 🐘 🦠 CDC (Ctr -> FTE) | Data Architect |

🐘 Tuscaloosa, Alabama, USA — If you know anything about the SEC , it's a rare thing for a UGA grad to find happiness in Crimson Tide country!

# CDC DATA HUB

OFFICE OF PUBLIC HEALTH DATA, SURVEILLANCE, AND TECHNOLOGY

# OUR JOURNEY IN BIG DATA VISUALIZATION IN PUBLIC HEALTH

**Core components of our big data visualization accelerator methodology**

## Top 5 Process Improvements

- Common Visualization Gallery
- Standardized Agile User Stories and Recipes
- Machine Readable Requirements
- Standardized Data Product Visualizations
- Standardized Project Management

## Top 5 Technology Advancements

- Common Data Models
- Standardized Data Conversion
- Data Quality Expectations and Synthea
- Data Product Catalog and Workflows
- Data Product Lifecycle

# CDC DATA HUB PROCESS MODERNIZATION

# CDC DATA HUB VISUALIZATION GALLERY

**Process Improvement #1:** Enhancing Discoverability and Engagement

- No Intranet Dashboards
- Unmonitored Usage
- Limited Offline Formats

**Engagement** Challenge

- Dashboard Intranet Listings
- Usage Monitoring & Optimization
- PDF, DOC, XLSX Export

**Engagement** Solution

- Higher User Adoption
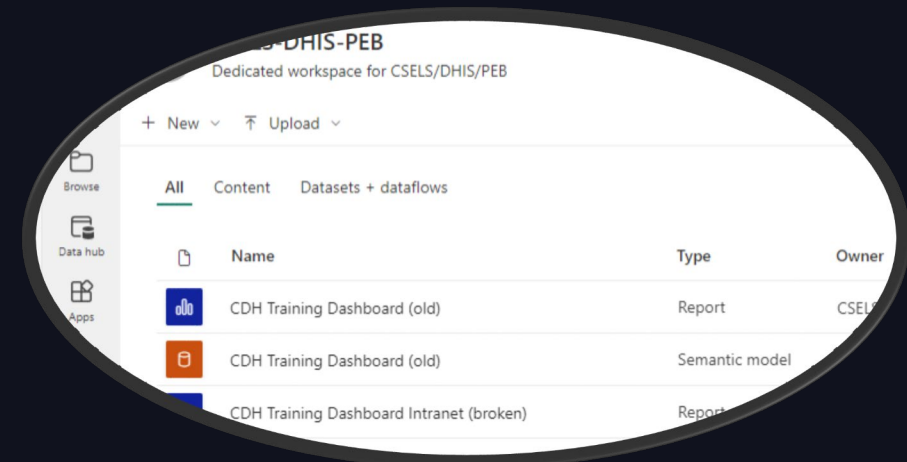- Better Data Accessibility
- Data-Driven Decisions

**User** Outcome

Before

After

Disconnected Power BI Silos

CSELS-DHIS-PEB
Dedicated workspace for CSELS/DHIS/PEB

+ New ∨    ⬆ Upload ∨

All    Content    Datasets + dataflows

| | Name | Type | Owner |
|---|---|---|---|
| | CDH Training Dashboard (old) | Report | CSEL |
| | CDH Training Dashboard (old) | Semantic model | |
| | CDH Training Dashboard Intranet (broken) | Report | |

Browse
Data hub
Apps

Discoverable, cost-effective, integrated Power BI

SAMPLE DATA

# STANDARDIZED AGILE USER STORIES AND RECIPES

## Process Improvement #2: Agile User Story and Recipe Authoring

- Old Method: 500-Page SOP
- Complex Information
- Poor Discoverability
- Expert/Non-expert Gap

**Communication** Challenge

- **Recipe Format:** Bite-sized, Visual
- **Training:** Regular Workshops

**Communication** Solution

- User-Crafted Stories
- Less Tech Dependency
- Better Process Ownership

**User** Impact

Before

After

500+ Page Document



Quick ~15-Minute Guided Recipes



Hide | View

# CROSSING THE CHASM: MACHINE READABLE REQUIREMENTS

**Process Improvement #3: Requirements Authoring Tools**

- **Old Methods:** Sticky Note Documentation
- **Issues:** Inaccessible and hardcoded

**Requirement** Challenge

- Machine-Readable Specs
- Accessibility: Non-technical Reader
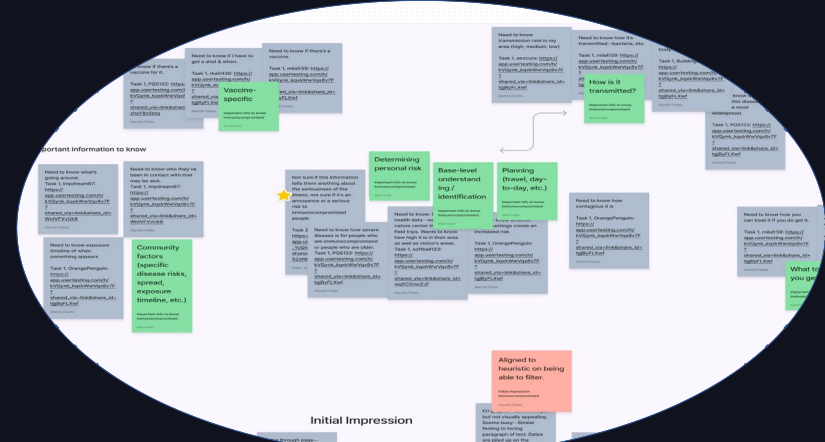- Version Control

**Requirement** Solutions

- **Empowerment:** Autonomy with Easier Tools

**User** Impact

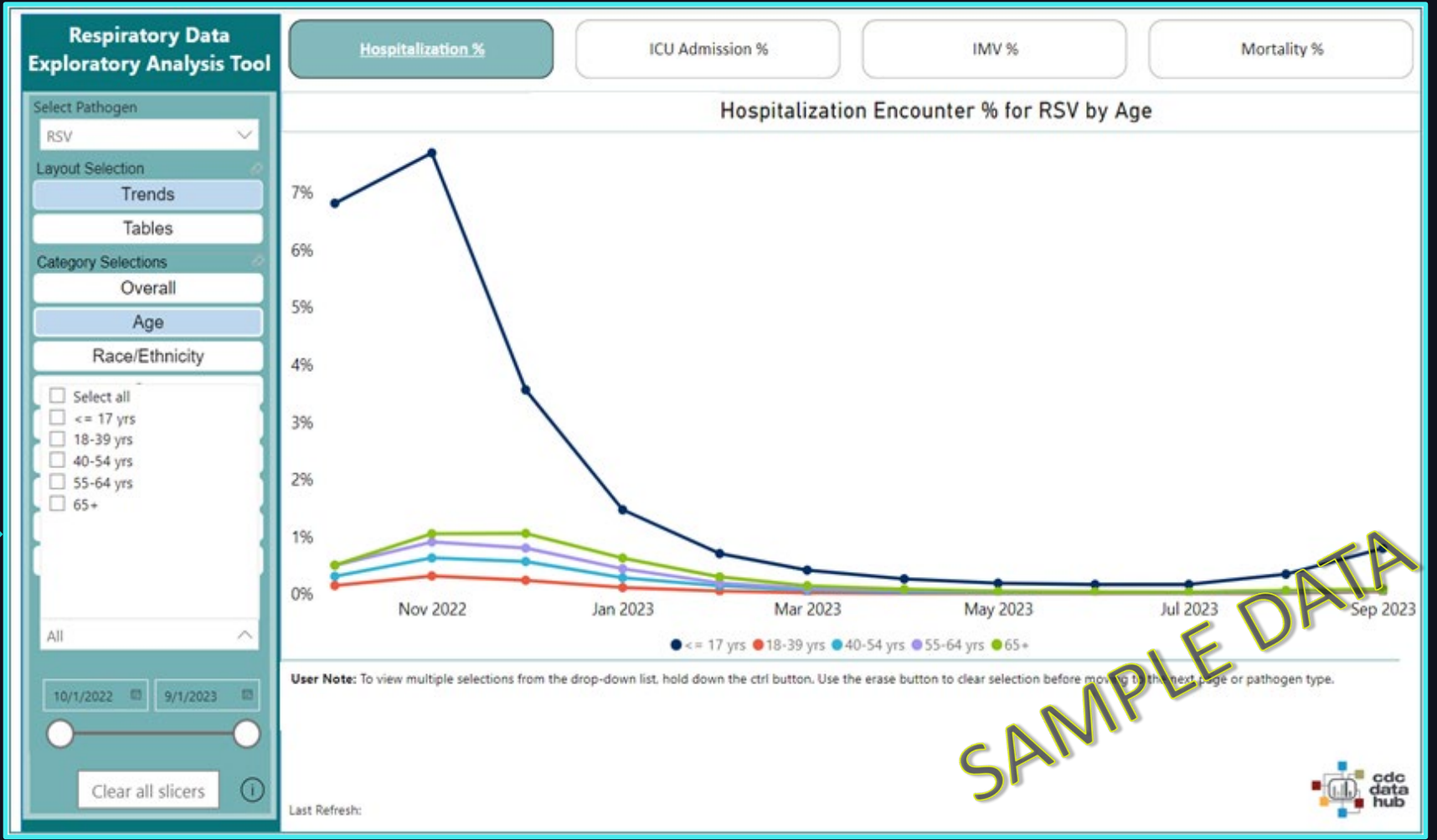Human-Readable Requirements
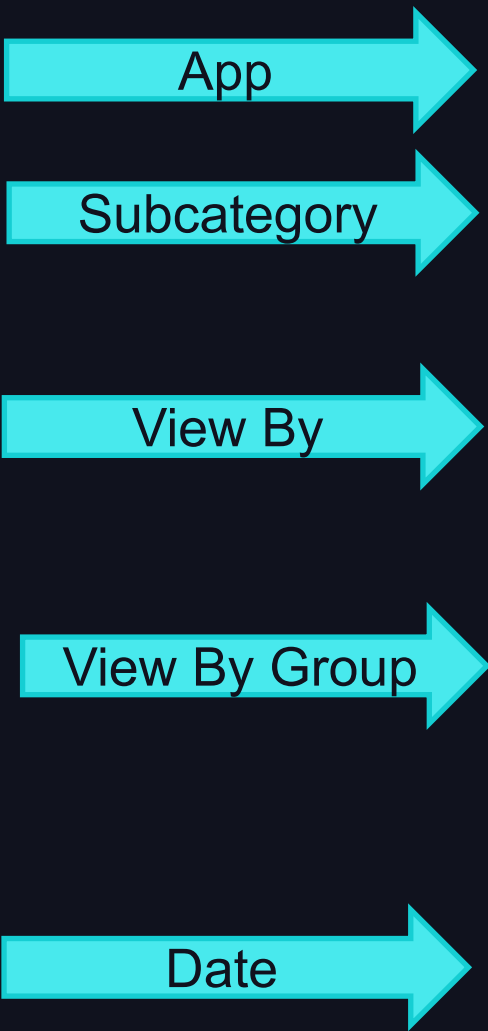
Human- and Machine-Readable Requirements

SAMPLE DATA

# STANDARDIZED DATA PRODUCT VISUALIZATIONS

**Process Improvement #4:** Streamlined and Consistent Dashboard Interface

Metric

App

Subcategory

View By

View By Group

Date



View

Hide

# STANDARDIZED PROJECT MANAGEMENT

● **Process Improvement #5: Introducing Agile Plan Authoring Templates and Tools**

- Details missing
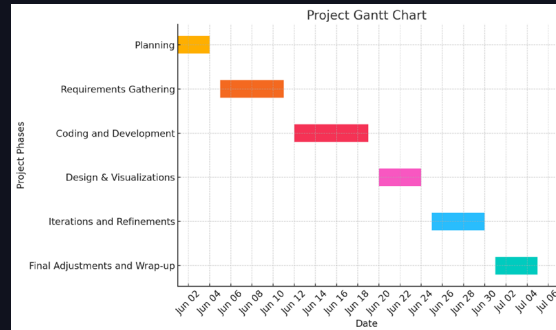- Started each project from

**Planning** Challenge

- Reused templates
- Reused naming
- Reused stories

**Planning** Solution

- Clarity
- Confidence
- Repeatability

**User** Impact

● Before

Inconsistent Project Reporting ≠



One Source of Truth – Office Linked to Jira

● After

| Tasks | Id | Jira Task | Status | 21-Mar |
|---|---|---|---|---|
| Document | LAVA-21 | LAVA-21 | ☑ | 100% |
| How to conduct Stakeholder Identification | LAVA-189 | LAVA-189 | ☑ | 100% |
| How to create Project Personas and Roles | LAVA-190 | LAVA-190 | ☑ | 100% |
| How to log and mitigate risks | LAVA-276 | LAVA-276 | ☑ | 25% |
| How to create a Sprint Plan | LAVA-191 | LAVA-191 | ☑ | 40% |
| How to adjust JIRA for Resource Allocation | LAVA-191 | LAVA-191 | ☑ | 30% |
| Charter Document | LAVA-275 | LAVA-275 | ☑ | 100% |
| How to implement Rules of Behavior | LAVA-271 | LAVA-271 | ☑ | 25% |

DATA+AI SUMMIT

# CDC DATA HUB TECHNOLOGY MODERNIZATION

# ACCELERATING COMMON DATA MODELS

**Technology Advancement #1:** Multiple Subject Areas using Standardized Schemas

- Inconsistent and Wide Data Schemas
- Data Accessibility

**Data Model** Challenges

- Adopting Common Schema
- Harmonizing Multiple Standards

**Data Model** Advancements

- Enhanced Data Usability
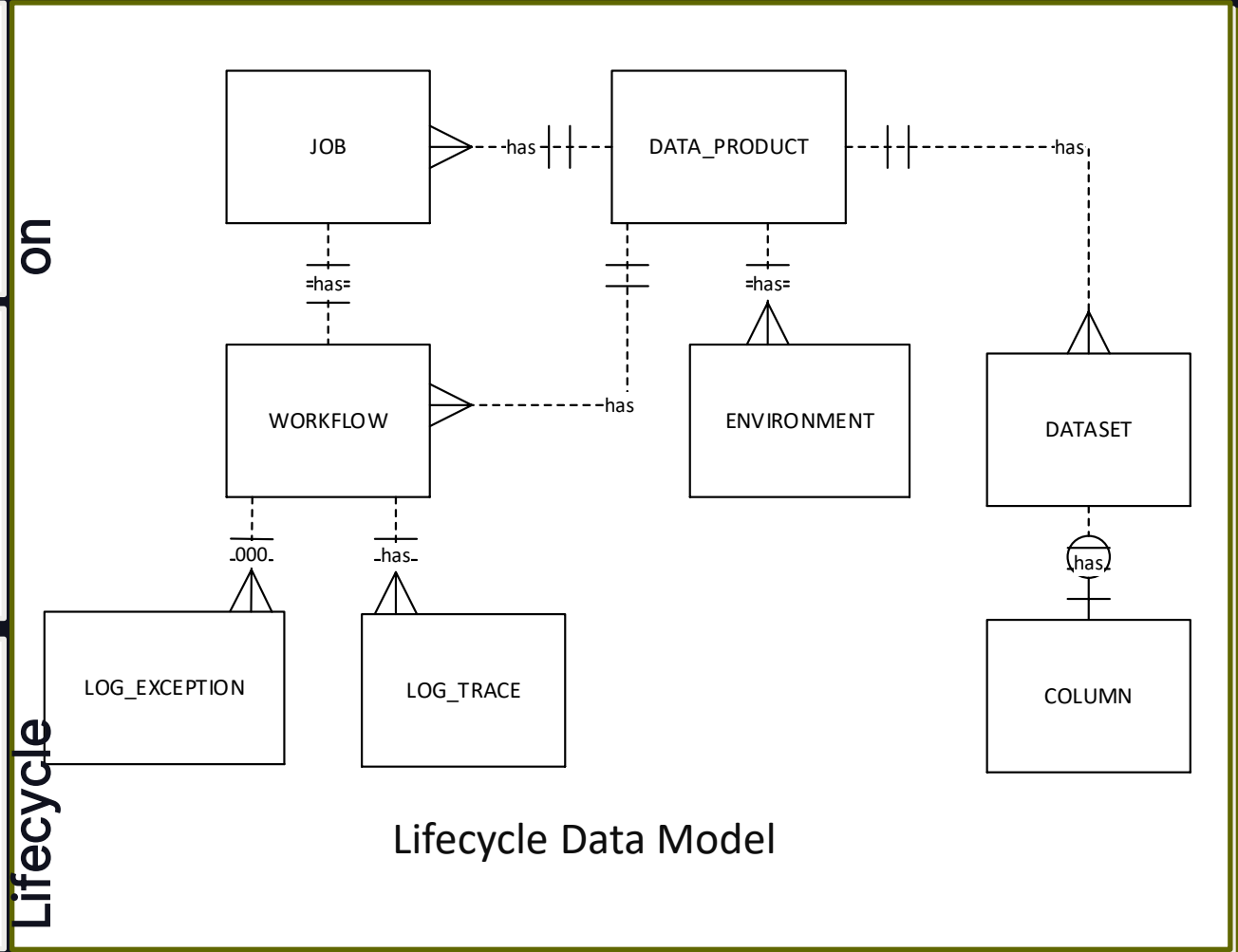- Streamlined Data Analysis
- Governance Compliance

**User** Impact



Lifecycle Data Model

# ACCELERATING STANDARDIZED DATA CONVERSION

Technology Advancement #2: Global Reference Data and Mapping Tools

- Custom Data Values
- Inconsistent and Custom Logic

**Conversion**
Challenges

- Universal Conversion Adapters
- Third Party Data Mapping Providers
- Parameterized SQL

**Conversion**
Advancements

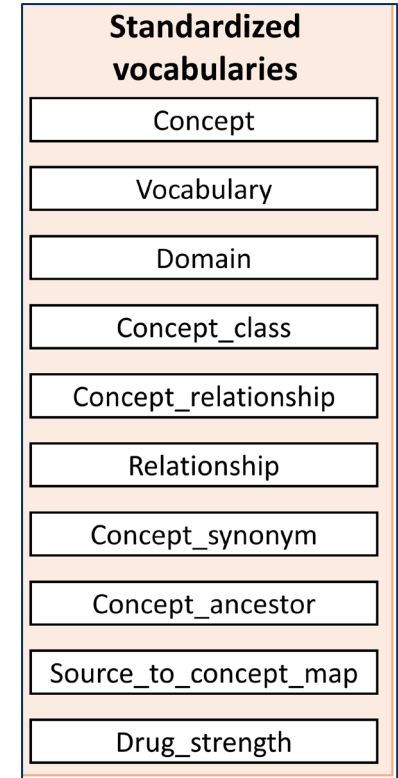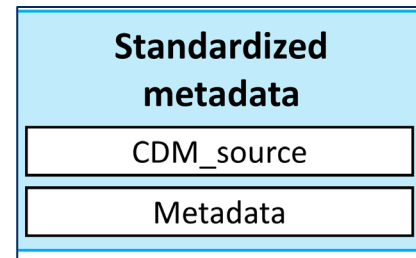- Increased Research Collaboration
- Enhanced Data Reliability

**User**
Impact

Concept

Parametrized

Model

**Standardized metadata**

CDM_source

Metadata

**Standardized vocabularies**

Concept

Vocabulary

Domain

Concept_class

Concept_relationship

Relationship

Concept_synonym

Concept_ancestor

Source_to_concept_map

Drug_strength

Global Reference Data Model

# DATA QUALITY EXPECTATION TRACKING WITH TEST DATA

**Technology Advancement #3:** Data Governance Expectation Tracking Tools

- Restrictive Data Access
- Large-scale Dataset Complexity
- Lack of visibility into data quality

**Quality** Challenge

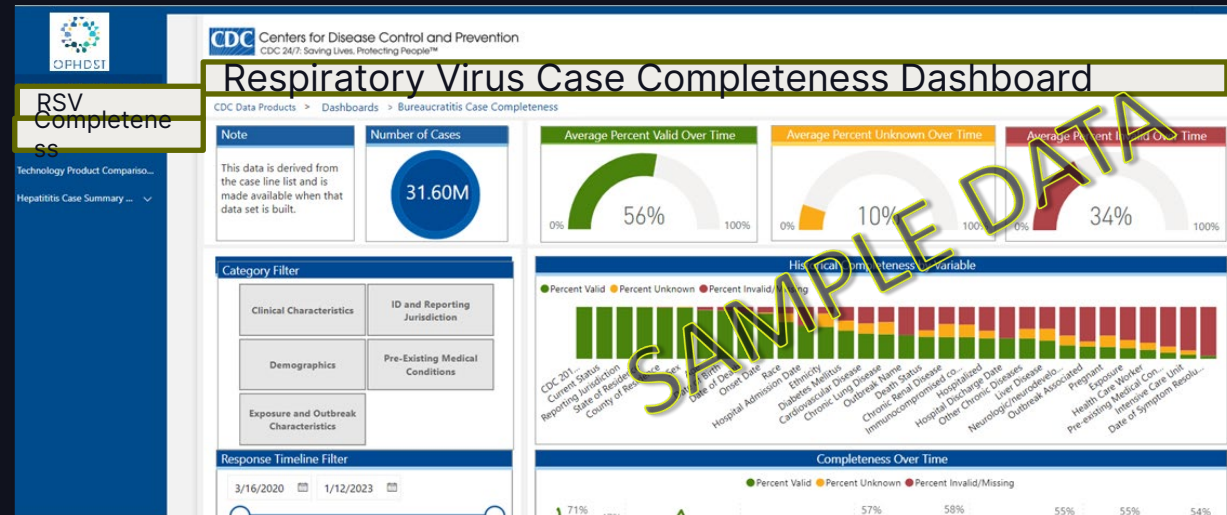- Synthetic & Real-World Data Testing
- Schema Monitoring with Unity Catalog

**Quality** Advancements

- Free, Immediate Data Access
- Enhanced Trust in Data
- Transparent Quality Insights

**User** Impact

Synthea Generated Test Data

Real-World Actual Test Data

Expectation Checks



CDC Centers for Disease Control and Prevention
CDC 24/7: Saving Lives, Protecting People™

Respiratory Virus Case Completeness Dashboard

RSV Completeness

CDC Data Products > Dashboards > Bureaucratitis Case Completeness

| Note | Number of Cases | Average Percent Valid Over Time | Average Percent Unknown Over Time | Average Percent Invalid Over Time |
|------|-----------------|----------------------------------|------------------------------------|------------------------------------|
| This data is derived from the case line list and is made available when that data set is built. | 31.60M | 56% | 10% | 34% |

SAMPLE DATA

Category Filter

Clinical Characteristics | ID and Reporting Jurisdiction
Demographics | Pre-Existing Medical Conditions
Exposure and Outbreak Characteristics

Response Timeline Filter
3/16/2020 — 1/12/2023

Quality Insights

# DATA PRODUCT CATALOG AND WORKFLOWS

## Technology Advancement #4: Prescriptive Guidance for Defining Data Products

- Lack of Guidance on data products
- No Analytical Cataloging
- Limited separation of billing and logging

**Data Product**
Challenges

- Established Governance
- API Calls: Pass data product ID
- Partitioned Logic: Shared and

**Data Product**
Advancements

- Findable & Accessible
- Interoperable & Reusable (FAIR)
- Multi-tenant Billing & Logging

**User**
Impact

**Data Product**

Contain

**Data Packages**

Have

**Package Datasets**

- Comprehensive set of data & metadata
- Includes multiple Data Packages
- Often maps 1-1 to a database

- Collection of datasets/tables
- Metadata + Actual data
- Often maps 1-1 to a schema

- Metadata in dictionary format
- Often maps 1-1 to a table or file

# STANDARDIZED DATA PRODUCT LIFECYCLE (SQL + PYTHON)

**Technology Advancement #5:** Data Lifecyle Management and Tracking

- Distinct processes per project for collecting data, grouping and basic analysis, and advanced analysis or ML techniques.

**Data Management**
**Challenges**

- Shared data ingestion framework
- Standardized 5 step IDEAS process

**Data Management**
**Advancements**

- Timely data for critical decision-making
- Reduced complexity in data

**User**
**Impact**

process_jobs_rsv_readmissi

run_job_name

parameters

config

job_name

parameters

get_config

get_job_name

get parameters

dbutils_widgets

process_env_metadata

process_job_metadata

process_standard_parameters

run_job_name

process_ingress

process_data_load

Bronze

process_enrichment

Silver

process_analytics

Gold

process_storytelling

Workflows and Orchestration

Medallion Architecture

DATA+AI SUMMIT

# CDC DATA HUB CONCLUSION AND FUTURE

# DATA PROCESS

## Data Hub Process Improvements

CDC Data Hub Results
- Increased data quality
- Easier user experience

Based on Simplifications
- Data Structure Reuse
- Code Reuse
- External Data Reuse

### Can Input Be Simplified?
- Reduce Hard Coding
- Use Global Reference Data

### Evaluate Visualization Tools
- Data-Driven Design
- Move >1 Buttons to 1 Filter

### Evaluate for Code Reuse
- Reduce Copy-Paste Actions
- Parameterize SQL

### Can Output Be Simplified?
- Make Output Pivot-Friendly

# DATA EFFICIENCY

## Data Hub Efficiency Gains

As you can see, the new dashboard performances captures the same data (and then some) but with less processing power, allowing the dashboard to run faster and smoother than before.

Over **10x** performance increase per metric as demonstrated in the demo

Despite expanding from 4 to over 20 metrics and expanding from 4 to over 10 dimensions, we achieved

### SMALLER: Up to 5x File Size Reduction



File Size (Mb): New Dashboard 1, Initial Dashboard 5

### SIMPLER: Up to 7x Less Components



Tab Count: New Dashboard 2, Initial Dashboard 5
Table Count: New Dashboard 8, Initial Dashboard 20
Bookmark Count: New Dashboard 6, Initial Dashboard 44

### FASTER: Up to 3x Time to Load Reduction



Slowest to Load: New Dashboard 8, Initial Dashboard 25
Seconds to Load (Avg): New Dashboard 10, Initial Dashboard 6

# CONCLUSION: KEY TAKEAWAYS

**Demonstrated an accelerator methodology to build faster and more reliable visualization data products**

**Our Current Status**

**Process** Modernization
- Highlighted modernization solutions in CDC Data Hub visualization processes enhancing data handling efficiency.

**Technology** Modernization
- Discussed advancements in CDC Data Hub's technology stack improving overall performance.

**Respiratory Virus** Example Dashboard
- Reviewed findings showing up to a 10x improvement in data visualization performance.

DATA+AI SUMMIT

# FUTURE IDEAS FOR BIG DATA VISUALIZATION IN PUBLIC HEALTH

**Core components of our big data visualization accelerator methodology**

Future Potential

## Future Opportunities

Prescriptive Analytics

Put Conclusion in Headlines

Animate It (Charts and Maps)

Model Explanations / Driver Analysis

Make It About Them (Personalization)

AI/ML and LLM with Forecasting

Data Sharing

DATA+AI SUMMIT

# DEMONSTRATION

OFFICE OF PUBLIC HEALTH DATA, SURVEILLANCE, AND
TECHNOLOGY

# LAVA Portal

Welcome – Introduction to What is LAVA - https://bit.ly/cdhlava

# LAVA Portal

## Portal Sidebar

# LAVA Portal

Data sources

# LAVA Portal

Upload

# LAVA Portal

**Run Jobs**
**Online Option**

# LAVA Portal

**Run Jobs
Notebook Option**

# LAVA Portal

# LAVA Portal

**Log Search
Application Insights**

# LAVA Flask

**Open API Specification**

## CDC Data Hub LAVA Flask API [1.0]

[ Base URL: /content/f14ebb0a-f6f7-4eda-be53-b9a117d28a6b ]

/content/f14ebb0a-f6f7-4eda-be53-b9a117d28a6b/swagger.json

## API Documentation

CDC Data Hub LAVA (CDH) provides shared resources, practices and guardrails for analysts to discover, access, link, and use agency data in a consistent way. CDH impr___ reduce the effort required to find, access, and trust data.

Back to Home

Config Upload Page

Config Download Page

EDC Upload Page

EDC Download Page

For detailed logs, please visit the Log File Page.

---

**welcome**    Welcome to the CDC Data Hub LAVA API

---

**cdc_security**    The security service manages security of the data products and associated services. The package contains datasets that provide critical informa___ availability of the data products and associated services.

---

**cdh_lava**    The CDC Data Hub Lifecycle, Analysis and Visualization Accelerator (CDC Data Hub LAVA)makes building and deploying data products faster and mor___ data processes and technology.

---

**cdc_admin**    The admin service manages and monitors data products and associated logs. This package contains datasets that provide critical information for e___ of the data products and related services.
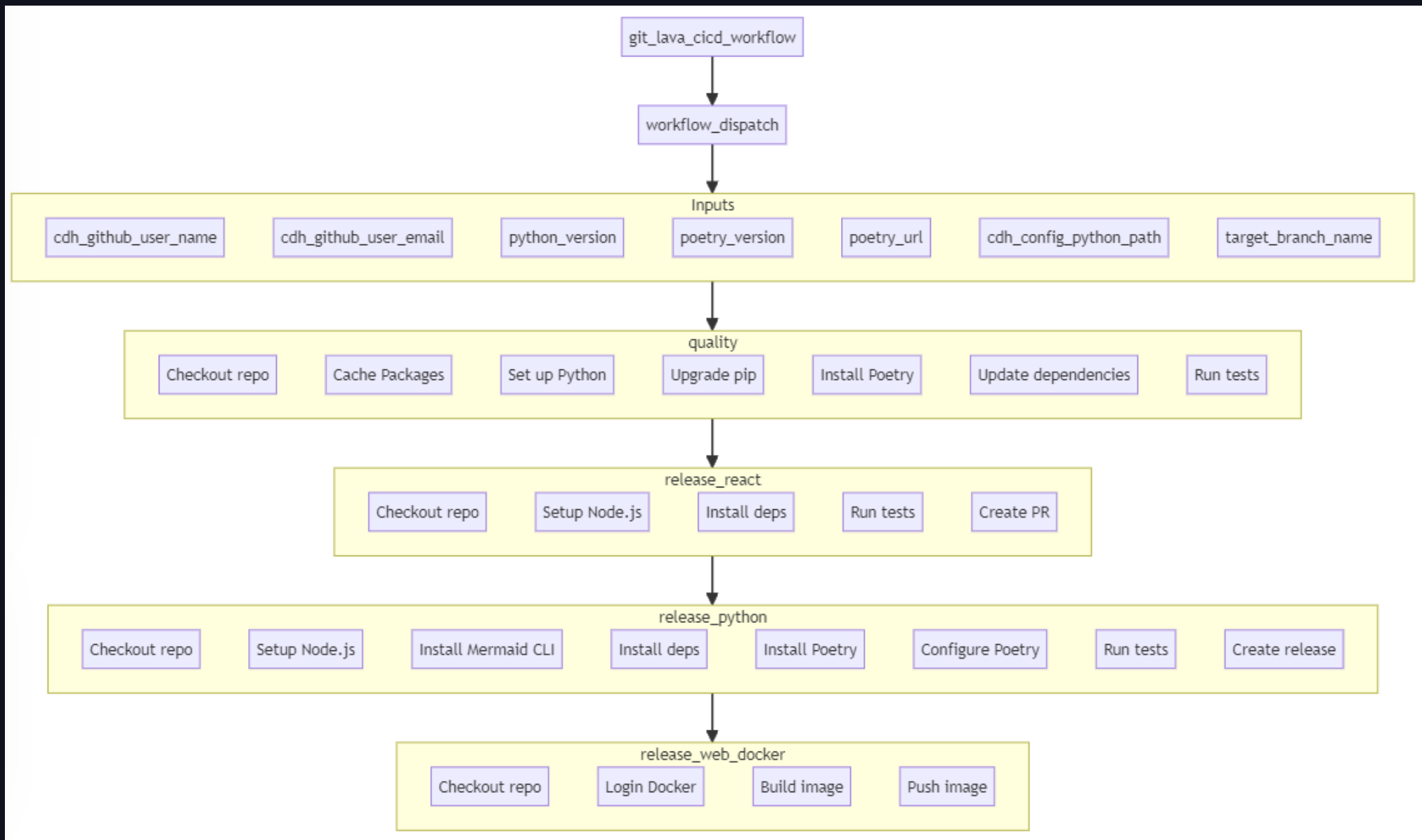
---

**GET**    `/cdc_admin/metadata_excel_file_download_codes/{schema_id}`    Retrieves the Excel metadata file from Alation based on the schema_id

# LAVA DEVSECOPS

# WHAT PROBLEMS ARE WE SOLVING

**Top 10 Pitfalls that challenge data product efficiency and quality**

### Dashboard Drifters
- Can't find dashboards.

### Death By Docs
- Documentation is too big and unwieldy

### Vague Vortex
- Requirements are not translatable to machines.

### Molasses Matrix
- Report navigation is slow and difficult

### Phantom Facts
- Planning summaries are not tied to details.

### The Pancake Stack
- Inconsistent and wide data schemas

### Conversion Chaos
- Inconsistent conversion logic.

### Opaque Oracle
- Undefined quality expectations

### Metadata Mystery
- No analytical cataloging.

### Reinvention Rocket
- Little to no ETL reuse.

# OUR JOURNEY IN BIG DATA VISUALIZATION IN PUBLIC HEALTH

**Core components of our big data visualization accelerator methodology**

## Top 5 Process Improvements

- Common Visualization Gallery
- Standardized Agile User Stories and Recipes
- Machine Readable Requirements
- Standardized Data Product Visualizations
- Standardized Project Management

## Top 5 Technology Advancements

- Common Data Models
- Standardized Data Conversion
- Data Quality Expectations and Synthea
- Data Product Catalog and Workflows
- Data Product Lifecycle

" We can complete our 20+ page Feedback Packet in 10 minutes....
"

-Matt Cole
Epidemiologist and Data Scientist
First Project Completed in 2021

" **The new process is 1000x faster for creating standard pivots.** "

-Stacey Adjei
Epidemiologist and Data Scientist
Latest Project Completed in 2024