

Multiomics Provider

Institute for Systems Biology

*Accelerating the creation of
high-value knowledge providers
from raw data*



Gustavo Glusman,
PhD



Jennifer Hadlock,
MD



Ilya Shmulevich,
PhD



Guangrong Qin,
PhD



Ryan Roper,
MS

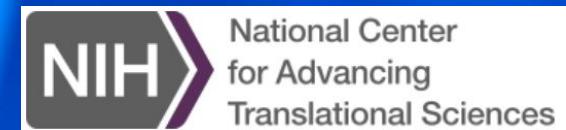


Denise Mauldin,
MS



Cheng Dai,
MEng

Thank you
Theo Knijnenburg
John Earls
Sui Huang
Jared Roach



Milestones

- M1.1 Ongoing support of Big GIM

- M1.2 Prototype DOCKET Overview

- M1.3 Prototype DOCKET Compare

- M1.4 Prototype DOCKET Integrate

- M1.5 Gather input on DOCKET features

- M1.6 Collaborate on Reasoner Standard API



Deliverables

DOCKET code, with Docker container

github.com/PriceLab/DOCKET

Instructions for running prototypes

drive.google.com/drive/u/1/folders/19CT2bu1kzVnXgORhgIQijJd7x8O8Ez6D

Translator Drive > Knowledge Providers

Try-it-yourself document

Sample data

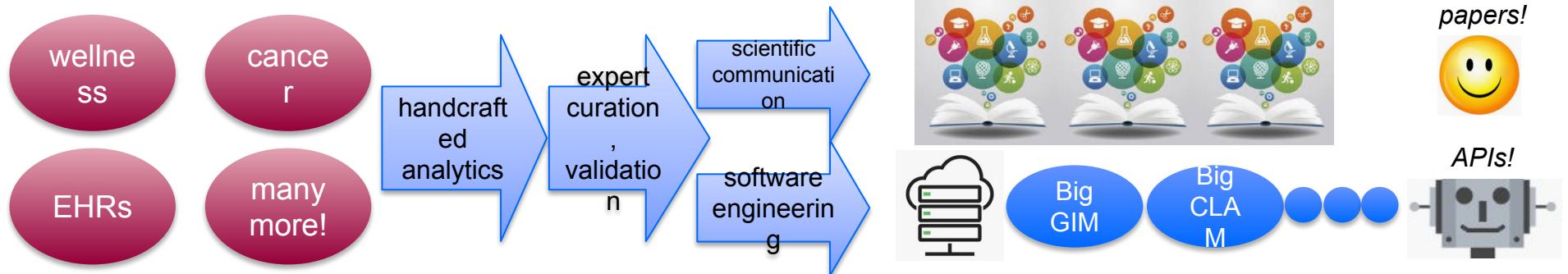
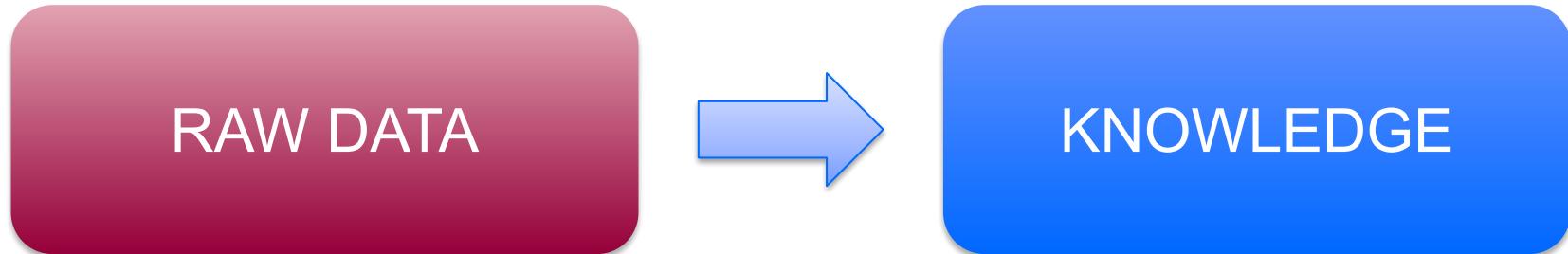
Big GIM: Gene Interaction Miner, and

Big CLAM: Cell Line Association Miner:

biggim.ncats.io/api

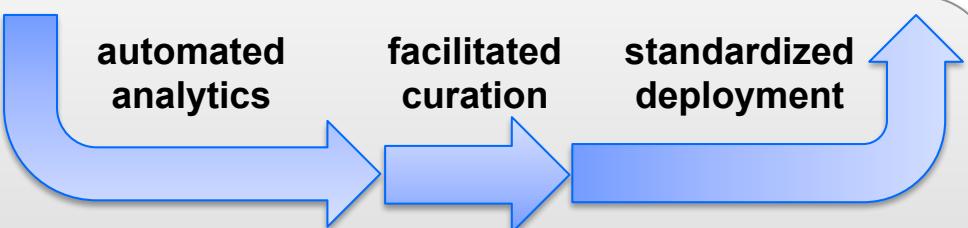
github.com/PriceLab/translator-bigquery-api





DOCKET:

- “Data Overview, Comparison, and Knowledge Extraction Tool”
- Utilities for extracting knowledge, evaluating, visualizing, serving...
- The resulting dossier, a resource readily available for ARAs



(1) Contribute multiomic Knowledge Providers

(2) Automate the creation of KPs via DOCKET



Big GIM: Gene Interaction Miner

Developed during Translator Phase 1 (Theo & John, Blue Team)

In use by workflows, collaborators, ARAs; M1.1: ongoing support

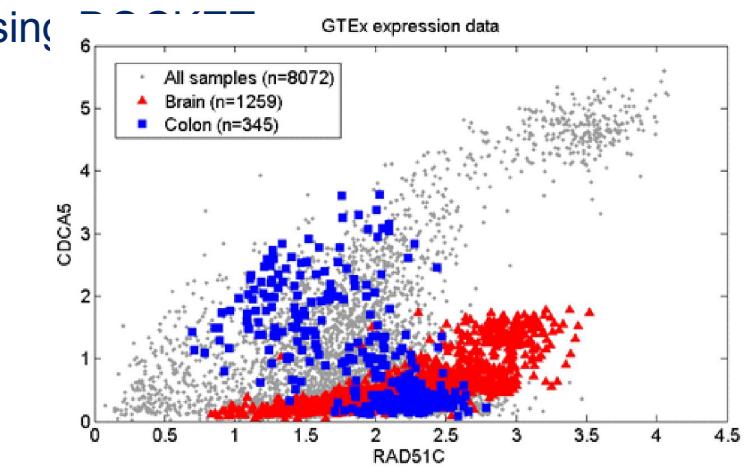
- **Precompute** all pairwise gene expression correlations, in tissue context
- **Curate** to ensure correctness
- **Store** results in Google Cloud Bigtable
- **Serve** via API as proto-KP on gene similarity



Gene 1	Gene 2	GTEx expression correlations		TCGA expression correlations		GIANT functional interaction		BioGRID direct interaction	
		Skin	Colon	SKCM	COAD	Skin	Colon		
tissue →	RAD51C	CDC6	0.65	-0.68	0.47	0.51	0.62	0.52	0
	KPNA3	BRCA2	0.70	-0.60	0.65	0.57	0.08	0.06	0
	PDS5B	BRCA2	-0.27	-0.57	0.66	0.69	0.12	0.10	1
							

410M rows

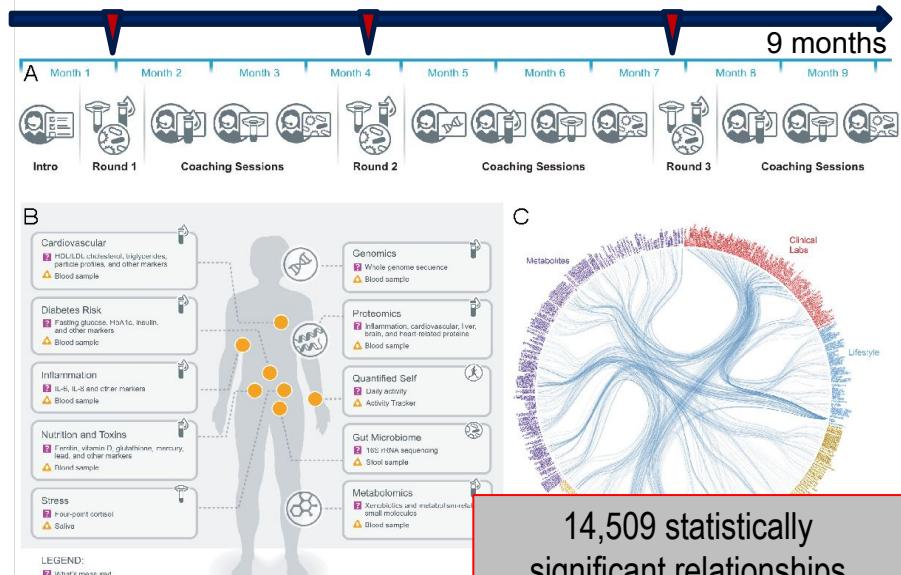
400 columns



Multiomics Data to Knowledge: Wellness

2016, Pioneer 100 Project: **N=108**

Whole genomes, 213 clinical lab tests, 643 metabolites, 392 proteins, gut microbiome



**nature
biotechnology**

A wellness study of 108 individuals using personal, dense, dynamic data clouds

Nathan D Price^{1,2,6,7}, Andrew T Magis^{2,6}, John C Earls^{2,6}, Gustavo Glusman¹, Roie Levy¹, Christopher Lausted¹, Daniel T McDonald^{1,5}, Ulrike Kusebauch¹, Christopher L Moss¹, Yong Zhou¹, Shizhen Qin¹, Robert L Moritz¹, Kristin Brogaard², Gilbert S Omenn^{1,3}, Jennifer C Lovejoy^{1,2} & Leroy Hood^{1,4,7}

Today: **N~5000**, 96 data snapshots

Genetics, clinical labs, metabolites, proteins, microbiome, medications, activity, sleep, weight, hormones, assessments of vitals, diet, digestive health, lifestyle, etc.



Journals of Gerontology: Biological Sciences
cite as: J Gerontol A Biol Sci Med Sci, 2019, Vol. S1, S52-S60
doi:10.1093/gerona/glz220

OXFORD

Healthy Longevity 2019: Supplement Article

Multi-Omic Biological Age Estimation and Its Correlation With Wellness and Disease Phenotypes: A Longitudinal Study of 3,558 Individuals

SCIENTIFIC REPORTS

OPEN Genetic Predisposition Impacts Clinical Changes in a Lifestyle Coaching Program

International Journal of Obesity

Original Article | Published: 12 September 2017

Habitual sleep duration and sleep duration variation are independently associated with body mass index

**nature
biotechnology**

Article | Published: 02 September 2019

Blood metabolome predicts gut microbiome α -diversity in humans

Cell Reports
Article

A Multi-omic Association Study of Trimethylamine N-Oxide

Ohad Manor,^{1,3,*} Niha Zubair,¹ Matthew P. Conomos,¹ Xiaojing Xu,¹ Jesse E. Rohwer,¹ Cynthia E. Krafft,¹ Jennifer C. Lovejoy,^{1,2} and Andrew T. Magis¹

¹Arivale, Inc., Seattle, WA 98104, USA

²Institute for Systems Biology, Seattle, WA 98109, USA

New Results

Gut Microbiome Pattern Reflects Healthy Aging and Predicts Extended Survival in Humans

Tomasz Wilmanski, Christian Diener, Noa Rappaport, Sushmita Patwardhan, Jack Wiedrick, Jodi Lapidus, John C. Earls, Anat Zimmer, Gustavo Glusman, Max Robinson, James T. Turkovich, Deborah M. Kado, Jane A. Cauley, Joseph Zmuda, Nancy E. Lane, Andrew T. Magis, Jennifer C. Lovejoy, Sean M. Gibbons, Leroy Hood, Eric S. Orwoll, Nathan D. Price

doi: <https://doi.org/10.1101/2020.02.26.966747>



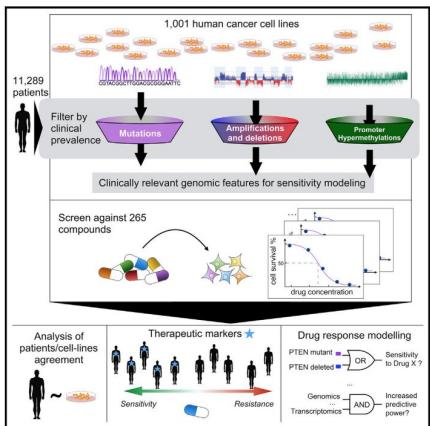
Multiomics Data to Knowledge: Cancer

Which **gene mutations** in any of **these pathways** are associated with sensitivity or resistance to any of **these drugs** in the cell lines from **these tumor types**?

Cell

A Landscape of Pharmacogenomic Interactions in Cancer

Graphical Abstract



Authors

Francesco Iorio, Theo A. Knijnenburg, Daniel J. Vis, ..., Julio Saez-Rodriguez, Ultan McDermott, Mathew J. Garnett

Correspondence

um1@sanger.ac.uk (U.M.), mg12@sanger.ac.uk (M.J.G.)

In Brief

A look at the pharmacogenomic landscape of 1,001 human cancer cell lines points to new treatment applications for hundreds of known anti-cancer drugs.

Accession Numbers

GSE68379
E-MTAB-3610

Highlights

- We integrate heterogeneous molecular data of 11,289 tumors and 1,001 cell lines
- We measure the response of 1,001 cancer cell lines to 265 anti-cancer drugs
- We uncover numerous oncogenic aberrations that sensitize to an anti-cancer drug
- Our study forms a resource to identify therapeutic options for cancer sub-populations

Mutation

+

Drug response

Expression

Methylation

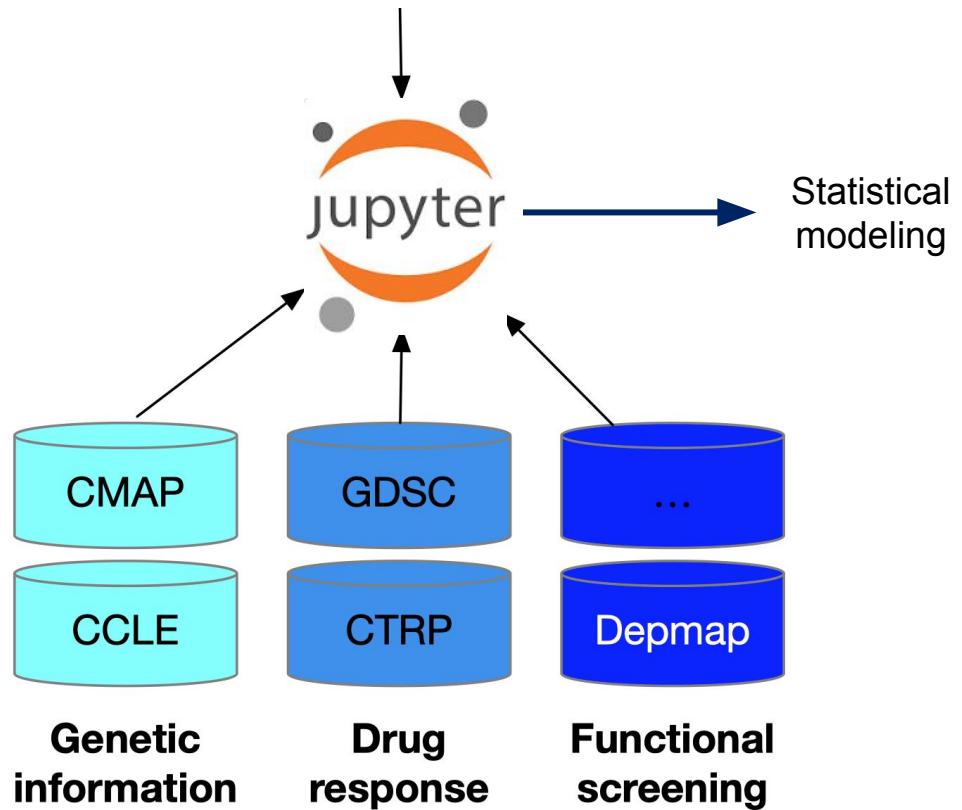


1001 cancer cell lines, 265 anti-cancer drugs
Over 30 tumor types



Multiomics Data to Knowledge: Cancer

Which **gene mutations** in any of **these pathways** are associated with relative increased or decreased viability of cells that have been exposed to any of **these gene knockdowns** in the cell lines from **these tumor types**?



Pathway List

Ribosome (n=153) × Cell cycle (n=124) × p53 signaling pathway (n=72) × Homologous recombination (n=41) × [View Pathway List](#) | [Clear All](#)

Enter Pathway Name

Gene Knockdown List

PARP1 × Enter Gene Name [Upload Gene List](#) | [Clear All](#)

Tumor Type List

Breast invasive carcinoma (n=51) × Ovarian serous cystadenocarcinoma (n=34) × Liver hepatocellular carcinoma (n=17) × [View Tumor Type List](#) | [Clear All](#)

Enter Tumor Type Name

Answer

There are no gene mutations in the selected pathways associated with increased viability to gene knockdown **KD_PARP1**.
Mutations in **BRCA2** is associated with decreased viability to gene knockdown **KD_PARP1**.

[View Details](#)

[View in Notebook](#)

Plot

[Download](#) | [View Details](#)

Tables

Decreased viability upon gene KD [Download](#) | [View Details](#)

Gene KD	Mutation	Effect Size	P-Value	FDR
KD_PARP1	BRCA2	-5.91e-01	3.25e-02	4.02e-01

Increased Viability upon gene KD [Download](#) | [View Details](#)

Gene KD	Mutation	Effect Size	P-Value	FDR
---------	----------	-------------	---------	-----

-log₁₀(P-value)

(Decreased viability) <----- Effect Size -----> (Increased viability)

Detailed description: This section shows a user interface for querying multiomics data. It includes dropdown menus for selecting pathways (Ribosome, Cell cycle, p53 signaling, Homologous recombination), gene knockdowns (PARP1), and tumor types (Breast invasive carcinoma, Ovarian serous cystadenocarcinoma, Liver hepatocellular carcinoma). The "Answer" section displays results for the query, mentioning that mutations in BRCA2 are associated with decreased viability to gene knockdown KD_PARP1. Below this is a scatter plot of -log10(P-value) versus Effect Size, with points for KD_PARP1 and BRCA2. Tables provide detailed statistics for each gene mutation.



EHR Phenotype Data to Knowledge

Clin Transl Sci. 2019 Jul;12(4):329-333.

Clinical Data: Sources and Types, Regulatory Constraints, Applications.

Ahalt SC¹, Chute CG², Fecho K¹, Glusman G³, Hadlock J³, Taylor CO², Pfaff ER⁴, Robinson PN⁵, Solbrig H², Ta C⁶, Tatonetti N⁶, Weng C⁶; Biomedical Data Translator Consortium.

Opportunity

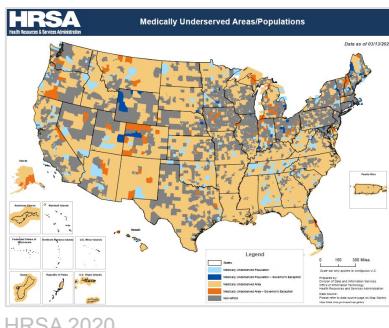
Private data to public knowledge*

4.3 million patients

- 7 states from Alaska to Texas
- 51 hospitals, 829 clinics
- Medically underserved areas
- Longitudinal records

30 billion data elements

- 330 M encounters
- 27 M diagnoses
- 7 M vital signs
- 260 M drug details
- 947 M lab results
- 11 M surgical logs
- Full text notes



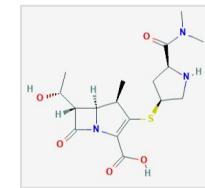
Challenge

Systemic biases are refractory to scale

Accurate phenotype requires higher fidelity than ICD-10 code
Biomedical logic supports substratification and classification

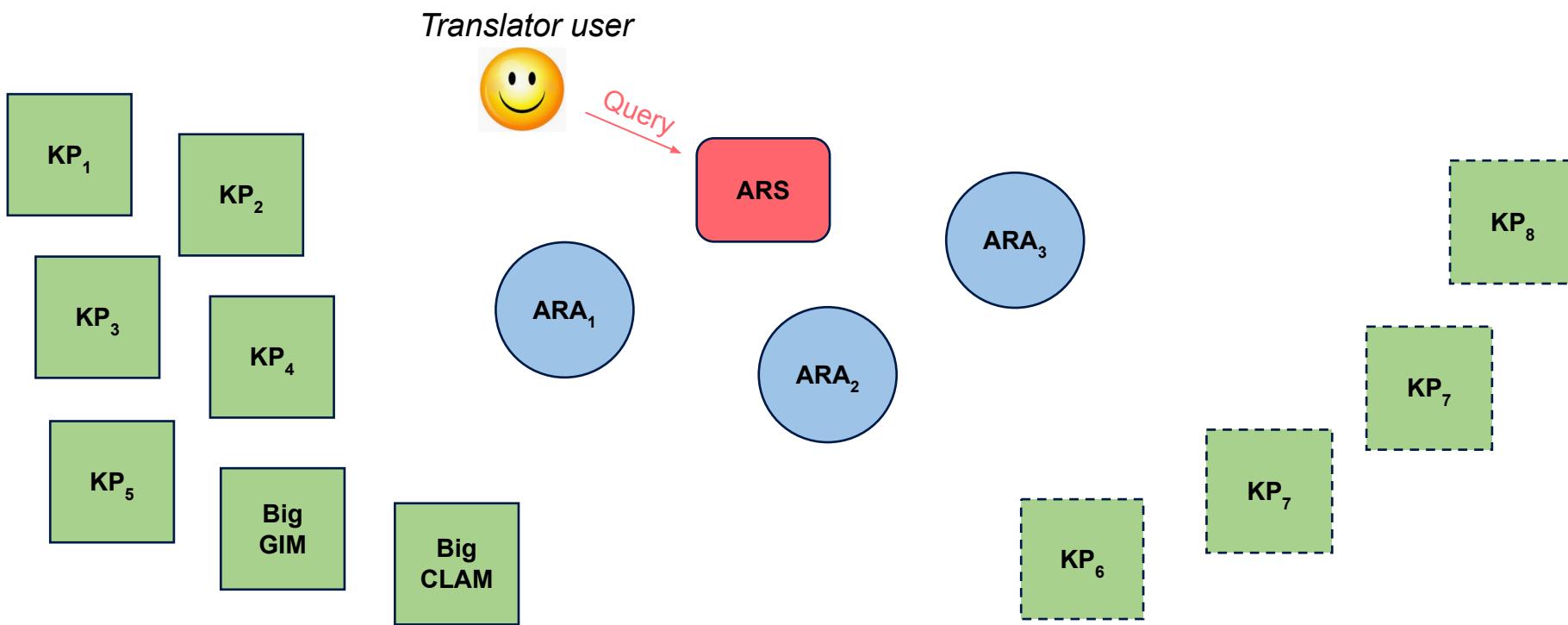
Examples:

- HbA1C correlations, going beyond billing label for diabetes
- Empirical treatment for suspected sepsis vs. clinical sepsis
- Care spread across large catchment areas
- Maternal care, links between maternal and fetal records



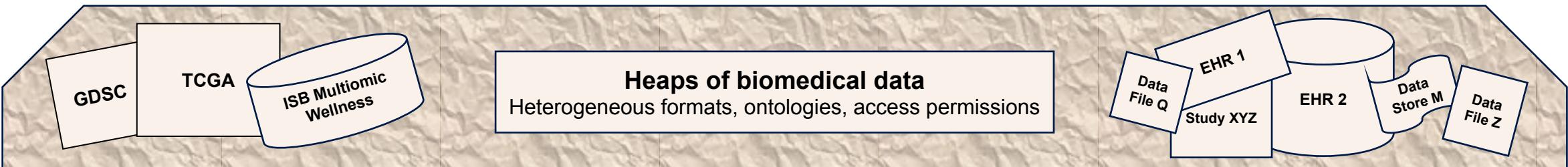
The following list provides a sense of the breadth of data types to be considered for the Biomedical Data Translator, but it is not exhaustive:

Diseases	Microbiome data
Signs of disease	Molecular mechanisms
Symptoms (from patient registries and natural history studies)	Signaling pathways
Patient-reported outcomes	Molecular and cellular networks
Electronic health records	Environmental factors
Clinical encounters	Disease etiologies
Prescription data	Proteomes
Health insurance claims data	Proteins
Diagnostic labs	Post translational modifications (PTMs)
Biomedical imaging data	Co-factors
Adverse event reports	Transcriptomes
Biomarkers	Epigenomes
Organ systems	Genes
Subanatomy	Genetic mutations
Tissue types	Functional polymorphisms
Cell types	Therapeutic interventions
Cell lineages	Intervention exposure
Cell processes	Pharmacokinetics/Pharmacodynamics
Organelles	Clinical trial data
Orthologs/Animal models	



The important question, motivating DOCKET:

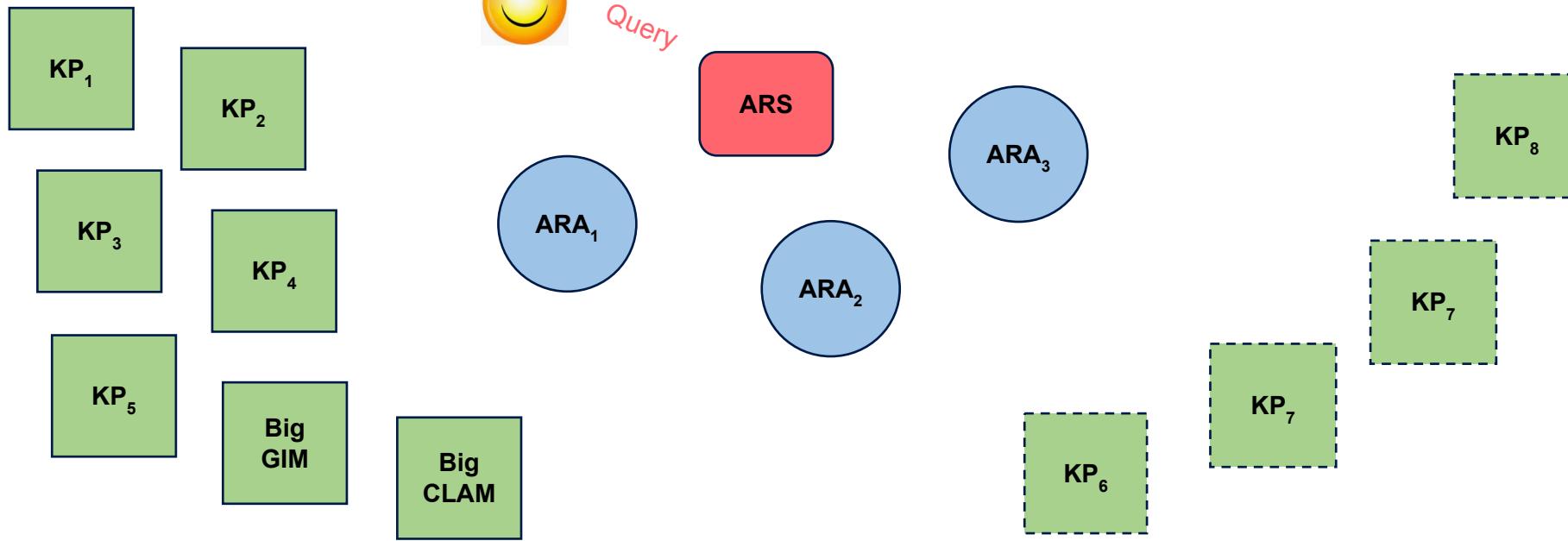
***How are we going to generate KPs to cover all desired data types?
(And how much will it cost?)***



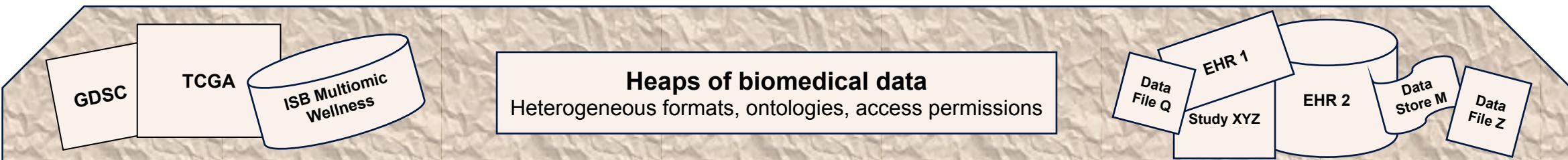
Translator user



Query



The first mile problem



Translator user

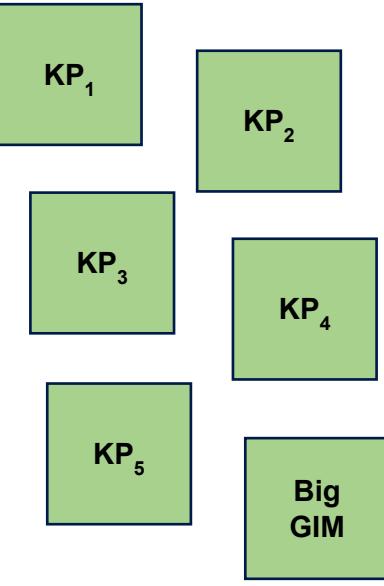


Query

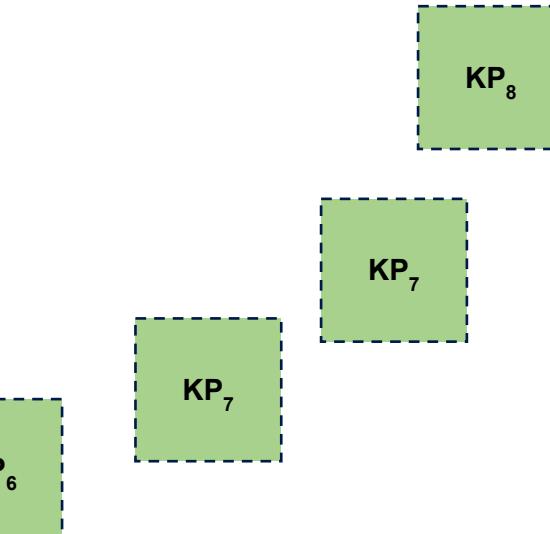
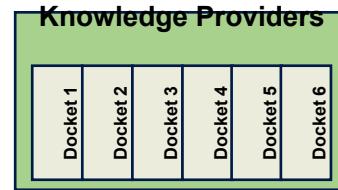
ARS

ARA₁

ARA₃



Big CLAM



Domain expert



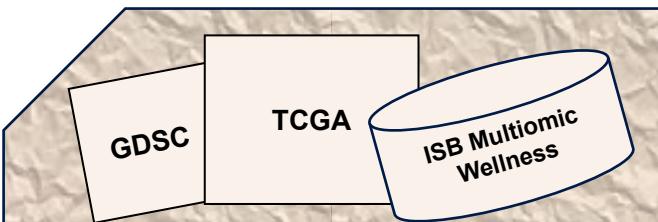
Interact

Our solution: DOCKET workbench

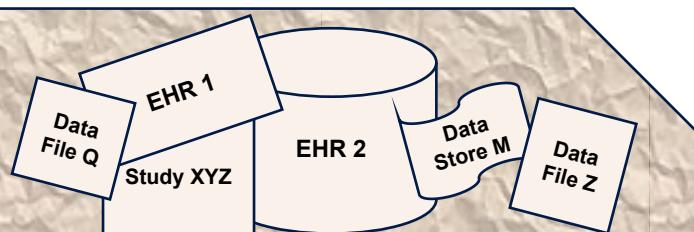
Rapidly study, compare and integrate data

Accelerate curation of knowledge

A TIDBIT generator, and a “compiler” of KPs



Heaps of biomedical data
Heterogeneous formats, ontologies, access permissions



Overview of DOCKET (prototype)

M1.2

DOCKET Overview:

a system for better understanding the data

M1.3

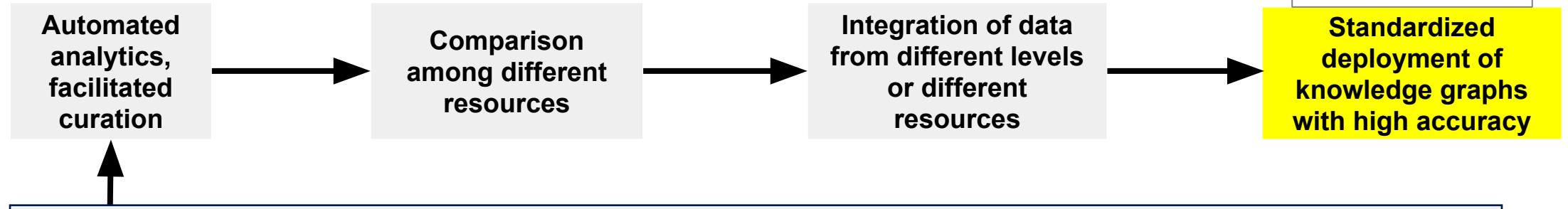
DOCKET Compare:

compare different datasets
compare different snapshots

M1.4

DOCKET Integrate:

statistical models to rationalize the connections among different entities or features

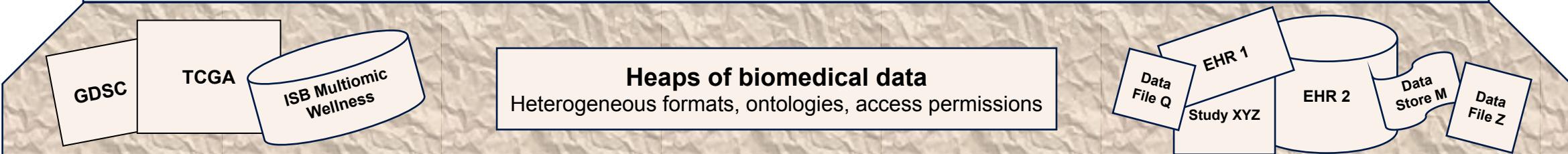


Our solution: DOCKET workbench

Rapidly study, compare and integrate data

Accelerate curation of knowledge

A TIDBIT generator, and a “compiler” of KPs



Features of the DOCKET workbench

What we mean by...

DOCKET: The utilities for extracting knowledge, evaluating, visualizing, serving...

A docket: The resulting dossier, a resource readily available for ARAs

M1.2

DOCKET Overview:
a system for better
understanding the data

DOCKET Study

Automated pipeline from data to docket

DOCKET Overview

Visualize contents of docket, interact to refine

M1.3

DOCKET Compare:
compare different datasets
compare different snapshots

DOCKET Match

Compare dockets to find common data types

DOCKET Compare

Compare versions of the data to assess
knowledge robustness

M1.4

DOCKET Integrate:
statistical models to rationalize
the connections among different
entities or features

DOCKET Merge

Merge compatible tables

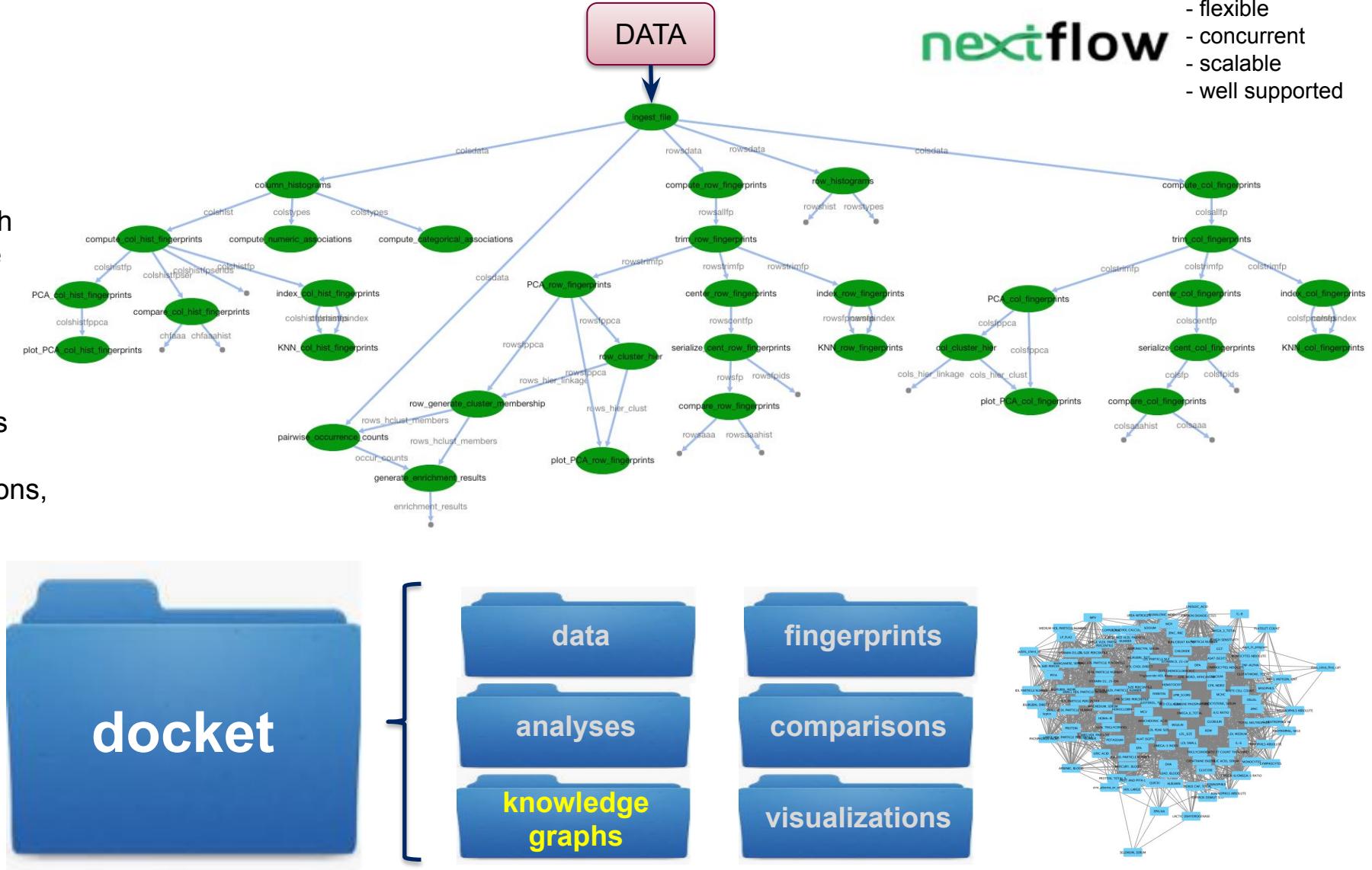
DOCKET Integrate

Statistical modeling to derive knowledge



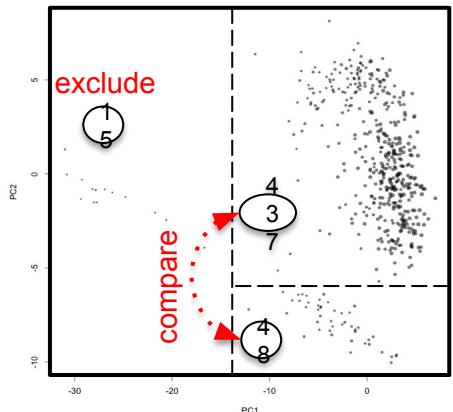
An automated, concurrent analysis pipeline

- Ingest data
 - Statistics, distributions
 - Dimensionality reduction (fingerprinting, PCA) – both row-wise and column-wise
 - Hierarchical clustering
 - Enrichment analysis
 - PCA visualizations
 - Data indexing
 - All-against-all comparisons
 - k-nearest neighbors
 - Pairwise column associations, numerical and categorical (bivariate statistics)

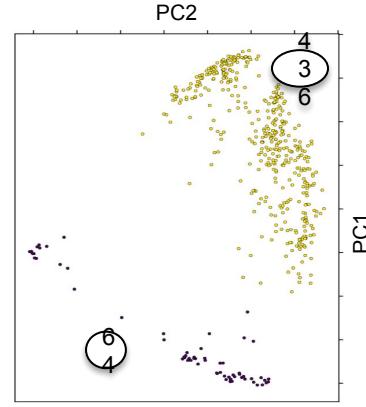


Automated analytics, curation support

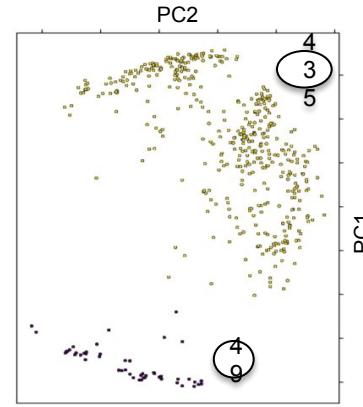
St. Jude Life cohort N=500



automation



curation,
refinement



Data types: demographics, body metrics, genetic tests, diagnoses (chronic), secondary primary cancers, chemotherapy + dosage, radiation + dosage

- 5.0x Adrenal insufficiency
- 4.9x Cerebellar dysfunction
- 4.6x Growth hormone deficiency
- 3.0x Vascular disease
- 2.6x Obstructive sleep apnea
- 2.5x Hypothyroidism
- 2.5x Cerebrovascular accident
- 1.8x Radiation to thyroid
- 1.7x GERD

Term	p-value
Hypothyroidism	1.7E-08
Adrenal insufficiency	9.9E-08
Childhood GHD	4.3E-06
Adult GHD	1.6E-05
Hearing loss	1.7E-05
Cerebellar dysfunction	3.7E-05

Enrichment analysis,
unclear statistics

Knowledge on pediatric oncology multi-morbidities:
association between disease entities (biolink)



DOCKET Overview

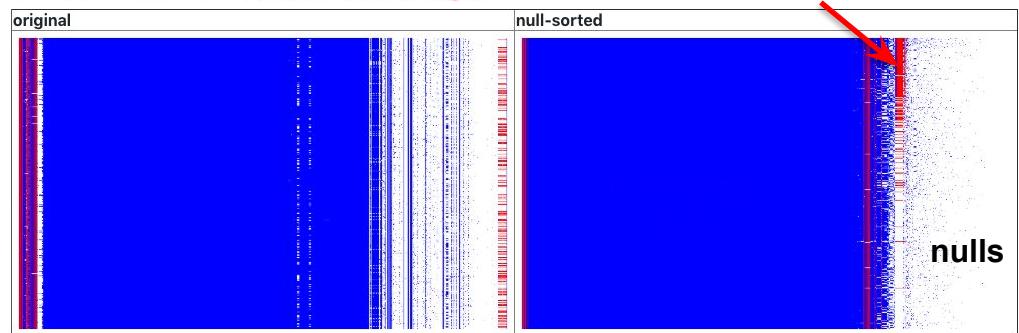
Interactive report □ discovery and curation

DOCKET Overview for SJL500 nonLabVars

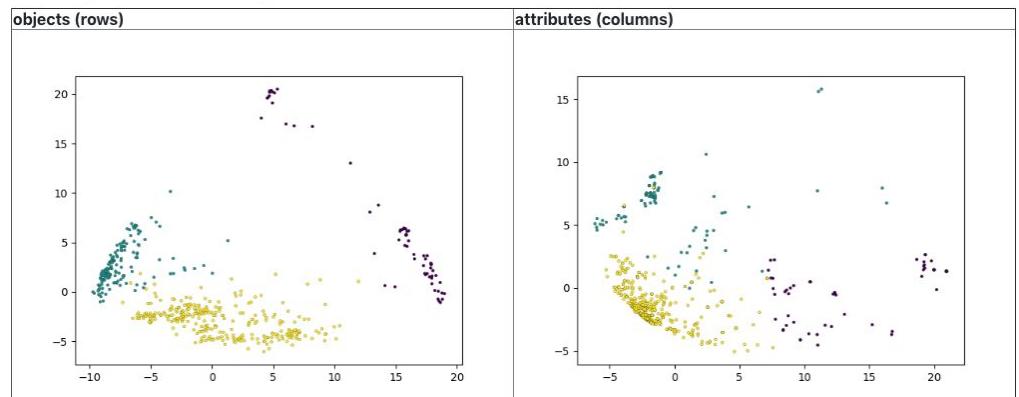
Docket path: demo_dockets/SJL500_nlv-200

Table size: 500 rows x 838 columns (41 of which are empty, 765 are numerical)

Data overviews (numbers, strings)



PCA (PC1 vs. PC2)



Numerical associations

Show 10 entries

variableA	variableB	N	rho	pval
chemodose.vinca_dose_prim	chemodose.vinca_dose_10	347	1.000	0.00e+00
chemodose.vinca_dose_prim	chemodose.vinca_dose_any	346	1.000	0.00e+00
chemodose.vincrist_dose_5	chemodose.vcr_dose_5	331	1.000	0.00e+00
chemodose.mtx_iv_dose_10	chemodose.mtx_iv_dose_prim	91	1.000	0.00e+00
chemodose.mtx_iv_dose_10	chemodose.mtx_iv_dose_any	91	1.000	0.00e+00
chemodose.mtx_dose_10	chemodose.mtx_dose_any	156	1.000	0.00e+00
chemodose.cortico_dose_10	chemodose.cortico_dose_any	15	1.000	0.00e+00
chemodose.cortico_dose_10	chemodose.cortico_dose_prim	15	1.000	0.00e+00
chemodose.cortico_dose_10	chemodose.cortico_dose_5	15	1.000	0.00e+00
chemodose.mercapto_po_dose_5	chemodose.mercapto_po_dose_10	9	1.000	0.00e+00

Showing 1 to 10 of 528 entries

Previous 1 2 3 4 5 ... 53 Next

Categorical associations

Show 10 entries

variableA	variableB	N	Theil_B_A	Theil_A_B
chemodose.mercapto_po_dose_10	chemodose.mercapto_po_dose_5	409	1.000	1.000
chemodose.mercapto_po_dose_10	chemodose.mercapto_po_dose_any	409	1.000	1.000
chemodose.mercapto_po_dose_10	chemodose.mercapto_po_dose_prim	409	1.000	1.000
chemodose.mercapto_po_dose_10	chronic.Cardiovascular_Thrombus	10	1.000	1.000
chemodose.mercapto_po_dose_5	chemodose.mercapto_po_dose_any	409	1.000	1.000
chemodose.mercapto_po_dose_5	chemodose.mercapto_po_dose_prim	409	1.000	1.000
chemodose.mercapto_po_dose_5	chronic.Cardiovascular_Thrombus	10	1.000	1.000
chemodose.mercapto_po_dose_any	chemodose.mercapto_po_dose_prim	409	1.000	1.000
chemodose.mercapto_po_dose_any	chronic.Cardiovascular_Thrombus	10	1.000	1.000
chemodose.mercapto_po_dose_prim	chronic.Cardiovascular_Thrombus	10	1.000	1.000

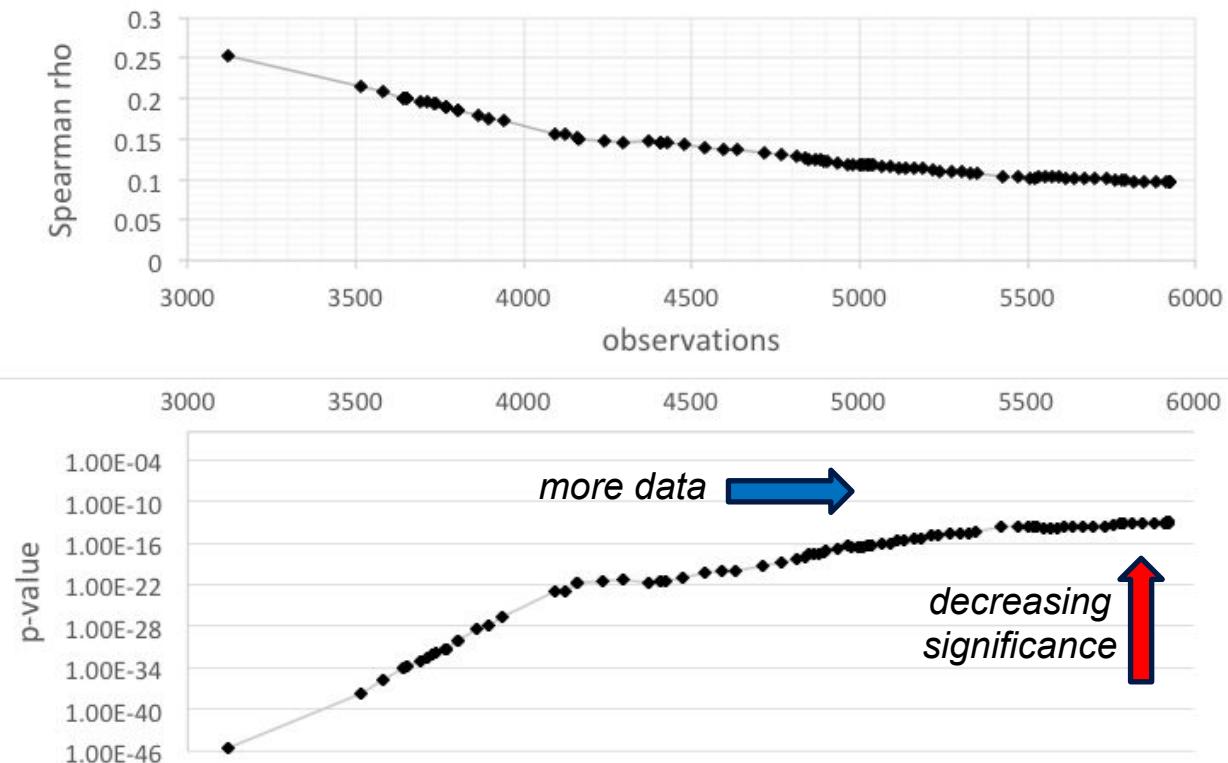
Showing 1 to 10 of 1,104 entries

Previous 1 2 3 4 5 ... 111 Next

Derive knowledge robustness statistics

ISB wellness data set, September 2017 correlations

- Glycohemoglobin A1C (LOINC 4548-4) and neutrophils (LOINC 26507-4): N=4970, rho=0.091, p-value=1.60E-10
- Glycohemoglobin A1C (LOINC 4548-4) and basophils (LOINC 706-2): N=3121, rho=0.252, p-value=2.48E-46



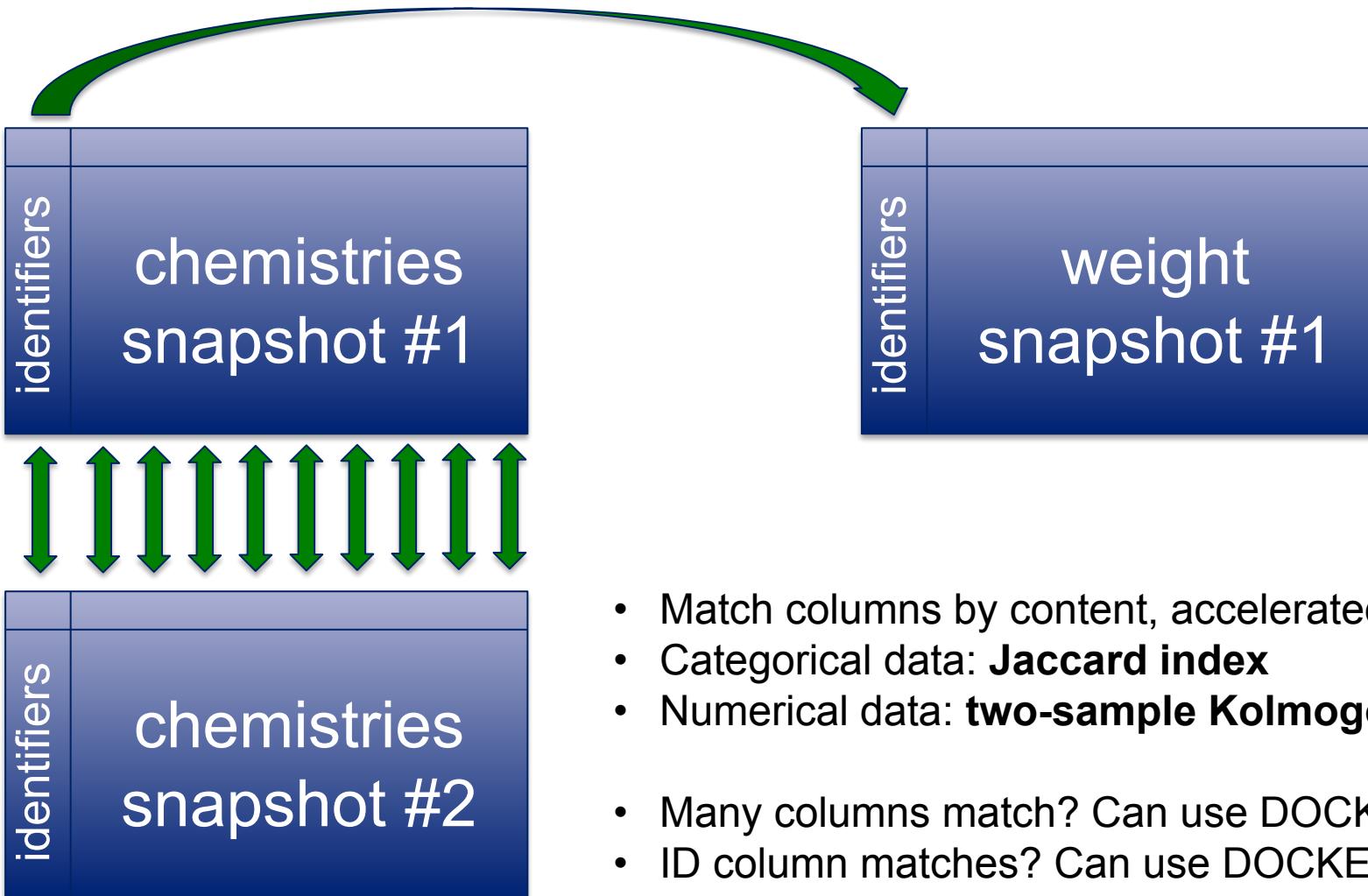
Confidence metric of knowledge robustness

- support for expert curation and validation
- metric of quality and validity of results returned to ARAs



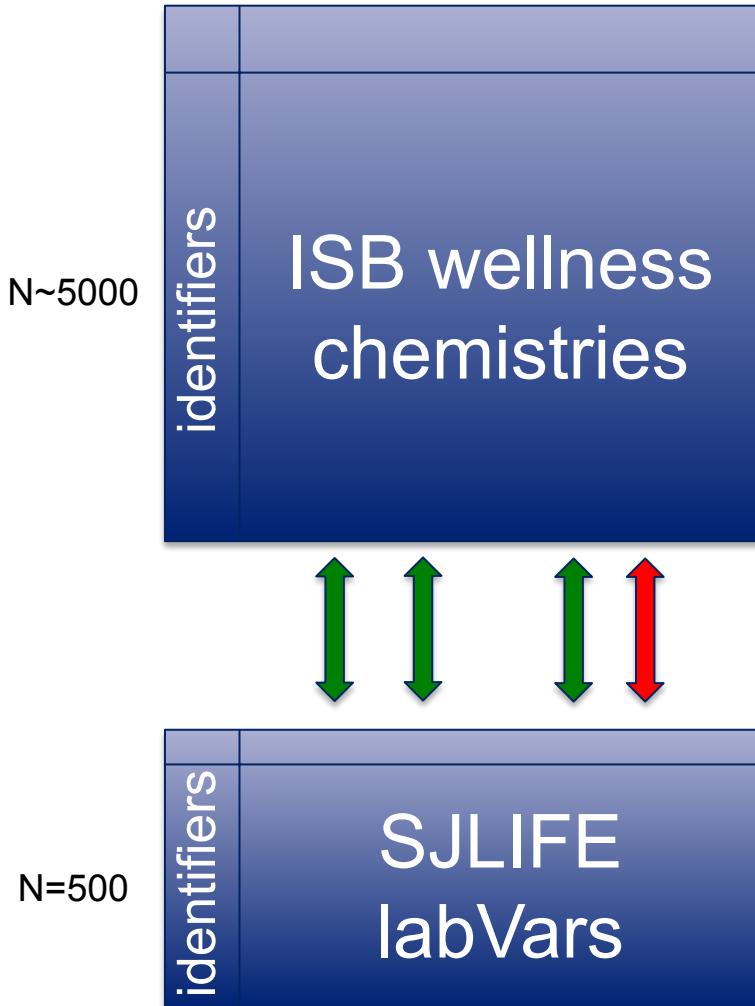
DOCKET Match

Compare dockets to find what they share in common



DOCKET Match

Can also match content from different data sets



Wellness “chemistries” vs. SJLIFE “labVars”:

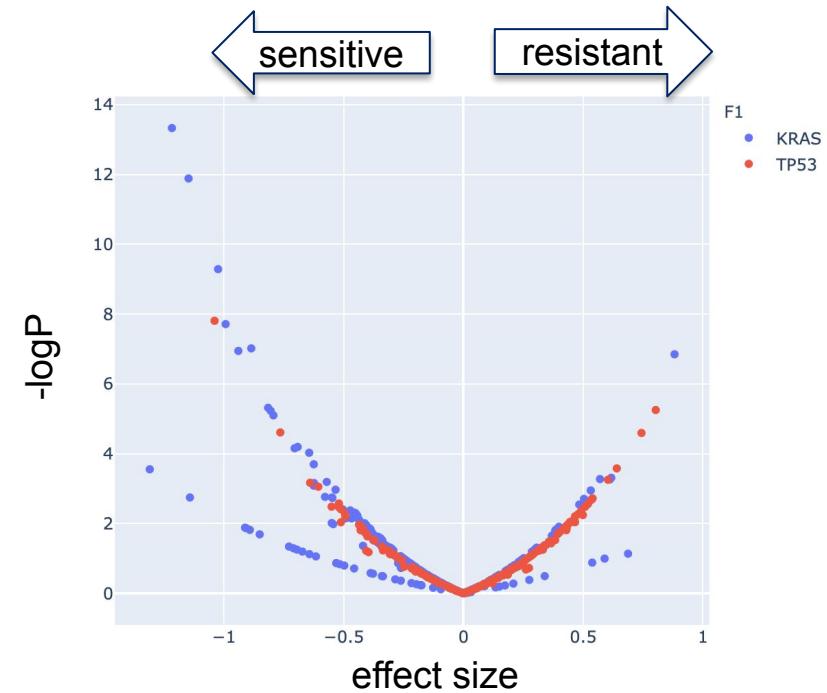
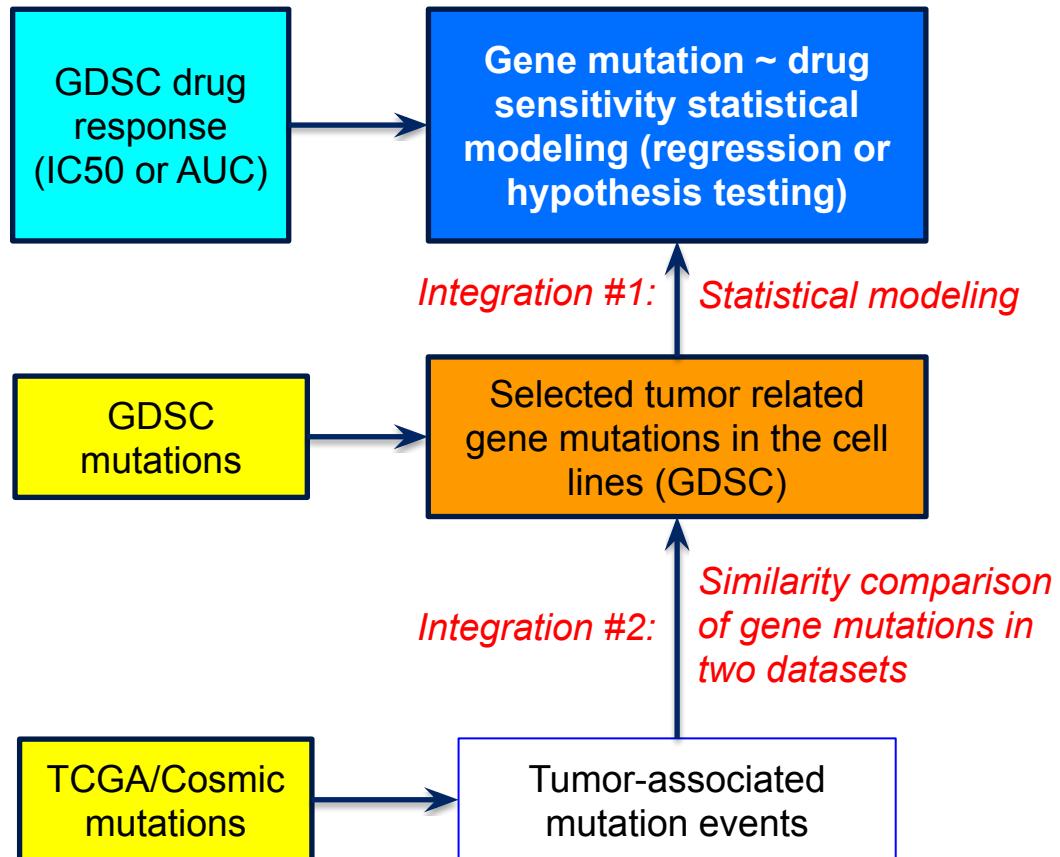
- OK “MCH” (785-6) ☐☐ “labs_numeric.MCH”
- OK “HEMOGLOBIN” (718-7) ☐☐ “labs_numeric.Hgb”
- info gain “RED CELL COUNT” (789-8) ☐☐ “labs_numeric.RBC_mm3”
- QC ? “HOMA-IR” ☐☐ “labs_numeric.Thyroid_Stimulating_Hormone_Analyzer”

Assistance to expert curation!



DOCKET Integrate

Statistical modeling (drug response vs. mutations, from multiple resources)



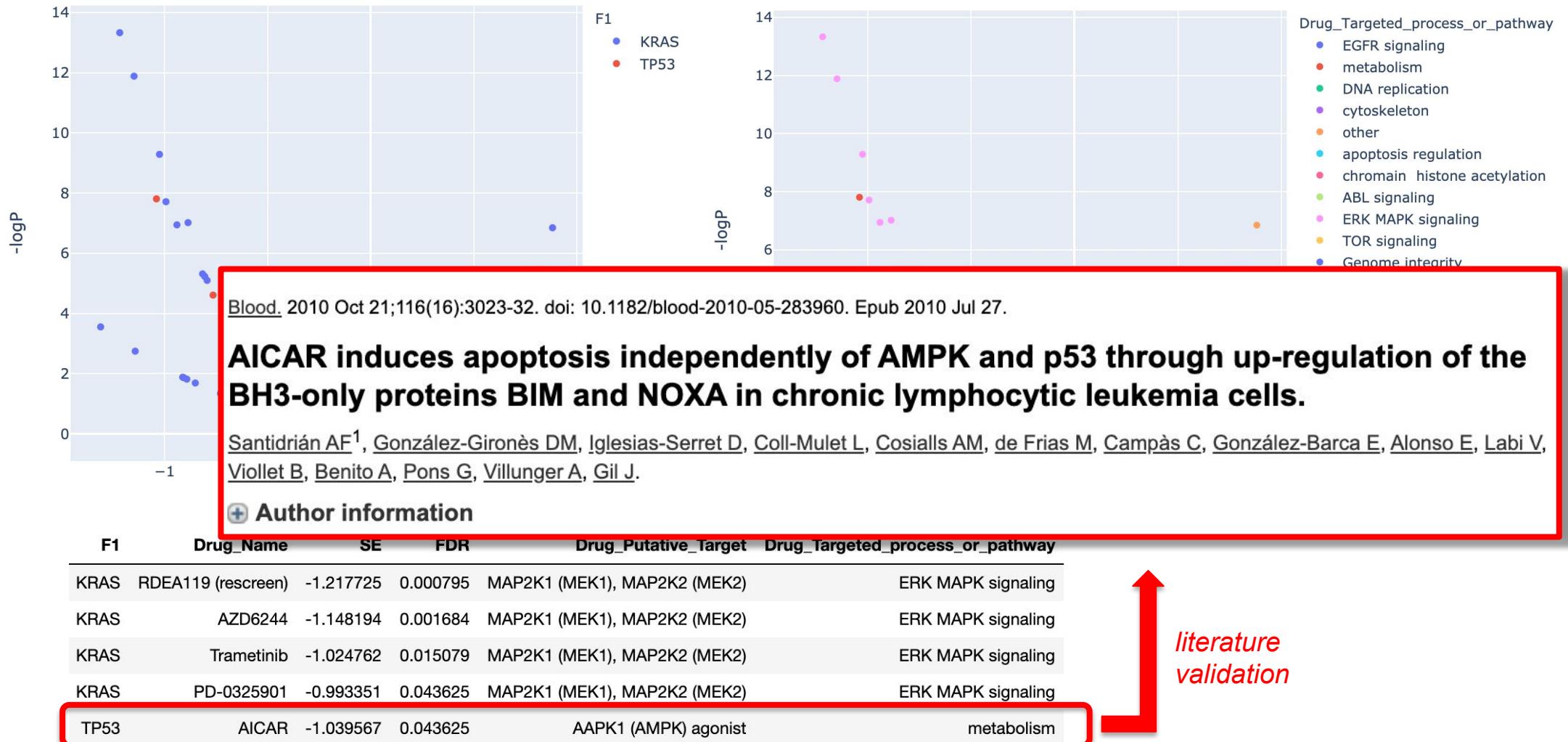
KP: Graphs and tables for users and ARAs

Biolink knowledge graph:

Disease --- Gene --- Drug (association: drug sensitivity)



Annotate knowledge graphs with external resources



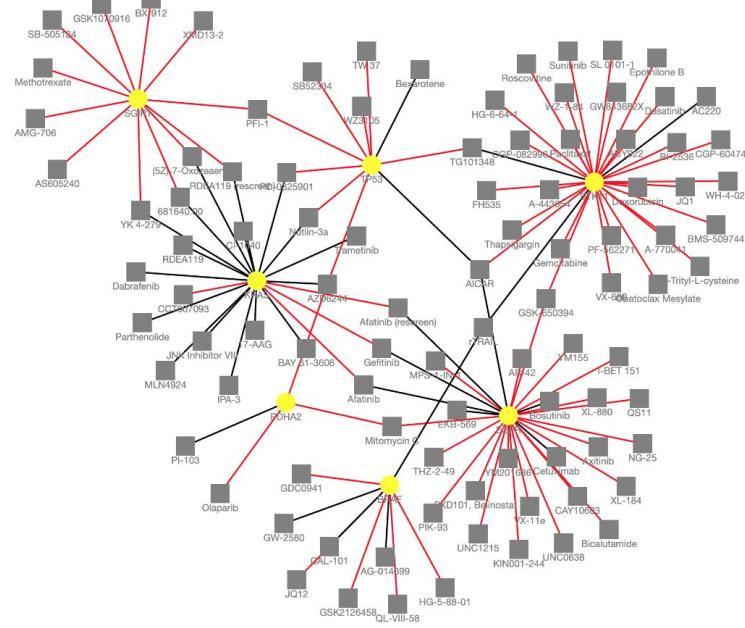
Enrichment of drug sensitivity or resistance:

Sensitivity pairs:

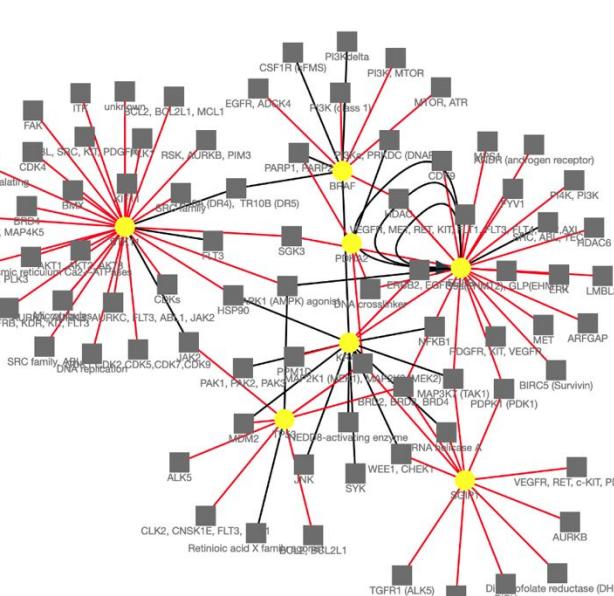
	Gene	mol_action	p	Resistance pairs:
0	KRAS	ERK MAPK signaling	0.000150	0 KRAS EGFR signaling 0.003023
1	TP53	metabolism	0.047059	

Visualization for knowledge graphs

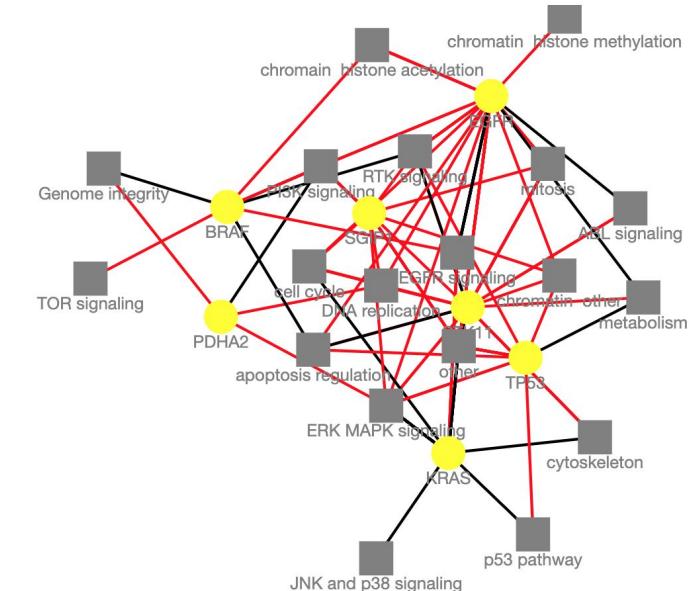
Association between drug and gene mutations



Association between Drug targets and gene mutations



Association between Targeted_process_or_pathway and gene mutations



Evidence level: GDSC data, LUAD, sample size: 64, P < 0.05, t-test

Provide visualizations of graphs



Milestones



M1.1 Ongoing support of Big GIM

- Migration to new Google Cloud Platform, connected with Translator users
- Deployed on NCATS server, BDD tests 91% code coverage



M1.2 Prototype DOCKET Overview

- Implemented **Study** (nextflow pipeline), **Overview** (interactive reports)
- Automatically studied wellness and cancer multiomics data
- Yielded visualizations and knowledge graphs



M1.3 Prototype DOCKET Compare

- Implemented **Match** (data types) and **Compare** (versions)
- Demonstrated ability to quantify knowledge stability



M1.4 Prototype DOCKET Integrate

- Implemented **Merge** (join tables with human expert validation)
- Implemented **Integrate** (statistical analysis, Jupyter notebooks)
- Demonstrated integrative analysis of cancer data



M1.5 Gather input on DOCKET features

- Designed survey, followed up with respondents
- CI, [C]D: Containerized DOCKET, established BDD test framework



M1.6 Collaborate on Reasoner Standard API

Deliverables

DOCKET code, with Docker container

github.com/PriceLab/DOCKET

Instructions for running prototypes

drive.google.com/drive/u/1/folders/19CT2bu1kzVnXgORhgIQijJd7x8O8Ez6D

Translator Drive > Knowledge Providers

Try-it-yourself document

Sample data

Big GIM: Gene Interaction Miner, and

Big CLAM: Cell Line Association Miner:

biggim.ncats.io/api

github.com/PriceLab/translator-bigquery-api



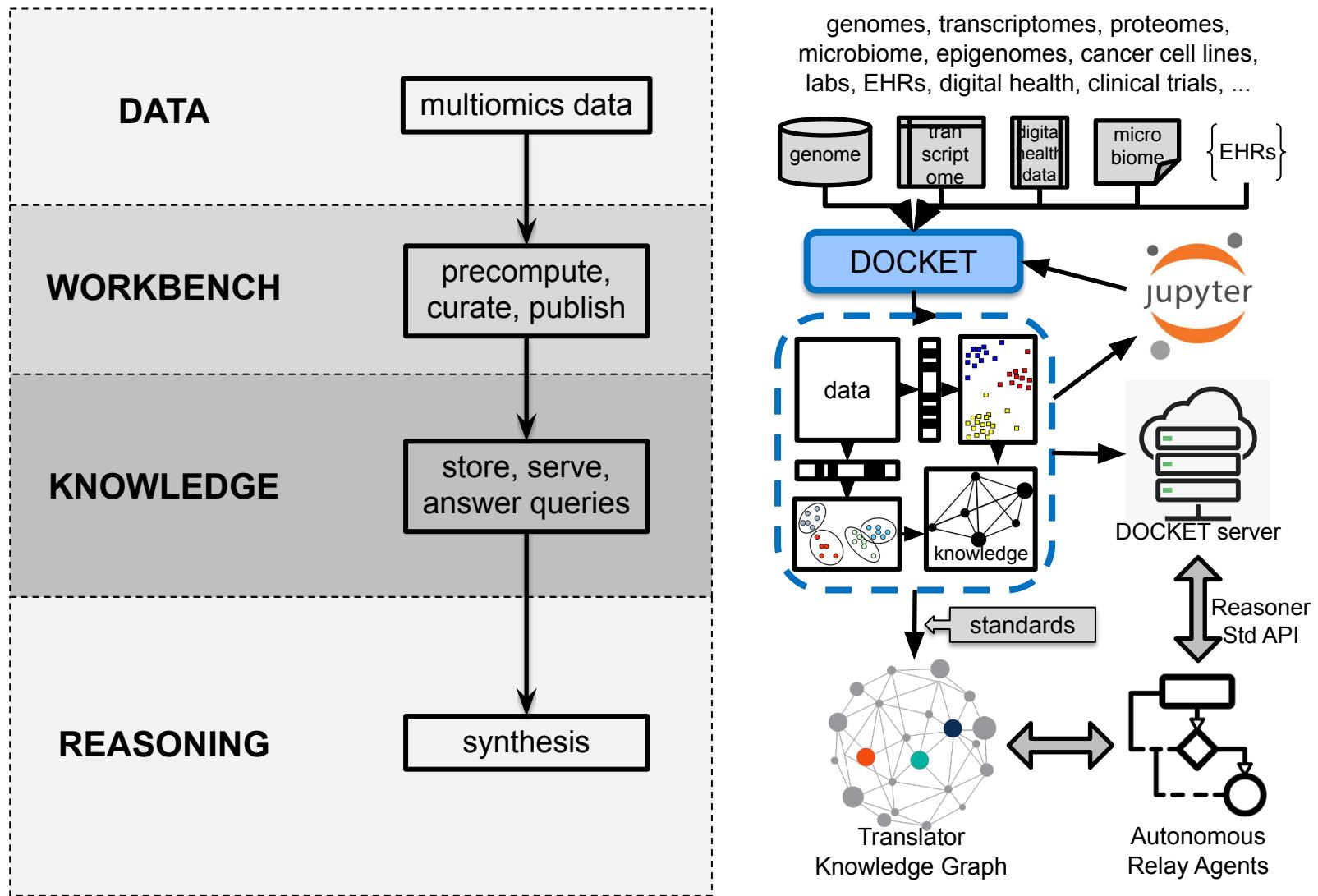
DEMO TIME

EXTRA SLIDES

Survey Highlights

High priority scenarios and features

Scenarios	
Allowing subject matter experts to curate and validate knowledge without needing to run code.	100%
Coordinating parallel workflow across groups of people curating and validating knowledge	83%
Features	
Support for rapid exploration, analysis, visualization and curation of raw data from heterogeneous locations, storage formats and versions.	100%
Automatically generated descriptive statistics and visualization to support understanding of data structure, content, missing data, and outliers.	100%
Automatically detected suggestions for observations in different data sources that are likely to be referring to the same concept.	100%
Automatically generated machine learning analyses, such as dimensionality reduction, clustering and enrichment.	100%
Automatically detected suggestions for mapping to Biolink entities and associations.	83%
A user interface for rapid human validation and curation of which correlations should be added as associations into your knowledge graph.	83%
Ability to incorporate existing provenance data from one or more data sets into your KP.	83%
Automatic generation of provenance information when you merge data from different sets, transform data, map to biolink entities and associations.	83%
Ability to ingest your graph and metadata and automatically deploy a knowledge graph.	83%



DOCKET features and information flow

