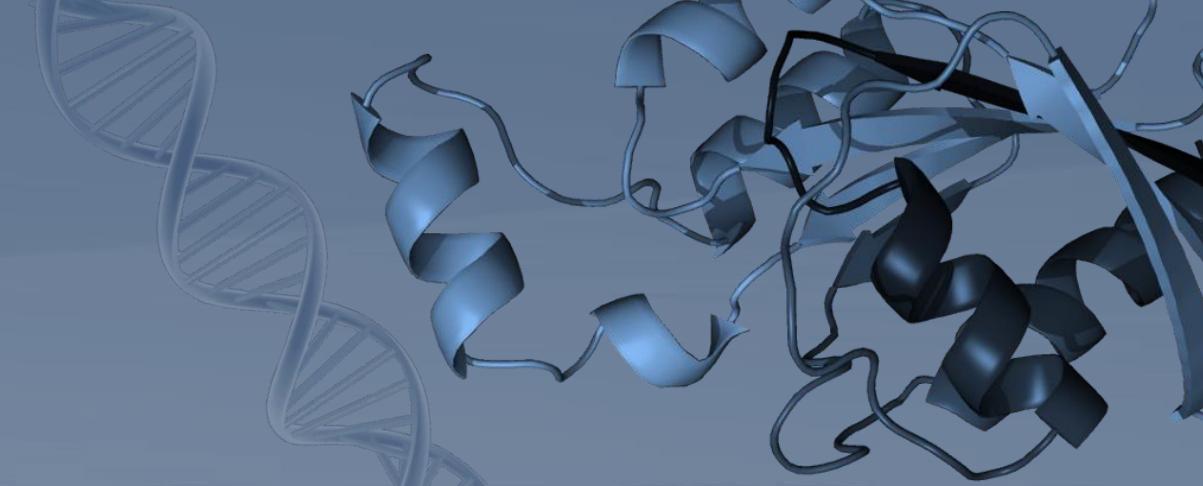




OFFICE OF CANCER CLINICAL  
PROTEOMICS RESEARCH



# Proteogenomic Approaches

**Chris Kinsinger, PhD**

**Program Director**

Office of Cancer Clinical Proteomics Research

Bring Your Own Bioinformatics

November 2019

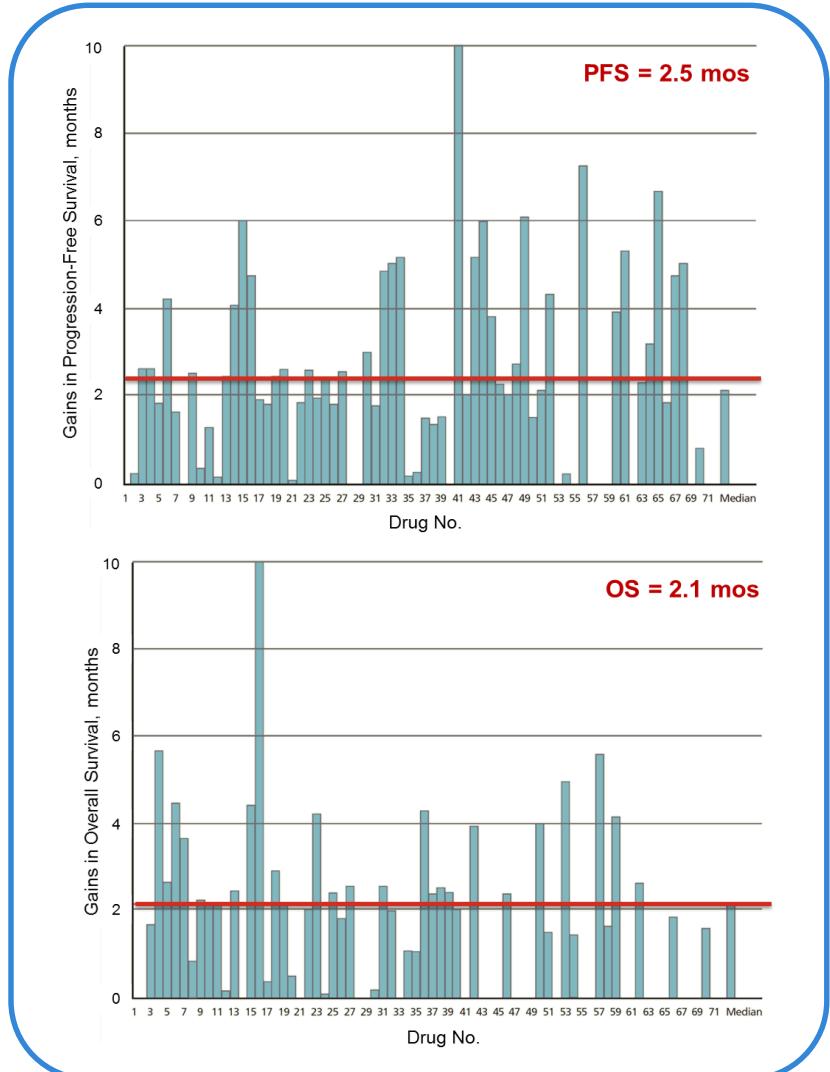


# Outline

- Proteogenomic approaches
  - Sequence-centric
  - Proteogenomic relationships
  - Integrative models
- Accessing the data
  - Cancer research data commons
  - Proteomic data commons
  - Genomic data commons
  - The Cancer Imaging Archive

# Drugs approved by FDA for advanced cancer

NIH NATIONAL CANCER INSTITUTE



Gains in overall survival  
for the 71 drugs approved  
by FDA from 2002 to  
2014 for advanced cancer

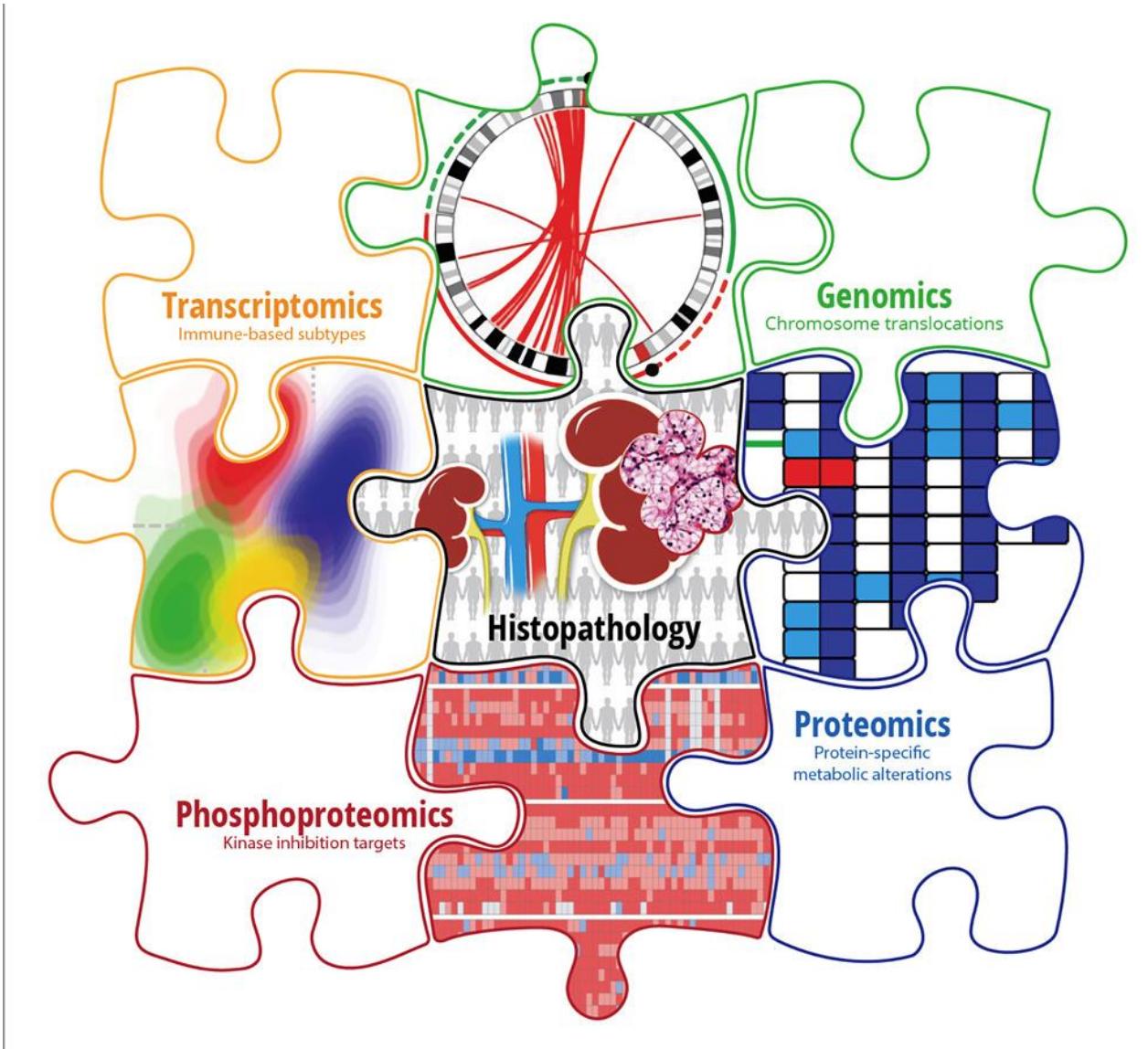


PFS: 2.5 mos

OS: 2.1 mos

# Proteogenomics

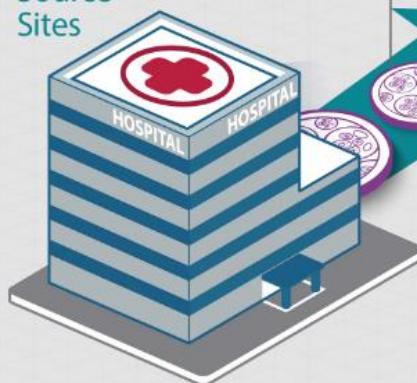
NIH NATIONAL CANCER INSTITUTE



# CPTAC PIPELINE

Tumor Characterization Program (Proteogenomic Tumor Atlas)

Tissue Source Sites



Tissue samples & clinical data



Biospecimen Core Resource  
(Pathology / Molecular QC)



Clinical, pathology & metadata

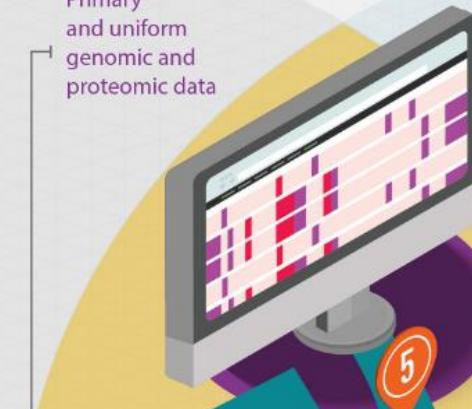
Protein, DNA, RNA analytes



Proteome Characterization Centers  
(DNA and RNA characterized in partnership with TCGA)

Primary genomic and proteomic data

Primary and uniform genomic and proteomic data



Proteogenomic Data Analysis Centers



CPTAC Data Coordinating Center

Primary genomic and proteomic data



Proteogenomic Translational Research Centers  
(Understanding drug responses and resistance)



Harmonized protected & public data

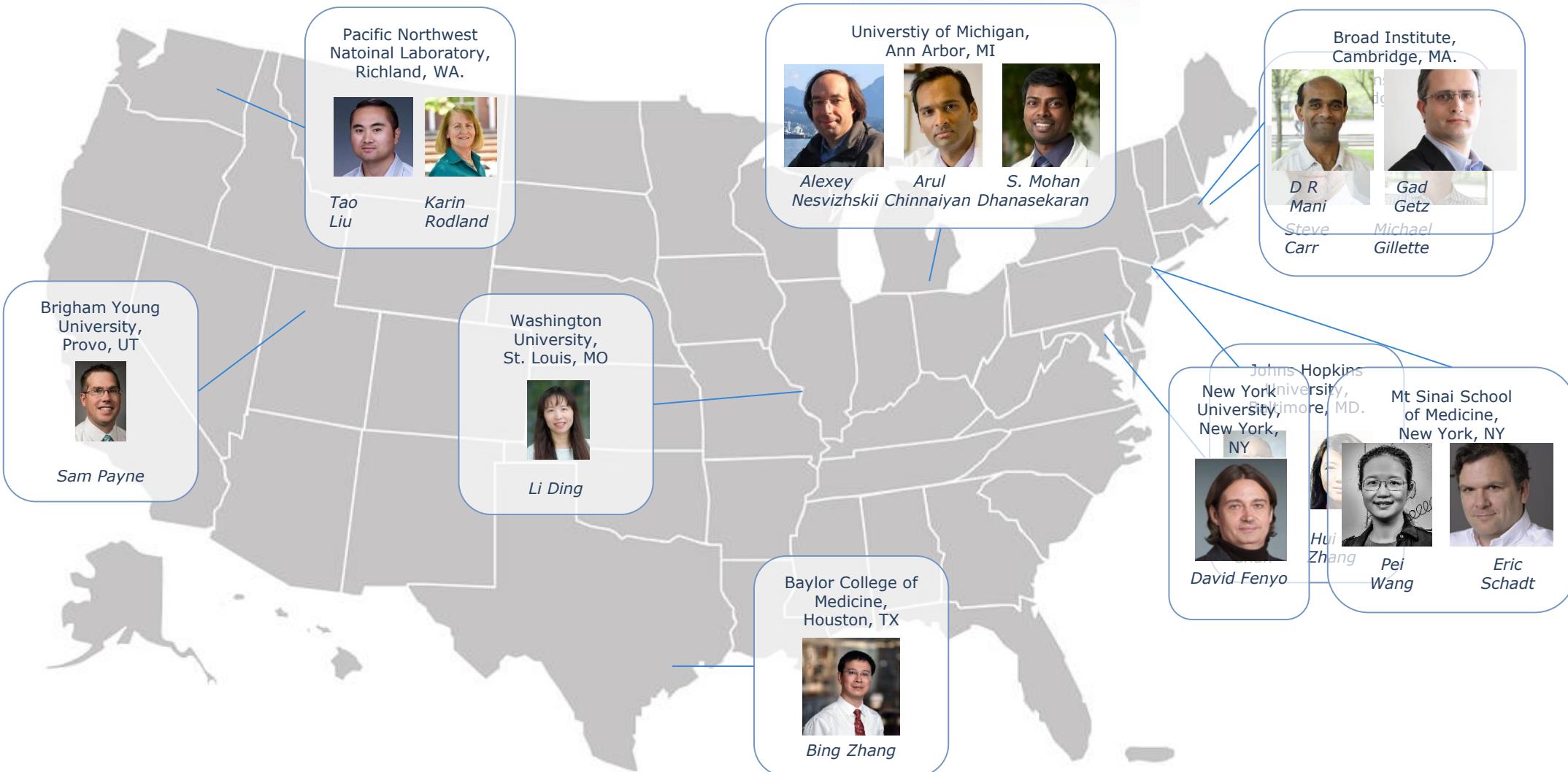


Research Community

Translational Research Program (in partnership with NCI-sponsored clinical trials)

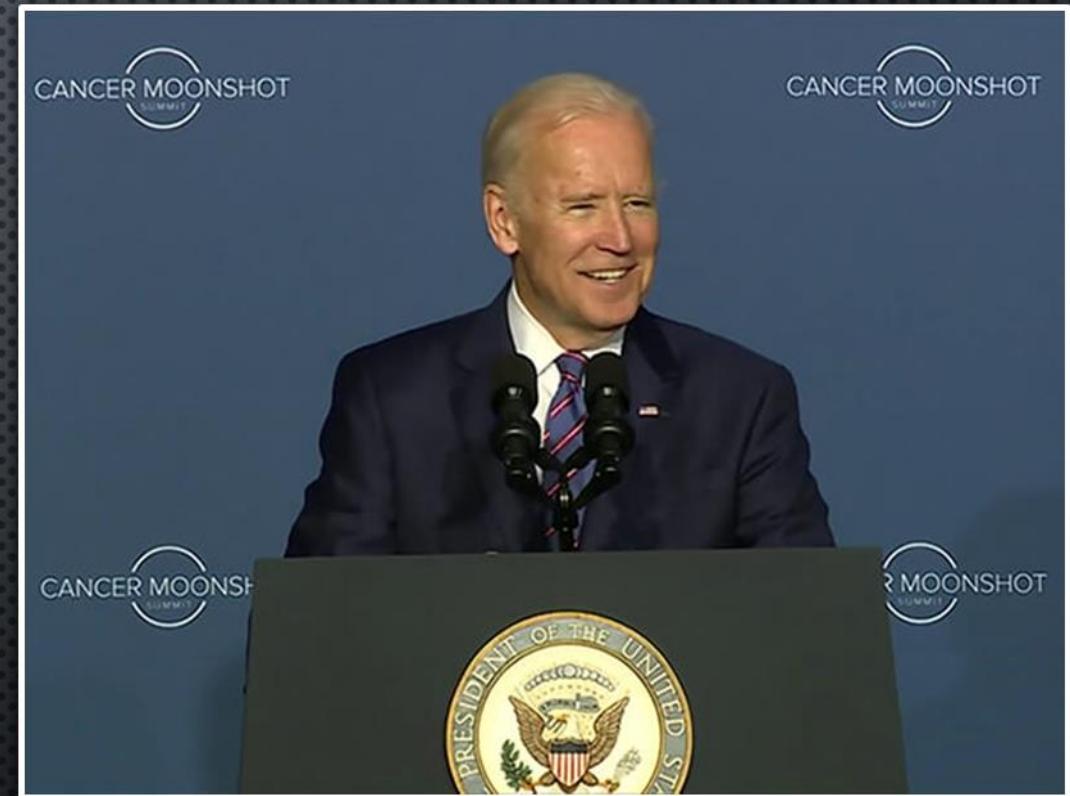
# The players

NIH NATIONAL CANCER INSTITUTE



# White House Cancer Moonshot

**GOAL:** Make 10 yrs  
of progress against  
cancer in 5 yrs.



Joe Biden  
Vice President of the United States

# APOLLO Clinical Network

## (Applied Proteogenomics Organizational Learning and Outcomes)



APOLLO Inaugural Leadership Meeting  
29 August 2016



Col. Craig Shriver  
(DoD)

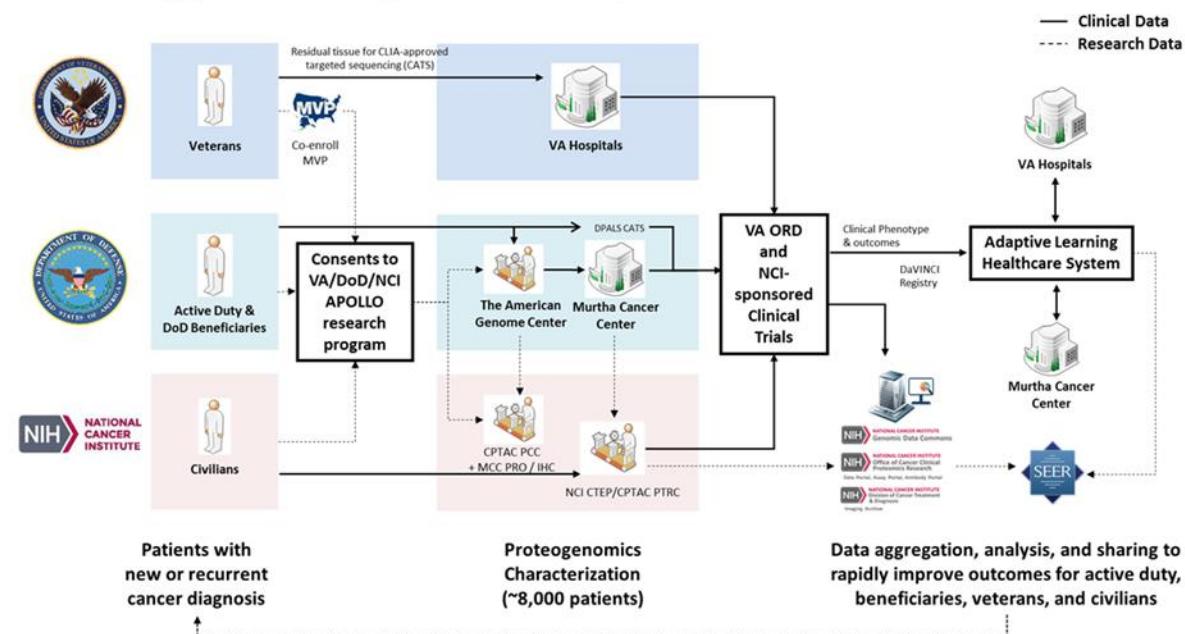


Jennifer Lee  
(VA)



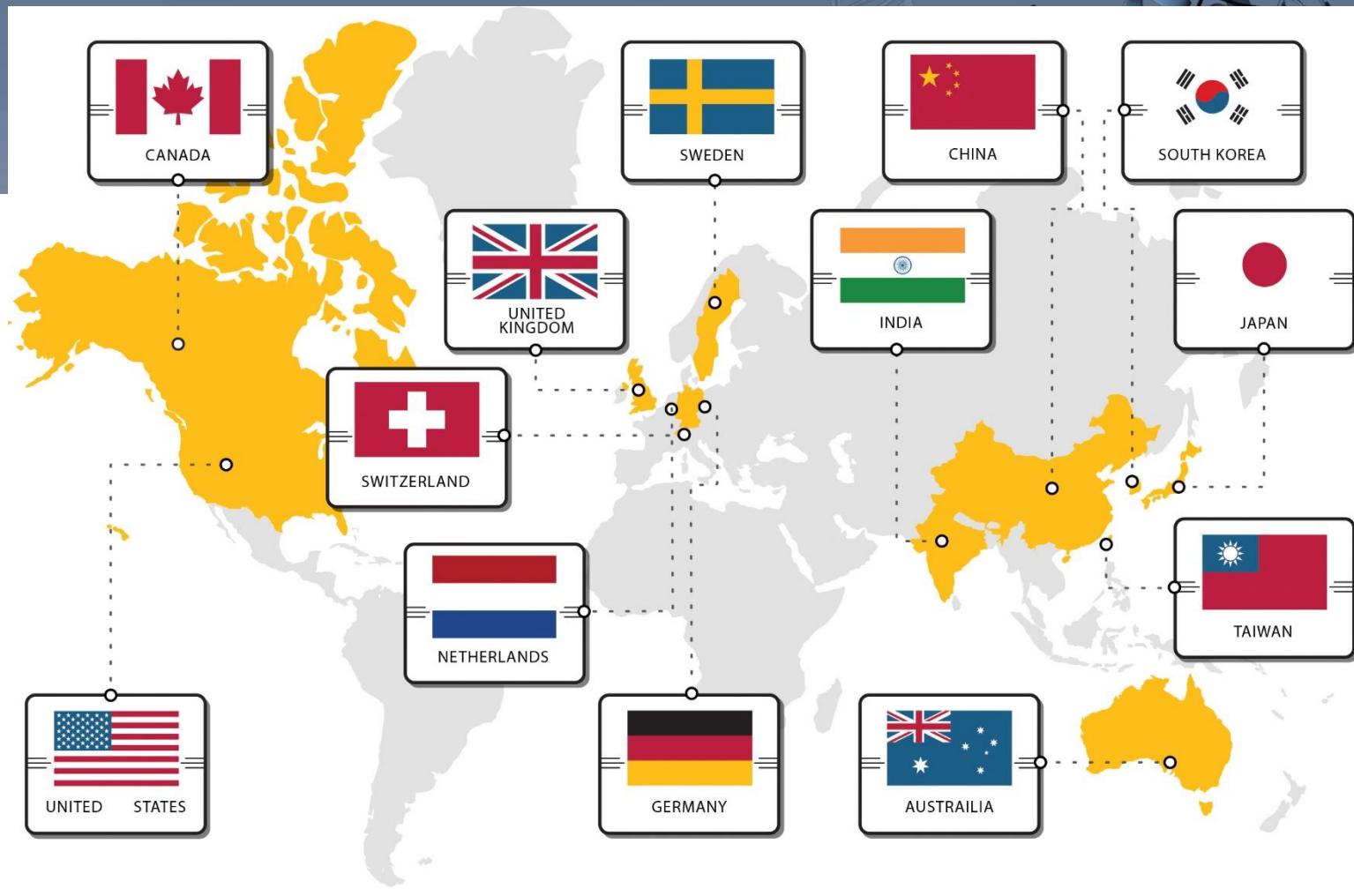
Henry Rodriguez  
(NCI)

- **Tri-lateral Collaboration:** NCI, Department of Defense, Veterans Administration
- Nation's first healthcare system where cancer patients routinely screened for genomic abnormalities and proteomic information to match tumor types to targeted therapies



13 countries  
from  
33 institutions

13 cancer types selected



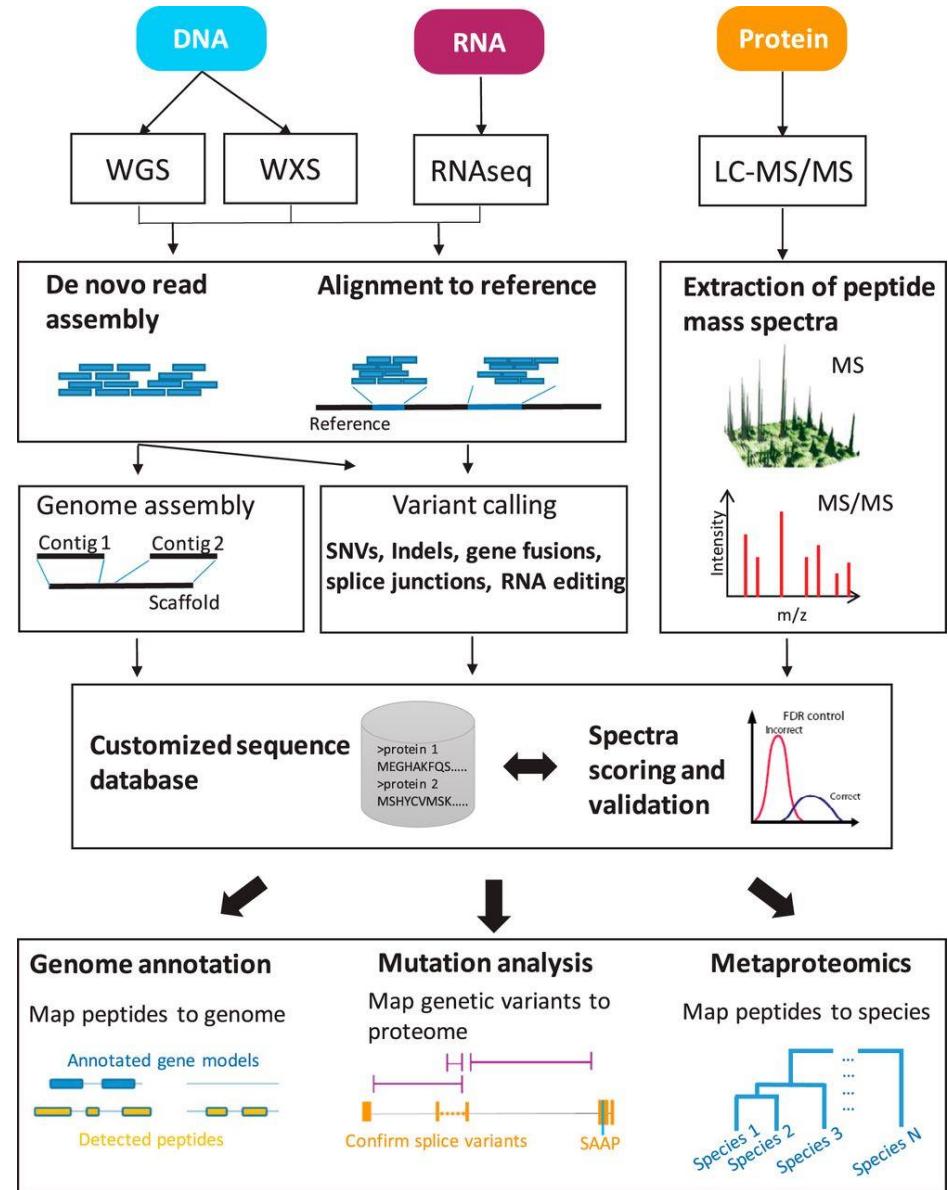
ICPC

INTERNATIONAL CANCER PROTEOGENOME CONSORTIUM

# Sequence-centric proteogenomics

NIH NATIONAL CANCER INSTITUTE

- QUILTS
- FDR caution



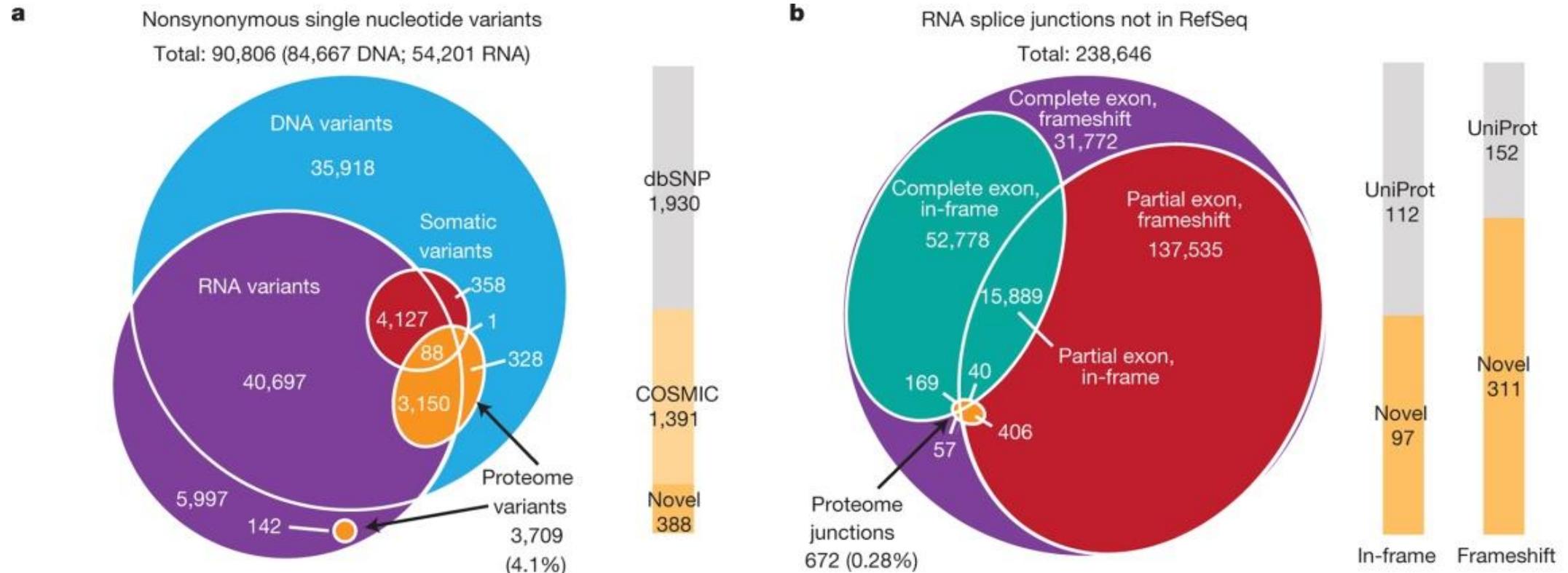
Ruggles, et al. MCP 16: 10.1074/mcp.MR117.000024, 959, 2017.

QUILTS – Ruggles et al. MCP 2017, 15.3 p. 1060

FDR cautions - Nesvizhskii, AI. "Proteogenomics: concepts, applications, and computational strategies." Nature Methods. 2014, 11, 1114.

# Proteogenomic mutation detection in breast cancer

NIH NATIONAL CANCER INSTITUTE



- Most single amino acid variants previously reported
- Few splice junctions detected, but many are novel

# Analysis of proteogenomic relationships

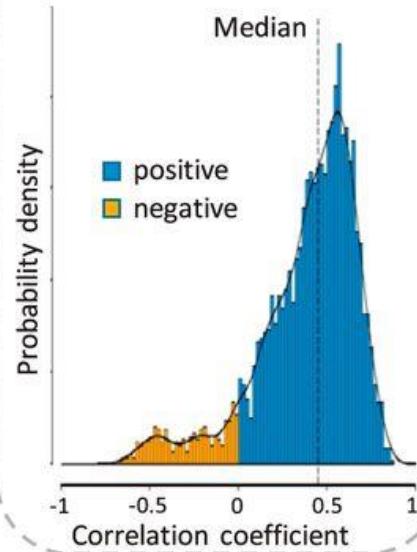
NIH NATIONAL CANCER INSTITUTE

A

mRNA-Protein correlation



mRNA to protein expression

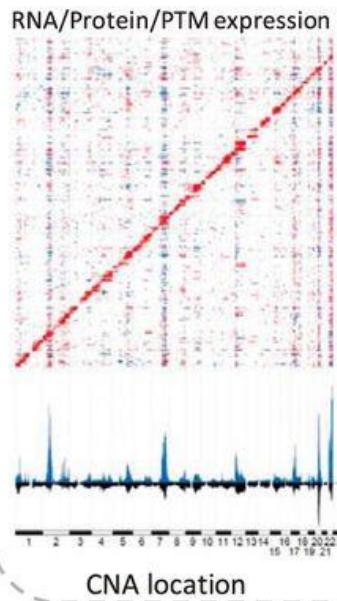


B

CNA



CNA to expression data

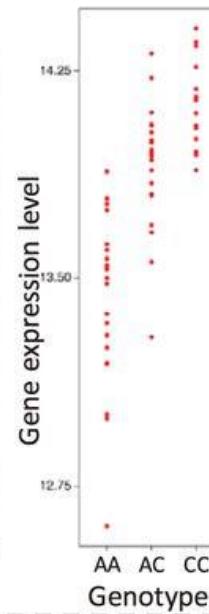


Regulatory correlation

eQTL



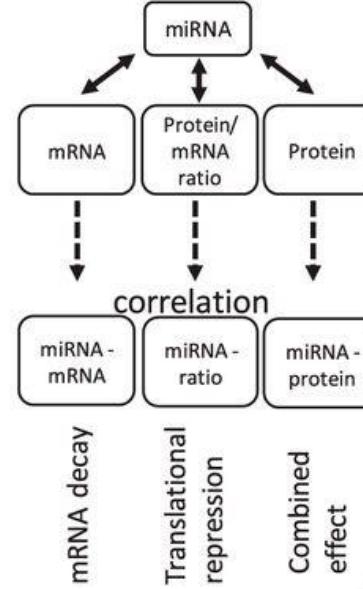
DNA variation to expression data



miRNA



miRNA regulation

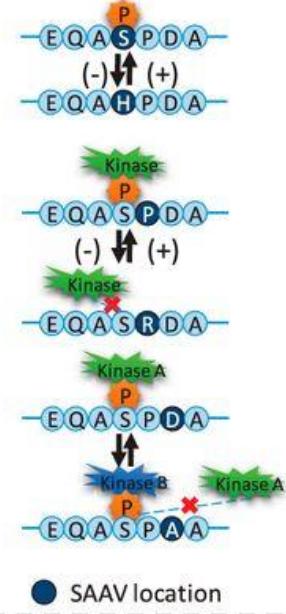


C

Mutations and signaling



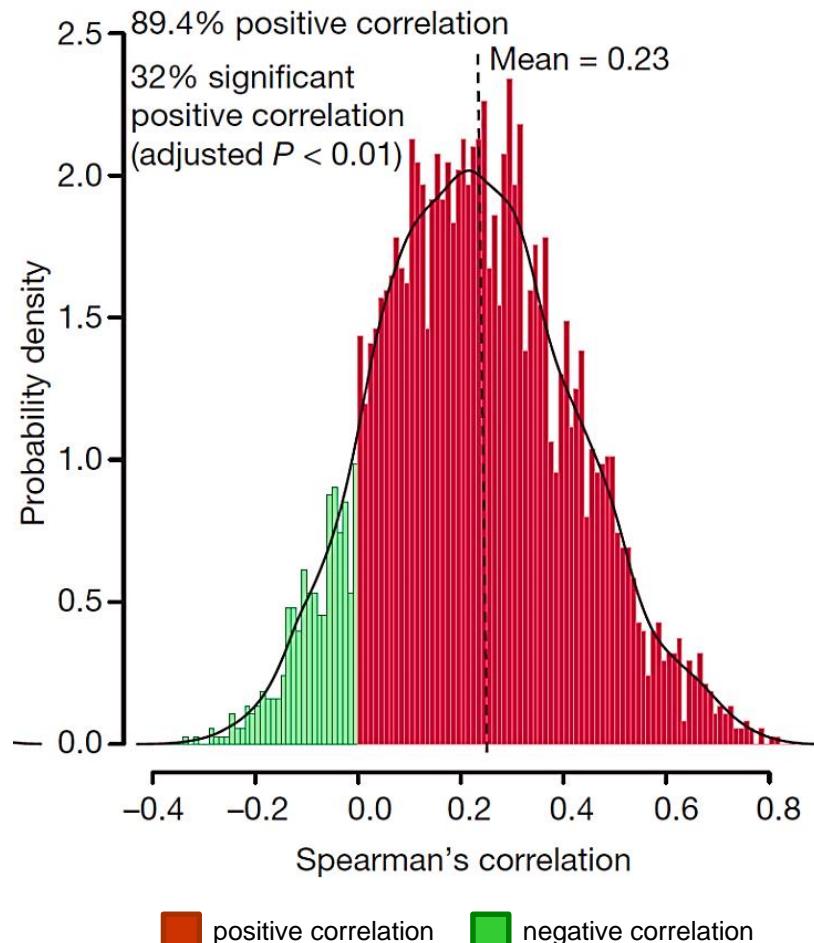
DNA variation to PTM sites



# Gene expression and protein abundance have fairly low correlation

NIH NATIONAL CANCER INSTITUTE

mRNA and protein abundance correlation for individual genes across all colon tumors (samples)



Mean Correlation:

- within 0.47
- across 0.23

- *Similar correlations have been shown for breast, ovarian, and gastric cancers*

# Cis and trans effects in ovarian cancer

NIH NATIONAL CANCER INSTITUTE

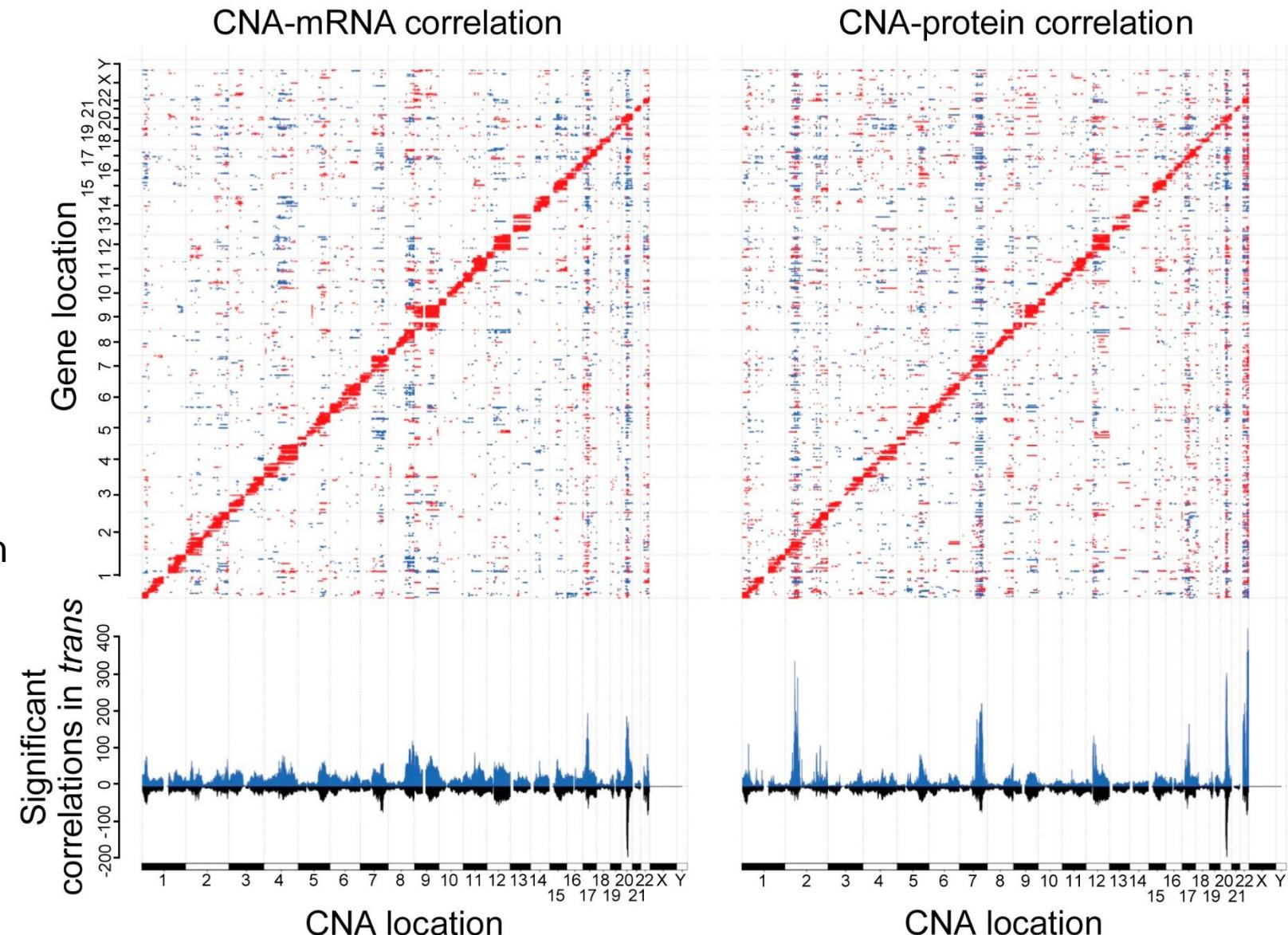
- **174 ovarian HGSC tumors**

- Selection criteria:
  - Overall Survival (OS)
  - Homologous Recombination Deficiency status (HRD)

- **5 proteomic subtypes**

(4 transcriptomic subtypes)

- mRNA subtypes translate at protein level
- New “stromal” subtype emerged



# Ovarian Cancer

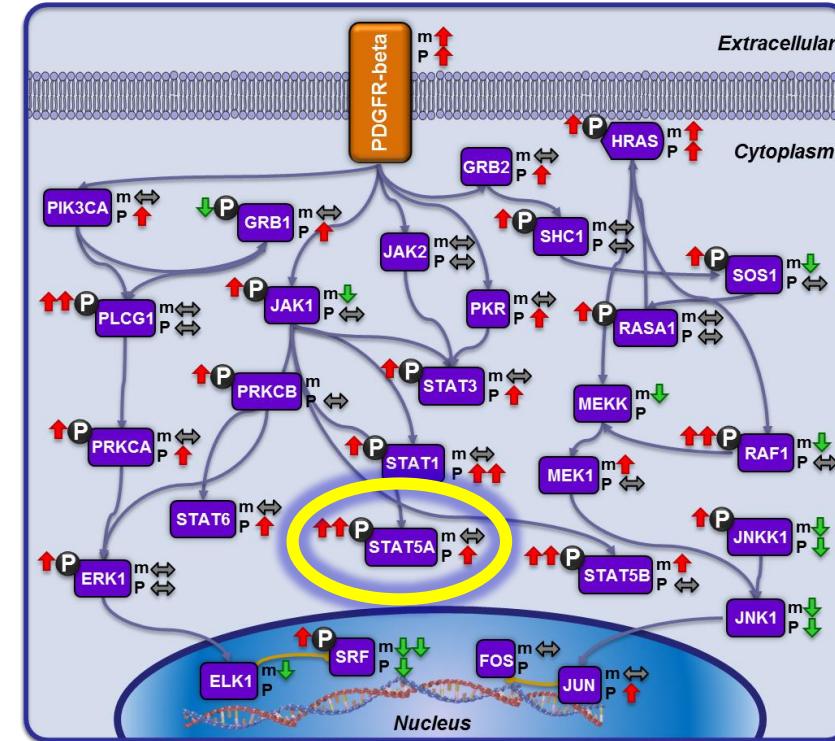
(pathway activation correlates with overall survival)

NIH NATIONAL CANCER INSTITUTE

## Network Data Exchange [NCI Pathway Interaction Database] (214 signaling pathways)

- Significantly upregulated pathways with short OS
  - Protein data ( $p<0.05$ )
  - Phosphorylation data ( $p<0.0001$ )
  - mRNA data ( $p<0.05$ )
- *Combining comprehensive proteomic, phosphoproteomic and transcriptomic analysis better elucidated the proteogenomic complexity of pathway activation not obtainable at the subtype level.*

PDGFR pathway upregulation in TCGA tumors with short OS



m = mRNA

P = protein abundance

P = phosphoprotein

↑ = upregulated

↑↑ = significantly upregulated

↓ = downregulated

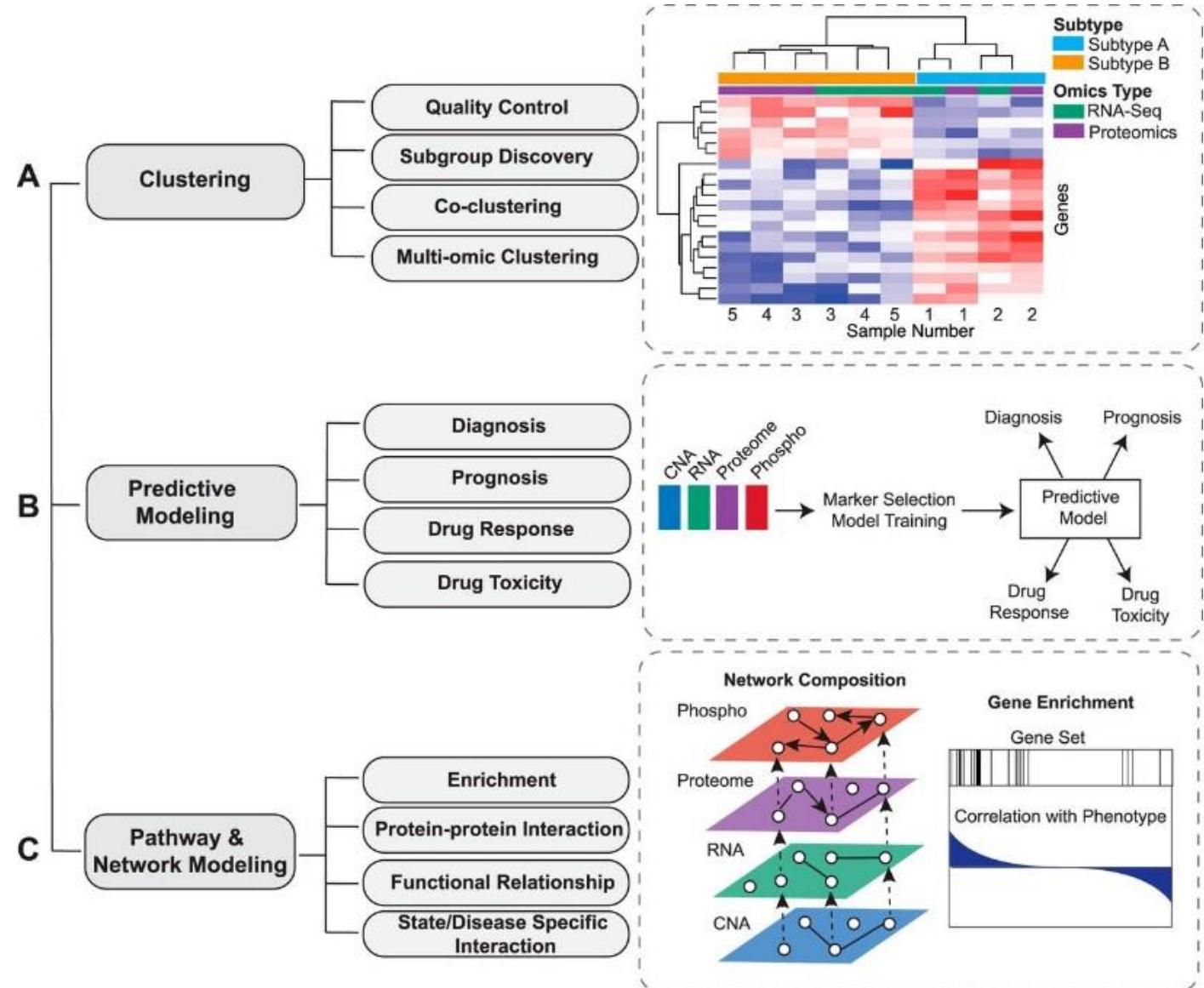
↓↓ = significantly downregulated

↔ = no difference

= not observed

# Integrative modeling of proteogenomic data

NIH NATIONAL CANCER INSTITUTE



# Colorectal Cancer: global protein abundance (new proteome subtypes identified)

NIH NATIONAL CANCER INSTITUTE

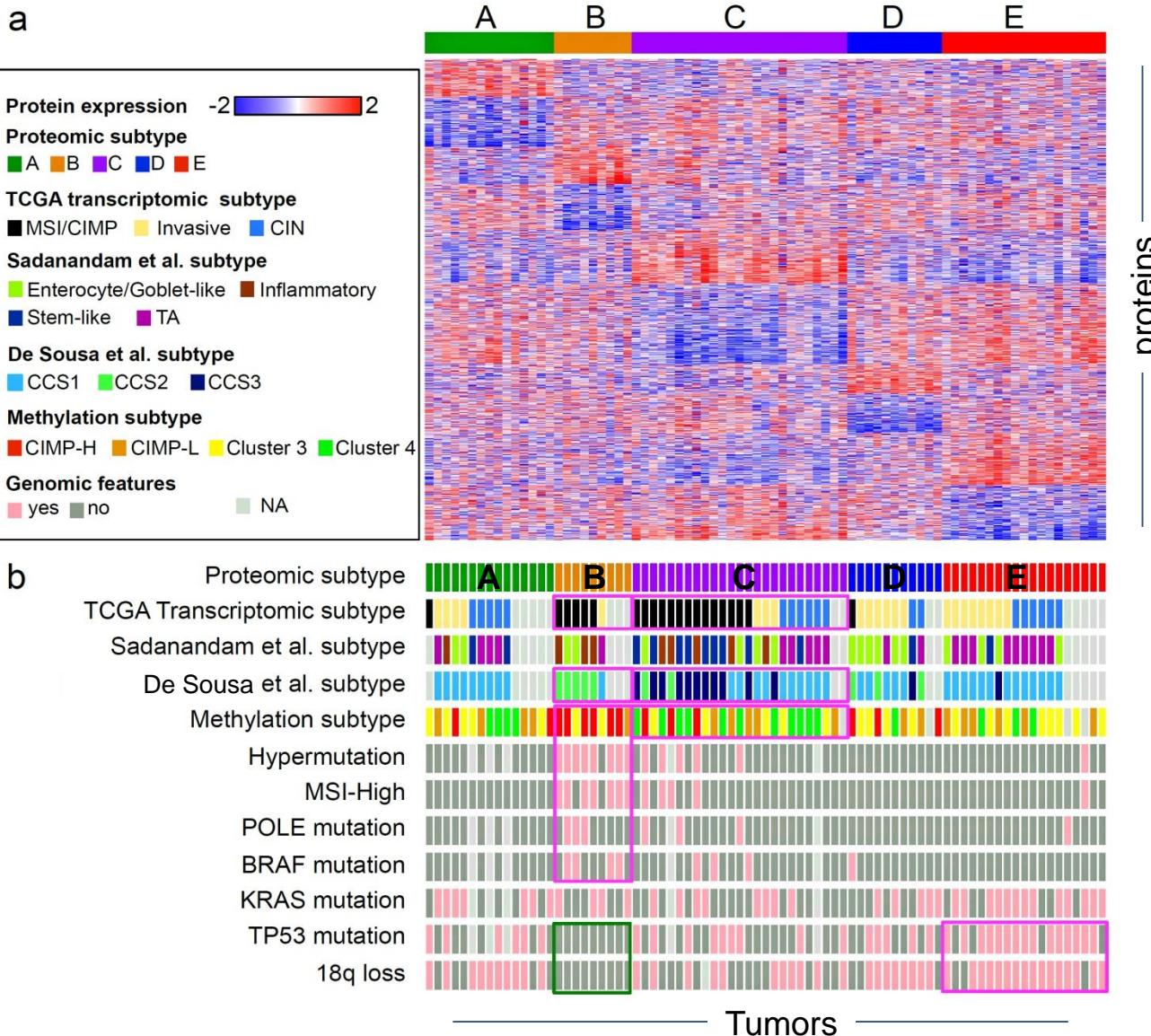
## Transcriptome Subtypes



## Proteome Subtypes



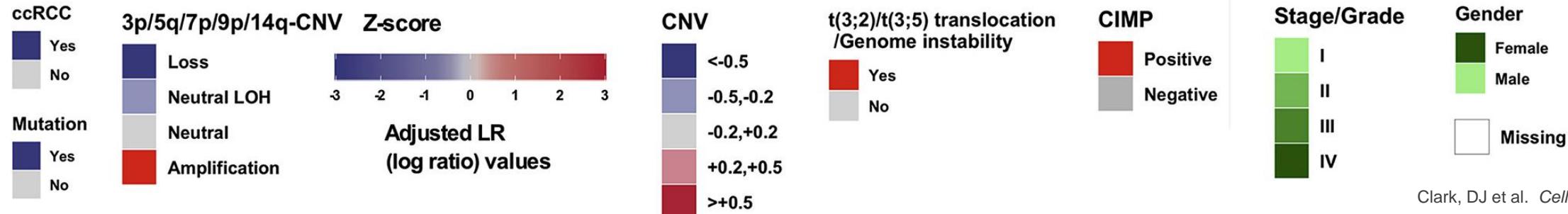
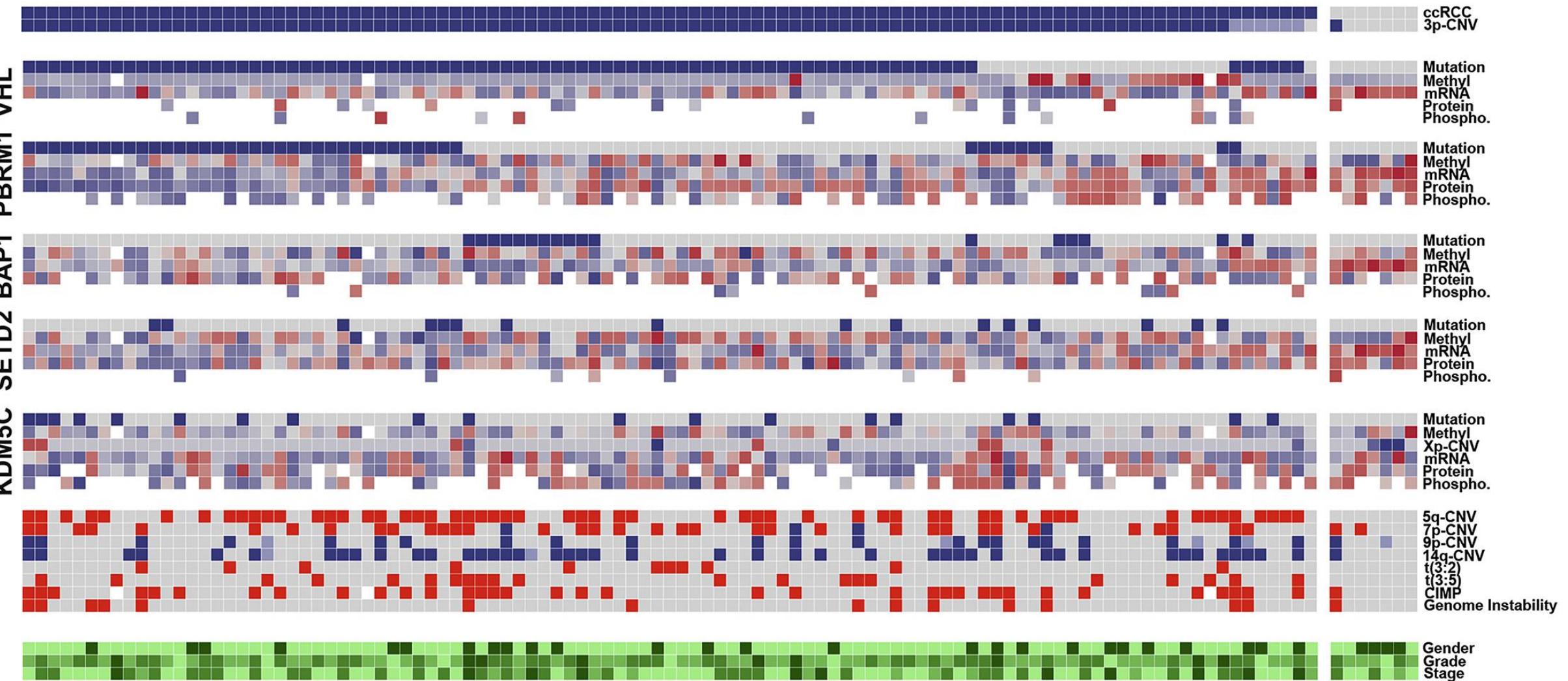
- MSI/CIMP transcriptome subtype split into two proteome subtypes
- Subtype C displayed protein network features characteristic of EMT, associated with rapid metastasis and overall poor survival



nature

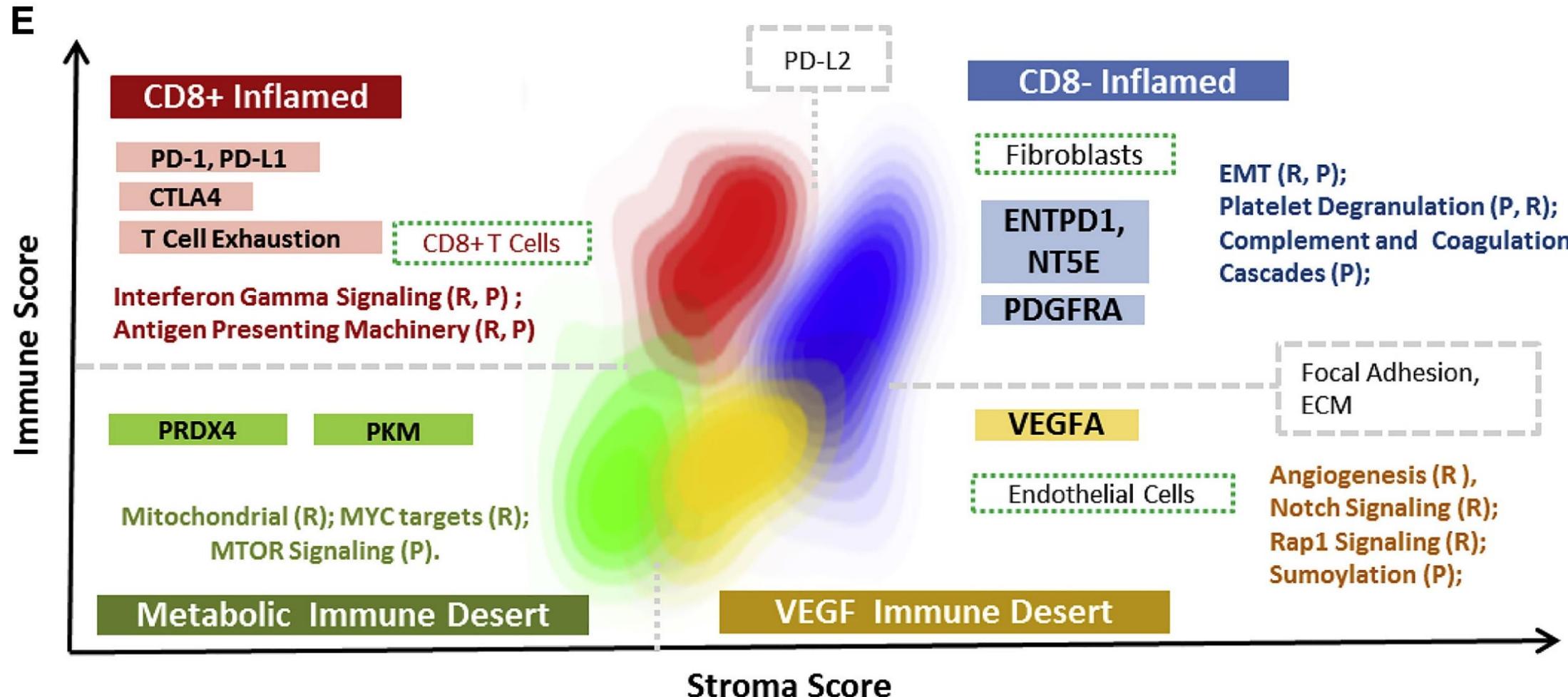
2014 Jul 20; doi:  
10.1038 / *Nature*  
13438

C



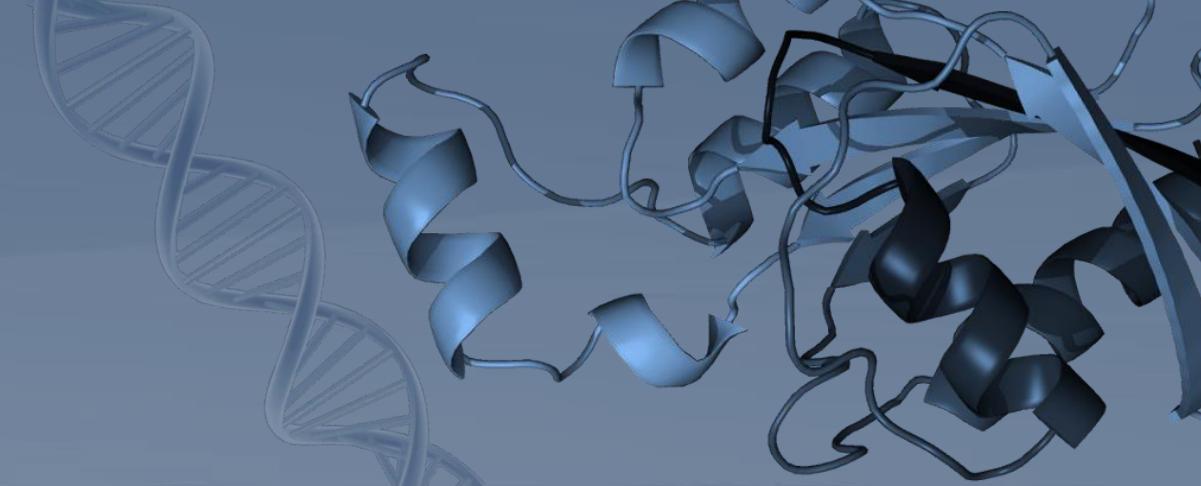
# Immune landscape of kidney cancer proteome

NIH NATIONAL CANCER INSTITUTE





OFFICE OF CANCER CLINICAL  
PROTEOMICS RESEARCH



# Accessing the data



NATIONAL  
CANCER  
INSTITUTE

# Two types of data repositories



- Repository by data type
  - Serves community specific to that data type
- Repository by program
  - Multiple types of data
  - Data generated by one program



# The Beau Biden Cancer Moonshot<sup>sm</sup>

NIH NATIONAL CANCER INSTITUTE

## Overarching goals

- Accelerate progress in cancer, including prevention & screening
  - From cutting edge basic research to wider uptake of standard of care
- Encourage greater cooperation and collaboration
  - Within and between academia, government, and private sector
- Enhance data sharing

## Blue Ribbon Panel – October, 2016

Recommendations include:

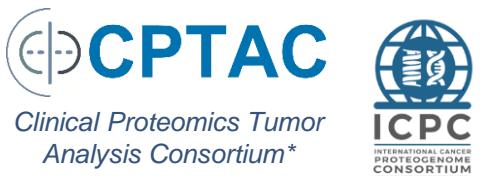
- Build a National Cancer Data Ecosystem
  - Enhanced cloud-computing platforms
  - Services that link disparate information, including clinical, image, and molecular data
  - Essential underlying data science infrastructure, standards, methods, and portals for the Cancer Data Ecosystem

# NCI Cancer Research Data Commons (CRDC)

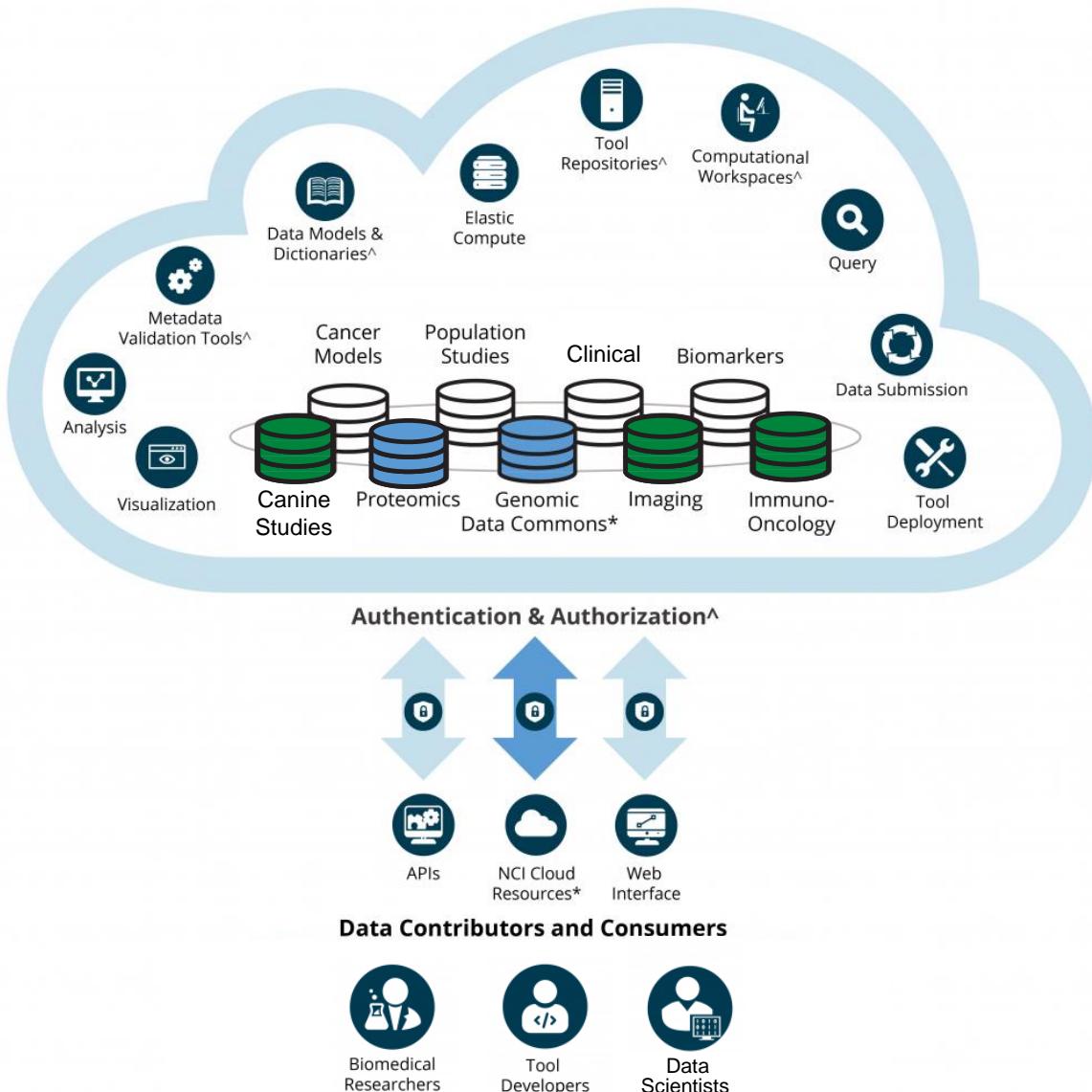
## Data Sources:



Canine Immuno-oncology studies



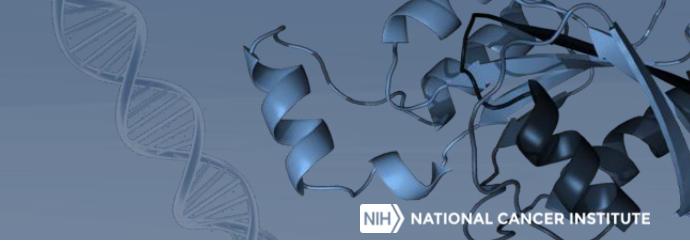
APOLLO Network  
A NCI-DoD-VA Proteogenomic Translational Initiative



## Components:

- Data Nodes
- Data Commons Framework
- Data Aggregators
- Cloud Resources
- APIs
- Elastic compute resources
- Portals
- Workspaces
- Analytic Tools
- Tool repositories

# Data Commons Framework Services



## Data Commons Core Framework Services

---

- Authentication and authorization
- Digital ID and associated metadata services



## Extensible, Interoperable Data Models

---

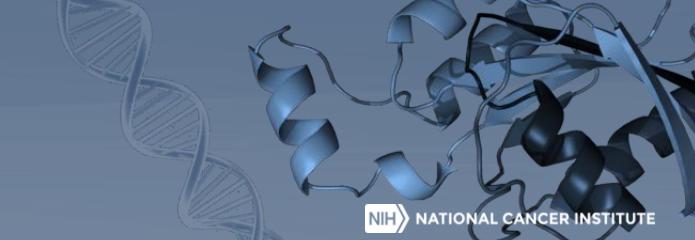
- Data models
- Data model based APIs

University of Chicago

---

- <https://gen3.org>

# NCI Cancer Genomics Cloud Pilots



Cloud Pilots provide:

- Access to large omic data sets without need to download
- Access to popular pipelines and visualization tools
- Ability for researchers to bring their own tools and pipelines to the data
- Ability for researchers to bring their own data and analyze in combination with existing genomic data
- Workspaces, for researchers to save and share their data and results of analyses

- Access and analyze 11,000 TCGA samples without having to download data
- Upload your own data for analysis

Data



- Perform large scale analysis using the elastic compute power of commercial cloud platforms

Compute



- dbGaP-authorized users can access controlled TCGA data
- Systems meet strict Federal security guidelines

Security



Democratize access to NCI-generated genomic and related data, and to create a cost-effective way to provide scalable computational capacity to the cancer research community.

# Proteomic Data Commons: High Level Goals

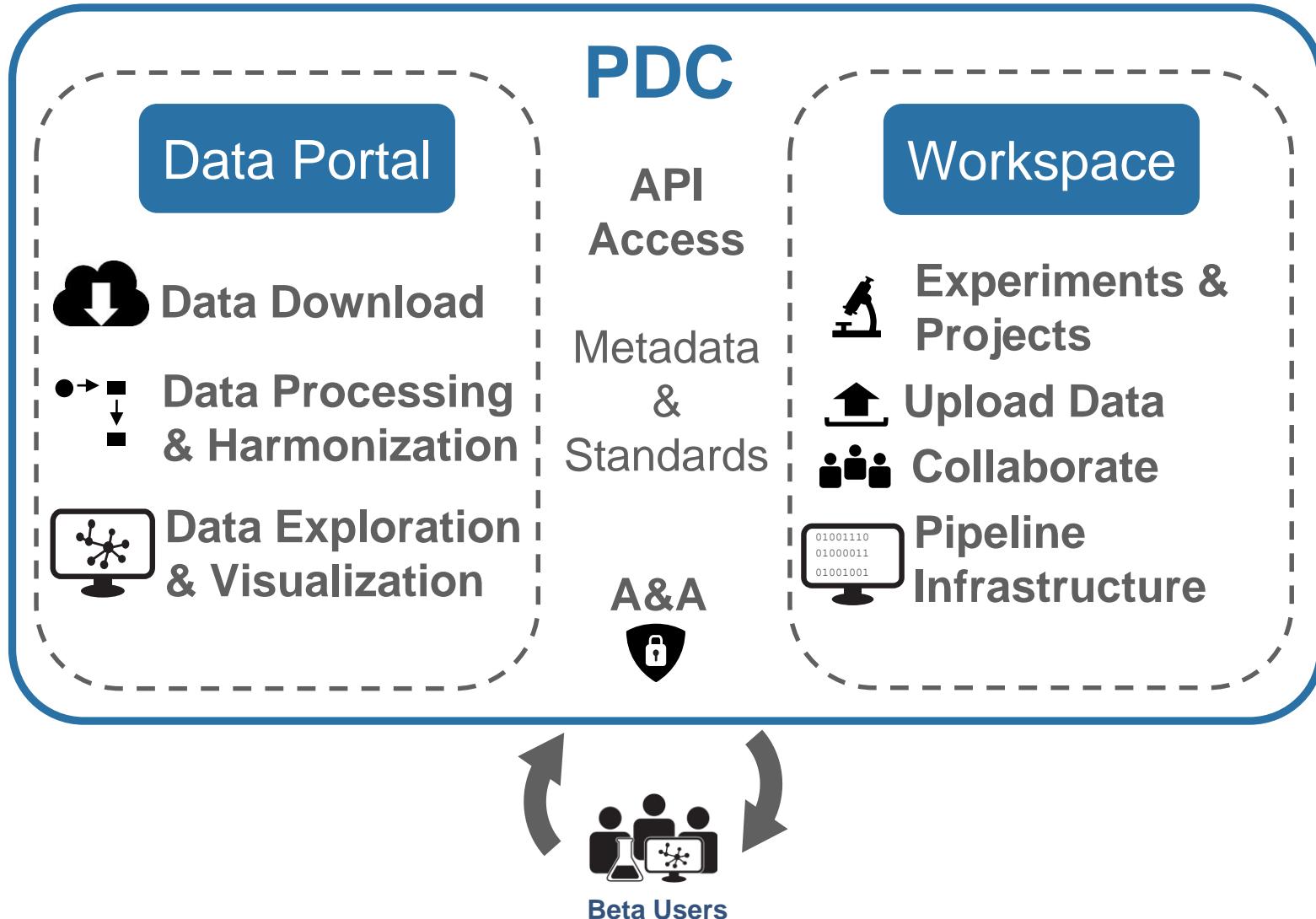


- 1 Connect proteomic data to other data types
- 2 Make data accessible to biologists and clinicians
- 3 Facilitate data submission and release
- 4 Co-locate data analysis with storage
- 5 Improve metadata annotation

# Proteome Data Commons – democratize access to cancer-related proteomic datasets

NIH NATIONAL CANCER INSTITUTE

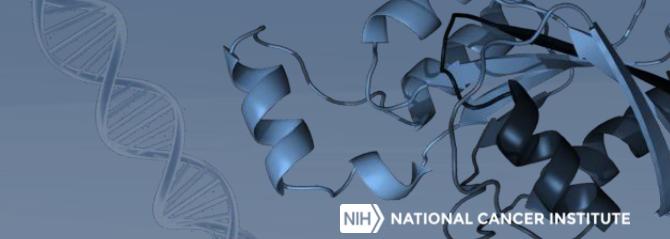
<https://pdc.esacinc.com/pdc/pdc>



# PDC fast facts



- Mass spectra
- TMT, label-free, and DIA data
- Data from CPTAC, CBTC
- In sync with GDC, KidsFirst, Seven Bridges, and TCIA



3  
Programs

9 TB  
Data volume

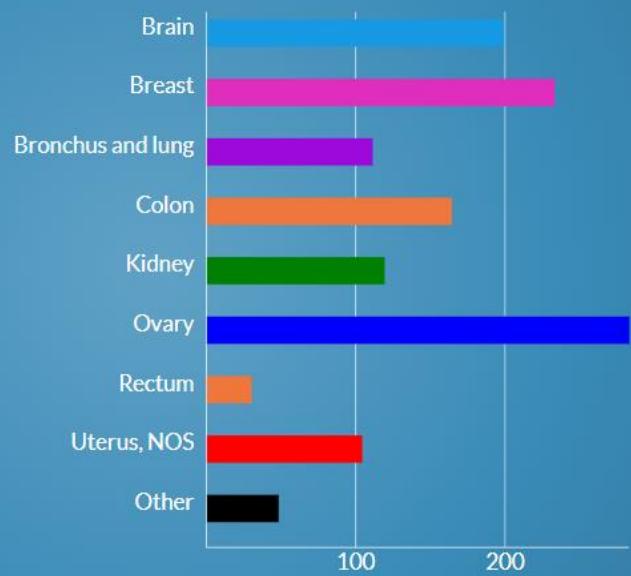
42,959  
Data files

64,720,831  
Spectra

991,578  
Peptides

13,713  
Proteins

## PRIMARY SITES



## DISEASE TYPES

Breast Invasive Carcinoma	233
Chromophobe Renal Cell Carcinoma	1
Clear Cell Renal Cell Carcinoma	116
Colon Adenocarcinoma	164
Lung Adenocarcinoma	111
Other	48
Ovarian Serous Cystadenocarcinoma	283
Papillary Renal Cell Carcinoma	2
Pediatric/AYA Brain Tumors	219
Rectum Adenocarcinoma	30
Uterine Corpus Endometrial Carcinoma	104

## RECENT RELEASES

- ✓ Pediatric Brain Tumor Atlas - CBTC program Pediatric/AYA brain tumor dataset released, November 2019
- ✓ CPTAC3 program Lung Adenocarcinoma (LUAD) dataset released, September 2019
- ✓ CRDC cloud resources piloted integration with Proteomic data commons for multiomic data analysis, August 2019
- ✓ MVP 0.5 Removed the requirement to register and login to the PDC Data Browser, June 2019

## NEWS

- 18th Human Proteome Organization World Congress, September 15-19, 2019, Adelaide, Australia
- 68th ASMS Conference on Mass Spec and Allied Topic June 2-6, 2019 Georgia World Congress Center, Atlanta, GA
- CPTAC Releases UCEC, ccRCC Discovery Data and Other Study Datasets



**FILTERS**

**GENERAL**

**Ethnicity**

- Hispanic or Latino
- Not Hispanic or Latino
- Not Reported
- Unknown

**Race**

- American Indian or Alaska Native
- Asian
- Black or African American
- Native Hawaiian or Other Pacific Islander
- Not Reported
- Unknown
- White

**Gender**

- Female
- Male
- Not Reported

**Tumor Grade**

**BIOSPECIMEN**

**CLINICAL**

**FILES**

**GENES**

**DISEASE TYPES**

Disease Type	Percentage
Ovarian Serous Cystadenocarcinoma	22.0%
Pediatric/AYA Brain Tumors	15.5%
Uterine Corpus Endometrial Carcinoma	8.1%
Breast Invasive Carcinoma	18.1%
Colon Adenocarcinoma	12.7%
Lung Adenocarcinoma	8.6%
Clear Cell Renal Cell Carcinoma	9.0%

**EXPERIMENT TYPES**

Experiment Type	Cases
iTRAQ4	~250
Label Free	~200
TMT10	~700
TMT11	~200

**ANALYTICAL FRACTIONS**

Analytical Fraction	Cases
Glycoproteome	~150
Phosphoproteome	~1100
Proteome	~1350

**Studies (24)**   **Biospecimens (1887)**   **Clinical (1311)**   **Files (42959)**   **Genes (13917)**

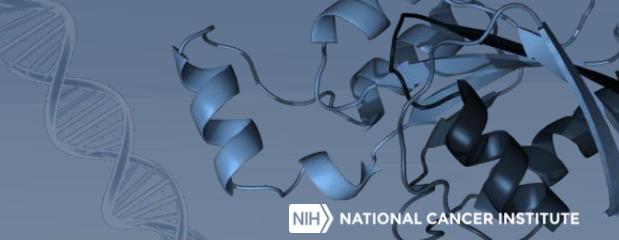
Total records: 1311

**Download All Manifests**   **Download Clinical Manifest**

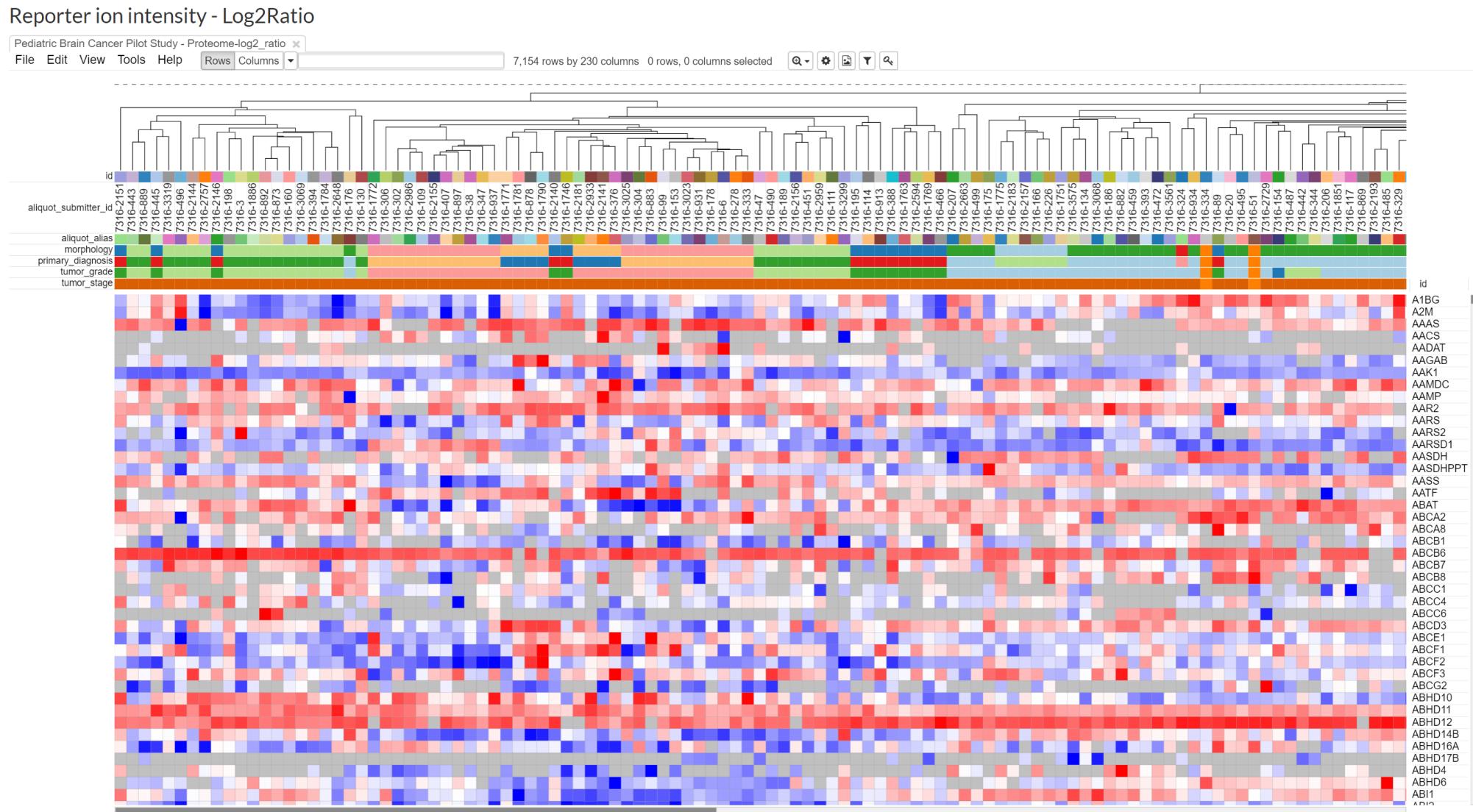
**Case Details**

#	Cases Submitter ID	External Data Resource	Ethnicity	Gender	Race	Morphology	Primary Diagnosis	Site of Resection or Biopsy	Tissue or Organ of Origin	Tumor Grade	Tumor Stage
1	C16974	Kids First Data Resource Center	Not Hispanic or Latino	Male	White	9505/1	Ganglioglioma, NOS	Temporal Lobe	Temporal Lobe	G1	Unknown
2	C94956	Kids First Data Resource Center	Not Hispanic or Latino	Male	White	9380/3	Glioma, malignant	Temporal Lobe	Temporal Lobe	High Grade	Unknown
3	C94956	Kids First Data Resource Center	Not Hispanic or Latino	Male	White	9380/3	Glioma, malignant	Temporal Lobe	Temporal Lobe	High Grade	Unknown

# Customizable heat maps of peptides



NIH NATIONAL CANCER INSTITUTE



# User flow to compute PDC data on Cloud Resources

NIH NATIONAL CANCER INSTITUTE

## NATIONAL CANCER INSTITUTE Proteomic Data Commons

1. User starts on PDC portal to identify cohort of files
2. User downloads **files manifest** of selected cohort
  - a. Note - there are multiple manifest options on the PDC, only the "files manifest" will work

<https://pdc.esacinc.com/pdc/>



## CANCER GENOMICS CLOUD

1. User moves to CGC, creates a project
  - a. Files → Add files → Import from PDC
2. User prompted to upload the manifest from the PDC
3. PDC files copied into user's project
  - a. FileID → file URL via Fence

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

# PDC: PepQuery

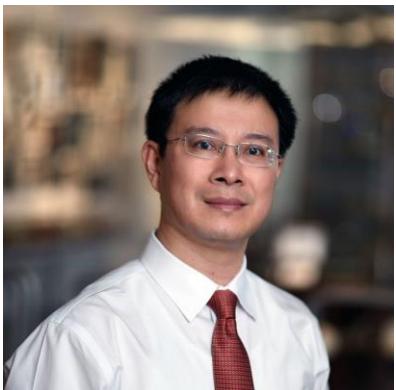
- Search data for specific variant peptide

## Resource

**PepQuery enables fast, accurate, and convenient proteomic validation of novel genomic alterations**

Bo Wen,<sup>1,2</sup> Xiaojing Wang,<sup>1,2</sup> and Bing Zhang<sup>1,2</sup>

<sup>1</sup>Lester and Sue Smith Breast Center, Baylor College of Medicine, Houston, Texas 77030, USA; <sup>2</sup>Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas 77030, USA



MS/MS dataset: (detailed information) i  
cptac\_colon\_2014\_nature ▼

Target event: i  
Peptide sequence ▼

Input peptide sequence i  
LVVVGADGVGK

Reference database: i  
RefSeq\_Hsapiens\_hg38 ▼

Scoring algorithm: i  
HyperScore ▼

Unrestricted modification filtering: i  
 Yes  No

Start Stop

# PDC: PepQuery

NIH NATIONAL CANCER INSTITUTE

## Identification overview

MS/MS searching parameters:

parameter	value
1 Enzyme	Trypsin
2 No. of missed cleavages	2
3 Fixed modifications	Carbamidomethylation of C
4 Variable modifications	Oxidation of M
5 Peptide tol. ±	20 ppm
6 MS/MS tol. ±	0.5 Da

## Identification result

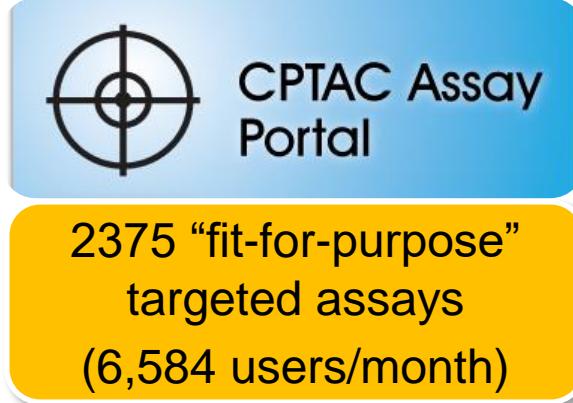
The identification result is presented in the below table:

Show 10 entries

Search:

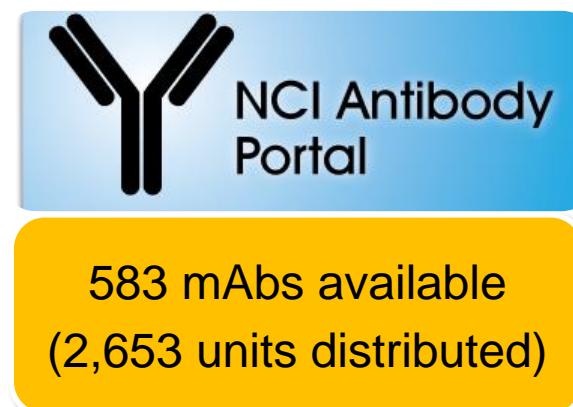
	peptide	modification	n	spectrum	sample_name	charge	exp_mass	ppm	pep_mass	mz	score	n_db	total_db	n_random	total_random
1	LVVVGADGVGK	-	84	23981	TCGA-AG-A00Y	2	1,012.587	4.748	1,012.592	507.301	32.602	0	249	0	974
2	LVVVGADGVGK	-	84	30908	TCGA-AA-A01R	2	1,012.595	-3.389	1,012.592	507.305	30.013	0	291	0	974
3	LVVVGADGVGK	-	84	20357	TCGA-AA-A02O	2	1,012.600	-7.789	1,012.592	507.307	23.684	0	252	1	974
4	LVVVGADGVGK	-	84	20415	TCGA-AA-A02O	2	1,012.592	-0.737	1,012.592	507.303	19.855	0	234	1	974

# Antibodies and targeted assays



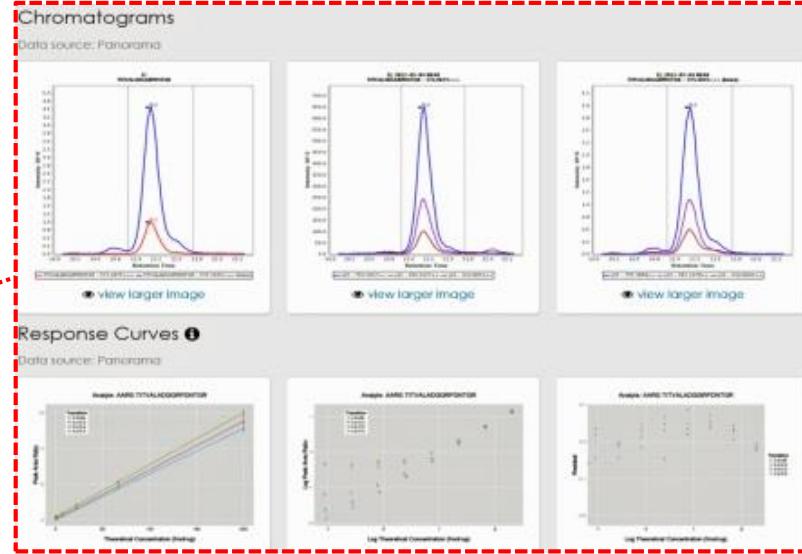
2375 “fit-for-purpose”  
targeted assays  
(6,584 users/month)

[assays.cancer.gov](http://assays.cancer.gov)

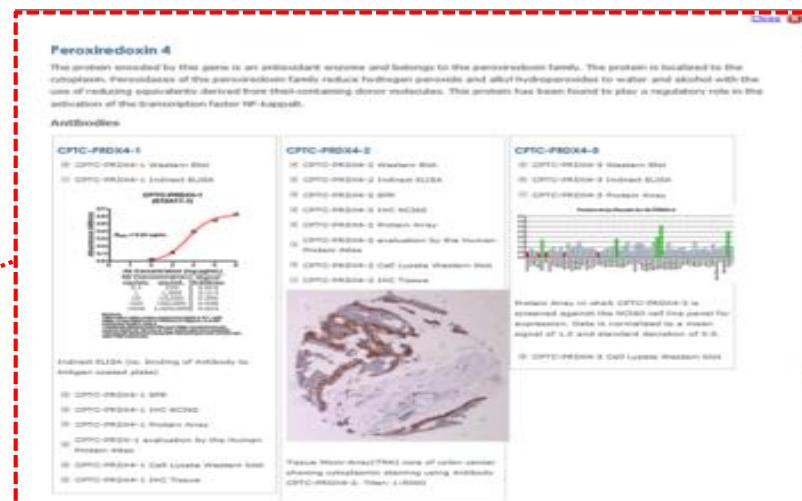


583 mAbs available  
(2,653 units distributed)

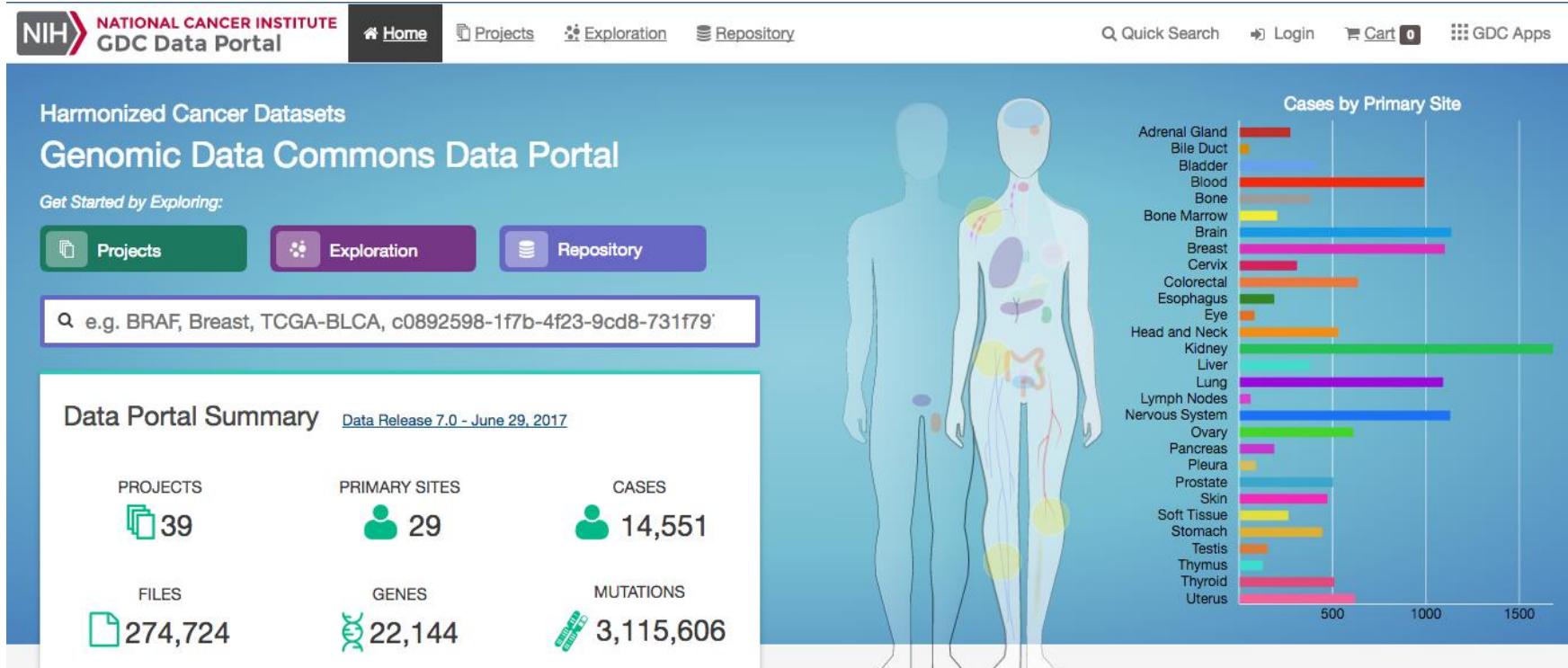
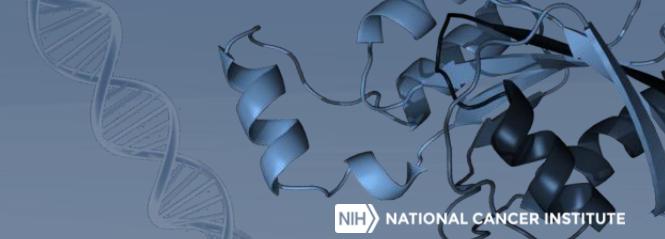
[antibodies.cancer.gov](http://antibodies.cancer.gov)



Available Antibodies		About	CPTAC Home	
Antibody Portal				
<input type="button" value="Browse By Antibodies"/>			<a href="#">Print This Page</a>	
Antigen Recognition		Show [25] entries	Showing 1 to 25 of 231 entries	Search:
<input type="checkbox"/> recombinant full-length				
<b>Antibody Images Type</b>				
<input type="checkbox"/> SBR	CPTC-AKR1B1-1	Aldo-keto Reductase Family 1 Member B1		
<input type="checkbox"/> DHC NCDB	CPTC-AKR1B1-2	Aldo-keto Reductase Family 1 Member B1		
<input type="checkbox"/> NAPPA	CPTC-AKR1B1-3	Aldo-keto Reductase Family 1 Member B1		
<input type="checkbox"/> Immuno-MS	CPTC-AKR1C1-1	Aldo-keto Reductase Family 1 Member C1		
<input type="checkbox"/> Western Blot	CPTC-AKR1C1-2	Aldo-keto Reductase Family 1 Member C1		
<input type="checkbox"/> Indirect ELISA	CPTC-AKR1C2-1	Aldo-keto reductase family 1 member C2		
<input type="checkbox"/> Immunoprecipitation Assay	CPTC-AKR1C2-2	Aldo-keto reductase family 1 member C2		
<input type="checkbox"/> IHC mAb	CPTC-AKR1C2-3	Aldo-keto reductase family 1 member C2		
<input type="checkbox"/> Cell Lysate Western Blot	CPTC-ANXA1-1	Anxa1 A1 (Annexin A1)		
<input type="checkbox"/> Cross Reactivity Data	CPTC-ANXA1-2	Anxa1 A1 (Annexin A1)		
<input type="checkbox"/> IHC Tissue	CPTC-ANXA1-3	Anxa1 A1 (Annexin A1)		
<b>Antibody Isotypes</b>				
<input type="checkbox"/> IgG1	CPTC-APOM1-2	APOM Nucleus 1		
<input type="checkbox"/> IgG2a	CPTC-BCL2L1-1	BCL2 like 1		
<input type="checkbox"/> IgG2b	CPTC-BCL2L1-2	BCL2 like 1		
<input type="checkbox"/> IgG3	CPTC-BCL2L1-3	BCL2 like 1		
<b>Monoclonal Source</b>				
<input type="checkbox"/> Mouse	CPTC-BCL2L2-1	BCL2 like 2		
<b>External Links</b>				
<input type="checkbox"/> Human Protein Atlas	CPTC-BCL2L2-2	BCL2 like 2		
<input type="checkbox"/> DSHB	CPTC-BCL2L2-3	Carboxic anhydrase VIII		
	CPTC-CAB-1	Carboxic anhydrase VIII		
	CPTC-CAB-2	Carboxic anhydrase VIII		
	CPTC-Caloydin-1	Caloydin (Prolectin Receptor Associated Protein)		
	CPTC-Caloydin-2	Caloydin (Prolectin Receptor Associated Protein)		
	CPTC-Caloydin-3	Caloydin (Prolectin Receptor Associated Protein)		
	CPTC-CD34-1	Cell division cycle 34 homolog (S...)		



# NCI Genomic Data Commons (GDC)



- Launched in 2016 with over 4 PB of data.
- Joint project with OICR.
- Used by 1000 -2000+ users per day.
- Based upon an open source software stack that can be used to build other data commons.

\*See: NCI Genomic Data Commons: Grossman, Robert L., et al. "Toward a shared vision for cancer genomic data." New England Journal of Medicine 375.12 (2016): 1109-1112.

# The Cancer Imaging Atlas (TCIA)

NIH NATIONAL CANCER INSTITUTE

CANCER IMAGING ARCHIVE

HOME ABOUT US SHARE YOUR DATA ACCESS THE ARCHIVE RESEARCH ACTIVITIES HELP

REQUEST DOI: Citation Linking for Datasets

TCIA has the ability to create Digital Object Identifier (DOI) linked to subsets of TCIA data, which authors may use as data citations in their scholarly papers. [Learn more...](#)

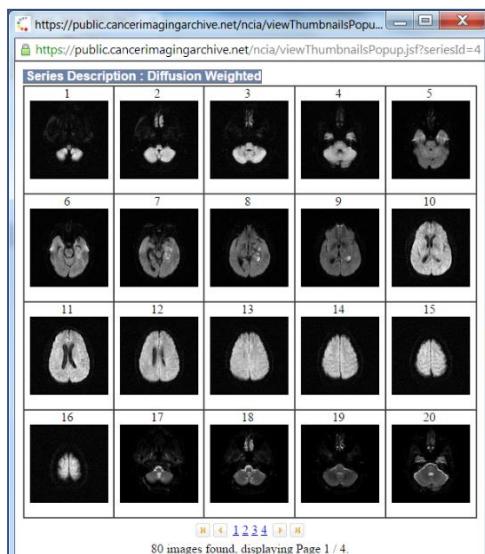
Search TCIA

Show 100 entries Filter table:

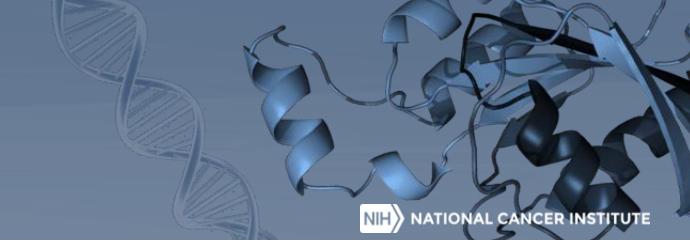
Cancer Type	Collection	Location	Subjects	Modalities
Ovarian Serous Cystadenocarcinoma	TCGA-OV	Ovary	111	CT, MR

<http://cancerimagingarchive.net>

- ~40,000 total subjects in the archive
- ~32 M images available
  - Quantitative Imaging Network
- TCGA - CPTAC
- ECOG-ACRIN and RTOG

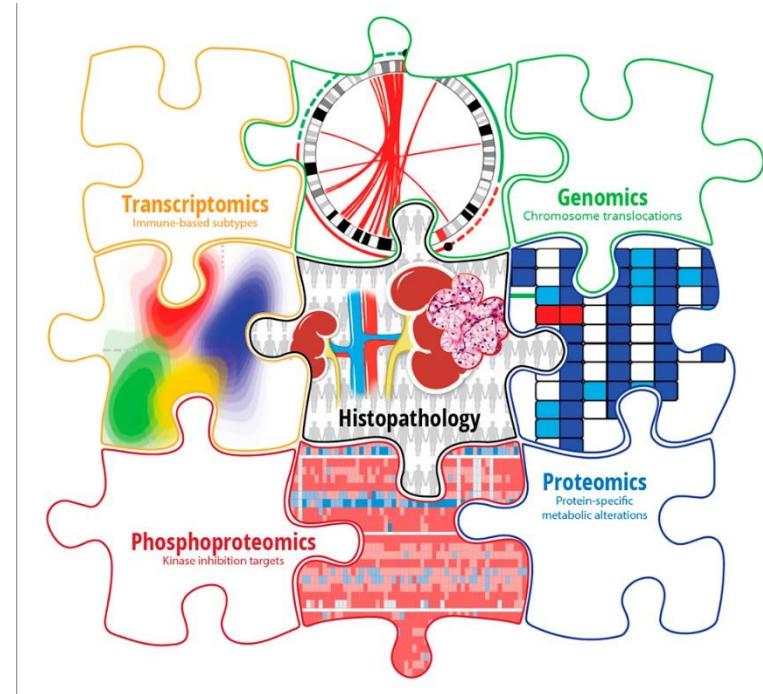


# Summary



## • Invitation

- Submit data to [\*\*https://pdc.esacinc.com\*\*](https://pdc.esacinc.com)
- Provide feedback [\*\*nci.pdc.help@esacinc.com\*\*](mailto:nci.pdc.help@esacinc.com)

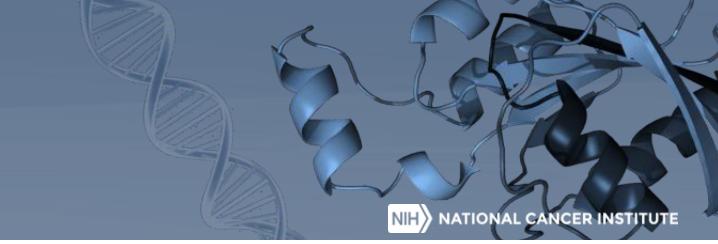


Thank You



# Backup slides

# Where to get the data



## ■ CPTAC data portal

<https://cptac-data-portal.georgetown.edu/cptacPublic/>

<https://proteomics.cancer.gov>  
and select “Data Portal”

Center for Strategic Scientific Initiatives

DATA PORTAL HOME ASSAY PORTAL ANTIBODY PORTAL ABOUT



### Data Portal

CPTAC (2006-2011)

CPTAC (2011-present)

External Studies

Query Data

Help

#### Latest Data Release and Publications:

March 2018

CPTAC ovarian cancer confirmatory study with global proteome and phosphoproteome analyses.

February 2018

Mass spectrometry-based proteomics reveals potential roles of NEK9 and MAP2K4 in resistance to PI3K inhibitors in triple negative breast cancer. Filip Mundt, Sandeep Rajput, Shunqiang Li, Kelly Ruggles, Arshag D. Mooradian, Philipp Mertins et al., *Cancer Res*, in press (2018).

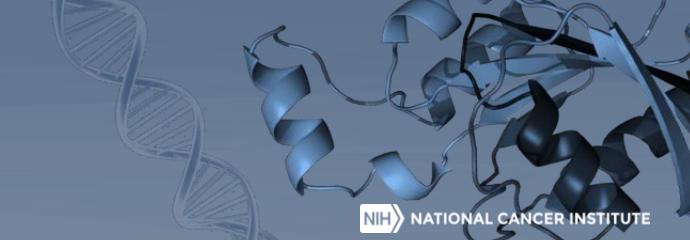
December 2017

CPTAC colon cancer confirmatory study with global proteome and phosphoproteome analyses.

Please install the latest Aspera Connect Client software (3.7) [here](#).

Study Name	Description	Publications
<a href="#">CPTAC Ovarian Cancer Confirmatory Study</a> <small>new</small>	Ovarian cancer tumors from a prospective CPTAC tissue collection were analyzed with liquid chromatography-tandem mass spectrometry (LC-MS/MS) global proteomic and phosphoproteomic profiling.	
<a href="#">Buparlisib treated xenograft tumors of TNBC</a> <small>new</small>	Intrinsic and adaptive mechanisms of resistance were analyzed in a panel of patient-derived xenograft models of triple negative breast cancer (TNBC) with varying responsiveness to buparlisib, a phosphoinositide 3-kinase (PI3K) inhibitor.	
<a href="#">CPTAC Colon Cancer Confirmatory Study</a> <small>new</small>	Colon cancer tumors from a prospective CPTAC tissue collection were analyzed with liquid chromatography-tandem mass spectrometry (LC-MS/MS) global proteomic and phosphoproteomic profiling.	

# Available datasets



Tumor type	TCGA residual		CPTAC confirmatory	
	Tumor	Normal	Tumor	Normal
Colorectal	90	0	104	96
Breast	105	0	100	
Ovarian	174	0	100	25



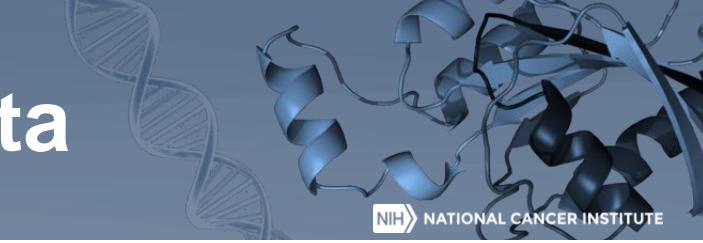
Genomic and imaging data also available for these samples at  
Genome Data Commons and The Cancer Imaging Archive

<https://gdc.cancer.gov/>

<http://www.cancerimagingarchive.net/>

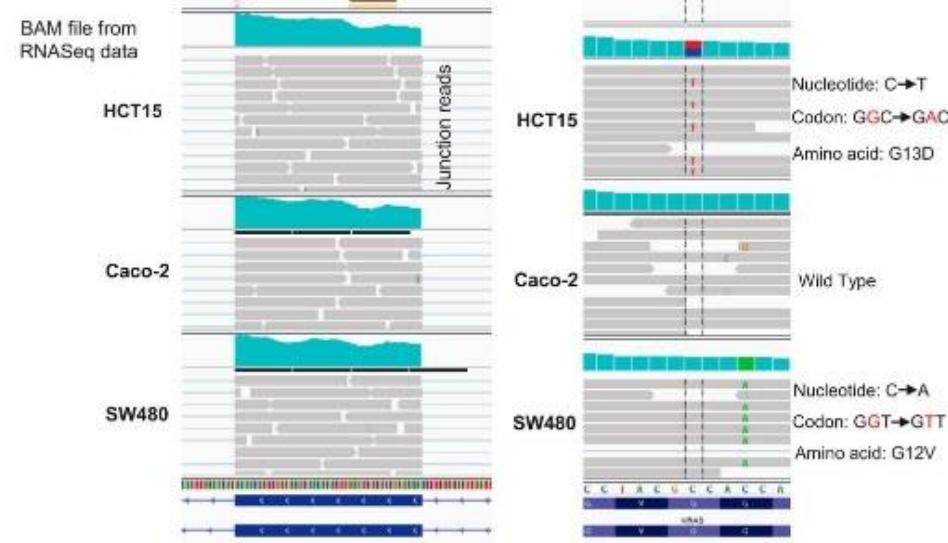
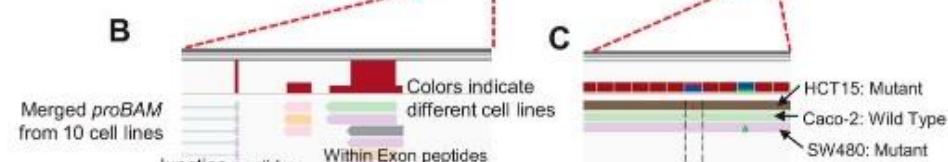
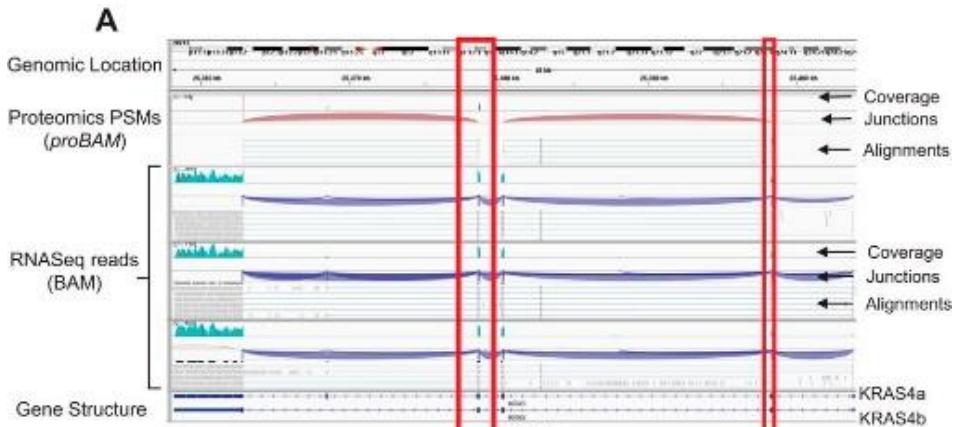


# 2 approaches to manage proteogenomic data



- Current approach
  - Local server storage
  - Users download data
- Future approach
  - Cloud storage
  - Users bring tools to cloud for data access and computation

# Data sharing and visualization



## Data sharing and visualization

# What problem will the PDC solve?

NIH NATIONAL CANCER INSTITUTE

- Simplify access to cancer research data



# What problem will the PDC solve?

NIH NATIONAL CANCER INSTITUTE

- Simplify access to cancer research data

Proteomic data



Genomic data



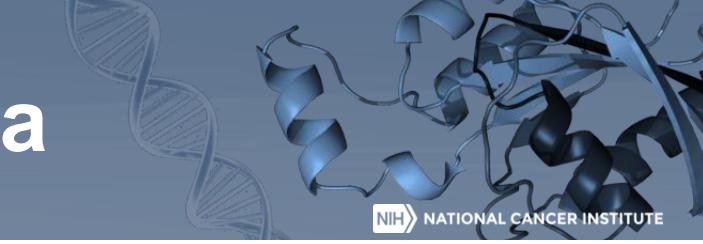
Imaging data

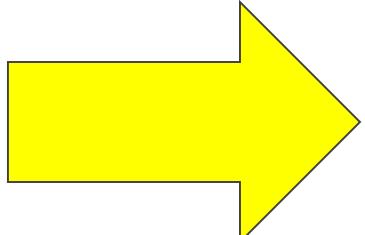


Clinical data

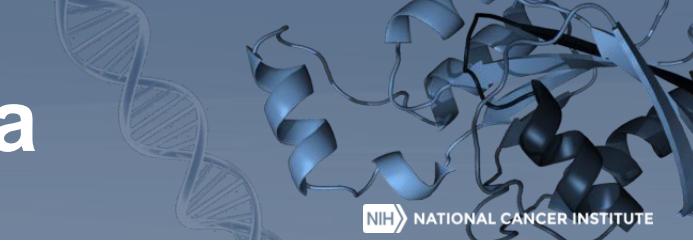


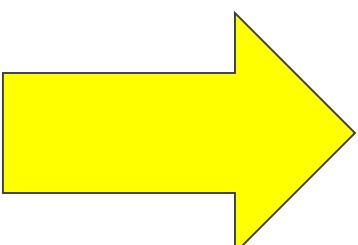
# Big Data V's applied to Proteogenomics data



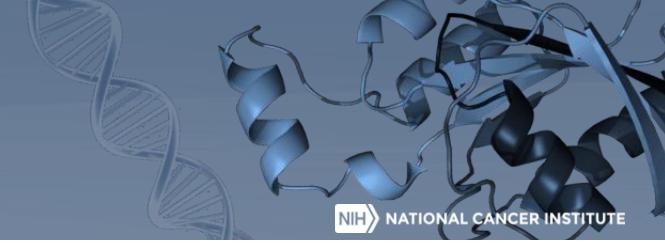
- Big Data V's
    - Volume
    - Variety
    - Velocity
    - Veracity
    - Valence
  - Proteogenomics data
    - P – TB; G - PB
    - Highly multidimensional
    - Low in research space
    - Analytical, clinical quality
    - Gene-protein, protein-protein, subject-sample
- 

# Big Data V's applied to Proteogenomics data



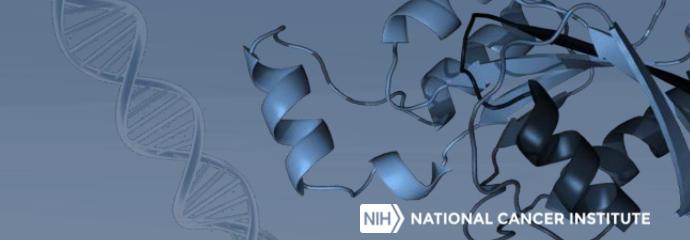
- Big Data V's
    - Volume
    - Variety
    - Velocity
    - Veracity
    - Valence
  - Proteogenomics data
    - P – TB; G - PB
    - **Highly multidimensional**
    - Low in research space
    - **Analytical, clinical quality**
    - Gene-protein, protein-protein, subject-sample
- 

# FAIR principles of data sharing



- Findable
- Accessible
- Interoperable
- Reusable

# Data Biosphere == Data Ecosystem



- Modular – composed of functional components with well-specified interfaces
- Community-driven – created by many groups to foster a diversity of ideas
- Open – developed under open-source licenses that enable extensibility and reuse
- Standards-based – consistent with standards developed by coalitions such as GA4GH

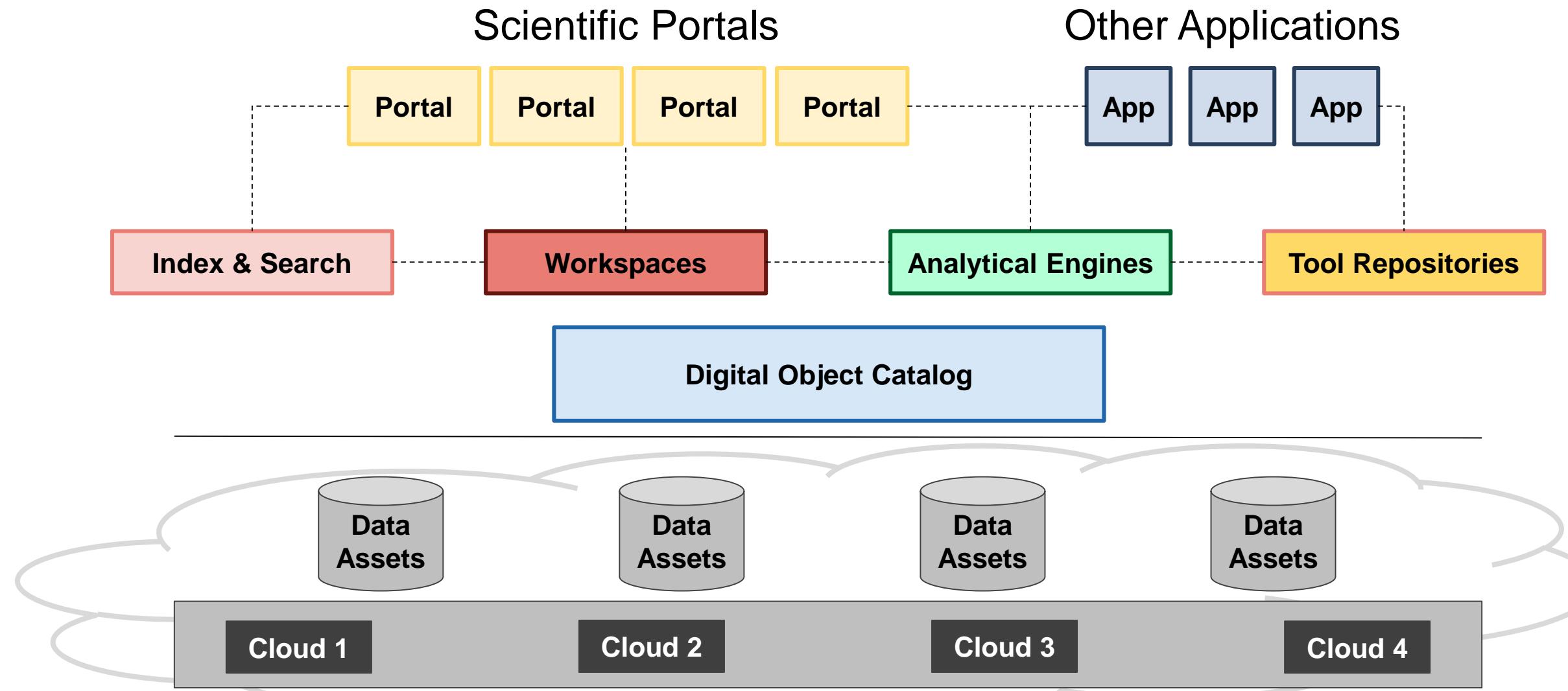
# Potential functions of a data commons

NIH NATIONAL CANCER INSTITUTE

- Data repository for submission, processing, and downloading
- System that applies a common set of bioinformatics pipelines to submitted data
- System that reanalyzes submitted data when new bioinformatics pipelines are developed
- System that allows users to build their own applications and systems that interoperate the commons through the API

# Data environment of modular components

NIH NATIONAL CANCER INSTITUTE



# What is a data commons?



- A data commons **colocates** data, storage, and computing infrastructure with commonly used services, tools, and apps for analyzing and sharing data to create an **interoperable** resource for the research community

# How is it different from a repository?

NIH NATIONAL CANCER INSTITUTE

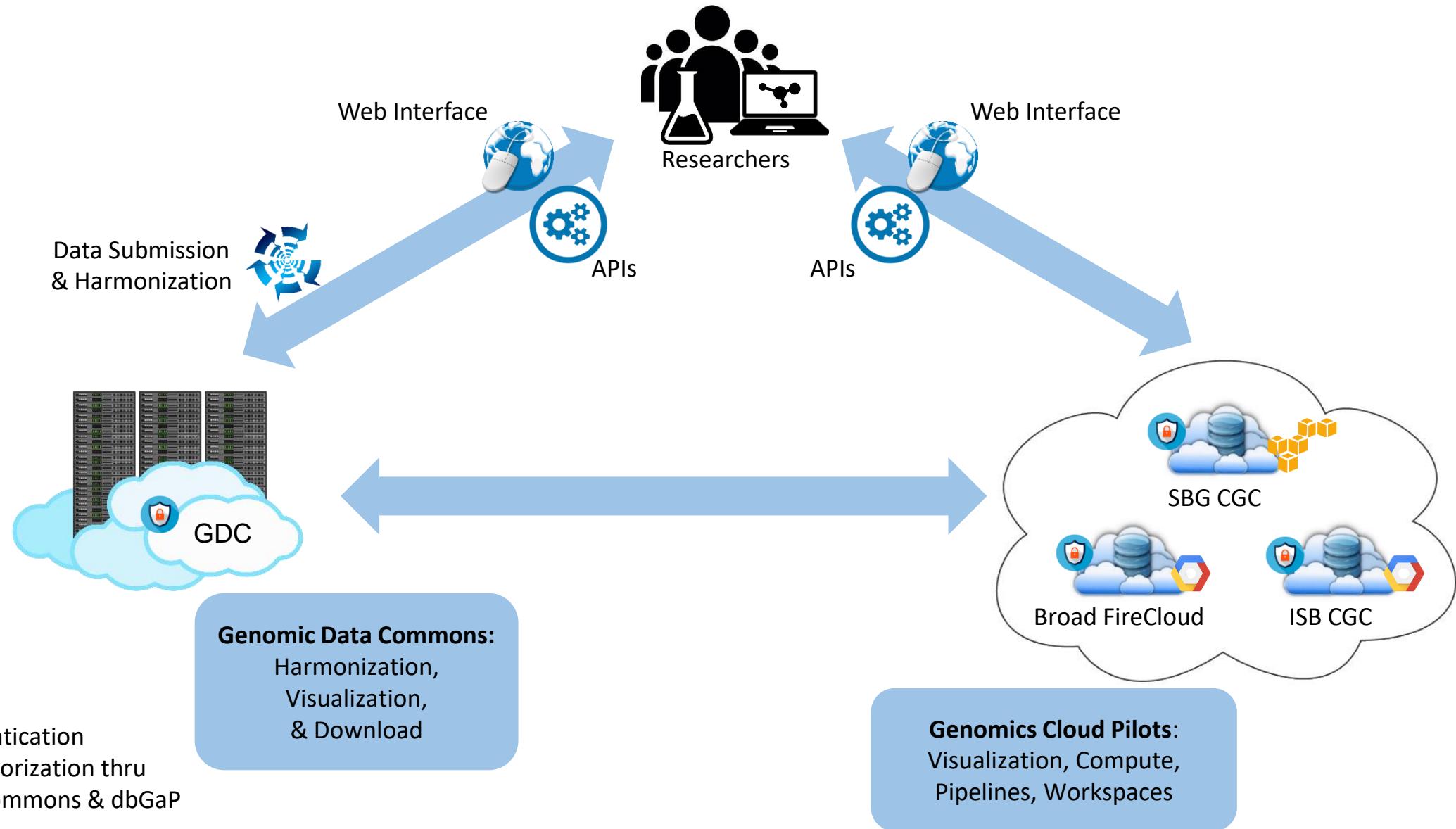
- Data repository
- Data commons

Data → Users

Data & Compute ← Users

# GDC and NCI Cloud Pilots/Framework Today

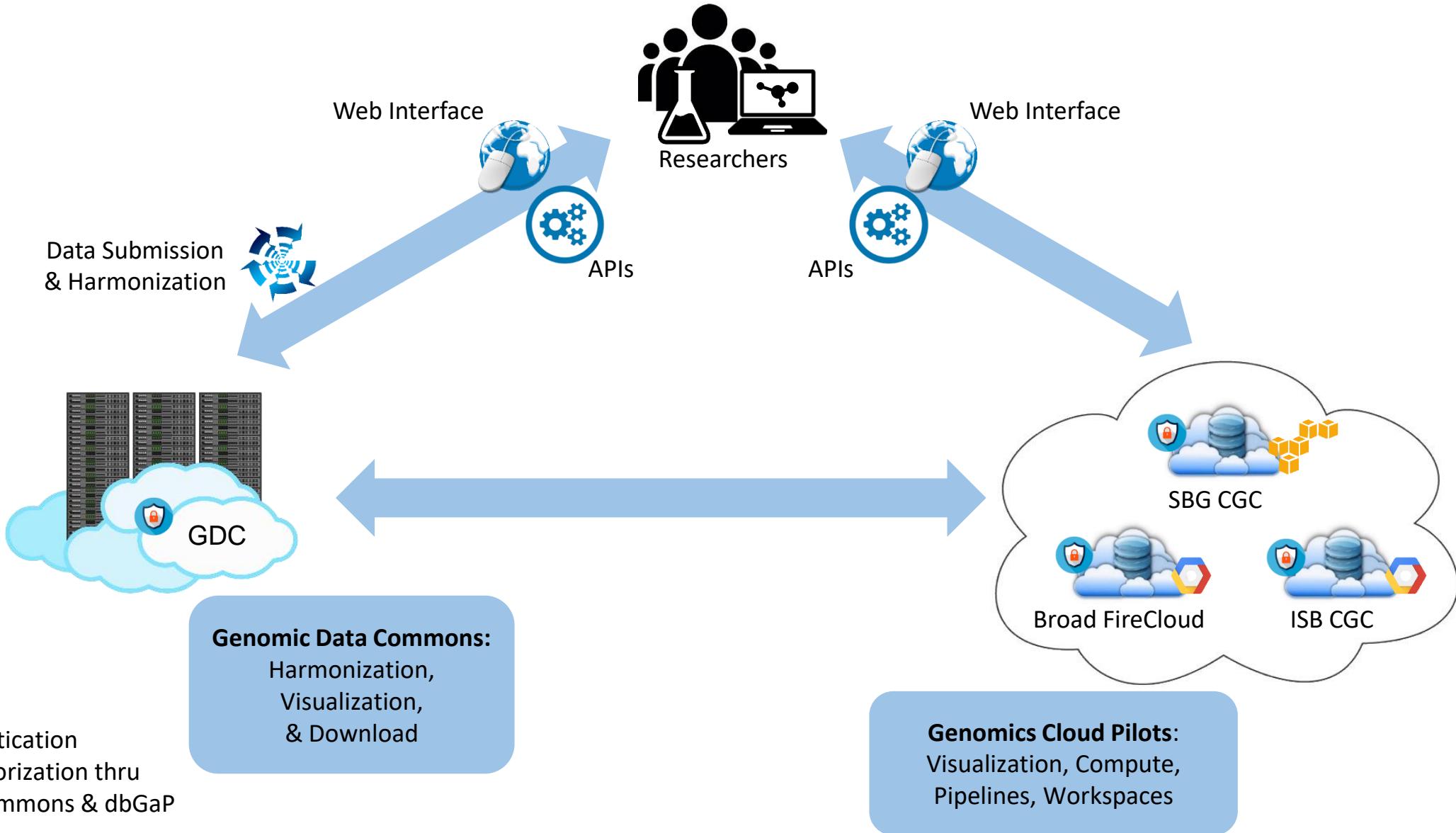
NIH NATIONAL CANCER INSTITUTE



# GDC and NCI Cloud ~~Pilots~~/Framework Today

Resources

NIH NATIONAL CANCER INSTITUTE



# NCI Cancer Research Data Commons (NCRDC)



Genomic  
node:  
GDC

Proteomic  
node:  
PDC

Imaging  
node:  
IDC

API's

Authentication  
Data Model

Framework – Core Services  
Workflows/worksheets  
Cross cloud services

Common architecture



## Data Acquisition and Curation

Mechanisms for the community to contribute and store data; ensuring the quality, consistency, and usability of the data

## Key Components of a Commons Node



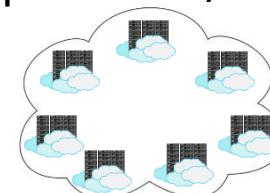
## Support / Help / Operations

Support for the community's needs to use and maximize impact of each Commons Node



## Infrastructure

Physical hardware, software, security policies and protocols, cloud computing

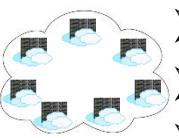


# The NCI Cancer Research Data Commons

## A virtual, expandable infrastructure



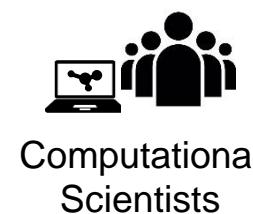
- Standardized data submission and Q/C
- Controlled vocabularies
- Harmonization by subject matter experts



- Secure data access through API or web UI
- Query across data domains
- Analytics, elastic compute, visualization



Tool /  
Algorithm  
Developers



Computational  
Scientists



Biologists / Clinical Researchers



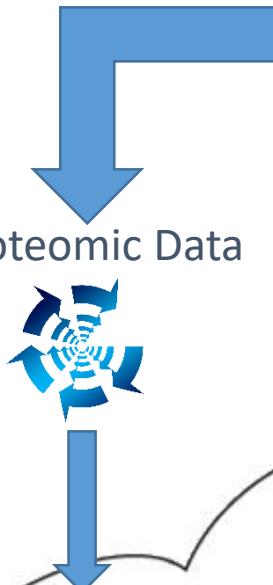
Authentication  
&  
Authorization

Clinicians and Patients

Data Contributors



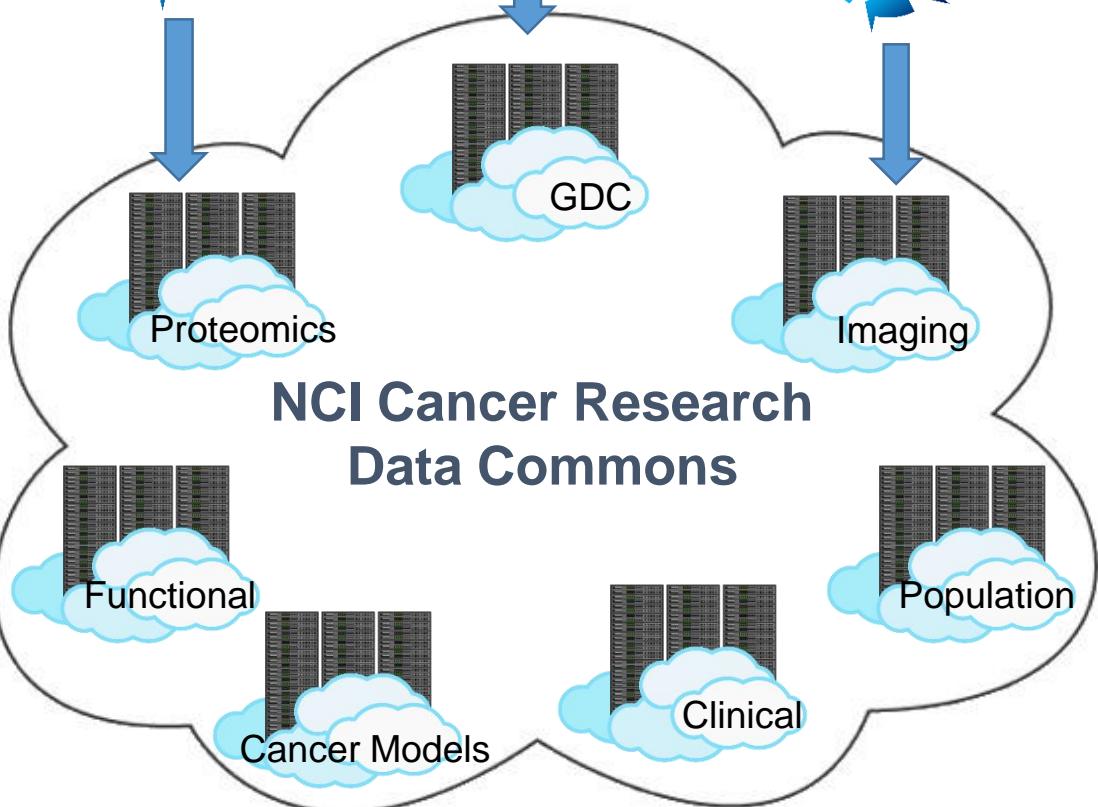
Genomic Data



Proteomic Data



Imaging Data



NCI Cancer Research  
Data Commons



GDC



Proteomics



Functional



Imaging



Population



Cancer Models



Clinical

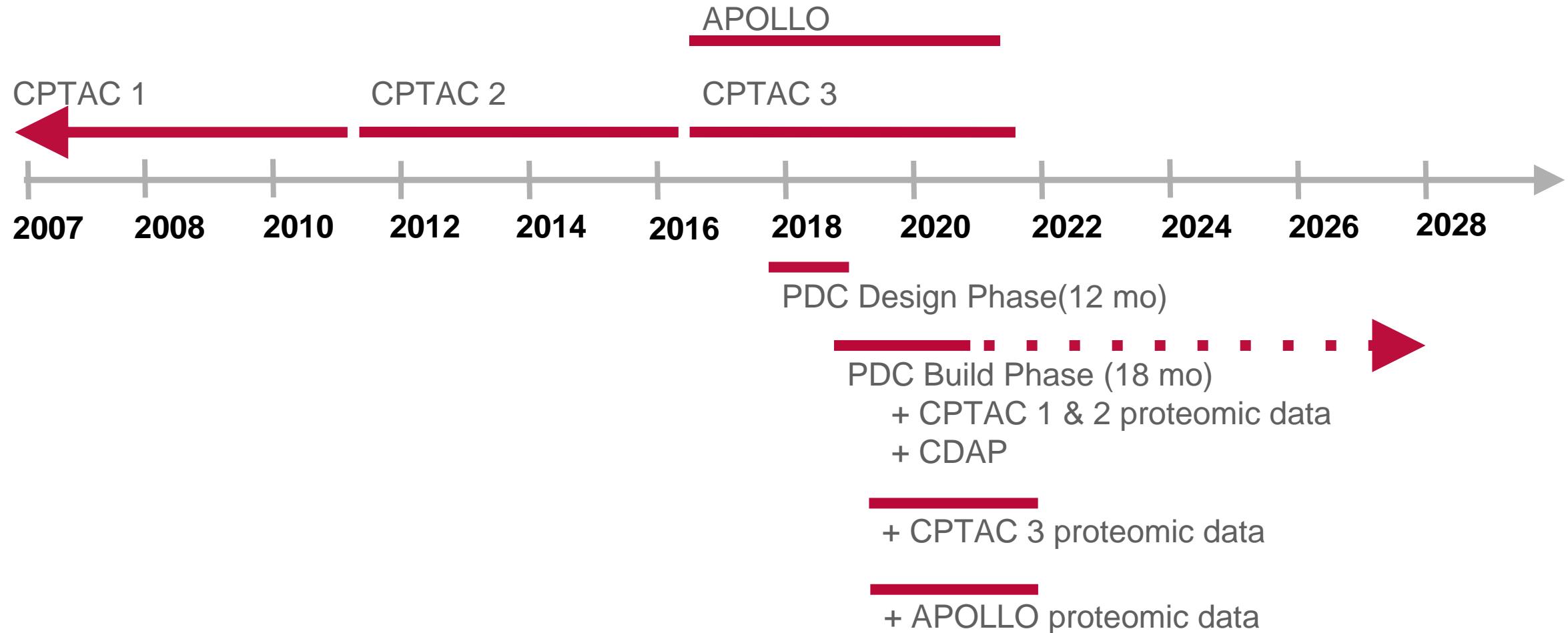
# What other problem will the PDC solve?



- Increase access and usage to researchers beyond proteomics domain

# PDC Timeline

NIH NATIONAL CANCER INSTITUTE



# What will the PDC do?



- Current award is to design and build a **prototype**
  - Gather requirements
  - Ensure interoperability with CRDC framework and other nodes
  - Design architecture, data model, analysis and visualization environments

# Who is the PDC?

NIH NATIONAL CANCER INSTITUTE

- ESAC
  - PI Anand Basu
  - Karen Ketchum, Rajesh Thangudu



# Acknowledgments



## NCI CBIIT

Izumi Hinkson

Tanja Davidsen

Tony Kerlavage

Juli Klemm

Allen Dearry

## NCI CSSI

Jerry Lee

Sean Hanlon

Henry Rodriguez

Steve Cole

## NCI OA

Andy May

Tim Crilley

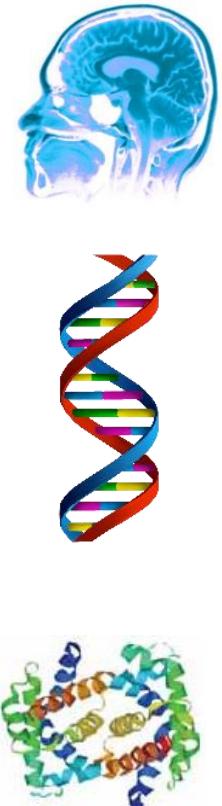
## Outside perspectives

Nuno Bandeira

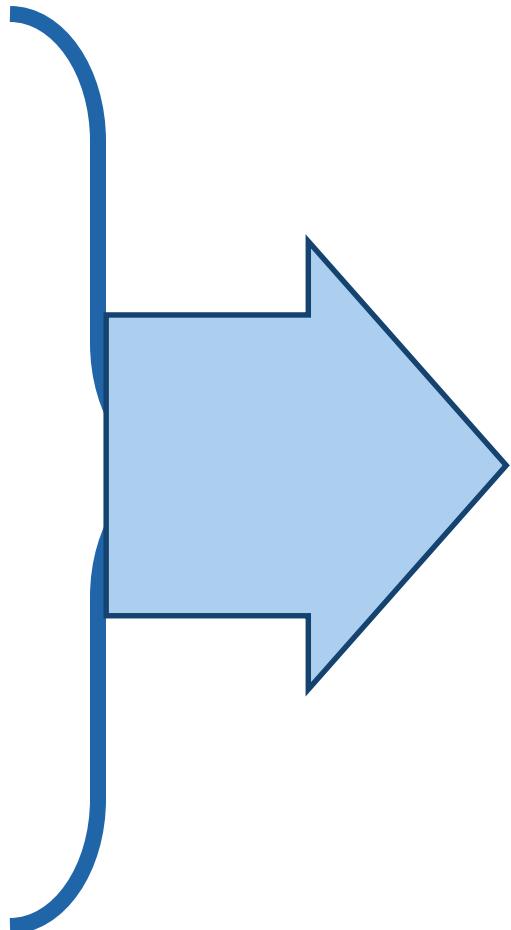
Jake Jaffe

# Opportunity for Big Data

NIH NATIONAL CANCER INSTITUTE



Imaging  
Genomics  
Proteomics



Improved  
Prevention,  
Diagnosis,  
Treatment

# Ovarian Cancer

(PROTEIN ABUNDANCE - new proteome subtype identified)



- 174 ovarian HGSC tumors

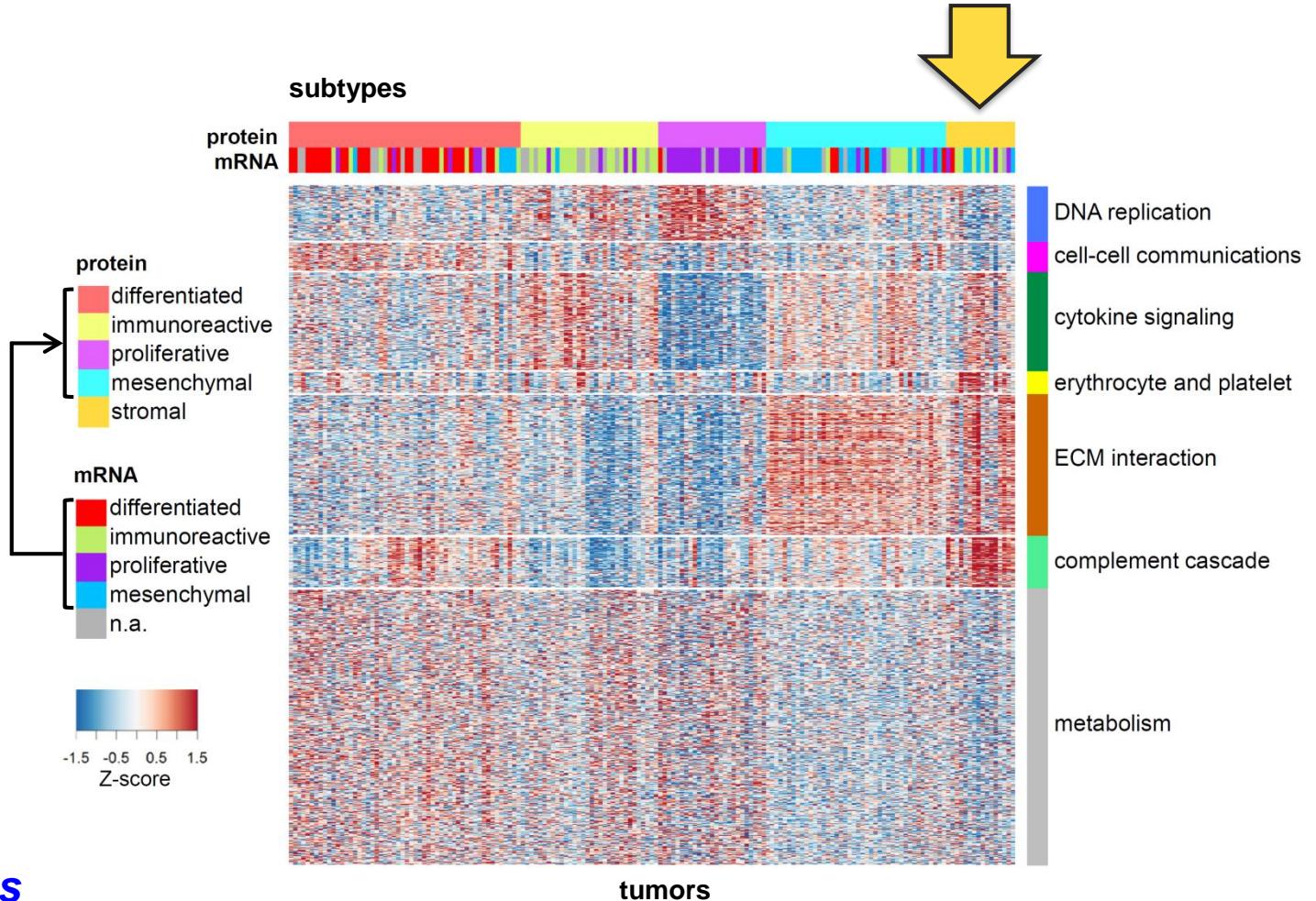
- Selection criteria:
  - Overall Survival (OS)
  - Homologous Recombination Deficiency status (HRD)

- **5 proteomic subtypes**

(4 transcriptomic subtypes)

- mRNA subtypes translate at protein level
- New “stromal” subtype emerged

• *While interesting observations, no strong separation of OS and HRD status*



# Clinical Proteomic Tumor Analysis Consortium (CPTAC)

NIH NATIONAL CANCER INSTITUTE

## Proteogenomics builds on TCGA

- **Tumor Characterization Program (treatment naïve – tumor and NAT):**

### Data Generation

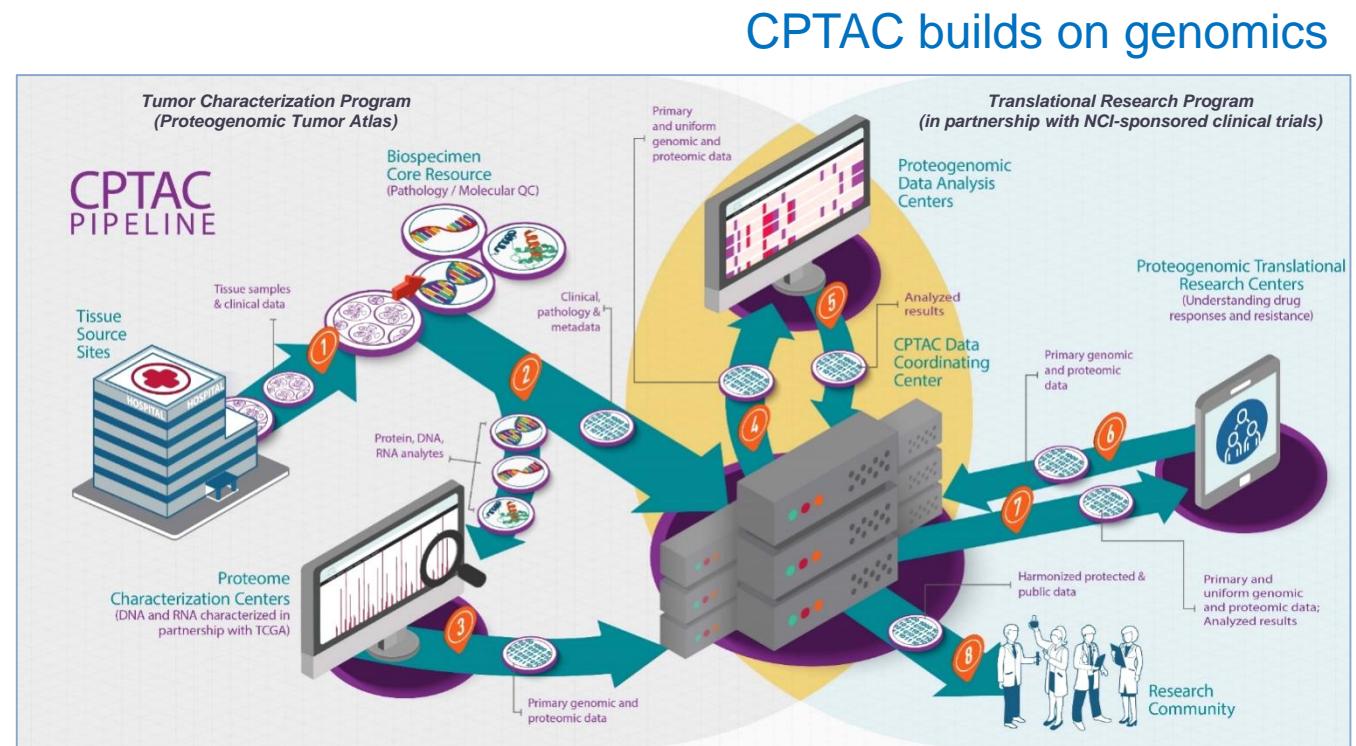
JHU  
PNNL  
Broad

### Data Analysis

BCM  
UMich  
NYU/WashU/BYU  
MtSinai  
Broad

- **Translational Research Program (pre-clinical and clinical trials):**

OHSU/PNNL  
BCM/Broad  
FHCRC/UAB



# Clinical Proteomic Tumor Analysis Consortium (CPTAC)

NIH NATIONAL CANCER INSTITUTE

## Proteogenomics builds on TCGA

- **Tumor Characterization Program (treatment naïve – tumor and NAT):**  
*characterize proteins and genes to better understand the molecular basis of cancer*
- **Translational Research Program (pre-clinical and clinical trials):**  
*understand [predict] drug response and resistance to therapies **in context of a clinical trial***

