| client | Labkey Users Conference |
| project | User-centric design for research tools: The CDS as a case study |
| date | September, 2012 |

# HIV COLLABORATIVE DATA SPACE

HIV VACCINE DATA CONNECTOR

Logout

Search    go

Plotting 2 variables & 2 groups ⊘

Each ... ...cipant visit    swap axes    ...sou...

Explore categories
**Plot data**
Chart by time
...
Vie...

Binding magnitude & breadth

**3**
**PARTICIPANTS**
5 participant visits
6 studies
3 vaccine regimens
7 assays
5 contributors
1,715 viruses
31 unique antibodies
2 of your saved groups

**CURRENT SELECTION**
● Range: x = .7 to 1, y = .65 to 1
keep overlap    keep all    exclude    save

**ACTIVE FILTERS**
Binding & neutralization (434)
save view    clear

**REFERENCE GROUPS**
CHAVI Broad Neutralizers
☐ Only show overlap with active filter

＋ add a reference group

Active filters (3)

CHAVI broad
neutralizers (1)

NO... ...
1
0

**Federated query using DCQL and credential delegation: input values**

Diagram

DCQL_Query

Port description

DCQL query to query molecular biospecimen across cagrid data services

Example value

<DCQLQuery xmlns:ns1="http://caGrid.caBIG/1.0/gov.nih.nci.cagrid.dcql">
<ns1:TargetObject name="edu.wustl.catissuecore.domain.CellSpecimen"
xmlns:ns1="http://caGrid.caBIG/1.0/gov.nih.nci.cagrid.dcql">

Delete    New value    Add file location(s)...    Add URL ...

<DCQLQuery...    <DCQLQuery xmlns:ns1="http://caGrid.caBIG/1.0/gov.nih.nci.c
<ns1:TargetObject name="edu.wustl.catissuecore.domain.C
<ns1:targetServiceURL>https://tissueinventory.cabig.upmc.ed

Workflow description

CDS_Activity issues an EPR of the delegated credential. FQP uses this EPR to fetch the actual delegated credential from CDS and uses it to invoke multiple data services (the query activity) on behalf of the invoker.
Need to install Taverna 2 caGrid integration suite from http://www.mcs.anl.gov/~wtan/t2/ and get a cagrid Dorian account (see http://wiki.cagrid.org/display/caGrid13/Home)

Workflow author

Wei Tan

---

**Lab Viewer** *Clinical Trials Object Data System (CTODS)*    CTODS LabViewer V1.5 || Help || Log out
Login ID: cdts@nih.gov

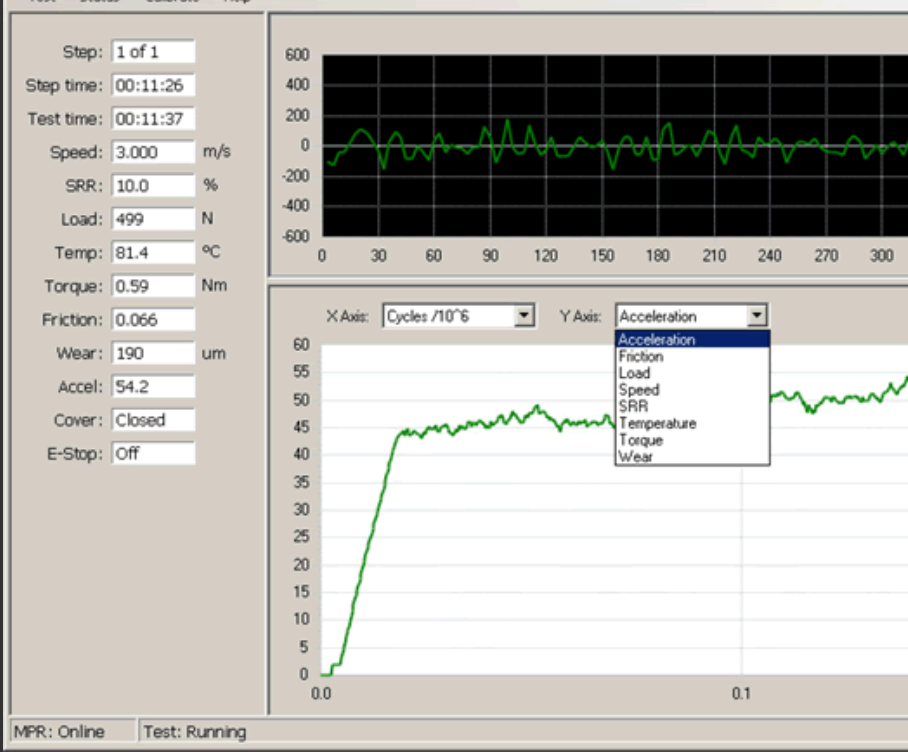Home    Study    Participant    Labs    Test    Administration

**Lab Activities - Search Results [26 record(s) found]**

View this patient in caAERS

| | Patient Id | Site | Date / Time | Lab Test | Text Result | Numeric Result | Unit Of Measure | Lower Limit | Upper Limit | Sent to CDMS | Sent to caAERS |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ | 69-98-14-2 | | 5/7/08 10:28 PM | ALK_PHOS | | 102.0 | U/L | 37.0 | 116.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | ANC | | 3.728 | mm3 | 1.0 | 7.0 | false | false |
| ☐ | 69-98-14-2 | | 4/29/08 10:28 PM | BANDS | with Polys | | | 0.0 | 4.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | BASO | | 0.1 | % | 0.0 | 3.0 | false | false |
| ☐ | 69-98-14-2 | | 4/29/08 10:28 PM | BASO | | 0.1 | % | 0.0 | 3.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | BASO_ABS | | 0.034 | mm3 | 0.0 | 0.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | BILIRUBIN_TOTAL | | 0.8 | mg/dL | 0.0 | 1.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | BUN | | 11.0 | mg/dL | 8.0 | 22.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | EOSIN_ABS | | 0.274 | mm3 | 0.0 | 0.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | GLUC_NONFASTING | | 110.0 | mg/dL | 70.0 | 115.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | HEMATOCRIT | | 32.5 | % | 31.0 | 43.0 | false | false |
| ☐ | 69-98-14-2 | | 4/29/08 10:28 PM | HEMATOCRIT | | 32.5 | % | 31.0 | 43.0 | false | false |
| ☐ | 69-98-14-2 | | 5/7/08 10:28 PM | INORG_PHOS | | 3.8 | mg/dL | 2.0 | 5.0 | false | false |
| | | | | | | 226.0 | | | | false | |

Send Labs to CDMS

---

**MPR-PC**

Test    Status    Calibrate    Help

Step: 1 of 1
Step time: 00:11:26
Test time: 00:11:37
Speed: 3.000    m/s
SRR: 10.0    %
Load: 499    N
Temp: 81.4    °C
Torque: 0.59    Nm
Friction: 0.066
Wear: 190    um
Accel: 54.2
Cover: Closed
E-Stop: Off

X Axis: Cycles /10^6    Y Axis: Acceleration
Acceleration
Friction
Load
Speed
SRR
Temperature
Torque
Wear

MPR: Online    Test: Running

---

**Microsoft Excel - xl-caBIG-smart-client**

File    Edit    View    Insert    Format    Tools    Data    Window    Help

Type a question for help

Arial    10

F7

# xl-caBIG Smart Client
Version: 0.0.1
http://xl-cabig-client.sourceforge.net

Query Parameters:
caBIG Data-Service:    Run Query
Service Data Element:    Run Query
Query Timestamp:    Never

0 Results

**Domain Object Details**

Domain Object's Model Name:    Nucleotide Sequence
Model Description:    Description Not Available

Model Identifiers
ID : 2223318    Short Name: C45374
Class Name:    NucleicAcidSequence
Package Name:    gov.nih.nci.cabio.domain

Object Attributes

| Name | Description | Val |
|---|---|---|
| id | Identifier | java |
| accessionNumber | Accession Number | java |
| accessionNumberVersion | Accession Number Version | java |
| type | Type | java |
| value | Value | java |
| length | Length | java |

EVS Concepts

| Name | Code | Definition |
|---|---|---|
| Nucleotide Sequence | C45374 | The sequence of nucleotide |

caBIG Data Browser

Ready

---

**Document Actions**

**xl-caBIG Smart Client**

Query Constructor | Connections Manager

**caBIG Query Constructor**

caBIG Data-Service: [caBIO] http://137.187.67.35:80/ogsa/s

Data-Service Details
Research Center:    NCICB
Research Center Type:    Cancer Research
Description:
National Cancer Institute Center for Bioinformatics

Comments:
The cancer Bioinformatics Infrastructure Objects (caBIO) architecture is the primary programming interface to caCORE. caBIO represents data as objects, and each object is part of a domain model that covers an area of biomedical information.

Contact Information:
Avinash Shanbhag
6116 Executive Blvd, Rockville, MD 20852
Phone Number: 301-594-9005

Data-Service's Available Objects:
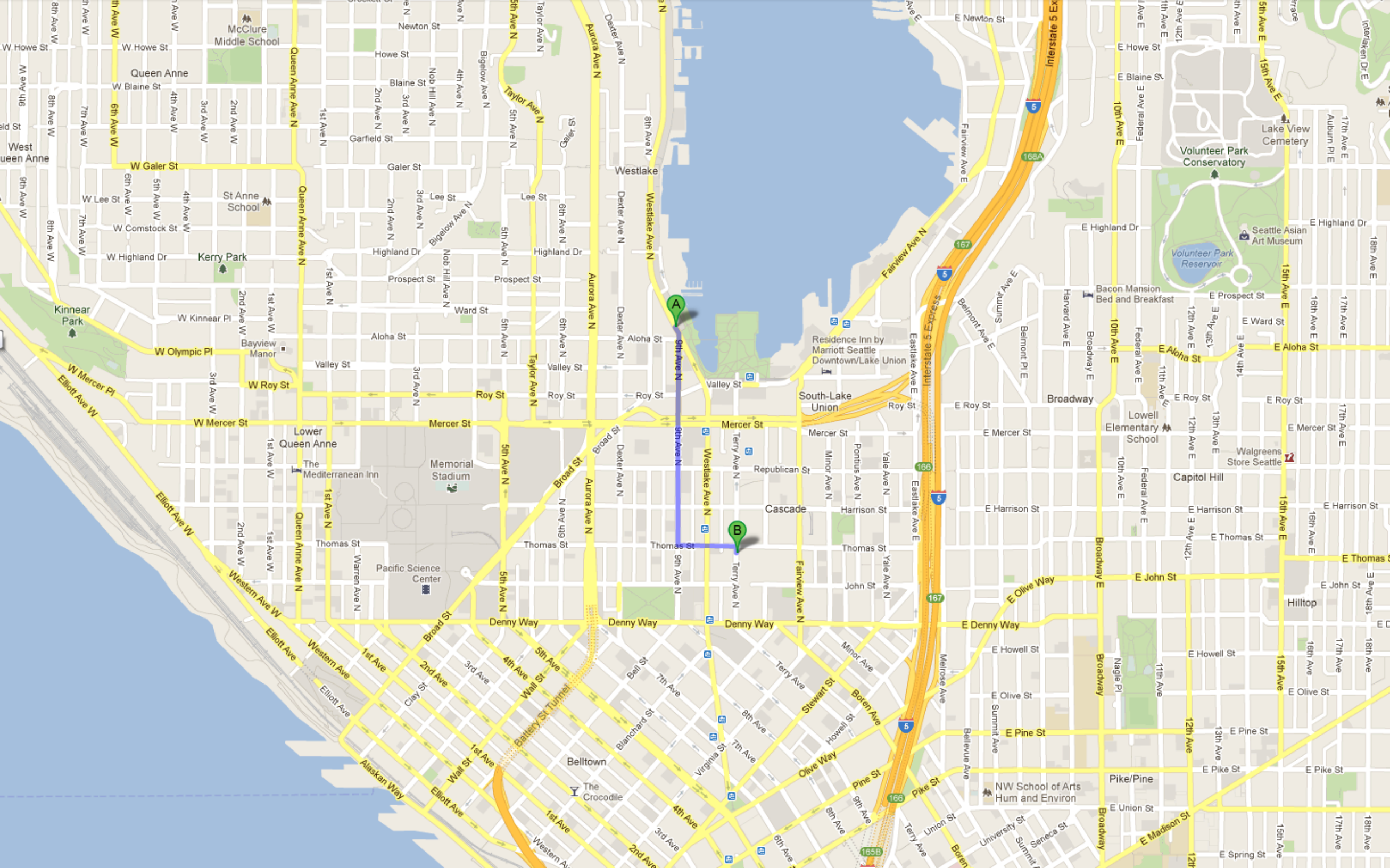Nucleotide Sequence

Import selected caBIG Data Service Objects

caBIG Connections Console
Successfully synced with Index Service.
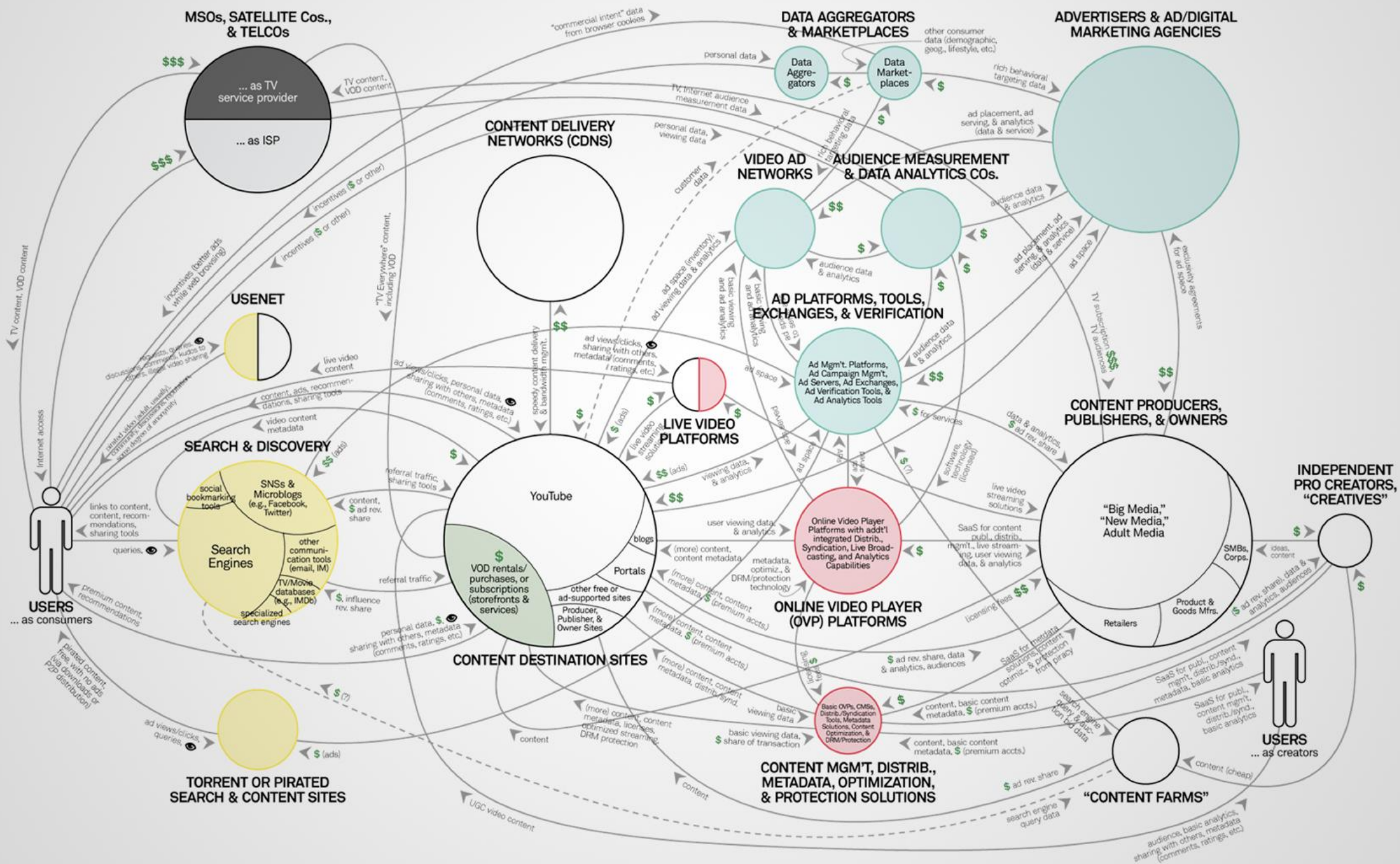
---

Start    Sent Items ...    caBIG_dotN...    2 Firefox    2 Microsof...    4 Window...    Command P...    iTunes    2 Micros...    06:18

1. Design is more about making the right thing than how it looks.

2. The process is available to you.

MSOs, SATELLITE Cos., & TELCOs

... as TV service provider

... as ISP

$$$

$$$

DATA AGGREGATORS & MARKETPLACES

Data Aggregators

Data Marketplaces

other consumer data (demographic, geog., lifestyle, etc.)

personal data

ADVERTISERS & AD/DIGITAL MARKETING AGENCIES

"commercial intent" data from browser cookies

TV content, VOD content

rich behavioral targeting data

ad placement, ad serving, & analytics (data & service)

CONTENT DELIVERY NETWORKS (CDNS)

TV, Internet audience measurement data

personal data, viewing data

VIDEO AD NETWORKS

AUDIENCE MEASUREMENT & DATA ANALYTICS COs.

$$

customer data

audience data & analytics

ad placement, ad serving, & analytics (data & service)

ad space

USENET

incentives ($ or other)

incentives ($ or other)

"TV Everywhere" content, including VOD

live video content

$$

ad views/clicks, sharing with others, metadata (comments, / ratings, etc.)

AD PLATFORMS, TOOLS, EXCHANGES, & VERIFICATION

Ad Mgm't. Platforms, Ad Campaign Mgm't, Ad Servers, Ad Exchanges, Ad Verification Tools, & Ad Analytics Tools

audience data & analytics

basic viewing and analytics

ad space (inventory), ad viewing data & analytics

speedy content delivery & bandwidth mgmt.

TV content, VOD content

Internet access

incentives (better ads while web browsing)

content, ads, recommendations, sharing tools

video content metadata

LIVE VIDEO PLATFORMS

live video streaming solutions

$$ (ads)

$$

exclusivity agreements for ad space

TV subscription $$$
TV audiences

$$

CONTENT PRODUCERS, PUBLISHERS, & OWNERS

data & analytics, $ ad rev. share

$ for services

software, technology (licensed)

SEARCH & DISCOVERY

SNSs & Microblogs (e.g., Facebook, Twitter)

social bookmarking tools

other communication tools (email, IM)

Search Engines

TV/Movie databases (e.g., IMDb)

specialized search engines

queries,

links to content, content, recommendations, sharing tools

$ (ads)

content, $ ad rev. share

referral traffic, sharing tools

YouTube

blogs

Portals

VOD rentals/ purchases, or subscriptions (storefronts & services)

$

other free or ad-supported sites

Producer, Publisher, & Owner Sites

CONTENT DESTINATION SITES

viewing data, & analytics

$$

user viewing data & analytics

(more) content, content metadata

metadata, optimiz., & DRM/protection technology

Online Video Player Platforms with add'l integrated Distrib., Syndication, Live Broadcasting, and Analytics Capabilities

ONLINE VIDEO PLAYER (OVP) PLATFORMS

live video streaming solutions

SaaS for content publ., distrib., mgm't., live streaming, user viewing data, & analytics

"Big Media," "New Media," Adult Media

SMBs, Corps.

INDEPENDENT PRO CREATORS, "CREATIVES"

$

ideas, content

$

Product & Goods Mfrs.

Retailers

licensing fees $$

USERS ... as consumers

premium content, recommendations

queries,

referral traffic

personal data, $, sharing with others, metadata (comments, ratings, etc.)

(more) content, content metadata, $ (premium accts.)

(more) content, content metadata, $ (premium accts.)

(more) content, content metadata, licenses, optimized streams, DRM protection

$ ad rev. share, data & analytics, audiences

SaaS for metadata solutions, content optimiz., & protection from piracy

SaaS for publ., content mgm't, distrib./synd., metadata, basic analytics

$ ad rev. share/ data & analytics, audiences

$

SaaS for publ., content mgm't, distrib./synd., basic analytics

USERS ... as creators

pirated content: free, with no ads (big downloads or P2P distribution)

TORRENT OR PIRATED SEARCH & CONTENT SITES

ad views/clicks, queries,

$ (ads)

$ (?)

content

basic viewing data

(more) content, content metadata, optimized streams, DRM protection

basic viewing data, $ share of transaction

Basic OVPs, CMSs, Distrib./Syndication Tools, Metadata Solutions, Content Optimization, & DRM/Protection

content, basic content metadata, $ (premium accts.)

content, basic content metadata, $ (premium accts.)

CONTENT MGM'T, DISTRIB., METADATA, OPTIMIZATION, & PROTECTION SOLUTIONS

$ ad rev. share, data & analytics

$ ad rev. share

UGC video content

content

search engine query data

search engine query & production data

"CONTENT FARMS"

content (cheap)

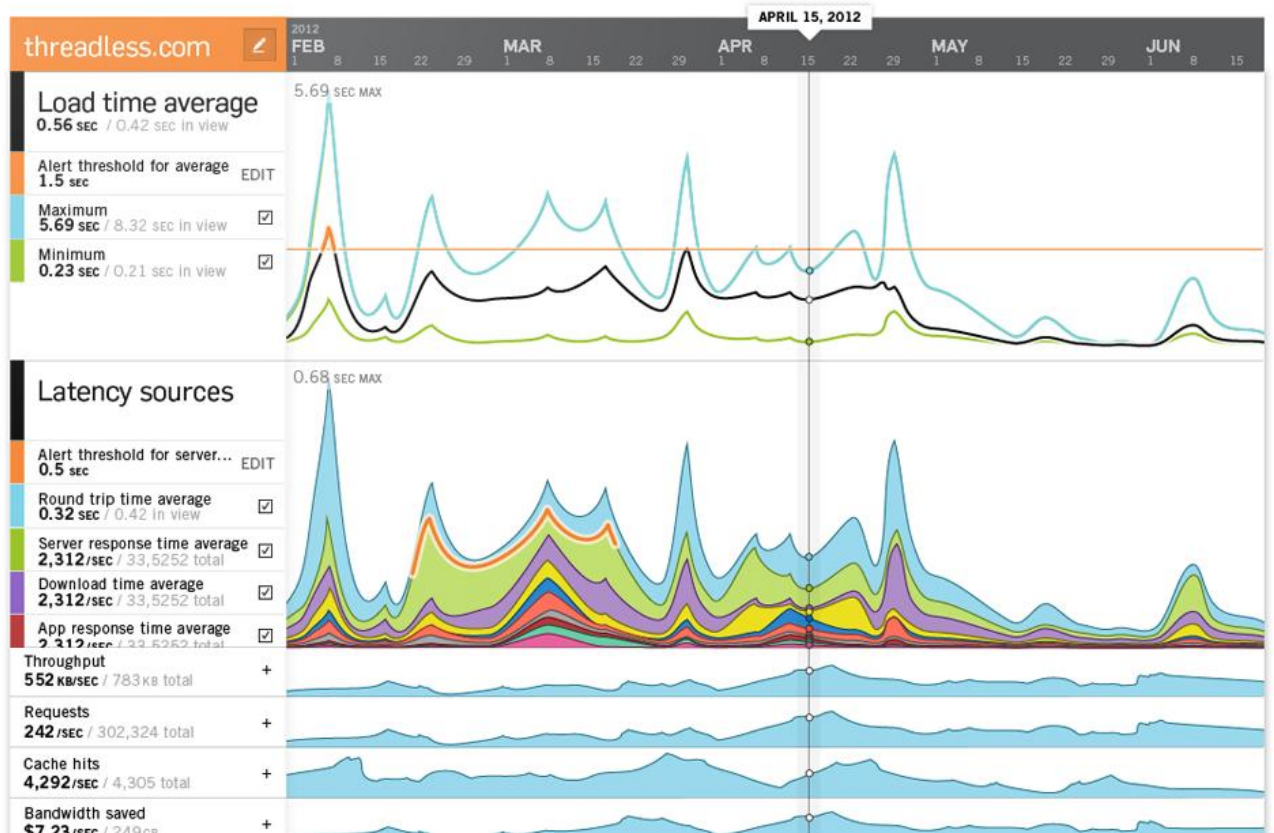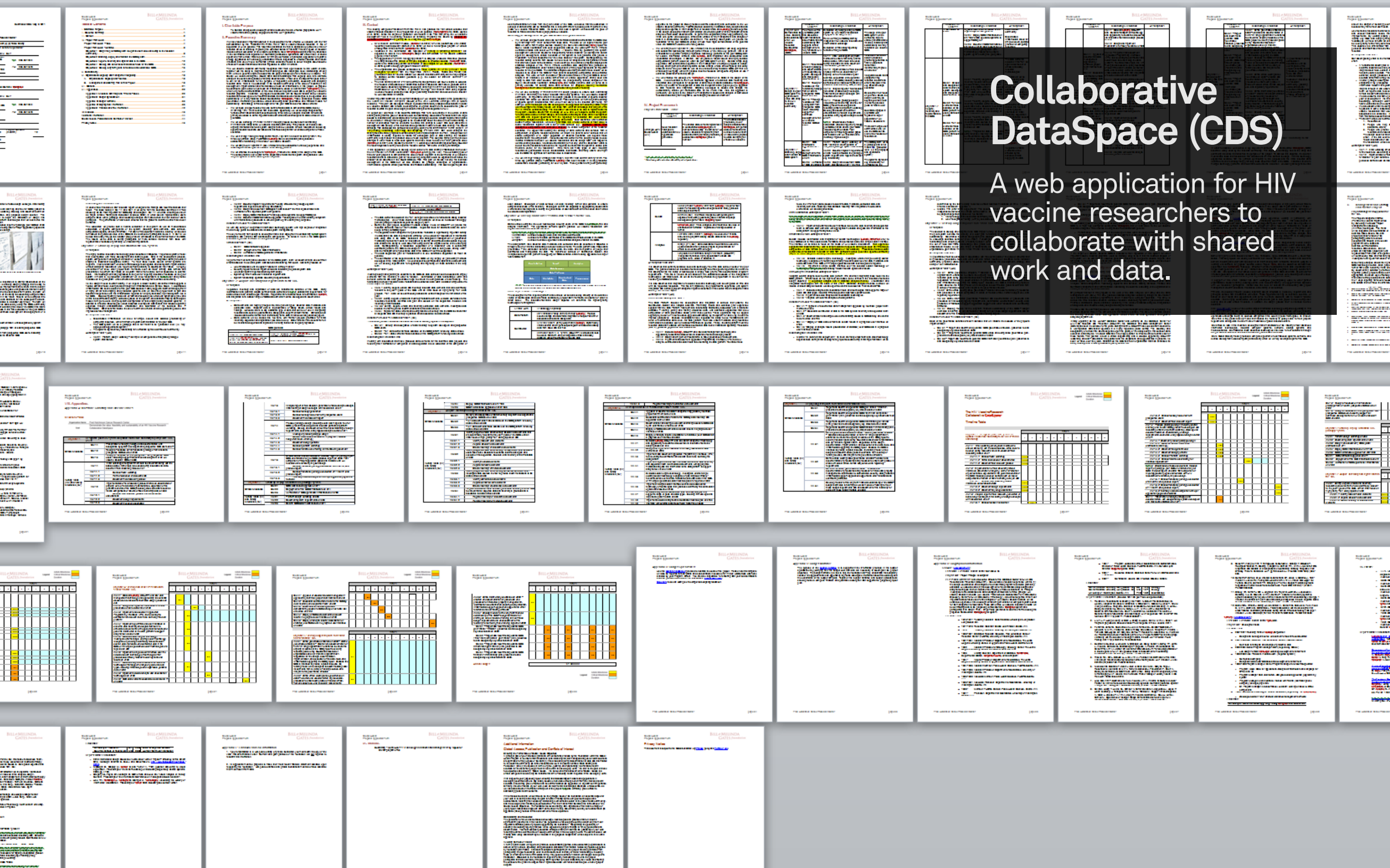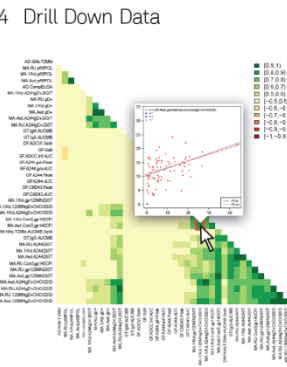audience, basic analytics, sharing with others, metadata (comments, ratings, etc.)

**Collaborative DataSpace (CDS)**

A web application for HIV vaccine researchers to collaborate with shared work and data.

# Drill Down Data

# Push Communications

# Drill Down Data

Inbox

## HIV/AIDS Research Monthly Digest
Compiled for you by the Collaborative Data Space

New Data This Month

New secondary analysis of RV 144

Lorem ipsum dolor sit amet, consectetur adipisicing 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

| Variable 1 | Variable 2 | Variable 3 |
|---|---|---|
| 5.1 | 1.4 | 0.2 |
| 4.9 | 1.4 | 0.1 |
| 4.7 | 1.3 | 0.3 |
| 4.6 | 1.4 | 0.2 |

# Automatically Identified Relationships For Your Data

The system found the following relationships. Please evaluate:

View underlying data

R=.8

View underlying data

R=.6

# Compare Data

| Data Set 1 | | |
|---|---|---|
| Variable 1 | Variable 2 | Variable 3 |
| 5.1 | 1.4 | 0.2 |
| 4.9 | 1.4 | 0.1 |
| 4.7 | 1.3 | 0.3 |
| 4.6 | 1.4 | 0.2 |
| 5.0 | 1.5 | 0.2 |

Similar Data Sets

98% similar

50% similar

50% similar

Upload Complete

# Interactive time visualizations

Patient 231a

Day 121

Day 1    Day 145

Day 121

# Subscriptions

PI: Gregory Thomas
University of Alabama at Birmingham, USA
Biology of Transmitted/Founder Viruses

SUBSCRIBE

New Data Sets

| Data Set 1 | | | | Data Set 2 | | |
|---|---|---|---|---|---|---|
| Variable 1 | Variable 2 | Variable 3 | | Variable 1 | Variable 2 | Variable 3 |
| 5.1 | 1.4 | 0.2 | | 5.1 | 1.4 | 0.2 |
| 4.9 | 1.4 | 0.1 | | 4.9 | 1.4 | 0.1 |
| 4.7 | 1.3 | 0.2 | | 4.7 | 1.3 | 0.3 |
| 4.6 | 1.4 | 0.2 | | 4.6 | 1.4 | 0.2 |
| 5.0 | 1.5 | 0.2 | | 5.0 | 1.5 | 0.2 |

New Visualizations

Inbox
RSS
CDS Home Page

# Remote syncronous session

Do you think this data supports our hypothesis?

Collaborative Data Space

I don't know. Can I have control so I can explore the data?

Collaborative Data Space

# Find Opportunities

Studies

| | CHAVI 001 | CHAVI 008 | CHAVI Broad Neut | RV 144 | RV217 | HVTN 304 | HVTN 068 | ... |
|---|---|---|---|---|---|---|---|---|
| | Less studied | Moderately studied | Less studied | Not studied | Not studied | Highly studied | Moderately studied | ... |
| | Highly studied | Not studied | Not studied | Less studied | Not studied | Less studied | Less studied | |
| | Less studied | Not studied | Less studied | Highly studied | Less studied | Not studied | Moderately studied | |
| | Moderately studied | Highly studied | Highly studied | Moderately studied | Moderately studied | Highly studied | Less studied | |
| | Highly studied | Not studied | Less studied | Less studied | Moderately studied | Highly studied | Highly studied | |
| ... | ... | ... | ... | ... | ... | ... | ... | |

# Physical Links To Live Data

Feel free to navigate this data using the collaborative data space.

Collaborative Data Space

Request Access

# Relationship Map

Antibodies

Epitopes

Your item of interest

Vaccines

Mutations

# Publicly Generated Analysis

Non-professionals generate ideas

What if...

Scientists review and follow-up

Filter the best ideas

What if...

# Discover Related Content

Collaborative Data Space

## RV 144 Study

| Data Set 1 | | |
|---|---|---|
| Variable 1 | Variable 2 | Variable 3 |
| 5.1 | 1.4 | 0.2 |
| 4.9 | 1.4 | 0.1 |
| 4.7 | 1.3 | 0.3 |
| 4.6 | 1.4 | 0.2 |
| 5.0 | 1.5 | 0.2 |

Related Materials

**Papers**
Lorem ipsum dolor
Consectetur adipiscing elit
Fusce et est
Tincidunt lorem
Dapibus lobortis

**Principal Investigators**
Lorem ipsum dolor
Consectetur adipiscing elit
Fusce et est
Tincidunt lorem
Dapibus lobortis

**Posters**
Lorem ipsum dolor
Consectetur adipiscing elit
Fusce et est
Tincidunt lorem
Dapibus lobortis

**Presentations**
Lorem ipsum dolor
Consectetur adipiscing elit
Fusce et est
Tincidunt lorem
Dapibus lobortis

**Abstracts**
Lorem ipsum dolor
Consectetur adipiscing elit
Fusce et est
Tincidunt lorem
Dapibus lobortis

**Annotations**
Lorem ipsum dolor
Consectetur adipiscing elit
Fusce et est
Tincidunt lorem
Dapibus lobortis

# Generate Citations

Copy to Clipboard

References

1. Vichey LA, Miller DA, Lindsay JM, et al. HIV prevalence and associated risks in young men. JAMA. 2000; 284: 198-204.

2. Nettles L, Han S, Bottler W, et al. Lifetime risk factors for HIV/sexually transmitted infections among male-to-female transgender persons. J Acquir Immune Defic Syndr. 2009; 52(3): 417-21.

3. Herbert JH, Jacoby ED, Finland TJ, et al. Estimating HIV prevalence and risk behaviors of transgender persons in the United States: a systematic review. AIDS Behav. 2008; 12(1): 1-17.

Create Citation

# What's Popular

Collaborative Data Space

## What's Hot this Week

Other Title

Viewed 500 times by 40 researchers.

Viewed 465 times by 34 researchers.

New

Viewed 237 times by 12 researchers.

Viewed 123 times by 10 researchers.

Viewed 50 times by 10 researchers.

Viewed 32 times by 4 researchers.

# Interactive Annotations

HIV Lab Experiment

Introduction
Lorem ipsum dolor sit amet, consectetur adipisicing elit. Fusce et est tincidunt lorem dapibus lobortis. Mauris ultricies ligula felis. Integer facilisis, felis at pretium felis.

Caveats
Lorem ipsum dolor sit amet, consectetur adipisicing elit. Fusce et est tincidunt lorem dapibus lobortis.

Conclusions

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Fusce et est tincidunt lorem facilisisdapibus lobortis. Mauris ultricies ligula felis. Integer facilisis, felis at consectetur iaculis facilisis.

## Super Watson

Imagine that you arrive at this new tool filled with data from multiple studies and labs working on HIV. There's a search box powered by a version of IBM's Watson from 10 years in the future. Watson will answer any question you have about all this data. That includes interpretations and judgments. But you only get 5 questions! Write out your questions.

# What could you do with aggregate data that you can't do now?

# Tools Inventory

$\checkmark$ = pro   $\otimes$ = con

| Sample processing | Data exploration and discovery | Visualization | Data analysis (descriptive, comparative, correlative) |
|---|---|---|---|

| Paper writing (text, figures) | Collaborating with other labs | Researching the field | Other |
|---|---|---|---|

**How should this tool relate to all the other tools in use?**

## Search Slices

Imagine that you arrive at this new tool filled with data from multiple studies and labs. There are only 3 ways that all this data is cataloged to help you find what you need, whether it's by browsing or searching or finding similarities. Circle the 3 that you would use the most frequently.

Studies (RV 144, CHAVI 008, HVTN 068, etc...)

People (participant IDs)

Participant attribute _____

Antibodies

Virus clades and subclades _____

Virus epitopes

ENVs

Vaccines _____

Geographic region _____

Assays

Visualization types

Analysis types

Other _____

Other _____

Other _____

# How should information be organized?

# Barriers and Benefits

In a few years, you will be able to upload your data and align it to what already there so that you can use it in combination with other's data, and others can use yours too. List the factors that influence whether you would choose to share.

## How likely are you to share your data?

○ definitely    ○ very likely    ○ somewhat likely    ○ don't know    ○ somewhat unlikely    ○ very unlikely

## Benefits

☐ CDS has a feature you need in order to analyze and interpret your own data

☐ Other researchers are uploading their data

☐ My grants require me to share in CDS

☐ My data are referenced in others' publications

☐ I am invited to be a coauthor because others are relying on my data

## Barriers

☐ Others will use my data and not credit me

☐ Others will publish something I'm planning before

☐ It's time consuming or difficult to upload and align

☐ There are too many caveats and nuances for other use of certain data

**What are specific privacy concerns and how can we overcome them?**

Do people really want to collaborate in here? How?

**1** Throughout the system there is **always a working set** of data.

Everywhere you go in the system, there an active set of data follows. "Everything" is selected at first, until filtered out by user actions.

---

**2** **Participants and visits** are the objects being viewed, sorted, filtered, and connected across data.

This overlap is needed for valid analysis. Every other noun (assay, vaccine, virus) exists primarily to help you find, analyze, and interpret. Getting information on these other nouns is secondary.

Communicating this model clearly is key.

---

**3** We prioritize **fast filtering** to intersecting data.

The interface is biased to filtering intersecting participants / visits (Boolean AND). Overlap leads to meaningful analysis and fast data reduction. Again, communicating this behavior is critical.

Communicating this model clearly is key.

---

**4** The current working set **persists across view** changes.

But each view is good at different kinds of filtering. E.g. Data Explorer for categorical filtering, Scatter Plot for interval filtering, Time for temporal...

---

**5** Information to help **interpret and define** is always near.

There is a persistent area devoted to this on the screen. Variables (columns) also show a description when selected.

---

**6** Data with **no or unknown value** are visible.

Seeing the missing pieces provides cues to filter changes or potential implications on conclusions. E.g. showing set-point viral load and seeing that there are some people in the current filter with none is important.

---

**7** Core tools **share consistent interactions** with each other.

Filling the role of each view with a different 'off-the-shelf' tool will lead to a disruptive user experience. Where divisions exist they must be made extremely clear and integrate with CDS core tools if data are returned.

# Can people use this?

# HIV VACCINE DATA
## CONNECTOR

Logout

Search  go

## All participants
Active filters (unsaved)
Search Results

**COMMUNITY GROUPS**

CHAVI broad neutralizers
CHAVI 20
All HLA-IIs
Male sex workers
RV144 from the 2012 NEJM paper

See all 19

**MY GROUPS**

My Contributed Data
Top 10% ADCVI samples
Infected female HLA-IIs
Broadest cytokine responders
Best NKTs and ADCC
Everyone tested on Ag CX29 Env

See all 8

## FIND PARTICIPANT GROUPS...

by **Studies** — CHAVI 08, CHAVI Broad Neutralizers, HVTN 204, HVTN 068, NSDP, MRK-AD5 — **6 total**

by **Antigens** — 7 clades, 3 tiers, 4 antigen sources (infection, reagents, IMCs, pseudoviruses) — **2 total**

by **Assays** — 23 Adaptive: humoral & B-Cell, 36 Adaptive: T-Cell, 4 Diagnostic & clinical, 6 Host genitic... — **63 total**

by **Vaccines** — 1 DNA primes, 2 Boosts, 2 Adjuvants — **3 total regimens**

by **Contributors** — Alpha, Beta, Carot, Donut, Eggplant — **38 total labs**

by **Demographics** — 9 races, 3 HLA types, 3 infection statuses, 31 locations, 2 genders — **2,131 total participants**

by **Antibodies** — 8,039 Env, 1,323 IgA, 812 IgE, 730 IgD, 1,400 IgG, 613 IgM — **14,321 total**

by **Saved views** — See all — **28 total**

**OR PASTE PARTICIPANT, VISIT, OR SAMPLE IDENTIFICATION NUMBER(S):**

PTID, PTID...  go

Logout

Search | go

## Assays ❯ by Choose category ❯

SORTED BY: **TYPE** ⌄

find assay | go

Showing number of: Participants ( hide empty ) ( export )

− Adaptive: humoral & B-Cell — 184

**Antibody dependent cellular cytotoxicity** / Last Name, Last Name, Last Name … 128

**Binding Antibody** Multiplex array / Last Name — 434 ⓘ view assay info

**Cytokine Multiplex Bead Array** / Last Name — 80

**Neutralizing Antibody** / Last Name — 140

Assay Name two / Last Name — 100

**Assay** Name three / Last Name — 180

Assay Name two / Last Name — 95

**Assay Name three** / Last Name — 204

+ Adaptive: T-cell — 200

+ Diagnostic & Clinical — 2,155

+ Host Genetic — 204

+ Innate — 690

+ Other & cross-category — 900

## 1,128

**PARTICIPANTS**

23,201 participant visits
6 studies
3 vaccine regimens
38 assays
22 contributors
1,715 viruses
31 unique antibodies
2 of your saved groups

**ACTIVE FILTERS**

All participants

( save view )

**REFERENCE GROUPS**

( + add a reference group )

to compare to your active filters

Logout

Search | go

# Assays ⊙

SORTED BY: **TYPE** ⌄

find assay | go

# by Studies ⊙

SORTED BY: **TYPE** ⌄

find studies | go

Showing number of: Participants   ( hide empty )  ( export )

| | CHAVI 08 | CHAVI Broad Neutralizers | HVTN 204 | RV 144 | |
|---|---|---|---|---|---|
| − Adaptive: humoral &... | 401 | 9 | 33 | 33 | |
| **Antibody dependent...** | 23 | 9 | 22 | 22 | |
| **Binding Antibody...** | 3 | 9 | 2 | 2 | |
| **Cytokine Multiplex...** | 23 | 9 | 22 | 22 | |
| **Neutralizing Antibody...** | 180 | 9 | 6 | 0 | |
| **Assay Name two / Last...** | 180 | 9 | 6 | 0 | |
| + Adaptive: T-cell | 400 | 9 | 19 | 19 | |
| + Diagnostic & Clinical | 281 | 9 | 0 | 0 | |
| + Host Genetic | 4 | 9 | 25 | 25 | |
| + Innate | 401 | 9 | 33 | 161 | |
| + Other & cross-category | 400 | 9 | 125 | 19 | |

# 1,128
**PARTICIPANTS**

23,201 participant visits
6 studies
3 vaccine regimens
38 assays
22 contributors
1,715 viruses
31 unique antibodies
2 of your saved groups

**ACTIVE FILTERS**

All participants

( save view )

**REFERENCE GROUPS**

( + add a reference group )

to compare to your active filters

Search | go

# Plotting 2 variables & 2 groups ❯

**Each dot represents:** Participant visit ⌄ | swap axes | view sources | export

r²=.65

Binding magnitude & breadth

NO OVERLAP

Neutralization magnitude & breadth

Active filters (3)

CHAVI broad neutralizers (1)

## 3
**PARTICIPANTS**
5 participant visits
6 studies
3 vaccine regimens
7 assays
5 contributors
1,715 viruses
31 unique antibodies
2 of your saved groups

**CURRENT SELECTION**
● Range: x = .7 to 1, y = .65 to 1
keep overlap | keep all | exclude | save

**ACTIVE FILTERS**
Binding & neutralization (434)
save view | clear

**REFERENCE GROUPS**
CHAVI Broad Neutralizers
☐ Only show overlap with active filter
+ add a reference group

# Charting 2 variables and 3 groups ⊙

**Each line represents:** Participant group ⌄    ⟨ view sources ⟩    ⟨ export ⟩

## 1

**PARTICIPANT**

5 participant visits
6 studies
3 vaccine regimens
7 assays
5 contributors
1,715 viruses
31 unique antibodies
2 of your saved groups

**CURRENT SELECTION**
● PTID: 3552623
⟨ keep all ⟩  ⟨ exclude ⟩  ⟨ save ⟩

**ACTIVE FILTERS**
● Binding & neutralization (434)    ⊞
⟨ save view ⟩  ⟨ clear ⟩

**REFERENCE GROUPS**
● Binding & neut high performers (3)    ⊟
● CHAVI Broad Neutralizers (9)    ⊠ ⊞
☐ Only show overlap with active filter

⟨ + add a reference group ⟩



Viral load — 200K, 100K, 0

Estimated infection date          12 weeks post infection

**Time since infection**

# 1. Design is more about making the right thing than how it looks.

# What is CDS not?

Dropbox / Atlas: directory and file-based sharing without added value

Completely public to 7 billion people

A specialized Wikipedia

The end of clarifying phone calls and emails

A replacement for statisticians

A replacement for new lab work

A source for HIV research news

"Shotgun science"

A new interactive paper format

A way to administer and evaluate study execution

Webex – live synchronous collaboration

# Dataset? Datacube.

## Data set-centric

Data set A + Data set B = Data set AB

## PTID- / visit-centric
(additional power of CDS)

Data attribute ? = Participants with data attribute

# Open? **Mixed**.

Private workspaces cannot go away. But fresh public data is critical.

**Envisioned levels of data access**



Whole HIV vaccine community

Research networks

Ad hoc groups

Me

# Collaborative? Communicative.

Ironically, the Collaborative DataSpace won't be a place for rich community collaboration after all.*

**\*But…**

# You **can't annotate completely.**

I needed to know what reactions participants had to the vaccine, which isn't part of their protocol to measure or publish but is very relevant to HIV diagnostics.
John Hural

I don't have a clue how to analyze others' data. There are so many vagaries to the process. "It's impossible to codify this well."
Mario Roederer

There might be 20 things per assay and it's really hard to standardize. It's even counter-productive; science moves on by the time you've standardized. You'll be out of date.
Rick Koup

Unless you're close to all the details you're going to screw up the analysis.
Peter Gilbert

Someone else looking at the same data might want something completely different. They might want to know how someone was infected. I don't care, though. I care about viral load and CD4.
Nicole Frahm

39

# **People don't consider** all the things that they should.

'My biggest concern isn't credit; it's misinterpretation.'
Shaunna Shen

"People... don't think about what region PTIDs come from or whether they have an STD. I help them find the right data considering all these factors."
Kelly Soderberg

You have to dig. It's only when I talked to them on the phone with these questions I realized they were f*#$ed.
Nathan Vandergrift

I would rather contact the person who posted the data. It's more efficient and faster.
Nicole Frahm

"Someone might make a claim using my data that I don't agree with. I want to know when they're presenting with it."
Georgia Tomaras

40

# It takes **too much effort** to annotate deeply. No one will do it.

It's "all about what I get in return. If it takes my time I have to have a reason... annotation is really hard."
Nicole Frahm

"It's theft of my time." Every login, every extra step, every administrative need.
Danny Douek

41

# Benefits of communication are the primary **carrots** to get future data.

Fresh ideas about my data, new connections to others' data, better context and interpretation of data, a greater likelihood of receiving credit, a chance to find collaboration opportunities

We feel the data is ours. If you share data that's in progress (even published data) it's especially important that people talk to you.'
Mario Roederer 1

If you host unpublished data and someone uses it without crediting the source, that's a problem you'll only have once because no one will share again.
Bart Haynes

42

# CDS should include **cues** about quality and linked **metadata.**

Is it peer reviewed and published? Link to the paper.

What were the key assay characteristics?

Is there an assay abstract? Is the assay experimental?

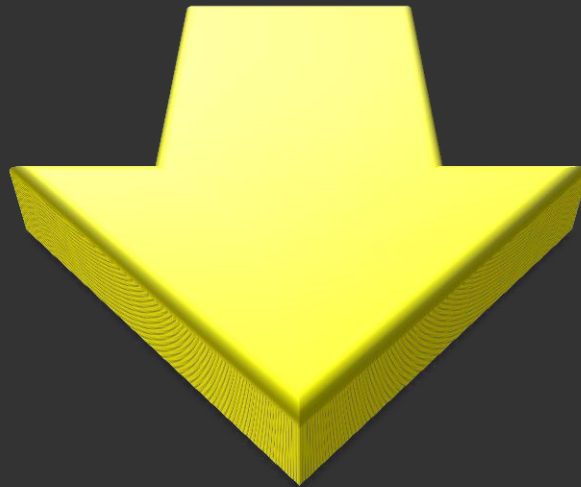Does it use GLP/GMP and provide metadata?

Whom should I contact about the data and how? Whom should I contact about the study?

43

# CDS will need **staff** to help annotate and align in the future.

44

# 2. The process is available to you.

Look for context: talk and observe before you start building

Validate assumptions: is the explicit ask really what's needed?

Iterate at low fidelity: Test, fail, learn early and often

Prioritize: optimize for key tasks rather than exposing everything

# Dialog