*Welcome to the...*

# NIH Cloud Platforms Interoperability Fall 2020 Workshop

*We'll be starting shortly!*

# Welcome & Introduction to Day 2

**Adam Resnick**
*Children's Hospital of Philadelphia*
**Valerie Cotton**
*Eunice Kennedy Shriver National Institute of Child Health and Human Development (NICHD), NIH*

# Draft Roadmap in FunRetro (from Day 1)

# CRDC Cloud Costs:
# Current Practices

**Tanja Davidsen**
NCI Center for Biomedical Informatics and
Information Technology (CBIIT)

Connecting NCI data and compute in the cloud

- Access to large cancer data sets without need to download
- Access to workspaces, analysis tools, and pipelines
- Ability for researchers to bring their own data/tools



**Data**



**Compute**



**Security**



Institute for
Systems Biology
isb-cgc.org

FireCloud
POWERED BY Terra
firecloud.terra.bio

CANCER GENOMICS CLOUD
SEVEN BRIDGES
cancergenomicscloud.org
aws

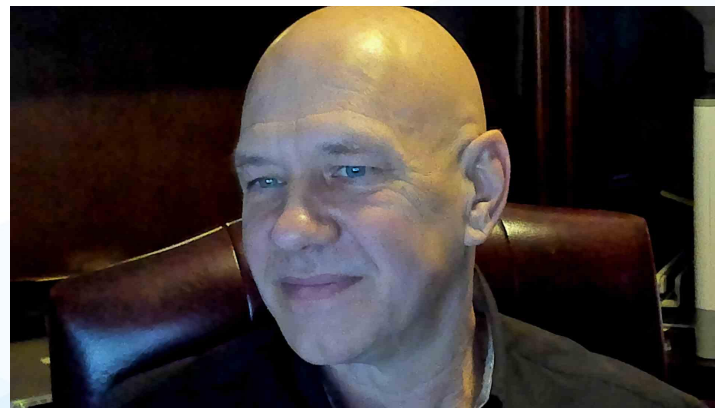**NCI Cloud Resources**

## CRDC Cloud Resources Compute

- Each Cloud Resource provides $300 Credit to every new user
  - Some additional free credits beyond $300 can be requested via an application
  - Beyond the free credits users can use a credit card or a billing account
- Links to free credit information
  - Seven Bridges CGC
  - Broad FireCloud
  - Institute for Systems Biology CGC

## CRDC Cloud Storage

- Copied on both Amazon and Google
- Mostly genomic, additional WGS, imaging, and proteomic data coming soon
- Size:
  - 3.5PB total and growing
  - 2-4PB additional data expected in FY21
- Costs:
  - 2PB free from STRIDES each: Amazon and Google
  - Additional FY21 budget for storage $1.5 million using STRIDES discount
- The future concerns:
  - Looking for funding opportunities as FY22 budgets will not cover storage needed
  - Looking for cost savings/sustainability: lower cost storage, compression, etc.

# BioData Catalyst Cloud Credit Overview

- Pilot Funding: The NHLBI currently provides **$500 in cloud credits** to new users of NHLBI BioData Catalyst via a billing group on either *BioData Catalyst Powered by Seven Bridges* or *BioData Catalyst Powered by Terra*.
- If the anticipated costs are in excess of $500:
  - Users can cover those costs using their own AWS and/or Google accounts which can be brought to BioData Catalyst
  - Users can apply for additional credits via the **NHLBI BioData Catalyst Cloud Credit Program**

More information at https://biodatacatalyst.nhlbi.nih.gov/resources/cloud-credits

# BioData Catalyst Cloud Credits

- BioData Catalyst has a lightweight form for requesting additional credits
- Requests are reviewed and approved based on projected costs and scientific merit

More information at
https://biodatacatalyst.nhlbi.nih.gov/resources/cloud-credits

# BioData Catalyst Cloud Credit Resources

- Currently, we offer these resources for **understanding cloud costs**
  - Controlling your cloud costs (*BioData Catalyst Powered by Terra*)
  - Cloud infrastructure pricing (*BioData Catalyst Powered by Seven Bridges*)
  - Comprehensive tips for reliable and efficient analysis set-up (*BioData Catalyst Powered by Seven Bridges*)
- Next steps:
  - White paper on estimating cloud cost is in the works

# BioData Catalyst Cloud Credit Lessons Learned

- Cloud costs needs to be more accurately tracked and managed
  - We need to develop **better reporting** on cloud cost disbursement and usage. STRIDES dashboard should help resolve many reporting issues (beta coming soon)
  - We need to **QA/QC user pipelines/workflows** to improve performance and ensure maximum cost efficiency
- Cloud costs create a lot of **anxiety** for users
  - They can present a significant perceived barrier to entry
  - Evaluating costs takes time and can return varied results
- **Shared data storage** facilitates research collaborations and ultimately reduces costs

# Kids First DRC Cloud Costs: Current Practices

**David Higgins, PhD**
Kids First Data Resource Center
Children's Hospital of Philadelphia

# Kids First DRC Cloud Credit Overview

- Pilot Funding: All users receive $100 in pilot funds upon making an account in CAVATICA (Seven Bridges).

- Cloud Credits: Users can receive for more funds through NIH Common Fund.
  - Phase 1 (FY2020): Up to $5,000 upon approval of an application.
  - Phase 2 (FY2021):
    - All Kids First X01 PIs: receive $1,000 when they receive their data
    - Can then apply for additional funds to complete their analysis
    - KFDRC can track disbursement and usage through an admin account

# Kids First DRC Cloud Credit Resources

[What costs are there for using Kids First and Cavatica?](#)

(FAQ on Kids First DRC Support Pages)

💵💲

## What costs are there for using Kids First and Cavatica?

### Are there fees for using the Kids First Data Resource Portal?

No. The portal itself is free to use. Anyone can make an account and browse available datasets using the Explore Data and File Repository tools. If you have been approved for access, you will have the ability to Send files to Cavatica for download and analysis.

### Are there fees for downloading Kids First files from Cavatica

Benchmarking Statistics for Kids First Workflows on CAVATICA

## Kids First DRC Joint Genotyping Workflow

Kids First Data Resource Center Joint Genotyping Workflow (cram-to-deNovoGVCF). Cohort sample variant calling and genotype refinement.

Using existing gVCFs, likely from GATK Haplotype Caller, we follow this workflow: Germline short variant discovery (SNPs + Indels), to create family joint calling and joint trios (typically mother-father-child) variant calls. Peddy is run to raise any potential issues in family relation definitions and sex assignment.

If you would like to run this workflow using the cavatica public app, a basic primer on running public apps can be found here. Alternatively, if you'd like to run it locally using cwltool, a basic primer on that can be found here and combined with app-specific info from the readme below. This workflow is the current production workflow, equivalent to this Cavatica public app.



### Runtime Estimates

- Single 5 GB gVCF Input: 90 Minutes & $2.25
- Trio of 6 GB gVCFs Input: 240 Minutes & $3.25

# Kids First DRC Cloud Credit Lessons Learned

- In brief: users don't have a strong sense of costs.

  - These might not be an issue on local clusters at their institutions.

  - Have access to spending, but no regular updates or alerts on balances.

- Adjustment from Pilot 1 to Pilot 2: give users more autonomy

  - Pilot 1: KFDRC manages financials as owners of the billing group

  - Pilot 2: User groups manage financials as owners of the billing group

# Kids First DRC Cloud Credit Lessons Learned



Individual tasks run in CAVATICA ➡

*"How much do these cost?"*
*"How much are we spending?"*

**Need two messages for two groups of people.**

PI

User

*"I'm writing a grant and want to include funds for cloud analysis, how much should I budget?"*

*"How do I track the usage of funds in the cloud?"*

*"How do I avoid recklessly spending down our grant before completing our analysis?"*

# AnVIL
# Cloud Costs:
# Current Practices

**Frederick Tan**
Carnegie Institution / JHU

# Barriers to Entry

Spending First Dollar

Estimating Costs

Reporting Charges

Writing Into Grants



MaGIC Jamboree
~100

BCC2020
~70

BiocCC2020
~40

Cost for NIH STRIDES to fund workshops

More information at datascience.nih.gov/strides

Default budget alert behavior

Budget alert → Billing admin

Automate Cost Control Responses
Example of reference architecture

Cloud Pub/Sub → Cloud Functions → Billing API → Cap spending

```python
# Disable billing
if cost_amount > budget_amount:
    body = {'billingAccountName': ''}
    projects.updateBillingInfo(
      name=project_name,
      body=body).execute()
```

★ **Note:** There is a delay of up to a few days between incurring costs and receiving budget notifications. Due to usage latency from the time that a resource is used to the

# AnVIL

More information at cloud.google.com/billing/docs/how-to/notify

# Reporting Charges - AnVILBilling



More information at bioconductor.org/packages/AnVILBilling

# Cloud Costs: Near Term Improvements

**Open Discussion**
*30 Minutes*

# Lunch Break

We will resume at 12:30 pm ET.

# Overview

- 10 min - Jiaqi and Brian will give an overview of RAS and Data Sharing work

- 5 min - Group will brainstorm on areas of RAS & Data Sharing to focus on in 2021

- 25 min - Group discussion, deep dive on ~3 of the most popular topics

# RAS Integration

**Jiaqi Liu** U. Chicago

# Commonality of Framework Services

*The commonality of AnVIL, BD Catalyst, CRDC, and Kids First framework services facilitates RAS adoption*

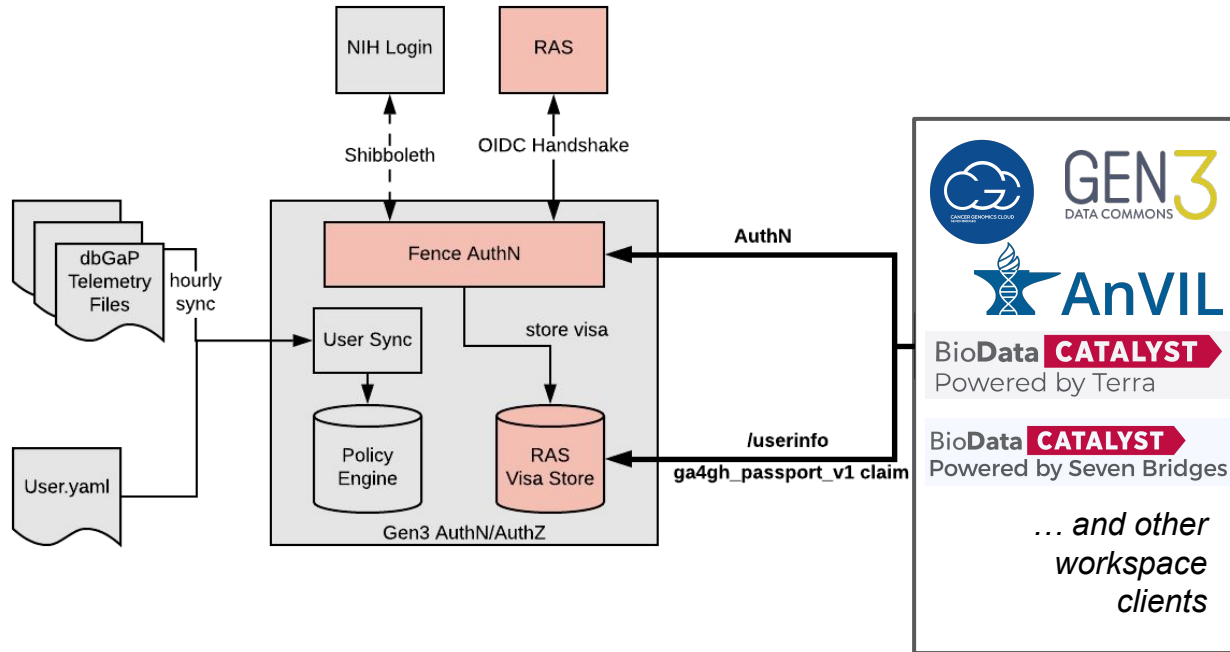# Shared Technical RAS Plan - Work Divided into 3 Milestones

The **U. Chicago**, **RAS**, **AnVIL, BDCat, CRDC, and Kids First**, and **GA4GH** teams have collaborated to bring RAS to our systems.

- We worked on a technical document to break down the work: [CRDC RAS Integration Proposal v1.4](#)
- CRDC-specific but highly applicable to **AnVIL**, **BD Catalyst**, **Kids First DRC**
- We continue to [coordinate](#) in ad hoc meetings, various GA4GH working groups, and the NCPI Systems Interoperation Working Group



Work is Broken into Milestones:
M1, M2, and M3

# RAS Integration - Milestone 1 (Now)



- Use RAS for user login (AuthN)

- Use Cases Enabled
  - Systems use Fence
  - SSO-like experience
  - Passport visas explored for securing derived results and data access UX

# RAS Integration - Milestone 2



- Use RAS for Login and Authorization
  - Replace dbGaP telemetry files
- Use Cases Enabled
  - *Previous Milestone 1 use cases*
  - Researcher benefit from *realtime* AuthZ instead of dbGaP whitelists which can lag

# RAS Integration - Milestone 3



**Workspaces:** Terra, SBG, Cavatica (SBG), ISB-CGC

Centralized Fence with OIDC connection to RAS

IndexD + DRS

IndexD + DRS

IndexD + DRS

IndexD + DRS

Google

Amazon & Google

Amazon & Google

Amazon & Google

- Fully distributed Authorization
  - A"central" Fence, users can access any dataset via single linking event
- Use Cases Enabled
  - *Previous Milestones*
  - SSO with single consent
  - Data across stacks in single linking
  - Provide non-dbGaP/other project access lists via single Fence

# RAS - Current Accomplishments

**Completed Work**:

- **Implemented RAS AuthN Support:** As part of RAS Phase 1, U. Chicago has completed the code changes to Fence to use RAS for user login
- **Coordinated Deployment:** U. Chicago and partner stacks have adapted their systems to use the RAS login via Fence, production deployment is happening ***now***
- **Future Work Design:** U. Chicago, Broad, SBG, RAS, and the GA4GH continue to work together on the future milestones for RAS and GA4GH Passports

**Continuing Work**:

- **Authorization:** U. Chicago will switch to RAS passport visas for Authorization info
- **Continued Technical Milestones:** Milestone 2 and 3, simplifying the user experience
- **GA4GH & Standards:** Laying foundation for future GA4GH Passports, brokers beyond RAS

# RAS Next Steps/Timeline

*Goal is to have all cloud stacks coordinate and use RAS identically across systems*

Time

Milestone 1 - Oct production

Milestone 2 - Q1 2021

Milestone 3 - Q2 2021

Users can **log in with RAS once**, multiple consents per system (e.g. AnVIL etc)

**Design for Milestone 2**

Users can log in with RAS, underlying workspace systems switch to GA4GH passports for authorization

Design for Milestone 3

Users can log in with RAS once, with a single consent, and workspace systems can use RAS or other systems' GA4GH passports to authorize (depending on dataset)

# Data Sharing

**Brian O'Connor** Broad

# What is Data Sharing?

There are many different facets of data sharing... many more than we could fit in a 40 minute conversation.

Some key areas of Data Sharing include:

- **Sharing data files between systems**
- **Exchanging data models between systems**
- Common data models across systems
- Search of data and data models
  - Specialty access/search for specific data types (think Beacon API, htsget)
- AuthN/AuthZ, SSO
- BYOD and sharing derived results
- And many more...

I'm going to dive into the first two aspects and report back progress from NCPI Systems Interoperation group over the last year.

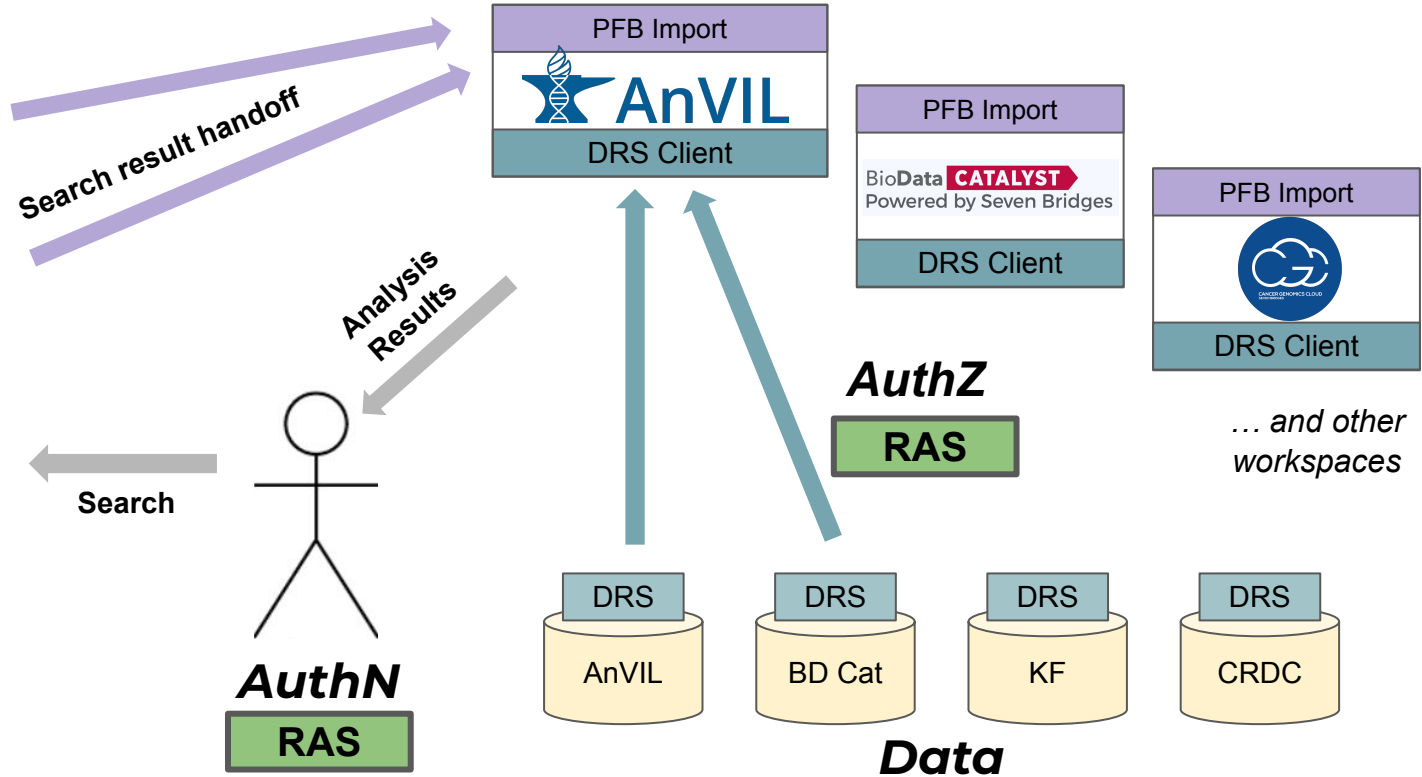*Open discussion at the end will help us frame what Data Sharing we should focus on in 2021*

# Data Sharing

*DRS* facilites
**data file sharing**,
*PFB* facilitates
**sharing data
model + DRS
URIs**,
*RAS* gives us a
common Auth
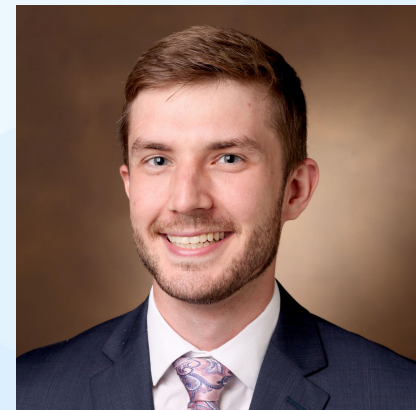system for **SSO**
and **data access
across systems**



**PFB Handoff**

Minimal set of data includes DRS IDs and can include clinical data

PFB File (*.avro)

AVRO™

**Workspaces**

Gen3 Workspaces  jupyter  R

External Workspace: Terra

External Workspace: Seven Bridges

External Workspace: ISB

Query for Cohorts

Access Cohort in Workspace

Data Portals

Phenotype Database

Federated Access

Access file objects in Cloud using DRS ID

**GA4GH Passports/Visas Ecosystem**

NIH RAS

Additional Passport Brokers

Authentication / Authorization

VISA

Centralized Fence

**GA4GH Data Respository Service (DRS) Servers**

KF DRS    CRDC DRS

BDC DRS    AnVIL DRS

From GA4GH Poster Session

# GA4GH DRS 1.1

- The **Data Repository Service (DRS)** API provides a generic interface to data repositories so data consumers, including workflow systems, can access data objects in a single, standard way regardless of where they are stored and how they are managed.
- **DRS 1.1** was released in 2020 and it added support for **compact identifiers** which was key to the DRS servers supporting AnVIL, BD Catalyst, CRDC, and Kids First
- See DRS 1.1 Transition within NCPI

**drs://dg.4DFC:0027045b-9ed6-45af-a68e-f55037b5184c**

```
{
  - access_methods: [
    - {
        access_id: "gs",
      - access_url: {
          url: "gs://gdc-tcga-phs000178-controlled/BRCA/RNA/RNA-Seq/UNC-LCCC/ILLU
        },
        region: "",
        type: "gs"
    },
    - {
        access_id: "s3",
      - access_url: {
          url: "s3://tcga-2-controlled/0027045b-9ed6-45af-a68e-f55037b5184c/UNCID
        },
        region: "",
        type: "s3"
    }
  ],
  aliases: [ ],
  - checksums: [
    - {
        checksum: "2edd5fdb4f1deac4ef2bdf969de9f8ad",
        type: "md5"
    }
  ],
  contents: [ ],
  created_time: "2018-06-27T10:28:06.398871",
  description: null,
  form: "object",
  id: "0027045b-9ed6-45af-a68e-f55037b5184c",
  mime_type: "application/json",
  name: null,
  self_uri: "drs://nci-crdc.datacommons.io/0027045b-9ed6-45af-a68e-f55037b5184c",
  size: 6703858793,
  updated_time: "2018-06-27T10:28:06.398882",
  version: "5eb15d8b"
}
```

# PFB

- Many portals, a smaller number of workspace environments
- DRS 1.1 useful for workspace environments to access data for compute
- But how do you find data in the multitude of portals and bring those results to a workspace environment?
- **Before PFB...**

# PFB

- Many portals, a smaller number of workspace environments
- DRS 1.1 useful for workspace environments to access data for compute
- But how do you find data in the multitude of portals and bring those results to a workspace environment?
- **After PFB...**
  - PFB handoff description
  - PFB "light"

# FunRetro for Discussion

- 5 min - Group will brainstorm on areas of RAS & Data Sharing to focus on in 2021
  - *Let's add discussion ideas/vote in this **FunRetro Board***

- ***Top Topics:***
  - ***Search for data across stacks***
  - ***Policies needed for sharing data between IC stacks***
  - ***Handoff of search results, PFB, FHIR, or both? Something else?***
  - ***BYOD: how does "bring your own data" work across platforms? Is there any difference between sharing data between the IC stacks if its from a canonical project or if it's provided by a researcher upload?***

- 25 min - Group discussion, deep dive on ~3 of the most popular topics
  - What is the challenge?
  - What solutions should be examined?
  - Who needs to work on this topic?
  - When does the topic need to be solved by?
  - ...

# Data Sharing

*DRS facilities **data file sharing**, PFB facilitates **sharing data model + DRS URIs**, RAS gives us a common Auth system for **SSO** and **data access across systems***

# Clinical Data Flow Landscape



Data submitters

Collect → Transform

Submit → Review → Ingest

Platform Staff

Release → Search → Export

Researchers

Transform → Analyze

Informed feedback loop for system change/growth

# Representation Experience with FHIR WG



NCPI Disease

ResearchStudy    Condition    Observation

ResearchSubject    Patient    Specimen    Task    DocumentReference

Observation

NCPI Family Relationship    NCPI Phenotype

NCPI DocumentReference (DRS)

Grey: Base FHIR Resources
Blue: Profiles on base FHIR resources

# MIS-C as Interoperability Use Case

# Ongoing Longitudinal data

| NIAID/PRISM | NICHD/POPS | NHLBI/MUSIC |
|---|---|---|
| • Determine the proportion with SARS-CoV-2 related death, rehospitalization or ongoing major complications at 12 months after presentation<br>• Determine immunologic mechanisms, immune signatures and predictive biomarkers associated with disease phenotypes | • Study the influence of genetic factors, metabolic, protein profiles on therapeutic exposure and response<br>• Evaluate PK/PD of understudied drugs in hospitalized children with SARS-CoV-2 related illness<br>• Establish drug safety profile and adverse events with specific cardiac or neurologic impact | • Characterize the occurrence and time course of coronary artery involvement and ventricular dysfunction<br>• Characterize the occurrence and time course of non-cardiac organ dysfunction, inflammation and major medical events |

**Longitudinal**

# MIS-C Interoperability Across the Landscape

- Collection
  - CRFs and CDEs
  - EHR
- "Platform staff" is across platforms
  - How to view incoming/provisional data?
- Feeding multiple downstream platforms/tools
  - Initial handoffs?
  - Awareness of new data availability?
- Intersecting with other existing datasets
  - Known or search?

# Quick Break

We will resume at 2:00 pm ET.

# Genomic Analysis Use Cases & Working Groups

**Jack DiGiovanna[1] & Michael Schatz[2,3]**
[1]Program Director - Seven Bridges
[2]Program Director - AnVIL
[3]Bloomberg Distinguished Associate Professor - Johns Hopkins

# Genomic Analysis

You have a great, testable hypothesis

You are authorized to use a large dataset

You have funding for yourself and cloud usage

You have access to a powerful cloud platform(s)

You already have the tools you need running locally

**You may already be a winner and are moments away from analyzing tens of thousands of samples!**

# Is your pipeline wrapped in a workflow description language that the platform understands?

This has been a blocker for many researchers

Some users are comfortable to wrap tools, other users are not

(*before NCPI; but still somewhat today*) Certain datasets were only available on one platform, users with tools in *workflow_language_a* could only analyze that data with *workflow_language_b.*

# How to get your research done

1. Search the available pipelines on *Platform_of_Choice* (or *Platform_of_Necessity*)
2. Search pipeline repos, e.g. dockstore.org
3. Reach out to the Support teams at Platform
4. Reach out to your IC Program Officer, they know magic

**We'd like the outcome of this session to be a concrete plan of how to leverage the NCPI work to get your research done. Some possibilities**

5. Bring data over to your platform to process with *your_workflow_language*
6. Use existing tools on other platform in *other_workflow_language*

**Do steps 5-6 work? What other strategies work? What other blockers do you face?**

**Johannes Köster**
@johanneskoester

If I am not missing something, #Snakemake is currently the #MostCited generic, discipline agnostic #workflow_engine! Data taken from dimensions.ai (Sep 2020), workflow engines considered are those with citeable articles from github.com/pditommaso/awe... (Mai 2020)

data source: dimensions.ai

citations in 2020

200

150

100

50

0

Snakemake
NextFlow
KNIME
toil
Bpipe
SCOOP
COMPSs
Ruffus
Popper
Anduril
Hyperloom
ClusterFlow
Cylc
BioQueue
BigDataScript
Tibanna
Jug
SciPipe
Pwrake

WMS

Not on the graph
Cromwell (**WDL**)
Most **CWL** executors
Galaxy

# CWL, WDL, Snakemake & Galaxy WF

*The Common Workflow Language (CWL)* is an open standard for describing analysis workflows and tools in a way that makes them portable and scalable across a variety of software and hardware environments, from workstations to cluster, cloud, and high performance computing (HPC) environments.                    https://www.commonwl.org/

*The Workflow Description Language (WDL)* is a way to specify data processing workflows with a human-readable and -writeable syntax. WDL makes it straightforward to define analysis tasks, chain them together in workflows, and parallelize their execution. The language makes common patterns simple to express, while also admitting uncommon or complicated behavior; and strives to achieve portability not only across execution platforms, but also different types of users.                    https://openwdl.org/

*The Snakemake workflow management system* is a tool to create reproducible and scalable data analyses. Workflows are described via a human readable, Python based language. They can be seamlessly scaled to server, cluster, grid and cloud environments, without the need to modify the workflow definition. Finally, Snakemake workflows can entail a description of required software, which will be automatically deployed to any execution environment.                    https://snakemake.readthedocs.io/en/stable/

*A Galaxy workflow* is a series of tools and dataset actions that run in sequence as a batch operation. Workflows can be generated quickly from the analysis already completed in a history. Workflow can be reused over and over, not only reducing tedious work, but enhancing reproducibility by applying the same exact methods to all of your data.
                    https://galaxyproject.org/learn/advanced-workflow/

# Workflow Interoperability

- Searching, storing, and publishing using multiple workflow languages (e.g. Dockstore)
- Single node solutions for launching a workflow written in language X within workflow engine Y?

**Geraldine Van der Auwera** 🏳️‍🌈 🧬 🌤️ @VdaGeraldine · Oct 28
Replying to @mike_schatz @infoecho and @jdidion
You should be able to wrap the snakemake job inside an individual WDL task — the command block would be whatever you'd type in the terminal to run snakemake, and you could either hardcode the path to the workflow script or feed it as an input file. Cromwell would run that.

💬 1          ♥ 6

- Converter from Workflow X to Workflow Y?

**Nils Homer** @nilshomer · Oct 28
This. Also, almost all Snakemake workflows have python code (not just the DSL), so this makes it hard to convert formats. Same for Nextflow. But that's the power of those two (support a first class PL versus roll your own). #Bioinformatics

- Better GUI/Lint/Debugger support

## Docker Hub Image Retention Policy Delayed, Subscription Updates



JEAN-LAURENT DE MORLHON
Oct 22 2020

Today we are announcing that **we are pausing enforcement of the changes to image retention until mid 2021.** Two months ago, we announced a change to Docker image retention policies to reduce overall resource consumption. As originally stated, this change, which was set to take effect on November 1, 2020, would result in the deletion of images for free Docker account users after six months of inactivity. After this announcement, we heard feedback from many members of the Docker community about challenges this posed, in terms of adjusting to the policy without visibility as well as tooling needed to manage an organization's Docker Hub images. Today's announcement means Docker **will not enforce** image expiration enforcement on November 1. Instead, Docker is focusing on consumption-based subscriptions that meet the needs of all of our customers. In this model, as the needs of a developer grow, they can upgrade to a subscription that meets their requirements without limits.

Post Tags

- docker hub
- docker subscription
- image retention
- subscription

Categories

- All
- Products
- Community
- Engineering
- Company

While our immediate issues have been pushed back until mid-2021, these issues could resurface at any time and on any product.

Should NCPI work towards developing an alternative resource for hosted binary & container management?

# Defining Best Practices



- Should NCPI organize DREAM-like challenges for interoperability technologies?
- Should NCPI organize DREAM-like challenges for genomic analysis?
- Where would we start?

# Planning for Obsolescence

# Planning for Obsolescence



For lncRNAs, circRNAs and mRNAs, paired-end reads were aligned to the human genome with Tophat.

# Planning for Obsolescence



galaxy workflow rnaseq

About 112,000 results (0.47 seconds)

usegalaxy.org › chmy › rna-seq-differential-analysis ▾
Published Workflow | RNA-seq differential expression ... - Galaxy
Is this single-end or paired-end data? Paired-end (as individual datasets). **RNA-Seq** FASTQ file, forward re...

Lior Pachter @lpachter · Dec 2, 2017
Please stop using Tophat scholar.google.com.mx/scholar?hl=es&... Cole and I developed the method in *2008*. It was greatly improved in TopHat2 then HISAT & HISAT2. There is no reason to use it anymore. I have been saying this for years yet it has more citations this year than last #methodsmatter

18          705          818

For lncRNAs, circRNAs and mRNAs, paired-end reads were aligned to the human genome with Tophat.

# Planning for Obsolescence

Google galaxy workflow rnaseq

About 112,000 results (0.47 seconds)

usegalaxy.org › chmy › rna-seq-differential-analysis

**Published Workflow | RNA-seq differential expression ... - Galaxy**
Is this single-end or paired-end data? Paired-end (as individual datasets). **RNA-Seq** FASTQ file, forward re...

RNA-seq differential expression analysis
Annotation: RNA-seq differential analysis

**Lior Pachter** @lpachter · Dec 2, 2017
Please stop using Tophat scholar.google.com.mx/scholar?hl=es&... Cole and I developed the method in *2008*. It was greatly improved in TopHat2 then HISAT & HISAT2. There is no reason to use it anymore. I have been saying this for years yet it has more citations this year than last #methodsmatter

18     705     818

Our metadata will be our longest-lasting artifacts

For lncRNAs, circRNAs and mRNAs, paired-end reads were aligned to the human genome with Tophat.

# Discussion

What **existing** workflow interop strategies do you know about today?

What **ideal** workflow interop strategy would be most impactful for your research today?

Is workflow obsolescence a topic we should address today or in the future?

# Imaging Data

**Ashok Krishnamurthy[1] & Steve Pieper[2]**
1: RENCI, UNC-Chapel Hill and BioData Catalyst
2: Isomics, Inc. and NCI Imaging Data Commons

# Imaging Data

You will have 40 minutes.

Notes:

https://docs.google.com/document/d/1TSnNqeP_FtQ2MqiEsZU4DnDFqjkf6_vyFeQnHNBodLw/edit

# Three Imaging Use Cases

- Image Data Ingestion Workflow
- Clinical/Research Radiologist Workflow
- Imaging/AI/Machine Learning Workflow

# BioData Catalyst Image Ingestion

- BDCatalyst to date has largely focused on ingestion of clinical and genomics data
- Images present unique challenges
  - Scale of image files
    - COPDGene Phase 1 has 22M DICOM files
  - Unique PHI issues
    - DICOM metadata
    - Embedded images (eg pacemaker with serial no.)
  - DICOM protocol standards

# Team Helium Rapid Image Ingestion Proposal

- **Image de-identification**
- **Images in bucket by dbGaP Consent groups**
- **Extraction of searchable metadata from DICOM files**

- **Ingest images into Google Health API**
- **Images are segregated by consent groups**

**New image data will be put into GC in buckets.**

**GCP Bucket**

**Google Health API**

# BDC Image Ingest, DICOM Viewer App, and Authorization Proposal

# Image Workflow Proposals

# Cancer Research Data Commons (CRDC) Imaging Data Commons (IDC)

*The NCI Imaging Data Commons will be a cloud-based resource that **connects** researchers with*

1. *cancer image collections*
2. *a robust infrastructure that contains imaging data and metadata*
3. *resources for searching, identifying and viewing images*
4. *links to other Cancer Research Data Commons nodes.*

* Available to the community.
^ Components of the Data Commons Framework

# IDC Uses Google Healthcare

- Scalable DICOM service
- BigQuery (SQL) for DICOM headers
- Authentication, data security, compute, GPU, notebooks ...

Goal is to use Google to make workable site, but use open standards so the same methodology can work anywhere

# Cloud Healthcare API <sup>BETA</sup>

Standards-based APIs powering actionable healthcare insights for security and compliance-focused environments.

**GO TO CONSOLE**   **VIEW DOCUMENTATION**

Cloud Healthcare API > Documentation

## The Cancer Imaging Archive (TCIA) datasets

☆☆☆☆☆
SEND FEEDBACK

★ **Beta**

This product or feature is in a pre-release state and might change or have limited support. For more information, see the product launch stages.

The Cancer Imaging Archive (TCIA) hosts collections of de-identified medical images, primarily in DICOM format. Collections are organized according to disease (such as lung cancer), image modality (such as MRI or CT), or research focus.

The Cloud Healthcare API provides access to these datasets via Google Cloud Platform (GCP), as described in GCP data access.

### DICOM

DICOM is the established standard for storing and exchanging medical images and their metadata across a wide range of modalities, including radiology, cardiology, ophthalmology, and dermatology. DICOMweb is a REST API used for storing, querying, and retrieving these images. The DICOMweb support in Cloud Healthcare API allows existing imaging devices, PACS solutions, and viewers to interact with the Cloud Healthcare API. This can be done either directly or via open source adapters designed to support existing DICOM DIMSE protocols. This allows customers to scalably store their medical imaging data and connect their data to powerful tools for analytics and machine learning.

76

# What's inside IDC pilot release

- Data in requester pays buckets
- Metadata in BigQuery tables
- Exploration portal
- Viewer (interfacing data via proxy)
- Documentation
- Video tutorials
- Forum (Discourse)
- Example integration with tools (Colab Notebooks, DataStudio)
- Analytics use cases - under development

# Image viewer: OHIF

IMAGING DATA COMMONS

Open Health Imaging Foundation

- Free open source,
  browser-based (zero install!)
  - Open source, modern Javascript
  - DICOM standard images,
    segmentations, annotations
  - Professional design
- DICOMweb supported by Google,
  Siemens, and open source
  servers
- VTK.js WebGL visualization
- Pathology plugin development
  - DICOM Whole Slide Imaging
  - Efficient DICOMweb pyramid access



https://github.com/OHIF/Viewers

Contributors: Steve Pieper, James Petts, Erik Ziegler, Trinity Urban, Gordon Harris (OHIF), Markus Herrmann (BWH/MGH CCDS)

# Kids First Imaging Pilots

1. Began piloting imaging data generation in response to user community requests for such data.

2. As the Kids First Datasets expanded, it became apparent imaging studies are core elements of clinical data collection efforts across X01 investigators in both congenital/birth defects and cancer contexts.

3. Many investigators have access to and/or collect imaging data, but have yet to fully leverage multimodal analysis that spans imaging/clinical/genomic data.

4. Growing and emerging needs for interoperating with MIS-C and INCLUDE data.

# Kids First Imaging Pilots



However, unlike bam/cram files -- getting to a DICOM file listing is often insufficient for investigators interested in using these data to make decisions on use.

The Cohort creation process for imaging datasets requires a dedicated imaging context, supporting both "human" and "machine" use/review/analysis setting.

# Kids First Imaging -- Data workflows/oportunities

Unlike genomics, imaging is a clinical standard with **imaging data "movement" as an inherent feature of most hospital system (even if not optimized).**

Unlike genomics, imaging is a clinical standard with imaging data "movement" as an inherent feature of most hospital system (even if not optimized).

De-identification workflows is a key need!

Like genomics - multi-cloud needs . . .

**We shouldn't "throw away" FHIR-based structuring defined for imaging**

# Kids First Imaging -- Flywheel Pilot

# Kids First Imaging -- Flywheel Pilot



"See" inside DICOMs

# Kids First Imaging -- Flywheel Pilot



Project-level management

# Kids First Imaging -- Flywheel Pilot

# Kids First Imaging -- Flywheel Pilot



**Shareable Workflows**

# Discussion

- Do the 3 use cases capture what is needed?
- How far are we from Imaging data being an important driver of interoperability
- Should we set up an Imaging Interoperability Workgroup?

# Quick Break

We will resume at 3:30 pm ET.

# What We See From Here

**Draft Roadmap in FunRetro
(from Day 1)**

**RAS and Data Sharing in FunRetro
(from Day 2)**

# Prioritizing Interoperability: Datasets and Initiatives

Can we leverage and define specific datasets and initiatives as the "canvas" against which we prioritize our interoperability efforts?

GTEx
👍 5 💬 1

Common Fund Data Ecosystem (CFDE) - goal is to build interoperability at the metadata and data level for 8 + Common Fund programs
👍 8 💬 0

Cancer Data Aggregator & CCDH for discovery and query across NCI datasets
👍 11 💬 1

Workflow interoperability?

Cloud vendor interoperability?

Imaging use case review?

Training and outreach?

FHIR Use Cases?

# 2021 Outline by End of November

## Governance WG

- Recap of current next steps

> Policies needed for sharing data between IC stacks
>
> 👍 23 💬 0

> BYOD: how does "bring your own data" work across platforms? Is there any difference between sharing data between the IC stacks if it's from a canonical project or if it's provided by a researcher upload?
>
> 👍 19 💬 0

Need to also consider intersections of governance and implementations, e.g.:
- Cross-stack workflow execution, meta-data association,
- Fully executable policies

## Outreach WG

- Recap of current next steps
- Tighter integration with FHIR group
- Exploring GA4GH Discovery API

> Outreach and Training: Code examples for interop. ✏️
>
> 👍 14 💬 0

# 2021 Outline by End of November

## System Interop WG

- Recap of current next steps

Sys Interop: working with GA4GH on DRS 1.2 -> including info on how to auth for that resource
👍 10  💬 1

Expanded "light" common metadata model across systems
👍 12  💬 0

Sys Interop + FHIR: PFB & Bulk FHIR -> one or both?
👍 14  💬 0

Sys Interop: GA4GH Discovery Search &/or FHIR for query -> researchers finding data across systems
👍 9  💬 0

## FHIR WG

- Document current best practices for NCPI FHIR Model
- Limited scope pilot of prod data, access, and  tools
- Community engagement
- Identify path to providing spanning set of data and metadata
- Continue to develop on new projects and data, especially emerging studies, eg MIS-C

# Cross WG Items

**Accurate workflow cost projections** ✏
👍 **13** 💬 **0**

Accurate workflow cost projections ✏
👍 **8** 💬 **0**

Tools/workflow portability ✏
👍 **13** 💬 **0**

Data search across systems ✏
👍 **14** 💬 **1**

Search for data across stacks (next group discussion)
👍 **26** 💬 **1**

# Other Template Slides

Feel Free to Copy/Paste as Needed

# This Is Where the Title or Headline Goes.

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

# This Is Where the Title or Headline Goes.

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

# This Is Where the Title or Headline Goes.

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

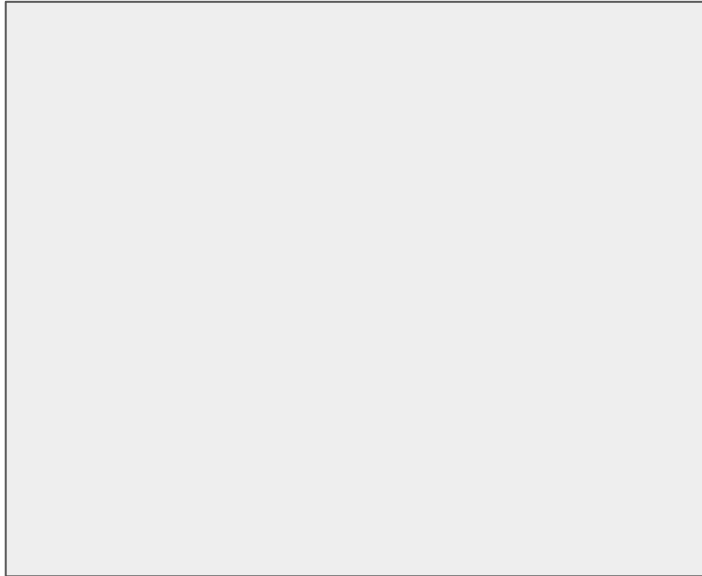# This Is Where the Title or Headline Goes.

## Compared Subject #1

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## Compared Subject #2

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

# This Is Where the Title or Headline Goes.

## Compared Subject #1

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## Compared Subject #2

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.