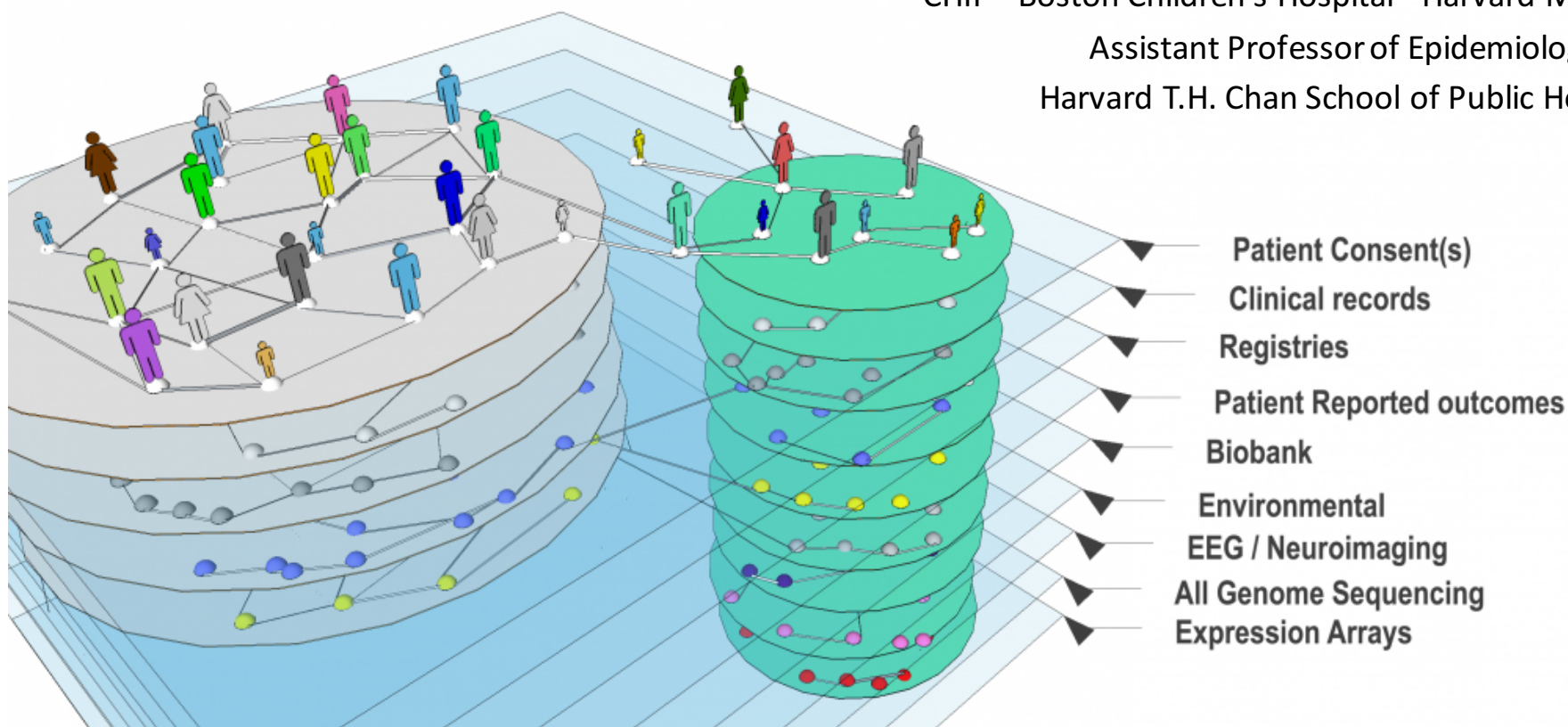


# I2b2/tranSMART BD2K PIC-SURE RESTful API

**Paul Avillach, MD, PhD**

Assistant Professor of Pediatrics and Biomedical Informatics  
CHIP - Boston Children's Hospital - Harvard Medical School

Assistant Professor of Epidemiology  
Harvard T.H. Chan School of Public Health



**HARVARD**  
MEDICAL SCHOOL

DEPARTMENT OF  
Biomedical Informatics



Advance statistical tools  
Biobank explorer  
Variant explorer



Advance cohort selection



Federated Advance cohort selection



**RESTful API**  
GET PUT POST DELETE

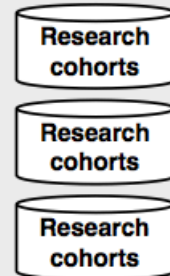
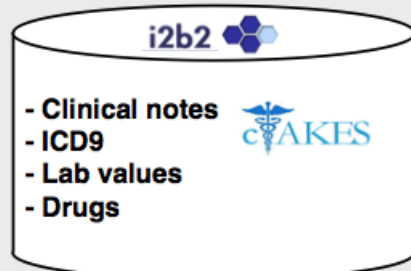


**SMART**<sup>®</sup>  
Patient level data lookup  
Interoperable tools





**hospital**  
**EHR**

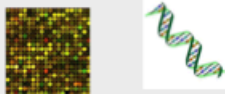
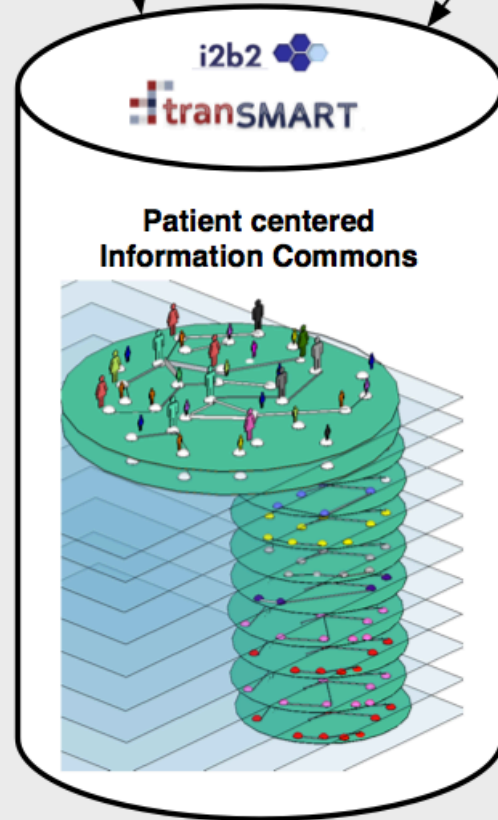


**Biobank Core**



**Patients' consent(s)**

Selected patients

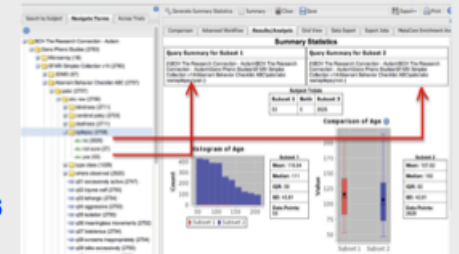


**'omics**  
Gene exp  
SNPs  
Whole Exome Seq  
RNA seq

**BCH Private network**



**Generation of Hypotheses by PIs**



**High throughput analysis by Biostatistician & Bioinformatician**

→ integration  
↔ analysis

- [-] Autism
  - [+] BCH Consent Information (2975)
  - [-] BCH EHR i2b2 (16587)
    - [+] Allergy (8184)
    - [+] Clinic (15547)
    - [+] Demographics (16522)
    - [+] Diagnoses (16426)
    - [+] Insurance Payors (13949)
    - [+] Labtests (12073)
    - [+] Medications (9371)
    - [+] Procedures (10140)
    - [+] Radiology (8736)
    - [+] Reports (4414)
    - [+] Service (16587)
    - [+] Vital Signs (11982)
  - [-] Expression Array (713)
    - [-] Peripheral Blood (713)
      - [-] Affymetrix Human Genome U133 Plus 2.0 Array (198)
        - [X] gcRNA Normalized Expression (198)
      - [-] Human Gene 1.0 ST Array (550)
        - [X] RMA Normalized Expression (550)
  - [+] GenoPheno (3189)
  - [+] SFARI Simplex Collection v14 (2760)
  - [-] WES\_Inventory (7373)
    - [+] Eichler (5349)
    - [+] State (2024)

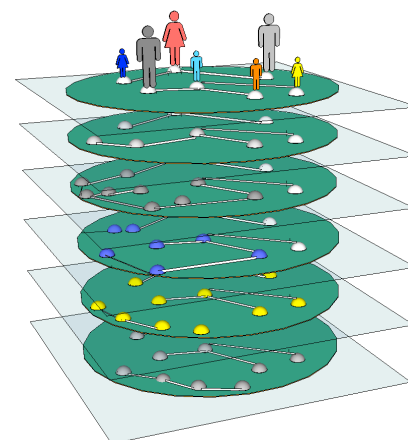
- **Patient consent(s)**

- **EHR longitudinal data**

- **Expression Arrays**

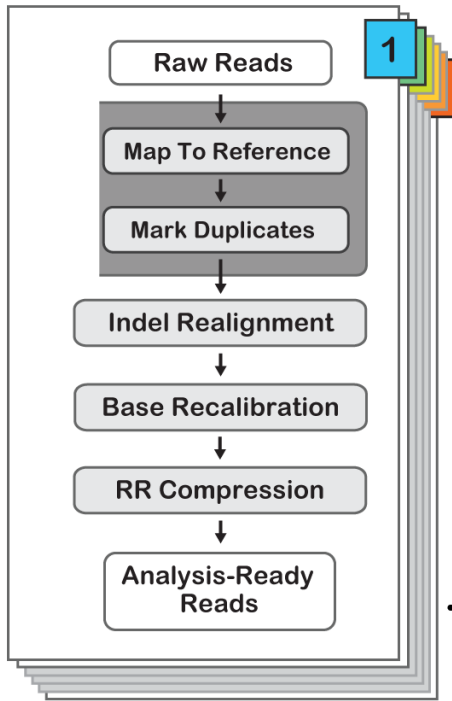
- **Clinical Cohorts**

- **WES data**

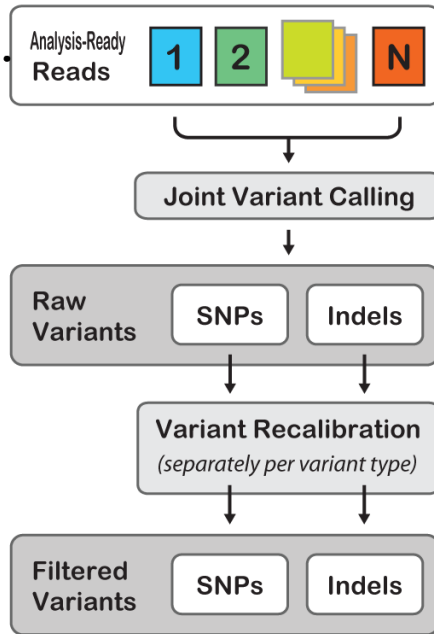




## Exome sequence data processing



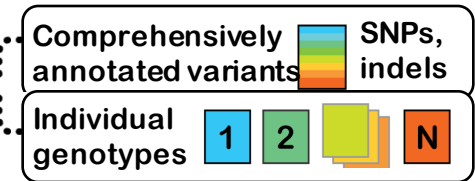
## Variant calling



## Variant annotation

- Physical location  
e.g. Chr:start-end  
Cytoband  
...
- Gene  
e.g. Gene name  
Variant function  
...
- Gene set  
e.g. Pathway  
Molecular process  
...
- Predicted variant impact  
e.g. SIFT  
PolyPhen  
...
- Conservation  
e.g. GERP  
PhyloP  
...
- Population frequency  
e.g. 1000 Genomes  
ESP 6500  
...
- Clinical significance  
e.g. ClinVar  
OMIM  
...
- Expression patterns  
e.g. GTEx  
BrainSpan  
...
- Transcriptional regulation  
e.g. ENCODE TFBS  
Histone modifications

## I2b2/tranSMART input



**ANNOVAR**

Search by Subject | **Navigate Terms** | Across

Generate Summary Statistics | Generate WES Statistics

**Comparison** | Advanced Workflow | Results/Analysis | Gr

**Subset 1**

Exclude | Enable Variant Panel

... \HLA-DQB1\ <0 **Phenotypic variable**

AND | Exclude | Disable Variant Panel

... \HLA-DQB1\ **Genomic variables**

AND | Exclude | Disable Variant Panel

... \0|1\  
... \1|0\  
... \1|1\  
AND | Exclude | Disable Variant Panel

... \nonsynonymous SNV\  
... \stopgain SNV\

AND | Exclude | Enable Variant Panel

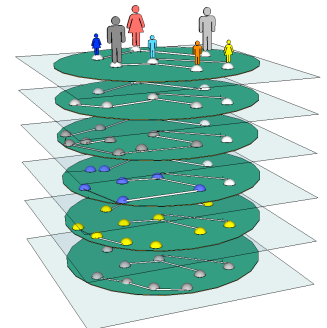
abc stopgain SNV (55)

1kG LCL Proteomics

- 01 Demographics (55)
  - Population (55)
  - Sex (55)
    - abc female (26)
    - abc male (29)
- 02 Whole Exome Variation (55)
  - 01 Physical location (55)
  - 02 Gene (55)
    - 01 Refseq (55)
      - 01 Gene symbol (55)
      - 02 Variant function (55)
      - 03 Exonic variant function (55)
        - abc frameshift deletion (55)
        - abc frameshift insertion (55)
        - abc frameshift substitution (39)
        - abc NA (55)
        - abc nonframeshift deletion (55)
        - abc nonframeshift insertion (55)
        - abc nonframeshift substitution (11)
        - abc nonsynonymous SNV (55)
        - abc stopgain SNV (55)**
        - abc stoploss SNV (55)
        - abc synonymous SNV (55)

Phenotypic variables

Exome variant annotations





### Harvard IRB security levels:

Level 5 - Extremely sensitive information

**Level 4 – Very sensitive information**

Level 3 – Sensitive, or Confidential information

Level 2 - Benign information to be held confidentially

Level 1 - Non-confidential research information





**HARVARD**  
MEDICAL SCHOOL



**AWS**

VPC

**PIC-SURE**



NIH Big Data to Knowledge (BD2K)

VPC

**CCD - PIC-SURE**



NIH Big Data to Knowledge (BD2K)

VPC

**CountEverything**



NIH Big Data to Knowledge (BD2K)

VPC

**GRDR**



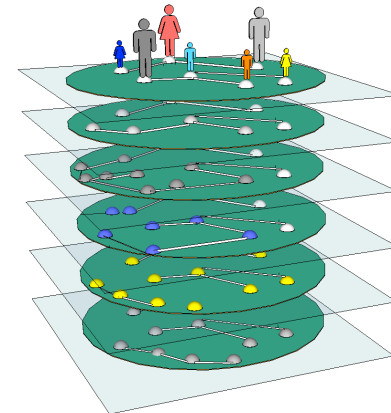
National Center for Advancing Translational Sciences

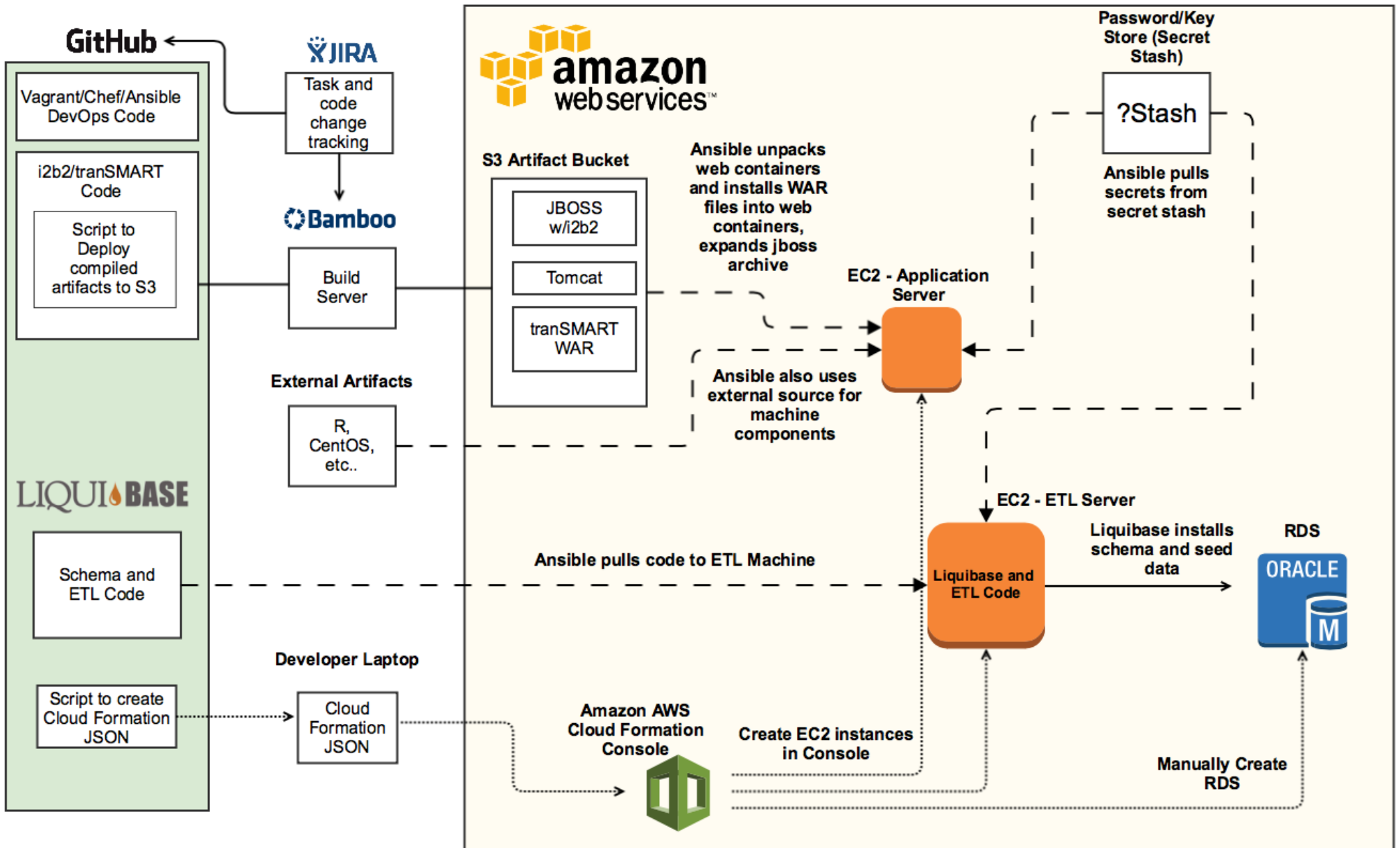
VPC



**PPRN  
PMS\_DN**

.....







# HIPPA Compliance on AWS



# Secured access control

## 1. Authentication

*Who are you?*

## 2. Authorization

*What are you allowed to do?*

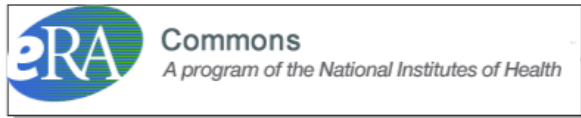
**Level 0** : *Authenticated BUT no access to data*

**Level 1** : *Aggregated data*

**Level 2** : *patient level data*

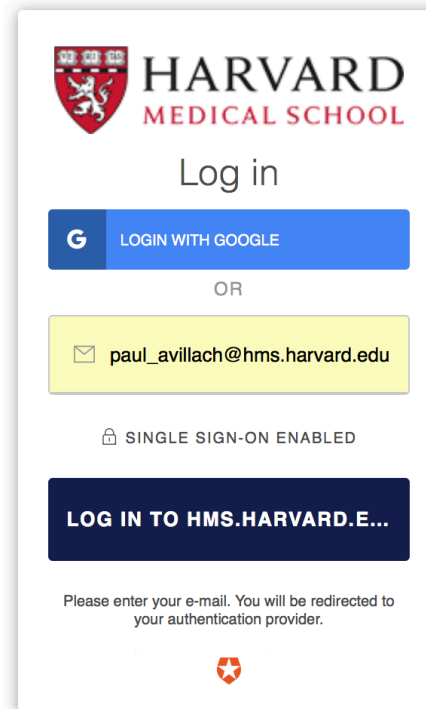
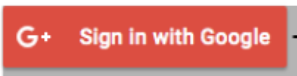
# 1. Authentication

## Enterprise Identity Providers



[.....]

## Public Identity Providers



## Service Providers

### Application User Interface



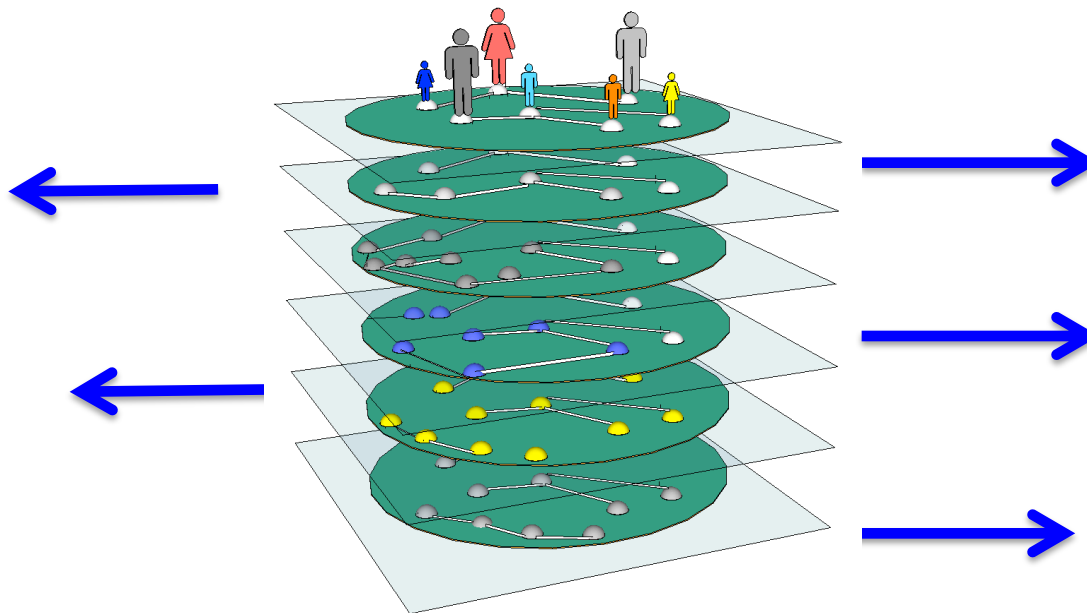
### Programmatic Interface

RESTful API



NIH Big Data to Knowledge (BD2K)

## Patient Centric Information Commons (PIC)



**RESTful API**

Play with API: <https://bd2k-picsure.hms.harvard.edu>

Source code: <https://github.com/hms-dbmi>



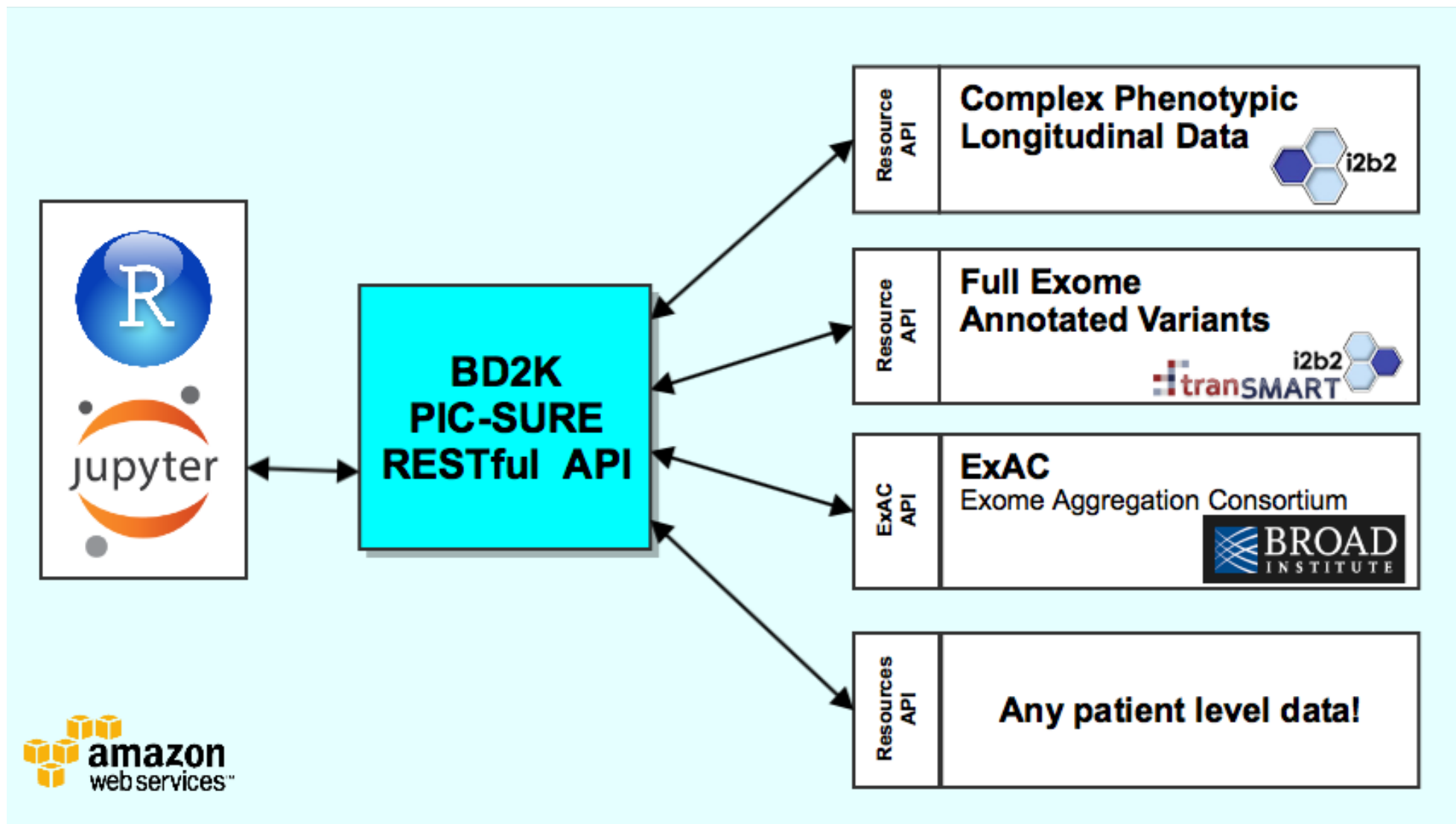
HARVARD  
MEDICAL SCHOOL

DEPARTMENT OF  
Biomedical Informatics



DEPARTMENT OF  
Biomedical Informatics





Play with PIC-SURE API: <https://bd2k-picsure.hms.harvard.edu>

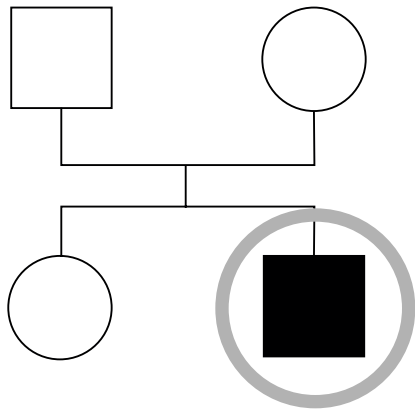
Play with ExAC API: <http://exac.hms.harvard.edu>



**HARVARD**  
MEDICAL SCHOOL

DEPARTMENT OF  
Biomedical Informatics

# Simons Simplex Collection: phenotypic data



2,800+ quad  
simplex families

27 standardized questionnaires:

- Neurodevelopmental assessment tools
- Background

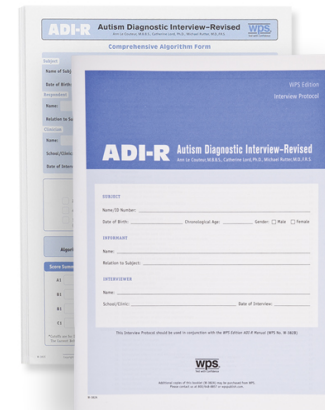
Quantitative, categorical and binary  
variables

6,000 clinical variables

**SFARI** SIMONS FOUNDATION  
AUTISM RESEARCH INITIATIVE

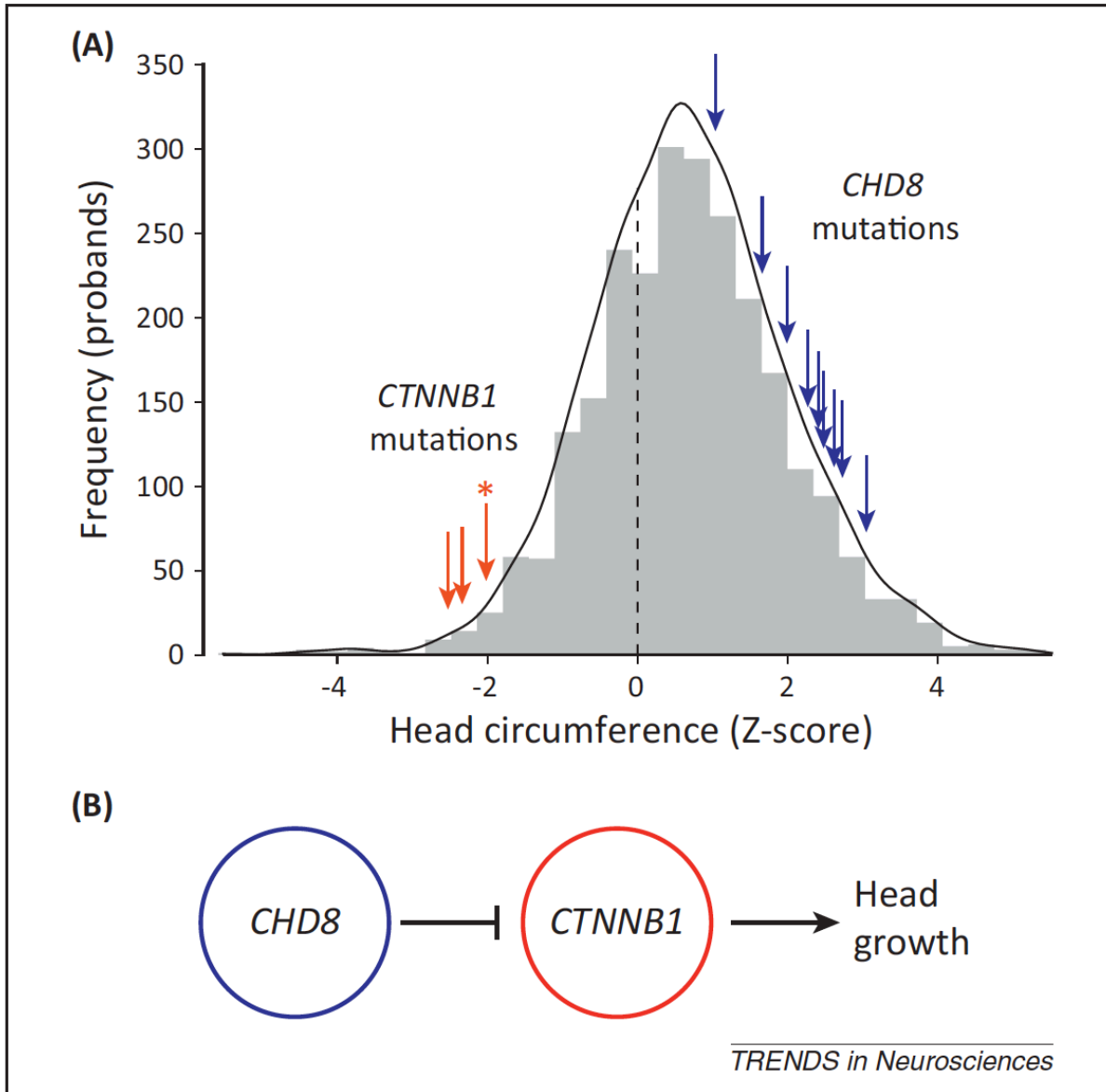
Simons Simplex collection

|                 |                |
|-----------------|----------------|
| ABC             | RBS-R          |
| ABCL            | SSC-Background |
| <b>ADI-R</b>    | SSC-Diagnosis  |
| ADOS            | SSC-HWHC       |
| BAPQ            | SSC-MHI        |
| CBCL            | SSC-Pedigree   |
| CTOPP-NR        | SSC-THF        |
| C-TRF           | SQC            |
| DAS-II          | SRS            |
| DCDQ            | SRS-ARV        |
| FHI-I           | TRF            |
| Mullen          | WISC-IV        |
| Purdue Pegboard | Vineland-II    |
|                 | WASI           |



**HARVARD**  
MEDICAL SCHOOL

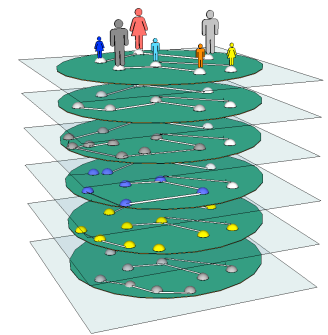
DEPARTMENT OF  
Biomedical Informatics



*CHD8* mutated =>  
Larger Head circumference

Phenotypes

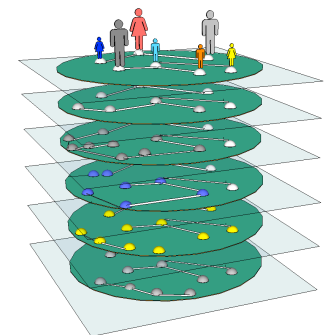
Genotypes



# Live demo

**Phenotypes**

**Genotypes**





**HARVARD**  
MEDICAL SCHOOL



**AWS**

VPC

**PIC-SURE**



NIH Big Data to Knowledge (BD2K)

VPC

**CCD - PIC-SURE**



NIH Big Data to Knowledge (BD2K)

VPC

**CountEverything**



NIH Big Data to Knowledge (BD2K)

VPC

**GRDR**



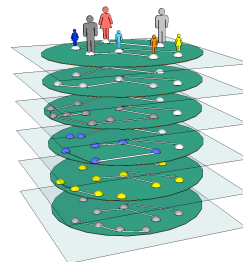
National Center for Advancing Translational Sciences

VPC



**PPRN  
PMS\_DN**

.....



**Boston Children's Hospital**



Until every child is well™

**AWS**

VPC

**Precision Link  
Biobank portal**

.....



**MASSACHUSETTS GENERAL HOSPITAL**



**NGRID**

National Institute of Mental Health

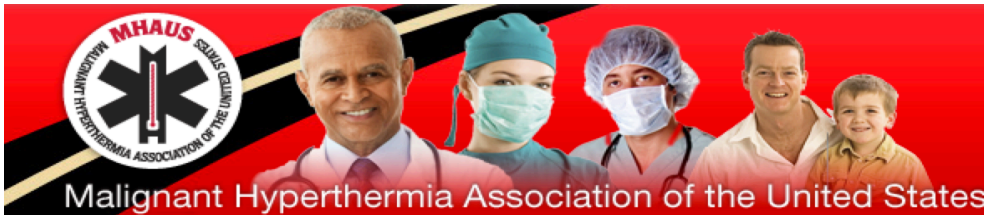






National Center  
for Advancing  
Translational Sciences

The NIH/NCATS GRDR® Program  
Global Rare Diseases Patient Registry  
Data Repository



Clinical Registry Investigating Bardet-Biedl Syndrome (CRIBBS)



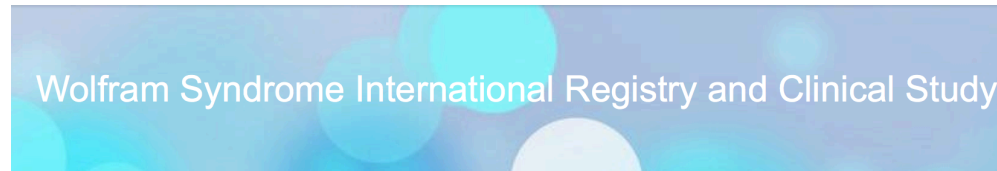
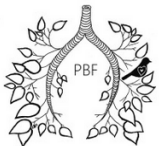
INTERNATIONAL WAGR  
SYNDROME ASSOCIATION

 Pachyonychia Congenita Project

*Fighting for a cure. Connecting & helping patients. Empowering research.*

**The Plastic Bronchitis Foundation**

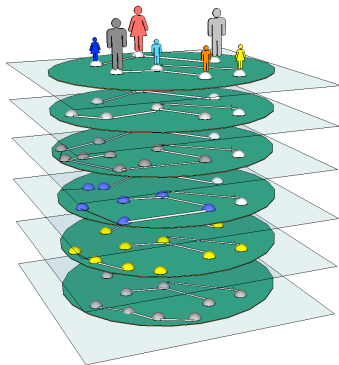
*Looking for a Cause, Working on a cure, Education and assisting*



<https://grdr.hms.harvard.edu>

| Rare disease registry                                  | Patients          | Variables                    |
|--|-------------------|------------------------------|
| Clinical Registry Investigating Bardet-Biedl Syndrome  | 180               | 708                          |
| International Pachyonychia Congenita Research Registry | 569               | 496                          |
| International Plastic Bronchitis registry              | 66                | 63                           |
| Intracranial Hypertension Registry                     | 1,349             | 91                           |
| North American Malignant Hyperthermia Registry         | 2,122             | 154                          |
| Wolfram Syndrome International Registry                | 124               | 580                          |
| Coordination of Rare Diseases at Sanford Registry      | 2,091             | 40                           |
| <i>including:</i>                                      | <i>including:</i> | <i>Additional variables:</i> |
| <i>National Ataxia Foundation</i>                      | 869               | 15                           |
| <i>International WAGR syndrome association</i>         | 52                | 379                          |
| <i>Cornelia De Lange Syndrome Registry</i>             | 67                | 485                          |

|              |              |              |
|--------------|--------------|--------------|
| <b>Total</b> | <b>6,501</b> | <b>2,132</b> |
|--------------|--------------|--------------|



<https://grdr.hms.harvard.edu>



# CDC National Health and Nutrition Examination Survey

41k patients  
2k patient level environmental variables

**I2b2/tranSMART user Interface:** <https://nhanes.hms.harvard.edu>  
**login: demo pass: demo**

**NHANES**

- demographics (41474)
- RACE (41474)
- SEX (41474)
- area (41474)
- 123 AGE (41474)**
- DMDBORN (41445)

Comparison | Advanced Workflow | Results/Analysis | Grid View

**Subset 1**      **Subset 2**

Exclude   Enable Variant Panel   X      Exclude   Enable Variant Panel   X

...AGE \ <25      ...AGE \ >=25

**NHANES**

- demographics (41474)
- examination (39274)
- laboratory (41474)
  - acrylamide (7535)
  - aging (7827)
  - allergen test (8339)
  - bacterial infection (41474)
  - biochemistry (35768)
  - blood (33718)
  - cotinine (31136)
  - diakyl (7540)
  - dioxins (5073)
    - 123 1,2,3,4,6,7,8,9-ocdd (fg/g) (4943)**
    - 123 1,2,3,4,6,7,8-hpcdd (fg/g) (4988)
    - 123 1,2,3,4,7,8-hxcdd (fg/g) (3100)
    - 123 1,2,3,6,7,8-hxcdd (fg/g) (4990)
    - 123 1,2,3,7,8,9-hxcdd (fg/g) (4977)
    - 123 1,2,3,7,8-pncdd (fg/g) (5029)
    - 123 2,3,7,8-todd (fg/g) (5002)
  - furans (5065)

Comparison | Advanced Workflow | **Results/Analysis** | Grid View

**Analysis of ...laboratory\dioxins\1,2,3,4,6,7,8,9-ocdd (fg/g) for subsets:**

Comparison of ...laboratory\dioxins\1, 2,3,4,6,7,8,9-ocdd (fg/g)

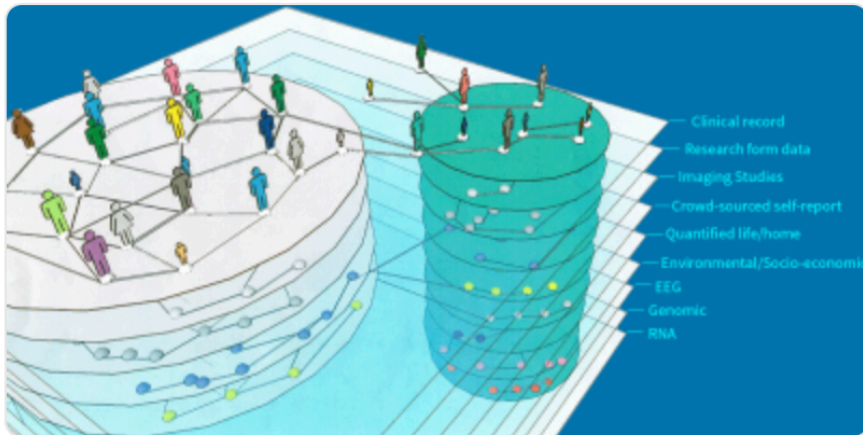
Histogram of ...laboratory\dioxins\1, 2,3,4,6,7,8,9-ocdd (fg/g)

|               | Subset 1      | Subset 2      |
|---------------|---------------|---------------|
| <b>NHANES</b> | <b>NHANES</b> | <b>NHANES</b> |
| Mean:         | 856.37        | 2,927.96      |
| Median:       | 594.89        | 2,012.15      |
| IQR:          | 436.65        | 2,469.5       |
| SD:           | 1,512.54      | 3,188.73      |
| Data Points:  | 1607          | 3336          |

|              |         |
|--------------|---------|
| t statistic: | -30.979 |
| p-value:     | 0.0000  |

The results are significant at a 95% confidence level.



## Creating an information commons for biomedical data centered on the patient.

*PIC-SURE combines genetic, environmental, imaging, behavioral, and clinical data on individual patients from multiple sources into integrated sets.*

[learn more](#) 

[www.pic-sure.org](http://www.pic-sure.org)

- **Mentors**

- Isaac Kohane (DBMI)
- Susanne Churchill (DBMI)

- **Collaborators**

- Ken Mandl (BCH)
- David Margulies (BCH)
- Sek Won Kong (BCH)
- Joel Hirschhorn (BCH)
- Florence Bourgeois (BCH)
- Jon Bickel (BCH)
- Ally Eran (BCH)
- Tim Yu (BCH)
- Savova Guergana (BCH)
- Doug McFaden (HMS)
- Shawn Murphy (Partners)
- Chirag Patel (DBMI)
- Nathan Palmer (DBMI)
- David Bernick (Broad)
- Finale Doshi-Velez (MIT)
- Peter Szolovits (MIT)
- Simon Lovestone (Oxford, UK)

- Michelle Williams (Harvard T Chan SPH)
- Tianxi Cai (Harvard T Chan SPH)
- Greg Cooper (U Pitt.)
- Lucila Ohno-Machado (UCSD)
- David Haussler (UCSD)
- Ida Sim (UCSF)
- Anthony J Brookes (U. Leicester, UK)
- Johan van Der Lei (Erasmus MC, NL)
- Anita Burgun (Paris, FR)
- Riccardo Bellazzi (U. Pavia, IT)
- Mehan O’Boyle (PMS foundation)
- Geraldine Bliss (PMS foundation)
- Thomas Bourgeron (Institut Pasteur, FR)

- **DBMI Administrative support**

- Katherine Flannery
- Sunny Alvear
- Jennifer Grandfield
- + all DBMI staff



<http://avillach-lab.hms.harvard.edu>

## Manager of Data Infrastructure

- Michael McDuffie, MSc

## Developers

- Jeremy Easton-Marks
- Gabor Korodi
- Thomas DeSain
- Sean Finan
- Ranjay Kumar
- Alexander Nikitin
- Ken Hoflen

## Graduate students

- Ombeline Dorval, MD
- Maxime Wack, MD
- Claire Hassen-Kodja, MD, MSc
- Emmanuelle Sylvestre, MD, MSc
- Yuri Ahuja, MD, PhD HMS candidate

## Funding:



THE NANCY LURIE MARKS FAMILY FOUNDATION



## Project Manager

Cassandra Perry, MS, CGC

## Research Associates - Postdocs

- Antoine Tran, MD, MSc
- Laurie Tran, MD, MSc
- Cartik Saravana, PhD
- Li Ly, PhD
- Joany Zachariasse, MD, MSc
- Antoine Neuraz, MD, MSc

## Previous members

- Samuel Finlayson, MD, PhD HMS candidate
- Ephi Sachs, MD
- Pei Chen
- Sushma Hanawal

## We are hiring now:

- Senior Software Developer \*2
- Postdocs \*3