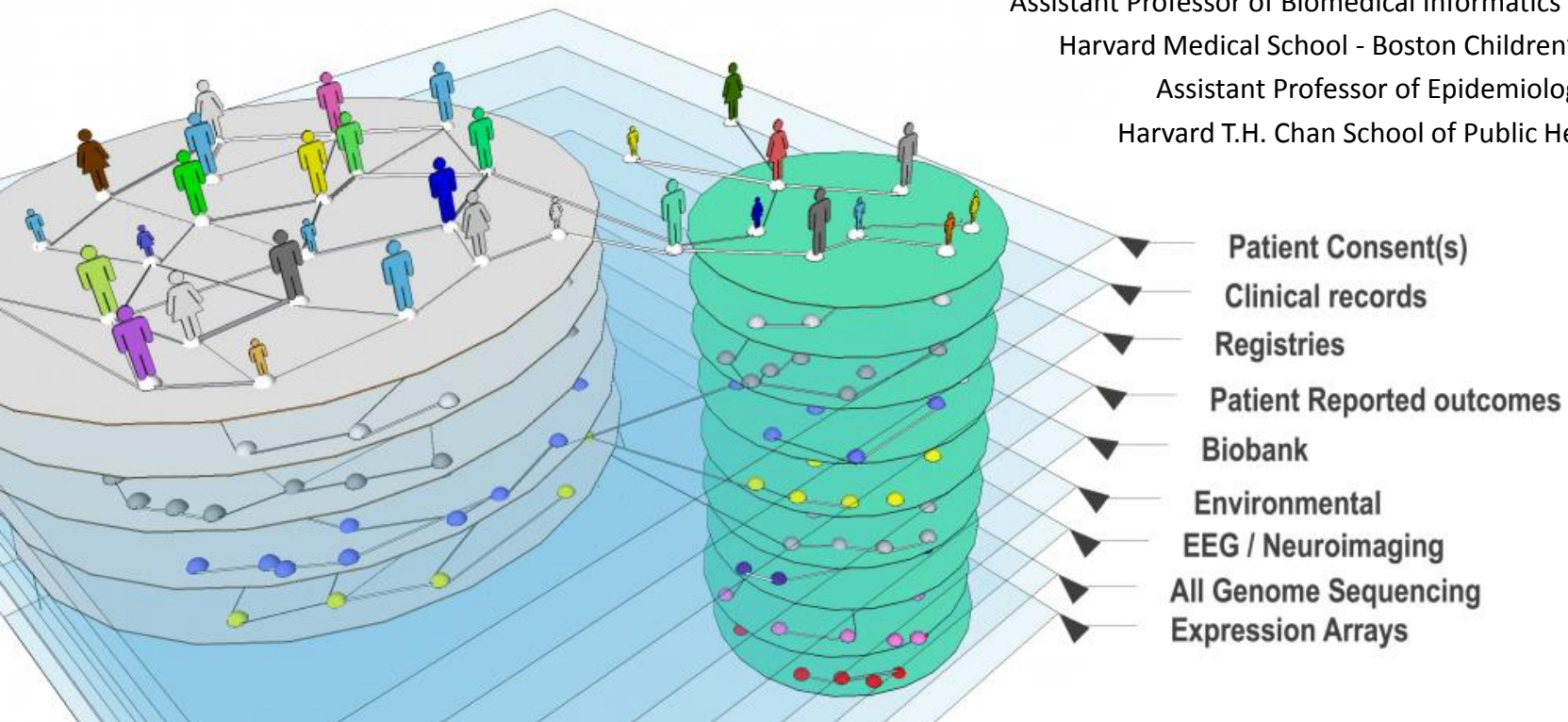


Creating scalable, secured, clinical and genomics platforms across clouds integrating i2b2 and tranSMART Platforms

Paul Avillach, MD, PhD

Assistant Professor of Biomedical Informatics and Pediatrics
Harvard Medical School - Boston Children's Hospital
Assistant Professor of Epidemiology
Harvard T.H. Chan School of Public Health



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

Genotypes

Biosamples

Baseline Genotypes

Burden, LOF, etc

Derived Genotypes

Phenotypes

Registries

EHR

Baseline Phenotypes

Custom algos, etc

Derived Phenotypes

Geno/Pheno integration

Insights & many iterations





Advance statistical tools
Biobank explorer
Variant explorer



Advance cohort selection



Federated Advance cohort selection



SMART®
Patient level data lookup
Interoperable tools



Driving Biology Projects

- Overview
- Current DBPs
 - Autoimmune/CV Diseases
 - Diabetes/CV Diseases
- Past DBPs
 - Airways Diseases
 - Hypertension
 - Type 2 Diabetes Mellitus
 - Huntington's Disease
 - Major Depressive Disorder
 - Rheumatoid Arthritis
 - Obesity

Overview

The disease-based driving biology projects (DBP's) serve as the testbed of i2b2. It is where we field, test and debug the methodologies and tools developed in Cores 1 and 3. For this reason we recognize that it is incumbent on us to use the DBP's to maximize interactions, oversight and corrections in the directions of Core 1 and 3. Consequently for each DBP, we have ensured the following:

- On each DBP there is an assigned computational/bioinformatics co-investigator to ensure that there is a very close collaboration between the clinical-genomic investigator and the methodologies and tools developed in Cores 1 and 3. Also this computational co-investigator provides a tighter loop for feedback and advice than would be ordinarily available.
- Frequent meetings between the clinical investigators of each DBP and the methodologists from Core 1 and 3 (see [Zak Kohane's blog](#) for details).



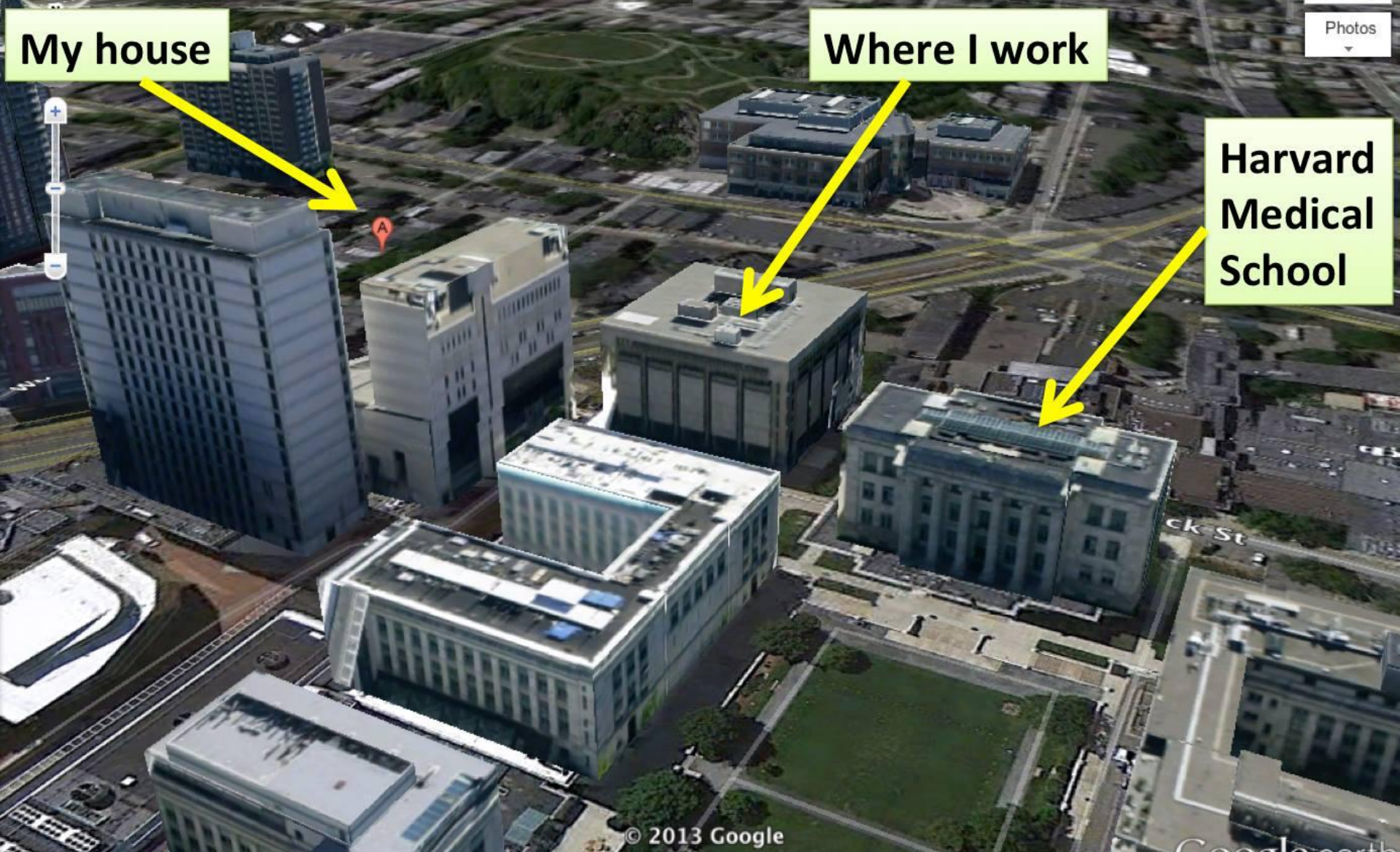


I2b2 database still there
But NOT i2b2 jboss app



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics



My house

Where I work

Harvard
Medical
School

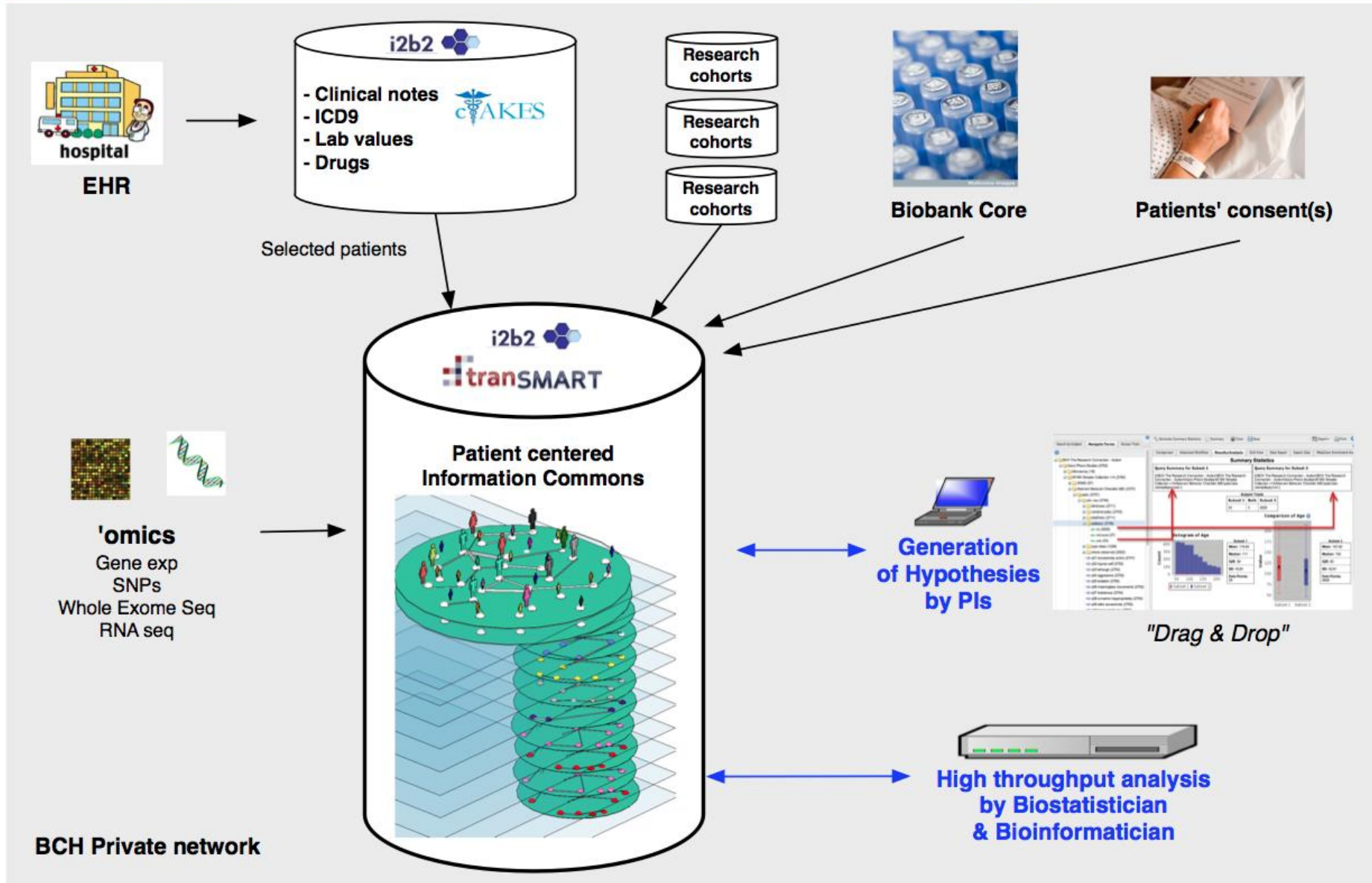
Photos

© 2013 Google



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics



→ integration
↔ analysis

Students / Postdocs



Cartik



Li



Romain



Qiu-Yue



Ombeline



Maxime



Antoine



Laurie



Laura



Haishuai



Mahdi



Niloofar



Alba



Joany



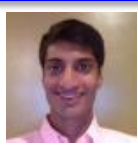
Carlos



Romina



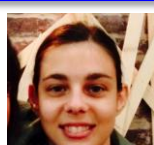
Emmanuelle



Yuri



Samuel

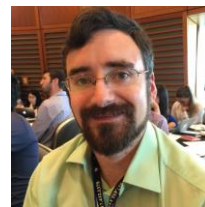


Claire



Antoine

Staff / Software developers



Jason



Ranjay



Gabor



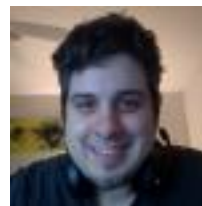
Thomas



Jaspreet



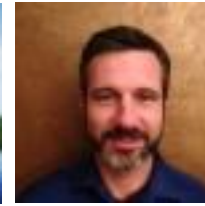
Alex



Andrew



Andre



Sean



Anoush



Cassandra



Libby



Sophia



Alyssa

Alumni



Michael



Sushma



Pei

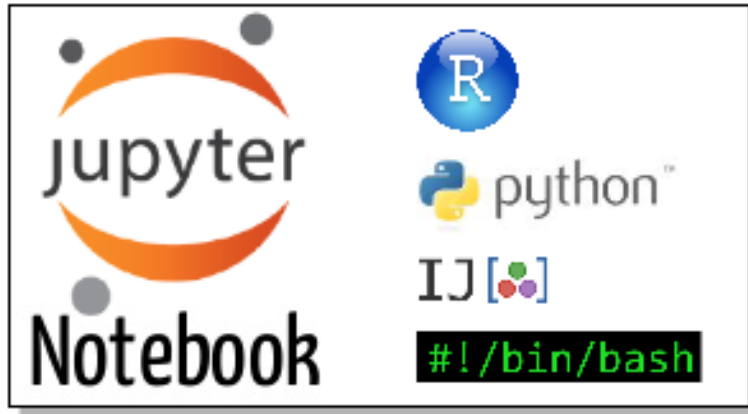


Ephi



Jeremy

Open Source Research Infrastructure





Harvard IRB security levels:

Level 5 - Extremely sensitive information

Level 4 – Very sensitive information

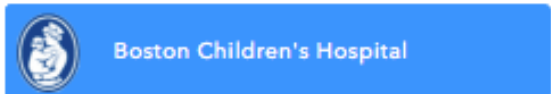
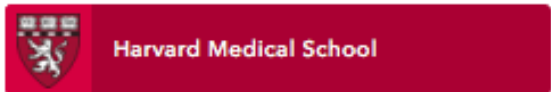
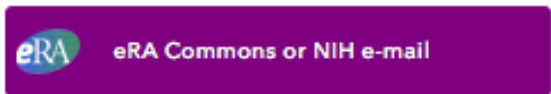
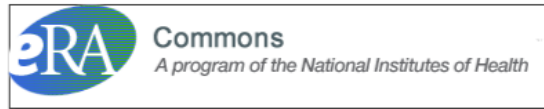
Level 3 – Sensitive, or Confidential information

Level 2 - Benign information to be held confidentially

Level 1 - Non-confidential research information

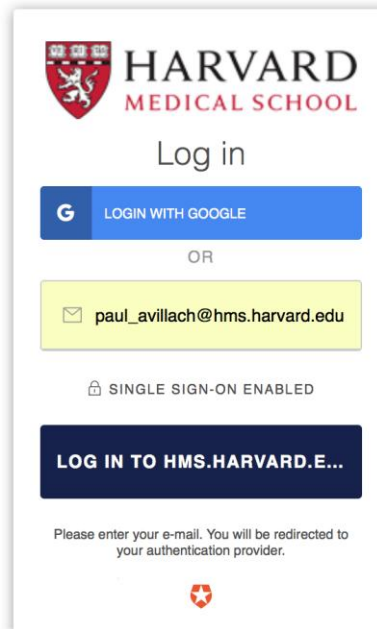
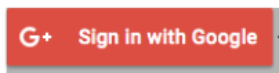
1. Authentication

Enterprise Identity Providers



[....]

Public Identity Providers



Service Providers

Application User Interface



Programmatic Interface



Secured access control

1. Authentication

Who are you?



2. Authorization

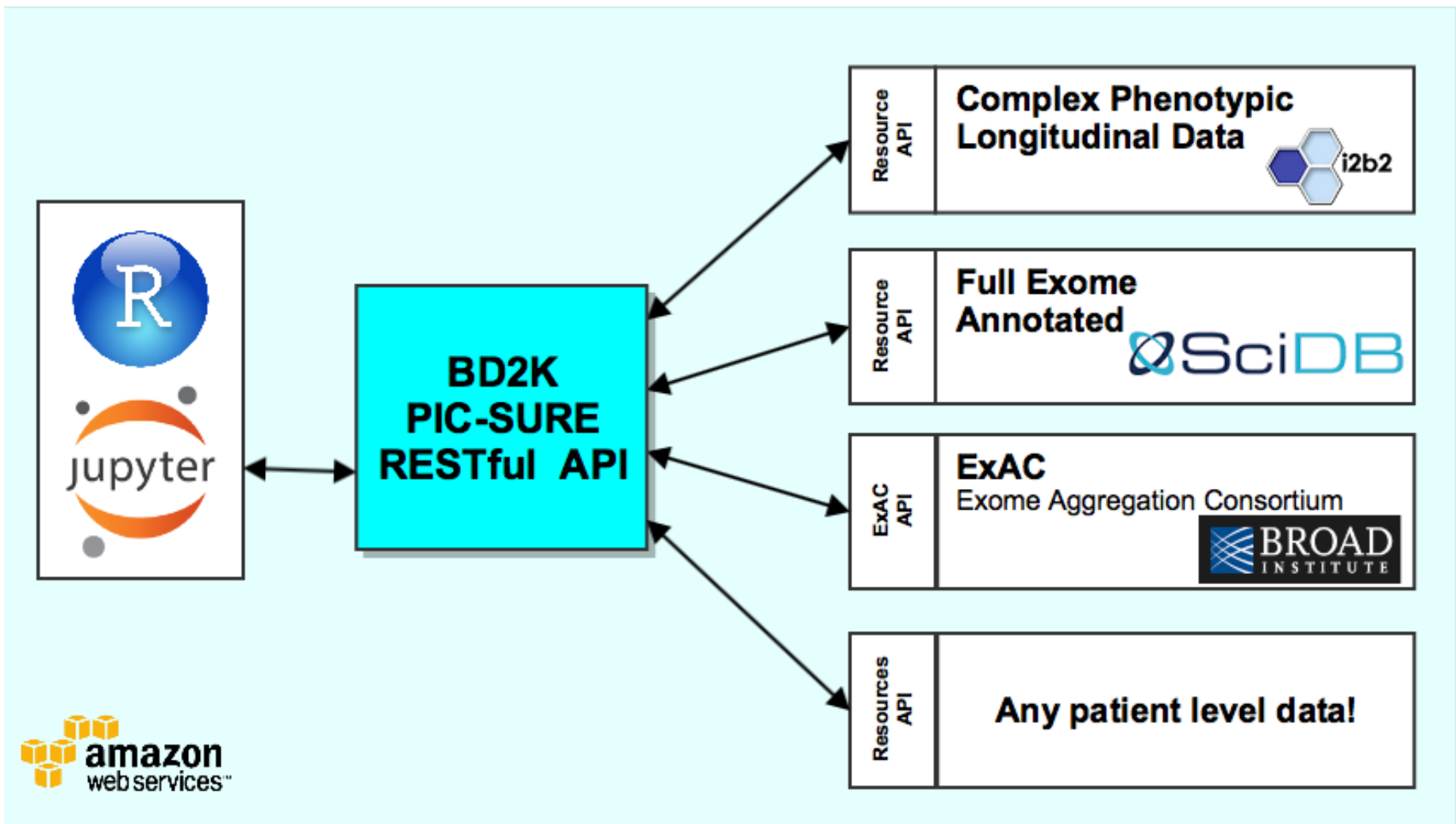
What are you allowed to do?

Level 0 : *Authenticated BUT no access to data*

Level 1 : *Aggregated data*

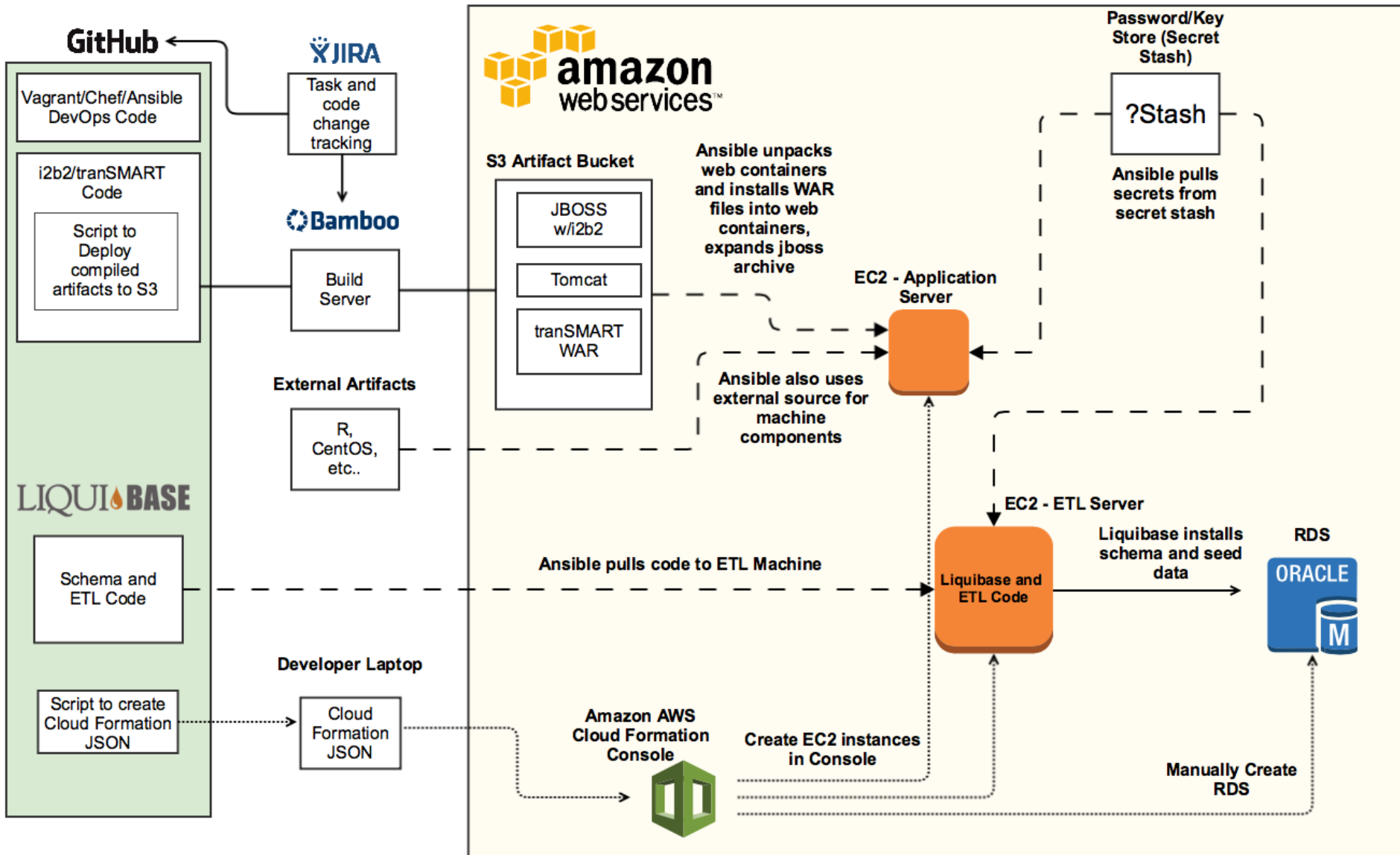
Level 2 : *patient level data*



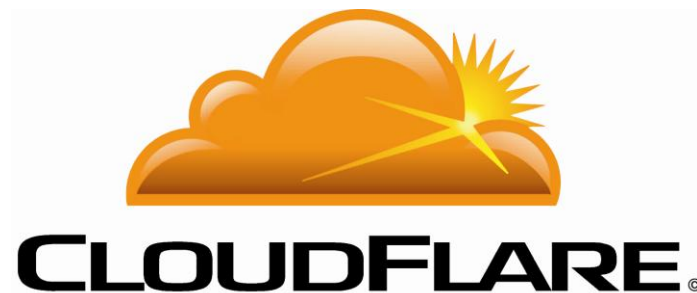


Play with PIC-SURE API: <https://bd2k-picsure.hms.harvard.edu>

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



HIPPA Compliance on AWS





Amazon mystery solved: A typo took down a big chunk of the Internet



1527

Elizabeth Weise, USATODAY Published 3:56 p.m. ET March 2, 2017 | Updated 5:57 p.m. ET March 2, 2017



(Photo: Amazon)

f 1527
CONNECT

TWEET

in 637
LINKEDIN

33
COMMENT

EMAIL

MORE

SAN FRANCISCO — The major outage that hit tens of thousands of websites using Amazon's AWS cloud computing service on Tuesday ends up having been the result of a simple typo — just one incorrectly-entered command.

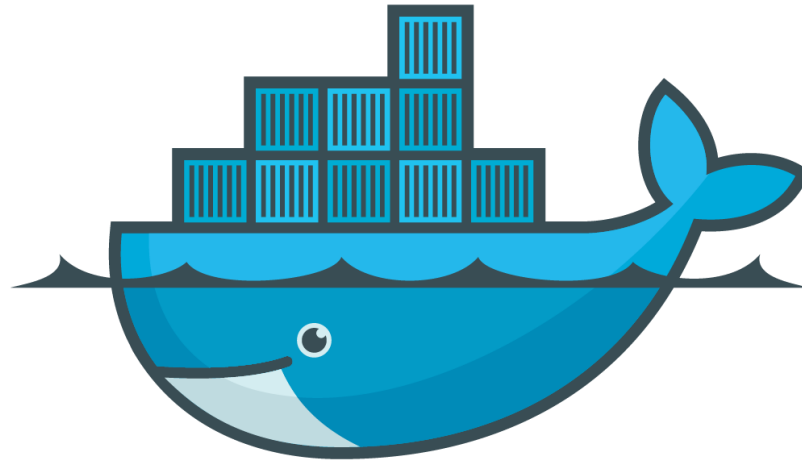
The four-hour outage at Amazon Web Services' S3 system, a giant provider of backend services for close to 150,000 websites, caused disruptions, slowdowns and failure-to-load errors across the United States.





HIPAA

Health Insurance Portability
and Accountability Act



docker



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics



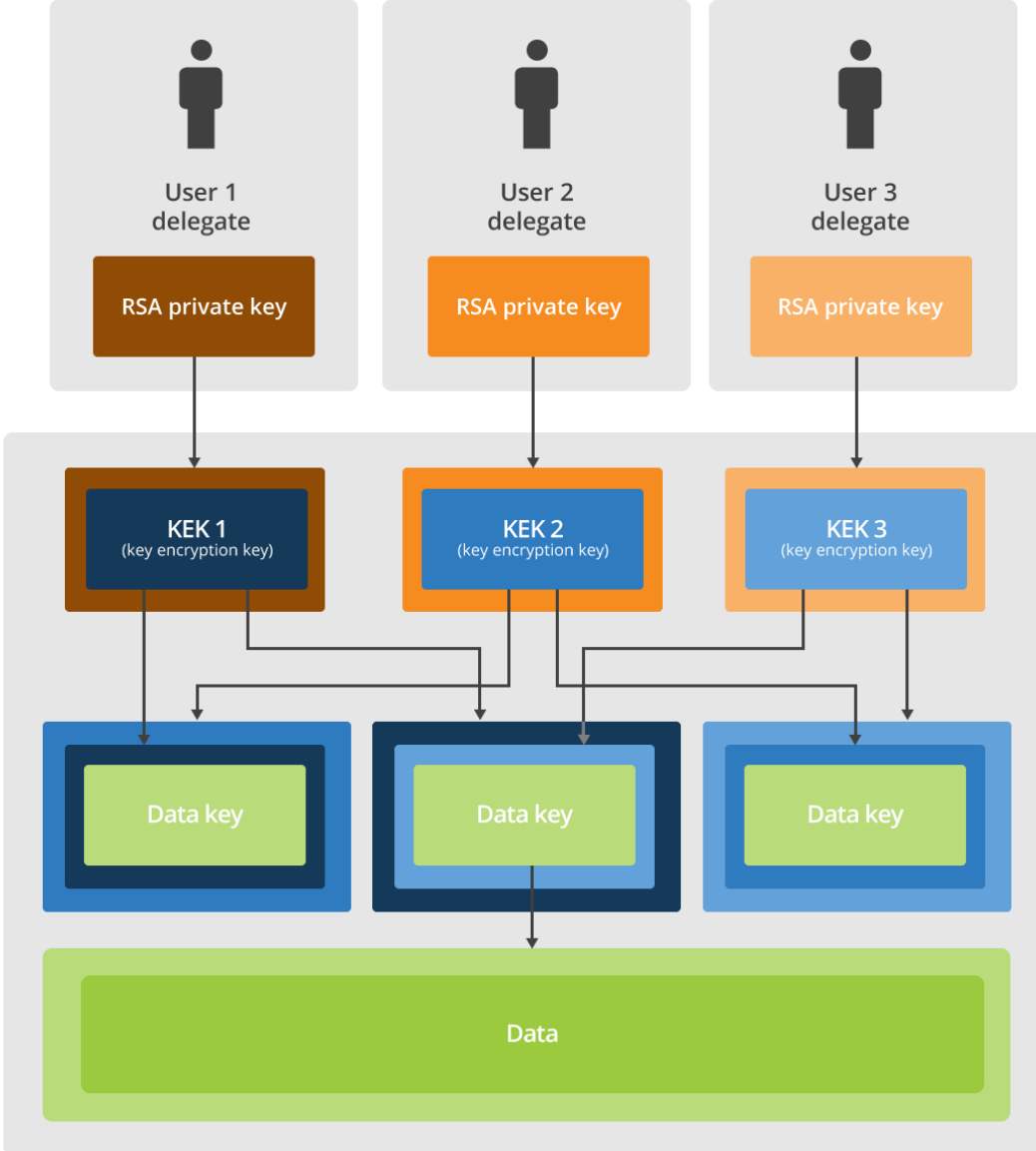
Commons Credits Portal



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

RedOctober

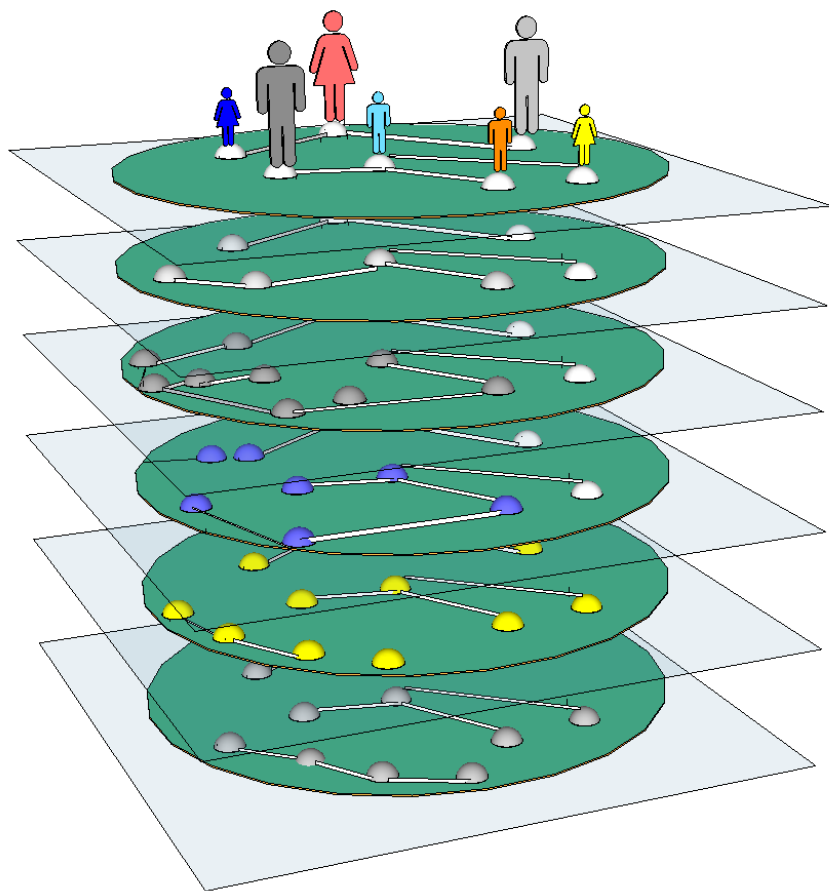


CLOUDFLARE



UDN Patient centric information commons

i2b2/tranSMART/gNOME

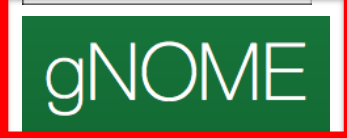
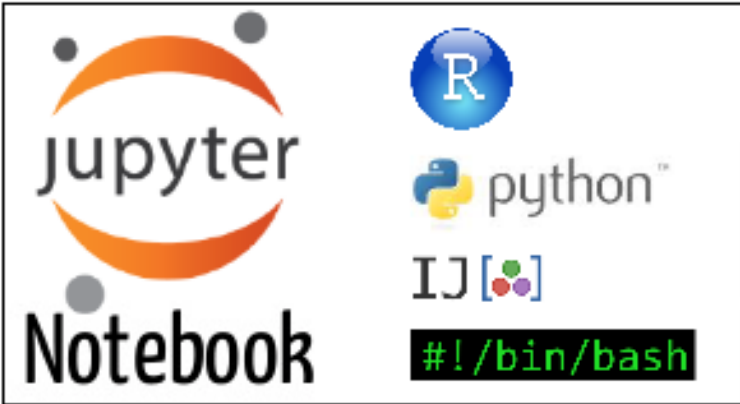


ALL UDN data available

- Demographics
- Online registration Form
- PhenoTips
 - HPO clinical terms
 - Ethnicity
 - Diagnosis OMIM
- Knowledge extracted from referral letters (via NLP)
- Type of sequencing
- Candidate genes
- Candidate variants
- **Full Annotated VCF**
 - **524 Whole genomes**
 - **616 Whole exomes**
 - **343 probands**
 - **259 trios**
- +



Modular Open Source Research platform



1



2

IP[y]: Notebook plotting_code (autosaved)

File Edit View Insert Cell Kernel Help this also runs code cell

Length and GC% script An instructions section

Instructions: Given a fasta DNA sequence file this prints the sequence length and GC content for each sequence.

```
In [11]: # this is a code cell
from Bio import SeqIO
from Bio.SeqUtils import GC
# specify the location/name of the fasta file below
for rec in SeqIO.parse("data/testseqs.fasta", "fasta"):
    length = len(rec)
    gc_value = (GC(rec.seq))
    print("Length: "+str(length)+" GC: "+str(gc_value))
# results print below
```

this specifies input file

This is a code cell

place cursor in cell then shift+enter to run code

output from the code

```
Length: 660 GC: 37.4242424242
Length: 661 GC: 37.2163388805
Length: 661 GC: 32.3751891074
Length: 661 GC: 33.2829046899
Length: 710 GC: 40.4225352113
Length: 708 GC: 40.5367231638
```

UDN338723

- 17 year old male with skeletal features suggestive of a **hereditary connective tissue disease**, also with proximal renal tubular acidosis, medullary nephrocalcinosis, hypophosphatemic rickets, polycythemia and chronic kidney disease, stage 3
- Marfan syndrome Confirmed
- FBN1 c.871G>T/p.E291X



Search by Subject

Navigate Terms

- + 000_UDN_ID
- + 00_Demographics
- + 01_Primary symptom category reported by patient or caregiver
- + 02_Type of sequencing
- + 03_UDN Clinical Site
- 04_Clinical symptoms and physical findings (in HPO, from P)
 - Phenotypic abnormality (350)
 - + Abnormal delivery (129)
 - + Abnormal eye (1)
 - + Abnormal growth (1)
 - + Abnormal thrombosis (1)
 - + Abnormality of blood and blood-forming tissues (251)
 - Abnormality of connective tissue (225)
 - + Abnormality of Sharpey fibers (34)
 - + Abnormality of adipose tissue (99)
 - + Abnormality of connective tissue (1)

Generate Summary Statistics | Generate WES Statistics | Exp

Comparison

Advanced Workflow

Results/Analysis

Grid View

Subset 1

Exclude

Enable Variant Panel

X

...\Abnormality of connective tissue\

AND

Exclude

Enable Variant Panel

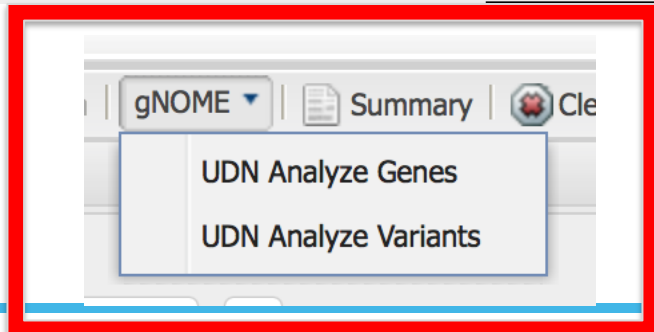
X

AND

Exclude

Enable Variant Panel

X



Analyze Variants

UDN_subset_de33b136c2a0c53 vs UDN_comple_subset_de33b136

Allele Frequency

Ancestry All

Allele Frequency ▶ ≤ 1% (rare)

Phenotype (powered by FindZebra)

Primary phenotype connective AND

Secondary phenotype AND

Tertiary phenotype

Use gene ranked in top 5 %

For each phenotype, please enter as many terms (separated by comma) describing the phenotype as possible. A single word (e.g., headache) or multiple words separated by space (e.g., neck pain) can be a single term. For each phenotype, only genes highly ranked by relevancy to the phenotype will be considered. Only the genes associated with all of the above phenotypes will be shown.

Functional Impact

Gene model RefSeq

Gene impact nonsynonymous

SIFT any

PolyPhen2 any

Condel any

Statistical Test Parameters

(not applied to per genome enrichment)

p-value threshold ≤ 1

N_{group A} genomes ≥ 1N_{group B} genomes ≤ 50

271 candidate variants (in 86 transcripts and 76 genes).

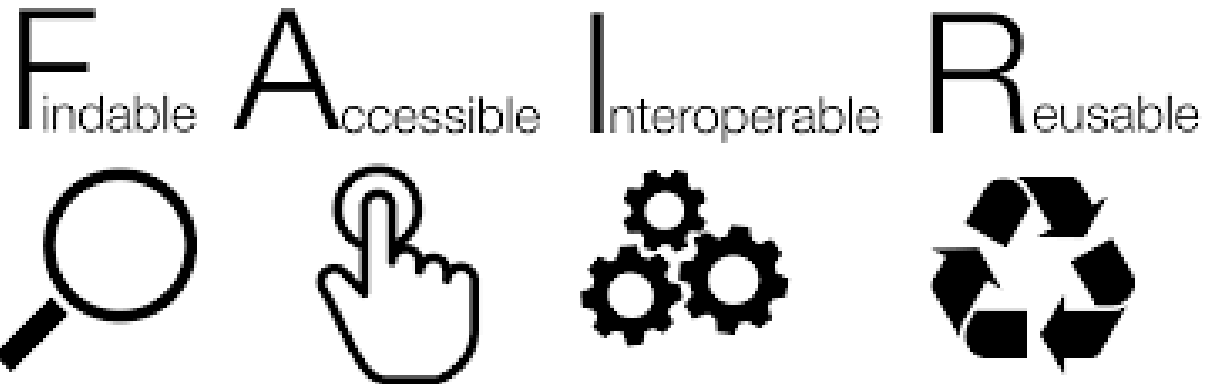
Transcript ID	Gene symbol (ID)	CDS size (bps)	Chromosome	Variant start	Variant end	Reference sequence		Wilcoxon p-value	number of case genomes	number of control genomes
NM_001127707	SERPINA1	1,257	chr14	94,845,944	94,845,944	C	T	0	34	46
NM_001127707	SERPINA1	1,257	chr14	94,845,914	94,845,914	T	A	0	32	40
NM_001127707	SERPINA1	1,257	chr14	94,845,886	94,845,886	-	G	0	23	30
NM_001127707	SERPINA1	1,257	chr14	94,844,968	94,844,968	T	C	0	20	30
NM_001127707	SERPINA1	1,257	chr14	94,845,885	94,845,885	-	TA	0.005	13	23
NM_032470	TNXB	2,022	chr6	32,010,523	32,010,523	A	G	0.262	12	40
NM_032470	TNXB	2,022	chr6	32,010,732	32,010,732	T	G	0.043	12	28
NM_000492	CFTR	4,443	chr7	117,188,877	117,188,877	G	T	0.004	11	17
NM_032470	TNXB	2,022	chr6	32,010,572	32,010,572	G	T	0.657	10	46
NM_000257	MYH7	5,808	chr14	23,889,445	23,889,445	-	G	0.779	9	47

[...]

NM_024312	GNPTAB	3,771	chr12	102,190,521	102,190,521	C	T	0.279	1	2
NM_000257	MYH7	5,808	chr14	23,884,281	23,884,281	C	T	0.023	1	0
NM_000257	MYH7	5,808	chr14	23,887,578	23,887,578	C	T	0.023	1	0
NM_000257	MYH7	5,808	chr14	23,892,910	23,892,910	A	G	0.023	1	0
NM_000257	MYH7	5,808	chr14	23,894,554	23,894,554	C	T	0.023	1	0
NM_001127707	SERPINA1	1,257	chr14	94,847,386	94,847,386	G	A	0.399	1	3
NM_000138	FBN1	8,616	chr15	48,704,813	48,704,813	C	T	0.023	1	0
NM_000138	FBN1	8,616	chr15	48,714,160	48,714,160	G	A	0.023	1	0
NM_000138	FBN1	8,616	chr15	48,782,072	48,782,072	T	C	0.143	1	1
NM_001009944	PKD1	12,912	chr16	2,140,294	2,140,294	C	T	0.497	1	4
NM_001009944	PKD1	12,912	chr16	2,153,345	2,153,345	C	T	0.399	1	3
NM_001009944	PKD1	12,912	chr16	2,157,984	2,157,984	G	A	0.143	1	1
NM_001009944	PKD1	12,912	chr16	2,158,419	2,158,419	G	A	0.279	1	2
NM_001009944	PKD1	12,912	chr16	2,158,869	2,158,869	G	A	0.023	1	0
NM_001009944	PKD1	12,912	chr16	2,159,557	2,159,557	C	T	0.497	1	4
NM_001009944	PKD1	12,912	chr16	2,159,757	2,159,757	C	A	0.005	1	1



OPEN DATA



CDC National Health and Nutrition Examination Survey



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



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

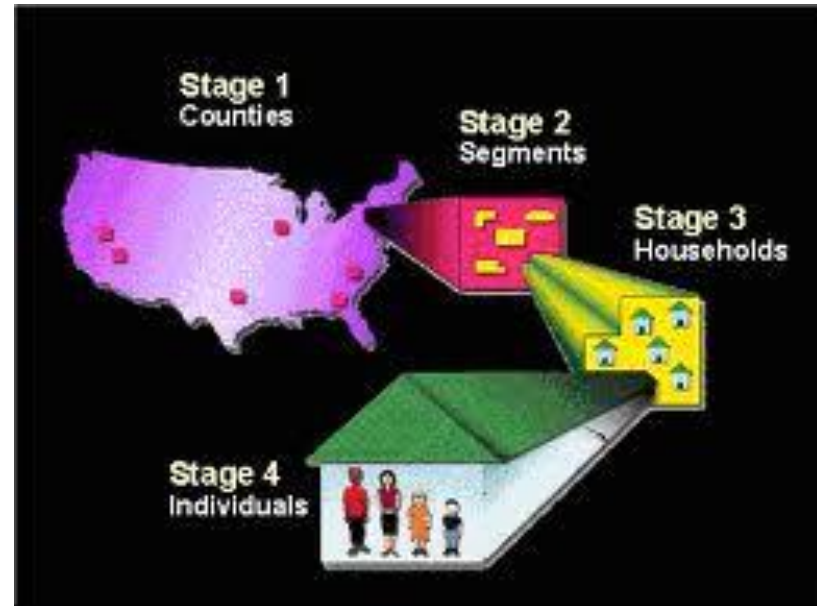
Gold standard resource for *phenotypes and exposure* information: National Health and Nutrition Examination Survey Examination Survey



since the 1960s
now biannual: 1999 onwards
5000 participants per year

1999-2006:
N=41,474
1,181 total measurements!

>250 exposures (serum + urine)
>85 quantitative clinical traits
(e.g., serum glucose, lipids,
body mass index)
nutritional and health status interviews





CDC National Health and Nutrition Examination Survey

41k patients
2k patient level environmental variables

I2b2/tranSMART user Interface:

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

10 min video <https://vimeo.com/182576739>

NHANES

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

Comparison Advanced Workflow Results/Analysis Grid View

Subset 1

Exclude Enable Variant Panel X

...AGE <25

Subset 2

Exclude Enable Variant Panel X

...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-tcdd (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)

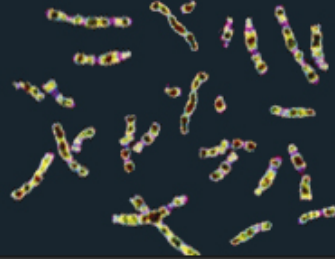
	Subset 1	Subset 2
NHANES		
Mean:	856.37	2,927.96
Median:	594.89	2,012.15
IQR:	436.65	2,469.5
SD:	1,512.54	3,188.73
Data Points:	1607	3336

t statistic:	-30.979
p-value:	0.0000

The results are significant at a 95% confidence level.

1000 Genomes

A Deep Catalog of Human Genetic Variation



- 56 full exomes with Phenotypic data from 1000 Genome project:
- *no registration at all*

<https://demo-ngs.hms.harvard.edu>



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

Exome sequence data processing



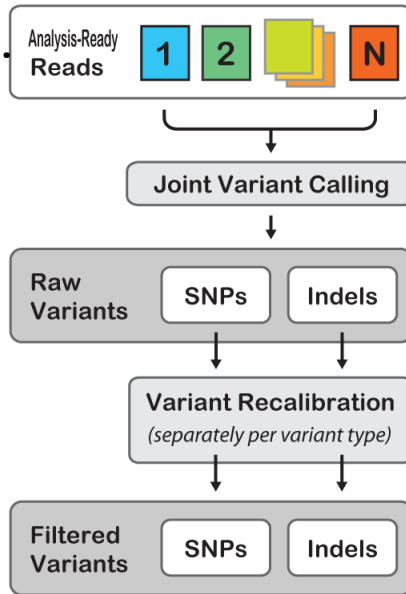
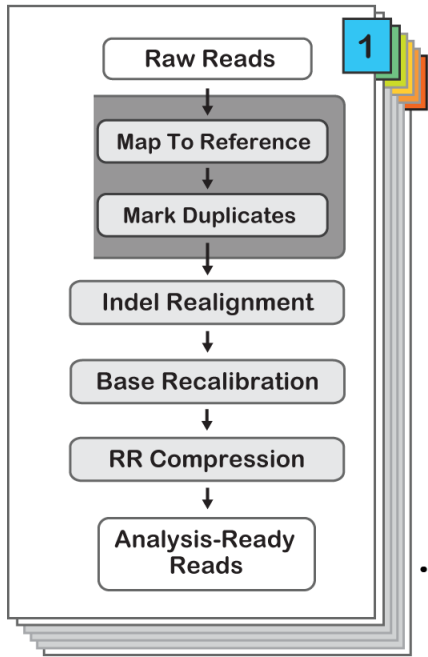
Variant calling



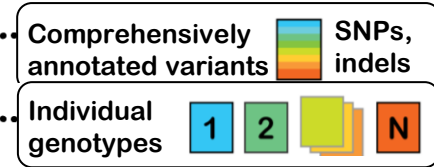
Variant annotation



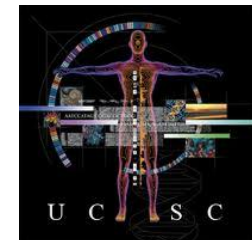
I2b2/tranSMART input



- Physical location
e.g. Chr:start-end
Cytoband
...
- Gene
e.g. Gene name
Variant function
...
- Gene set
e.g. Pathway
Molecular processes
...
- 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



<https://demo-ngs.hms.harvard.edu>



ANNOVAR

Generate Summary Statistics | Generate WES Statistics

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

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\
Genomic variables

AND | Exclude | Disable Variant Panel

... \nonsynonymous SNV\
... \stoppain SNV

AND | Exclude | Enable Variant Panel

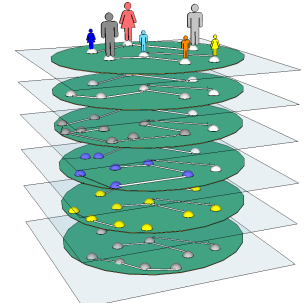
abc stopgain SNV (55)

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





National Center
for Advancing
Translational Sciences

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



Malignant Hyperthermia Association of the United States

Marshfield Clinic
Research Foundation

Clinical Registry Investigating Bardet-Biedl Syndrome (CRIBBS)



INTRACRANIAL HYPERTENSION
RESEARCH FOUNDATION

CONTACT

CoRDS Registry
Coordination of Rare Diseases
at Sanford



National Ataxia
Foundation



INTERNATIONAL WAGR
SYNDROME ASSOCIATION

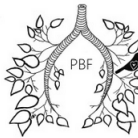


Pachyonychia Congenita Project

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

The Plastic Bronchitis Foundation

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



Wolfram Syndrome International Registry and Clinical Study



CdLS Foundation
Cornelia de Lange Syndrome Foundation, Inc.

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



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

Rare disease registry	Patient count	Variables	Version
Clinical Registry Investigating Bardet-Biedl Syndrome	180	708	V4
International Pachyonychia Congenita Research Registry	569	496	V3
International Plastic Bronchitis registry	66	63	V2
Intracranial Hypertension Registry	963	77	V4
North American Malignant Hyperthermia Registry	2,107	162	V4
Wolfram Syndrome International Registry	124	580	V2
Coordination of Rare Diseases at Sanford Registry	1,218	46	V2
<i>including:</i>	<i>including:</i>	<i>Additional data:</i>	
<i>National Ataxia Foundation</i>	487	16	V1
<i>International WAGR syndrome association</i>	50	380	V1
<i>Cornelia De Lange Syndrome Registry</i>	61	486	V1
<i>Total</i>	<i>5,227</i>	<i>3,014</i>	

Rare Cancer Registry (soon)

800

1,800

Centronuclear and myotubular myopathy (soon)

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



HARVARD
MEDICAL SCHOOL

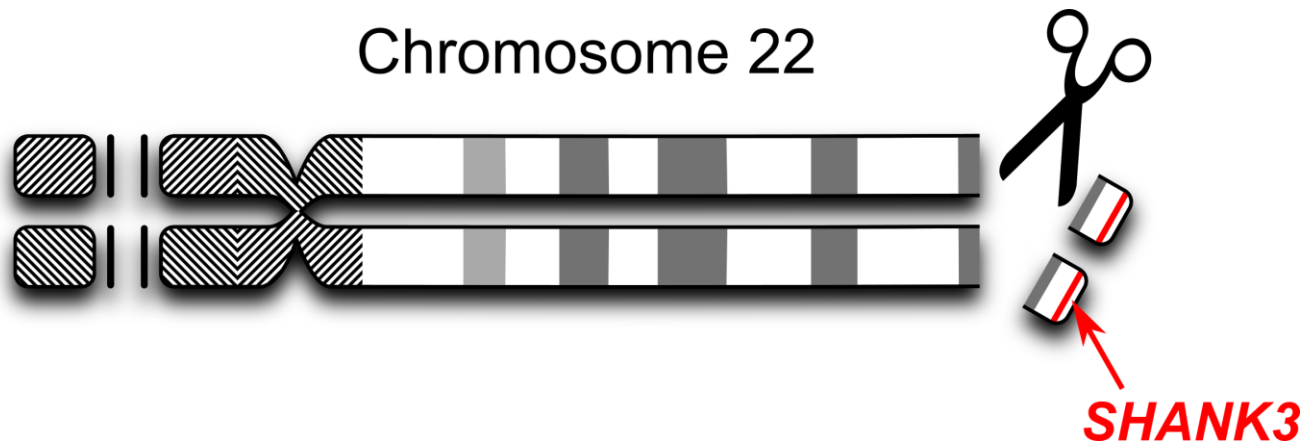
DEPARTMENT OF
Biomedical Informatics

- + 📁 Clinical Registry Investigating Bardet-Biedl Syndrome
- + 📁 Coordination of Rare Diseases at Sanford Registry
- + 📁 International Pachyonychia Congenita Research Registry
- + 📁 International Plastic Bronchitis Registry
- + 📁 Intracranial Hypertension Registry
- + 📁 North American Malignant Hyperthermia Registry
- 📁 Wolfram Syndrome International Registry
 - + 📁 1 Socio-demographic data (124) 📄
 - + 📁 2 Clinical data (124) 📄
 - + 📁 3 Administrative data (124) 📄
 - 📁 __International Ontologies (124) 📄
 - + 📁 Clinical Terms Version 3 (CTV3) (Read Codes) (124) 📄
 - + 📁 DSM-IV (62) 📄
 - + 📁 Gene Ontology (13) 📄
 - + 📁 Human Phenotype Ontology (124) 📄
 - + 📁 International Classification for Nursing Practice (68) 📄
 - + 📁 International Classification of Diseases, 10th Edition, Clinical Modification (124) 📄
 - + 📁 International Classification of Diseases, Ninth Revision, Clinical Modification (124) 📄
 - + 📁 LOINC (124) 📄
 - + 📁 MEDCIN (124) 📄
 - + 📁 Medical Dictionary for Regulatory Activities Terminology (MedDRA) (124) 📄
 - + 📁 Medical Subject Headings (124) 📄
 - + 📁 Metathesaurus Source Terminology Names (124) 📄
 - + 📁 NCI Thesaurus (124) 📄
 - + 📁 National Drug File - Reference Terminology (124) 📄
 - + 📁 Online Mendelian Inheritance in Man (124) 📄
 - + 📁 US Edition of SNOMED CT (124) 📄
- + 📁 __GRDR Common Data Elements

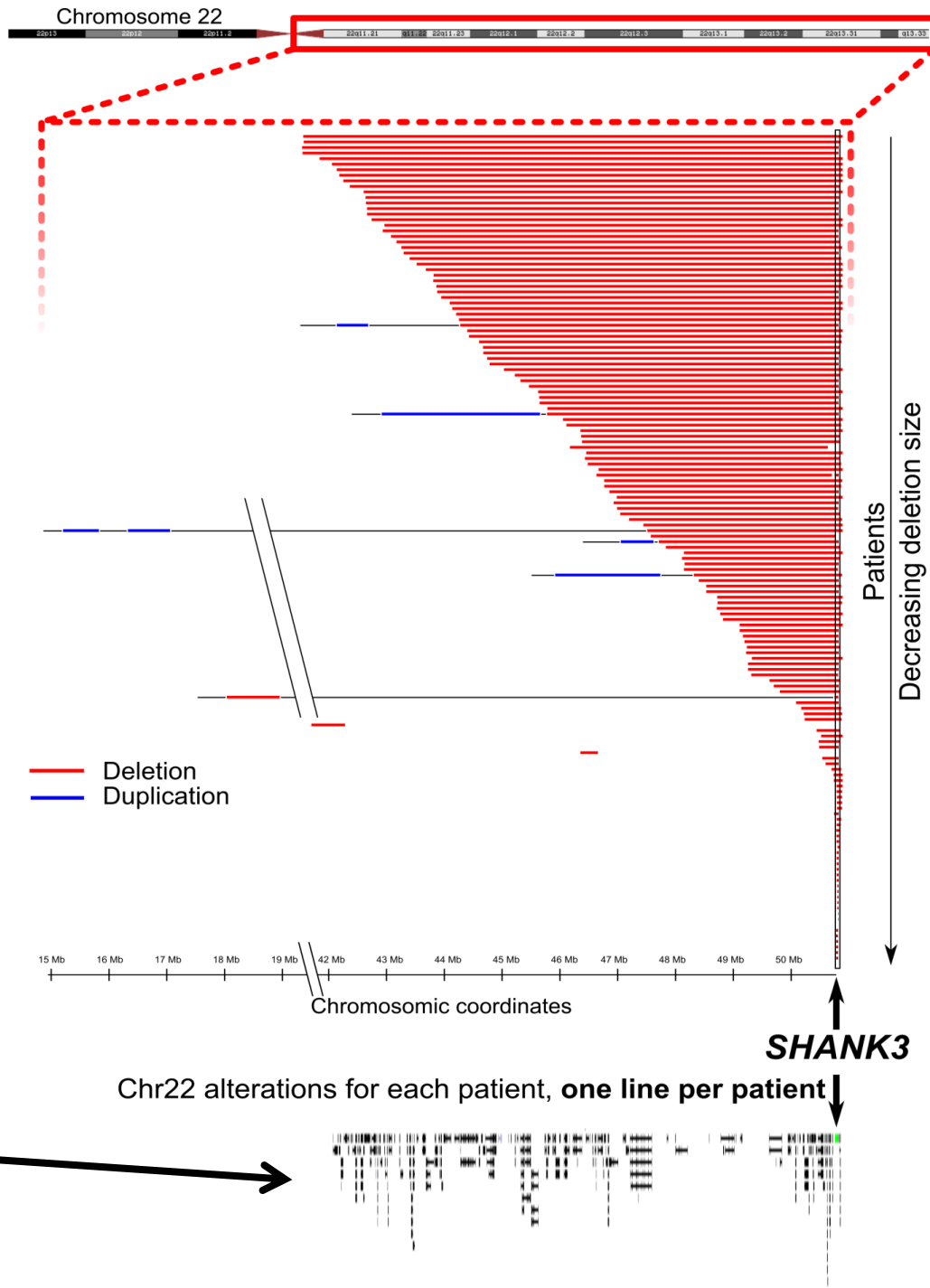
<https://grdr.hms.harvard.edu>

Phelan McDermid Syndrome

- Extremely rare genetic disease: **~1100 diagnosed patients worldwide today**
- Autistic traits, **intellectual deficiency**, slight dysmorphic features
- Also called **deletion 22q13 syndrome**
- Caused by deletions of the terminus of chromosome 22



Heterogeneity of the genetic material alterations



Heterogeneity of the phenotypes



All organs can be affected:

- Neuro-developmental
- Facial dysmorphic features
- GERD (Gastro-Esophagal Reflux)
- Renal problems
- Lax joints
- Dysplastic toenails
- Congenital cardiac diseases
- ...

Sources :

globalgenes.org

autismspeaks.org

sfari.org



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

Deep phenotyping

Knowledge from
Clinical Notes





The Apache Software Foundation

<http://www.apache.org/>

Box 1 | Natural language processing

Boundary detection] [Fx of obesity but no fx of coronary artery diseases.] [... ...															
Tokenization	Fx of obesity but no fx of coronary artery diseases .															
Normalization	- - - - - - - - - - disease_															
Part-of-speech tagging	NN IN NN CC DT NN IN JJ NN NNS															
Shallow parsing	NP PP NP NN NP															
Entity recognition	<table border="0" style="width: 100%;"> <tr> <td style="width: 33%;">Obesity</td> <td style="width: 33%;">Coronary artery disease</td> <td style="width: 33%;">Coronary artery</td> </tr> <tr> <td>Disease or disorder</td> <td>Disease or disorder</td> <td>Anatomy</td> </tr> <tr> <td>UMLS ID: C0028754</td> <td>UMLS ID: C0010054</td> <td>UMLS ID: C0205042</td> </tr> <tr> <td>Status: family history</td> <td>Status: family history</td> <td></td> </tr> <tr> <td>Negated: no</td> <td>Negated: yes</td> <td></td> </tr> </table>	Obesity	Coronary artery disease	Coronary artery	Disease or disorder	Disease or disorder	Anatomy	UMLS ID: C0028754	UMLS ID: C0010054	UMLS ID: C0205042	Status: family history	Status: family history		Negated: no	Negated: yes	
Obesity	Coronary artery disease	Coronary artery														
Disease or disorder	Disease or disorder	Anatomy														
UMLS ID: C0028754	UMLS ID: C0010054	UMLS ID: C0205042														
Status: family history	Status: family history															
Negated: no	Negated: yes															

Peter B. Jensen, Lars J. Jensen and Søren Brunak, Nat Rev Genet. 2012

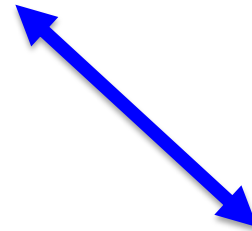


HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics



 caresync

The logo for caresync, featuring a stylized green apple icon to the left of the word "caresync" in a grey, lowercase, sans-serif font.

HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics



Health care providers



clinical notes from their Health care providers



Patients / Parents



PPRN: Phelan-McDermid Syndrome Data Network (PMS_DN)

➡ Already in place
➡ PCORI - PPRN project

Individual patient data entry including clinical notes



Patient ownership and governance of data

Aggregated or individual patient data consultation

Researchers



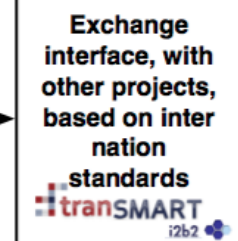
Omics' data

Collaboration with Clinical Data Research Networks (CDRN) - PCORI

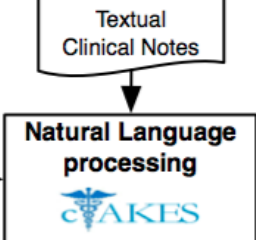
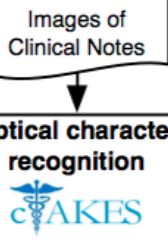
For example:
Scalable Collaborative Infrastructure for a Learning Health System (SCILHS) to find new patients with Phelan-McDermid Syndrome across all their network of 9 Hospitals

Firewall

Clinical data from Registry



Anonymized curated Clinical data from Clinical notes



Harvard Medical School Research Computing Private Orchestra: Phelan-McDermid syndrome research environment



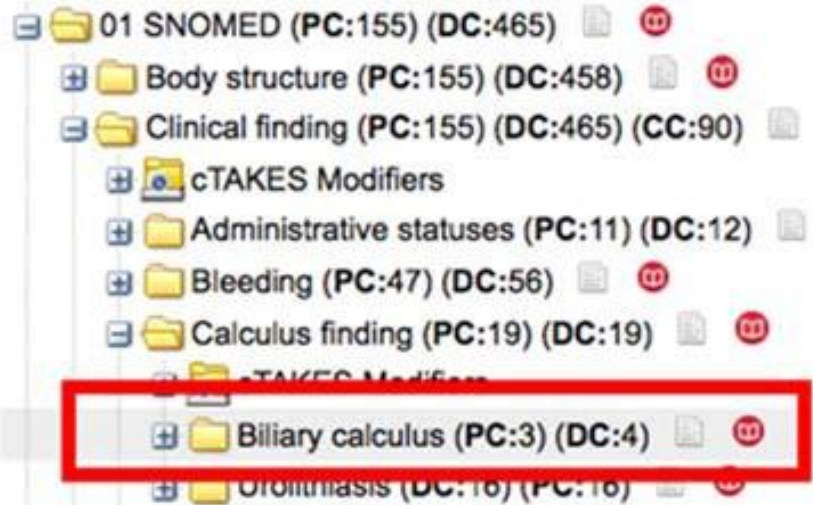
HARVARD MEDICAL SCHOOL

DEPARTMENT OF Biomedical Informatics

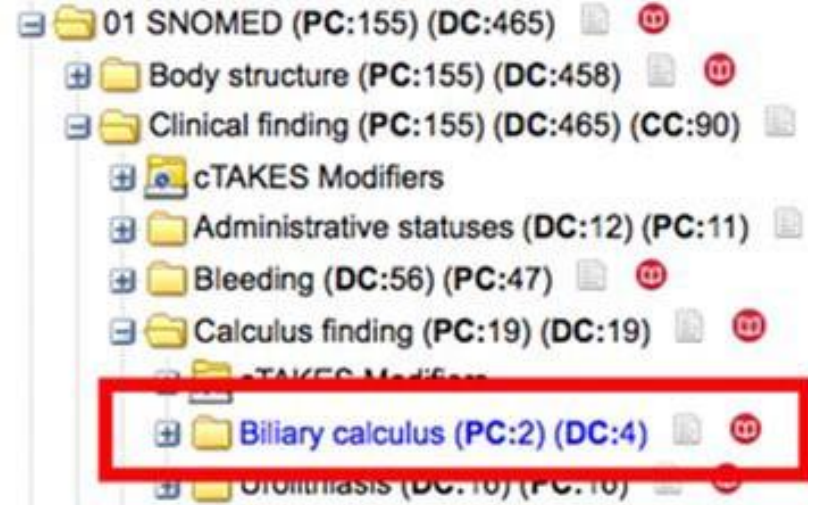
- [-] 465 Clinical Notes (cTAKES NLP)
 - [-] 01 SNOMED (PC:161) (DC:481) [document icon] [red circle with 'w']
 - [+] Body structure (PC:161) (DC:472) [document icon] [red circle with 'w']
 - [+] Clinical finding (CC:94) (PC:161) (DC:481) [document icon] [red circle with 'w']
 - [+] Event (PC:25) (DC:27) [document icon] [red circle with 'w']
 - [+] Observable entity (PC:154) (DC:431) [document icon] [red circle with 'w']
 - [+] Pharmaceutical / biologic product (PC:37) (DC:42) [document icon] [red circle with 'w']
 - [+] Procedure (CC:126) (PC:160) (DC:472) [document icon] [red circle with 'w']
 - [+] Qualifier value (PC:160) (DC:466) [document icon] [red circle with 'w']
 - [+] SNOMED CT Model Component (PC:78) (DC:108) [document icon] [red circle with 'w']
 - [+] Situation with explicit context (PC:96) (DC:145) [document icon] [red circle with 'w']
 - [+] Social context (PC:21) (DC:22) [document icon] [red circle with 'w']
 - [+] Special concept (PC:102) (DC:157) [document icon] [red circle with 'w']
 - [+] Specimen (DC:6) (PC:6) [document icon] [red circle with 'w']
 - [+] Staging and scales (DC:1) (PC:1) [document icon] [red circle with 'w']
 - [+] Substance (PC:129) (DC:251) [document icon] [red circle with 'w']
 - [+] 02 HPO (PC:153) (DC:398) [document icon] [red circle with 'w']
 - [+] 03 ICD9CM (PC:152) (DC:387) [document icon] [red circle with 'w']
 - [+] 04 ICD10CM (PC:147) (DC:355) [document icon] [red circle with 'w']
 - [+] 05 MeSH (PC:160) (DC:477) [document icon] [red circle with 'w']
 - [+] 06 NDFRT (PC:159) (DC:461) [document icon] [red circle with 'w']



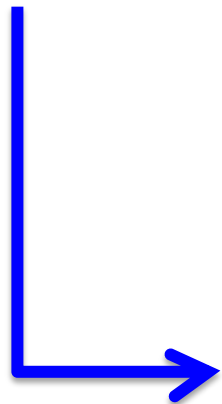
BEFORE validation



AFTER validation



Pop-up validation window

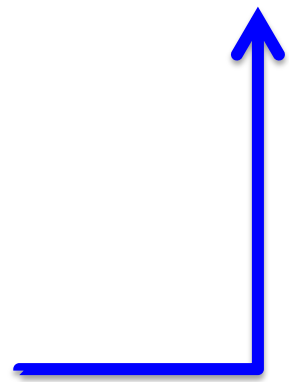


Node Metadata and Statistics

Patient 3 Patient automatically included by natural language processing

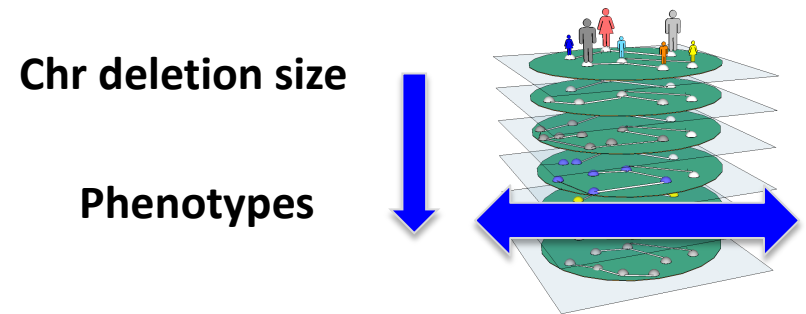
1 - Please note evaluation of the abdominal organs is secondary to the lack of intravenous contrast material. **Gallstones** are seen within the gallbladder lumen.

Patient 4 Patient automatically included by natural language processing



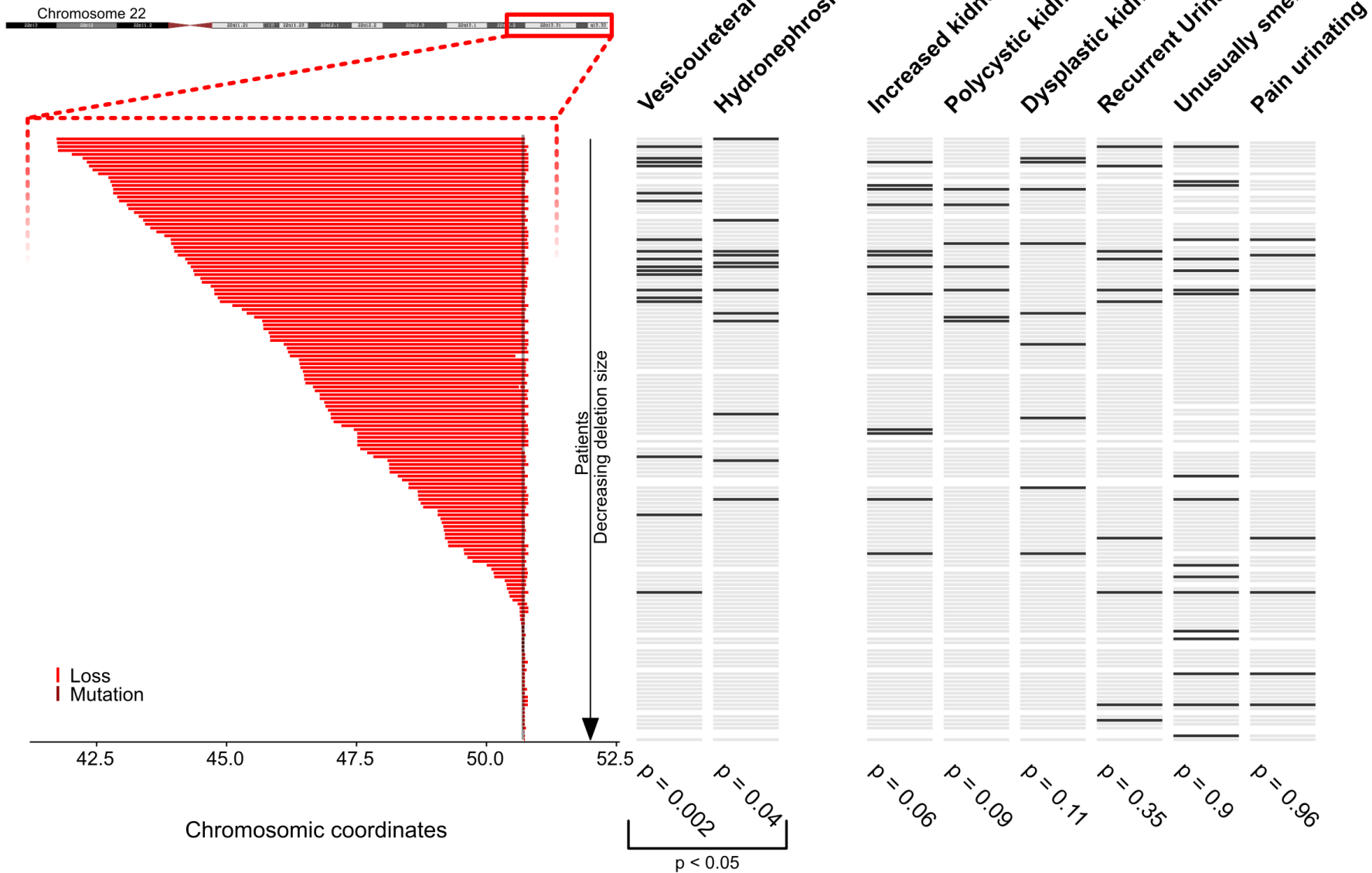
Objective:

Identify phenotypes linked
with the **deletion of other genes** than *SHANK3*

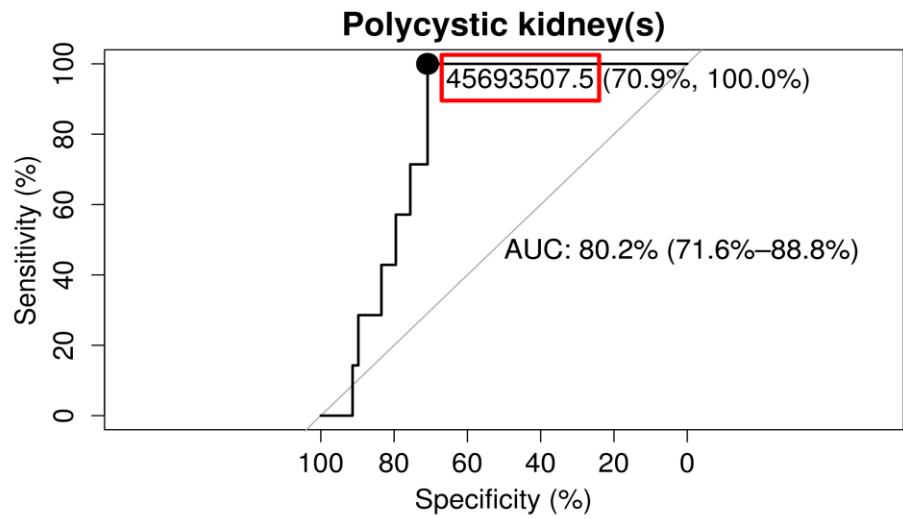
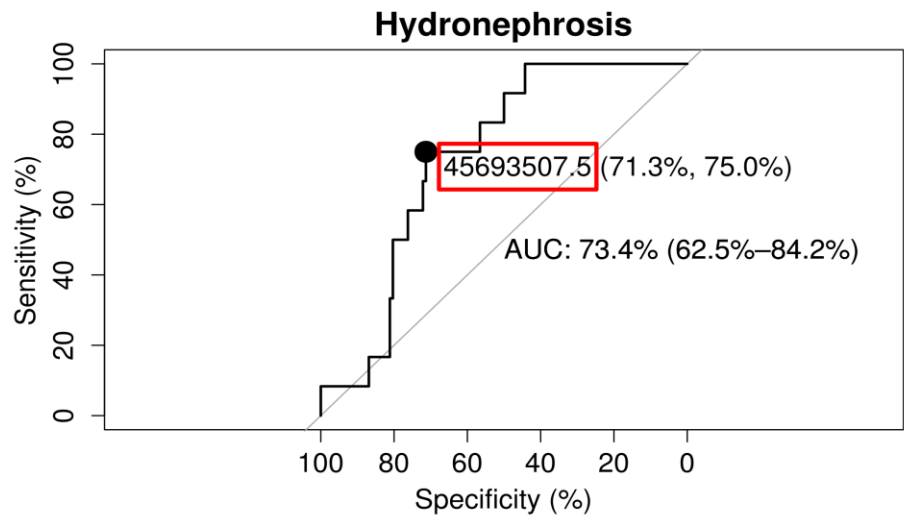
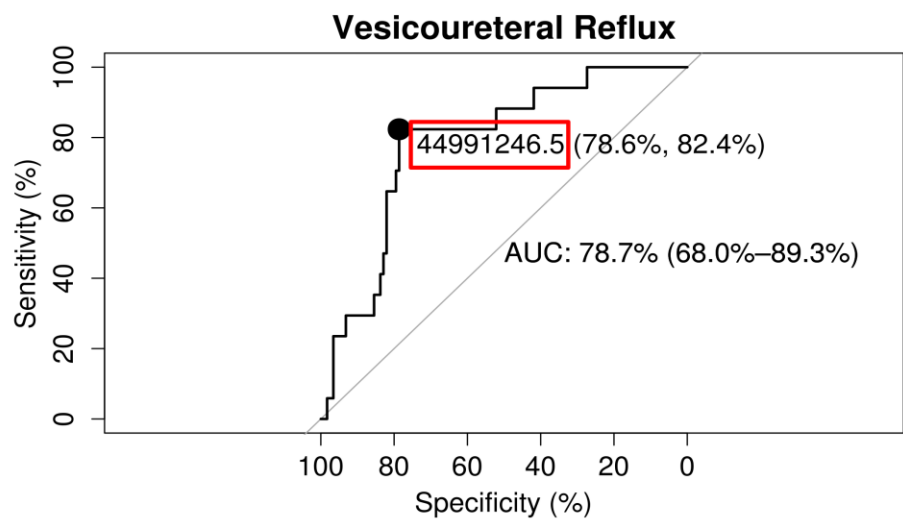
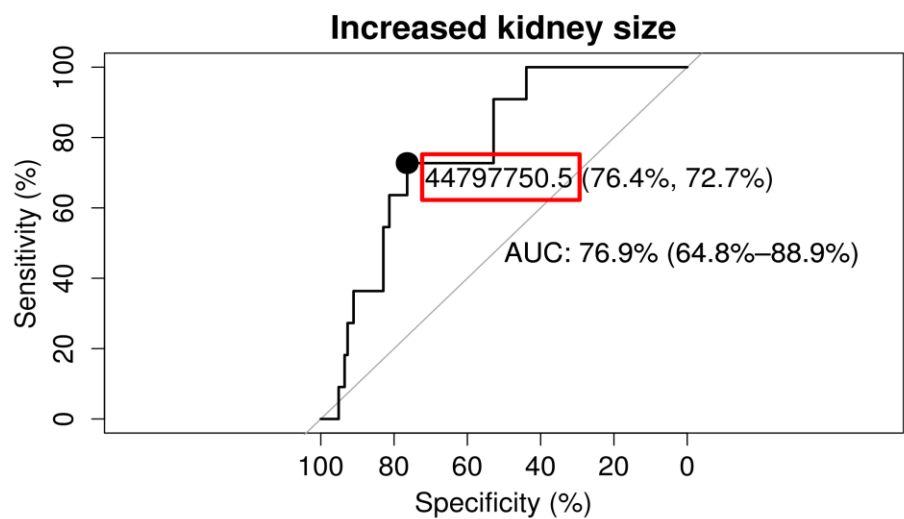


Results — Kidney malformations associated with other genes that *SHANK3*

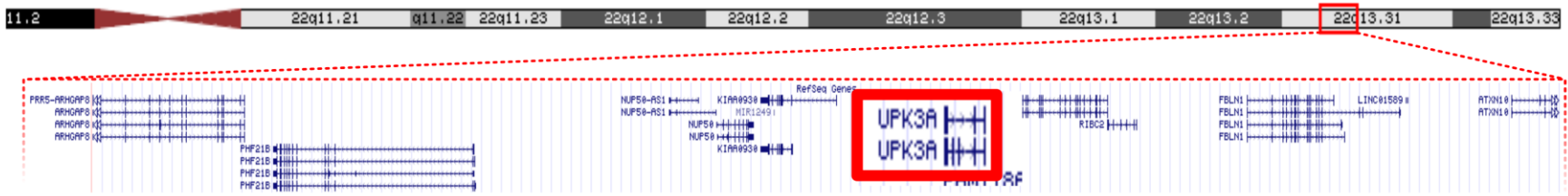
Phenotype Absent Present Missing



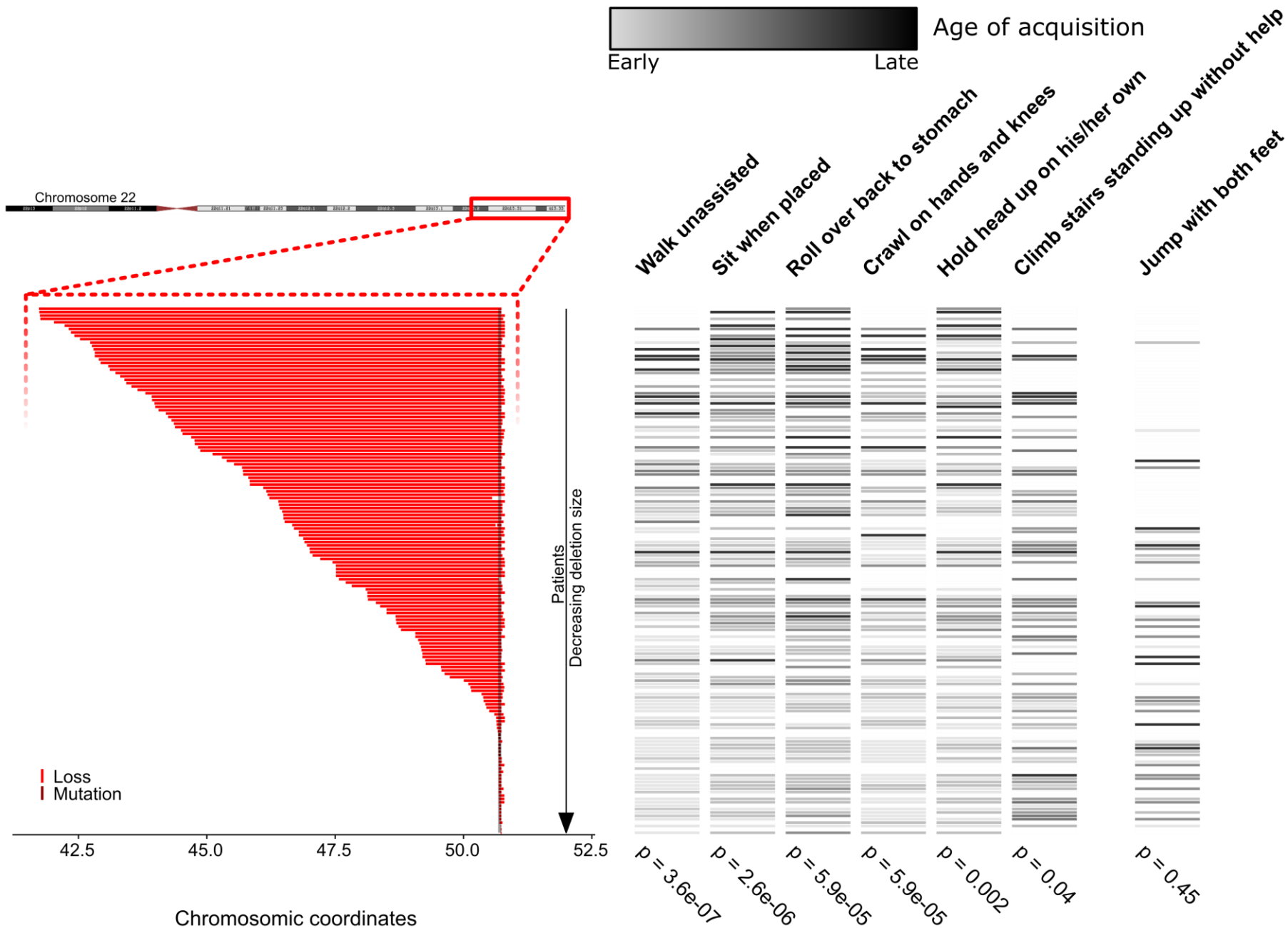
Panel A



Panel B

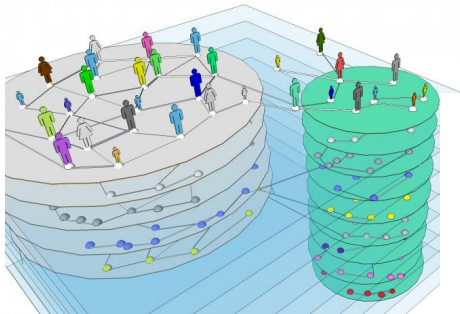


Results — Gradient of gross motor delays showing a cumulative effect



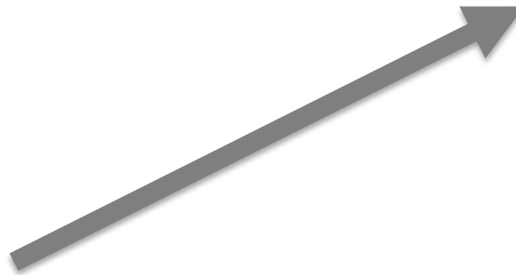
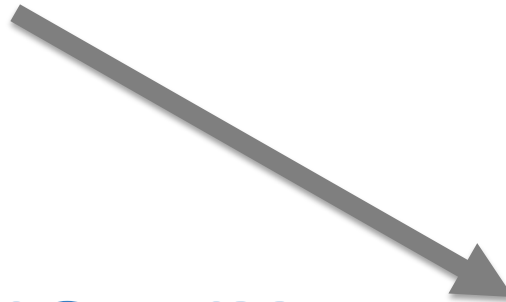


1.0_HOTFIX



Harvard – DBMI

AvillachLab improvements



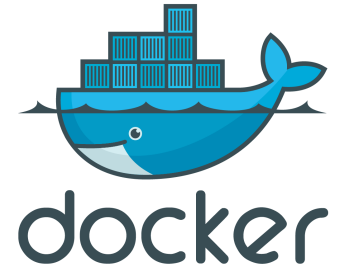
HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

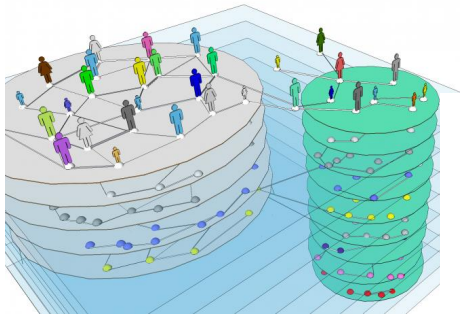


HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics



16.2 (~~17.12~~)



Harvard – DBMI

AvillachLab improvements



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

An i2b2/tranSMART project management committee

If interested to join, email us:

Diane:

dkeogh@i2b2foundation.org

Paul:

paul_avillach@hms.harvard.edu



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

- **Play with i2b2/tranSMART UI:**
 - <https://grdr.hms.harvard.edu>
 - <https://nhanes.hms.harvard.edu>
 - <https://demo-ngs.hms.harvard.edu>
 - <https://pmsdn.hms.harvard.edu> (*ASD researcher only*)
-
- **Play with BD2K PIC-SURE RESTful API**
 - <http://bd2k-picsure.hms.harvard.edu>
 - <http://exac.hms.harvard.edu>



Students / Postdocs



Cartik



Li



Romain



Qiu-Yue



Ombeline



Maxime



Antoine



Laurie



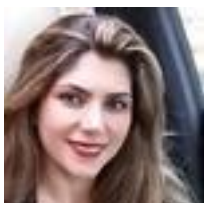
Laura



Haishuai



Mahdi



Niloofar



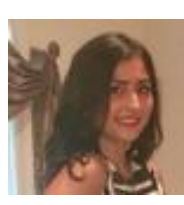
Alba



Joany



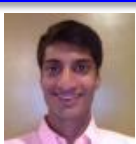
Carlos



Romina



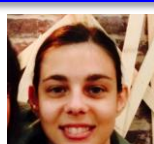
Emmanuelle



Yuri



Samuel



Claire



Antoine

Staff / Software developers



Jason



Ranjay



Gabor



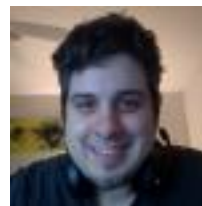
Thomas



Jaspreet



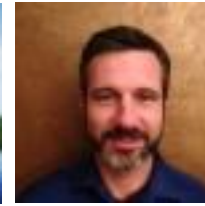
Alex



Andrew



Andre



Sean



Anoush



Cassandra



Libby



Sophia



Alyssa

Alumni



Michael



Sushma



Pei



Ephi



Jeremy

Postdoc? / Mobilité?



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics

www.avillach-lab.hms.harvard.edu

Paul_Avillach.hms.harvard.edu



HARVARD
MEDICAL SCHOOL

DEPARTMENT OF
Biomedical Informatics