



CPTAC Proteomic Data Processing

Olivier Gevaert
Department Medicine
Department Biomedical Data Science
Stanford University



Stanford
MEDICINE

Overview

- CPTAC data download and preprocessing issues
 - Quality controls
 - Batch correction
 - Different types of proteomics
- Integration of CPTAC proteomic data with other omics
 - ProteoMix: Example of integration with DNA methylation
- Proposed radioproteomics maps
 - Comparison with **traditional** radiogenomics maps.



CTPAC data download & preprocessing issues

Technical steps



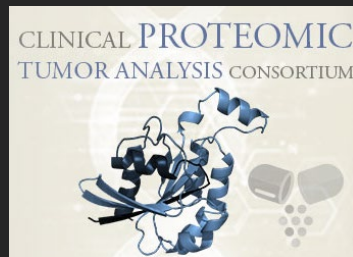
Stanford
MEDICINE

Clinical Proteomic Tumor Analysis Consortium (CPTAC)

Objective: Understand the molecular basis of cancer that is not fully elucidated or not possible through genomics by adding **complementary functional layer of protein biology** and to accelerate the translation of molecular findings into the clinical.

Timeline:

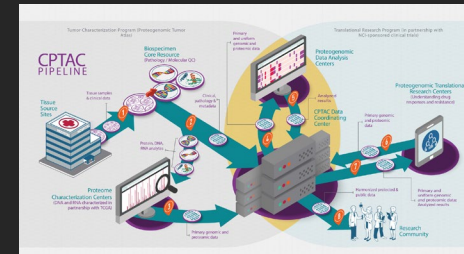
2006-2011 Phase I:
Process Development
and Reproducibility



2011-2016 Phase II:
Proteogenomic
Discovery

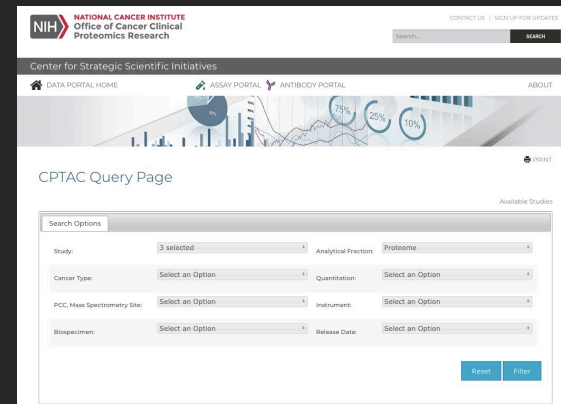


2016-Present Phase III:
Expansion & Clinical
Translation



Data Download:

- Method 1: Data portal
 - Hosted by Georgetown University
 - Requires IBM Aspera Launcher
 - <https://cptac-data-portal.georgetown.edu>
- Method 2: NCI Proteomic Data Commons (in beta)
 - Similar to genomic data commons
 - <https://pdc.esacinc.com/pdc/pdc>
- Method 3: Programmatically
 - Direct access using Linux and command line or Python executable script
 - <https://proteomics.cancer.gov/data-portal/about/faqs>



CPTAC Data Overview: Phase 2 data

Generated using Common Data Analysis Pipeline
(CDAP)

Dataset	Number Samples	Number Genes
BRCA – BROAD Institute	105	10624
COADREAD – Vanderbilt Uni.	95	5561
OV – Johns Hopkins Uni.	115*	8597
OV – PNNL	75*	7480

* 32 samples in common

Proteogenomics connects somatic mutations to signalling in breast cancer

Philipp Mertins^{1*}, D. R. Mani^{1*}, Kelly V. Ruggles^{2,3*}, Michael A. Gillette^{1,3*}, Karl R. Clauser¹, Pei Wang⁴, Xianlong Wang⁵, Jata W. Qiao⁶, Song Cao⁷, Francesca Petralia⁸, Emily Kawaler⁹, Filip Mundt¹⁰, Karsten Krug¹¹, Zhidong Tu¹², Jonathan T. Lei¹³, Michael L. Gatzra¹⁴, Matthew Wilkerson¹⁵, Charles M. Perou¹⁶, Venkata Yellapantula¹⁷, Kuan-lin Huang¹⁸, Chenwei Lin¹⁹, Michael D. McLellan¹, Ping Yan¹, Sherri R. Davies¹⁰, R. Reid Townsend¹⁰, Steven J. Skates¹¹, Jing Wang², Bing Zhang², Christopher R. Kinsinger¹¹, Mehdi Mesri¹¹, Henry Rodriguez¹¹, Li Ding¹², Amanda G. Paulovich¹³, David Fenyö⁷, Matthew J. Ellis¹⁴, Steven A. Carr¹ & the NCI CPTAC[†]

Somatic mutations have been extensively characterized in breast cancer, but the effects of these genetic alterations on

Proteogenomic characterization of human colon and rectal cancer

Bing Zhang^{1,2}, Jing Wang¹, Xiaojing Wang¹, Jing Zhu¹, Qi Liu¹, Zhao Shi^{3,4}, Matthew C. Chambers⁵, Lisa J. Zimmerman^{3,6}, Kent F. Shaddox⁶, Sangtae Kim¹, Sherri R. Davies⁸, Sean Wang⁹, Pei Wang¹⁰, Christopher R. Kinsinger¹¹, Robert C. Rivers¹¹, Henry Rodriguez¹¹, R. Reid Townsend⁸, Matthew J. C. Ellis¹², Steven A. Carr¹², David L. Tabb¹³, Robert J. Coffey¹³, Robert J. C. Slebos¹⁴, Daniel C. Liebler¹⁵ & the NCI CPTAC^{*}

Integrated Proteogenomic Characterization of Human High-Grade Serous Ovarian Cancer

Hui Zhang,^{1,19} Tao Liu,^{2,19} Zhen Zhang,^{1,19} Samuel H. Payne,^{2,19} Bai Zhang,¹ Jason E. McDermott,² Jian-Ying Zhou,¹ Vladislav A. Petyuk,² Li Chen,² Debit Ray,² Shisheng Sun,¹ Feng Yang,² Lijun Chen,¹ Jing Wang,³ Punit Shah,¹ Seong Won Cha,⁴ Paul Aiyetan,¹ Sunghee Woo,¹ Yuan Tian,¹ Marina A. Gritsenko,² Therese R. Clauss,² Caitlin Choi,¹ Matthew E. Monroe,² Stefani Thomas,² Song Nie,² Chaochao Wu,² Ronald J. Moore,² Kun-Hsing Yu,¹⁰ David L. Tabb,¹³ David Fenyö,⁷ Vineet Bafna,⁸ Yue Wang,¹ Henry Rodriguez,¹¹ Emily S. Boja,¹² Tara Hillke,¹⁰ Robert C. Rivers,¹¹ Lori Sokoll,¹ Heng Zhu,¹ Le-Ming Shi,¹ Leslie Cope,¹ Akhilesh Pandey,¹³ Bing Zhang,² Michael P. Snyder,² Douglas A. Levine,¹⁴ Richard D. Smith,² Daniel W. Chan,^{1,16} Karin D. Rodland,^{2,16} and the CPTAC Investigators

¹Department of Pathology, Johns Hopkins Medical Institutions, Baltimore, MD 21231, USA

²Biological Sciences Division, Pacific Northwest National Laboratory, Richland, WA 99352, USA

³Department of Biomedical Informatics, Vanderbilt University School of Medicine, Nashville, TN 37203, USA

⁴Department of Electrical and Computer Engineering, University of California, San Diego, La Jolla, CA 92093, USA

⁵Biomedical Informatics Training Program, Stanford University School of Medicine, Stanford, CA 94305, USA

⁶Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305, USA

⁷Center for Health Informatics and Bioinformatics and Department of Biochemistry and Molecular Pharmacology, New York University School of Medicine, New York, NY 10016, USA

⁸Department of Computer Science and Engineering, University of California, San Diego, La Jolla, CA 92093, USA

⁹Bradley Department of Electrical and Computer Engineering, Virginia Polytechnic Institute and State University, Arlington, VA 22203, USA

¹⁰Office of Cancer Clinical Proteomics Research, National Cancer Institute, Bethesda, MD 20892, USA

¹¹Department of Gynecology and Obstetrics, Johns Hopkins Medical Institutions, Baltimore, MD 21231, USA

¹²Department of Oncology, Johns Hopkins Medical Institutions, Baltimore, MD 21231, USA

¹³McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins Medical Institutions, Baltimore, MD 21287, USA

¹⁴Department of Gynecologic Oncology, Laura and Isaac Perlmutter Cancer Centre, NYU Langone Medical Center, New York, NY 10016, USA

¹⁵Co-first author

¹⁶Co-senior author

*Correspondence: dchan@hmi.edu (D.W.C.), karin.rodland@pnnl.gov (K.D.R.)

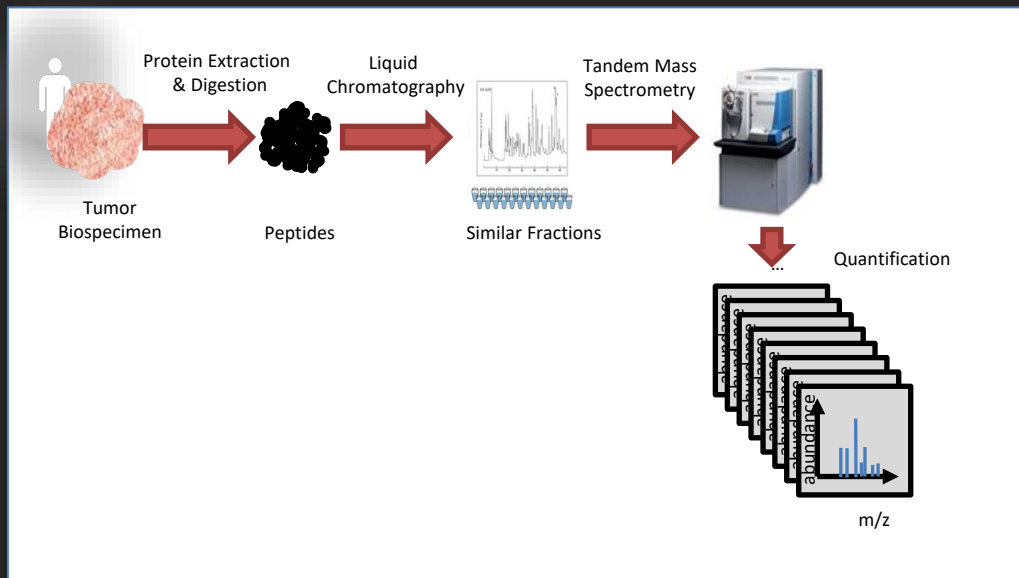
<http://dx.doi.org/10.1016/j.cell.2016.05.069>

Ongoing phase 3

- 10 new cancer sites
 - AML
 - CCRCC
 - Cutaneous Melanoma
 - GBM
 - HNSCC
 - LSCC
 - LUAD
 - Ductal Adenocarcinoma
 - Sarcomas
 - UCEC

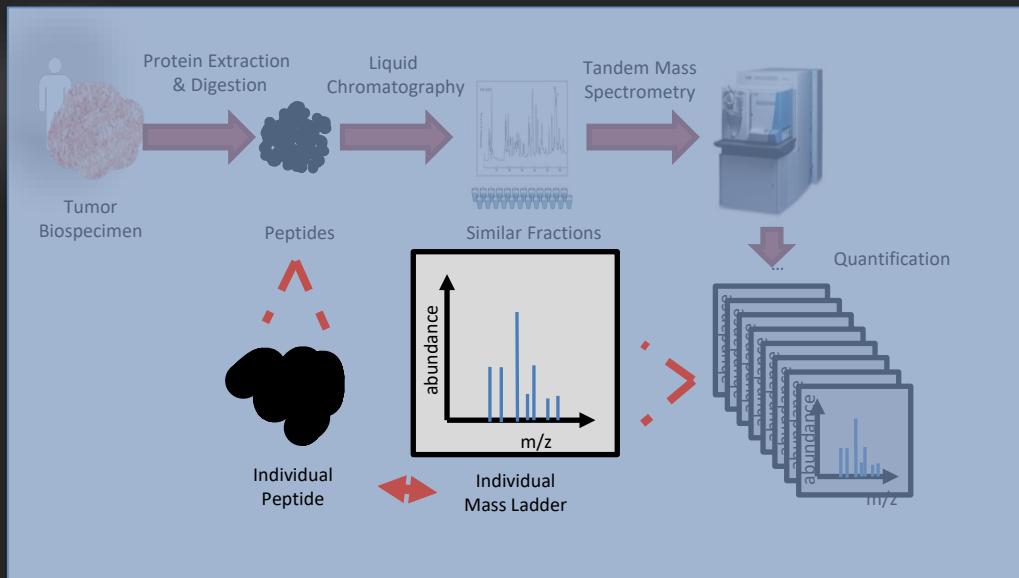
Proteomic workflow

Experimental Quantification



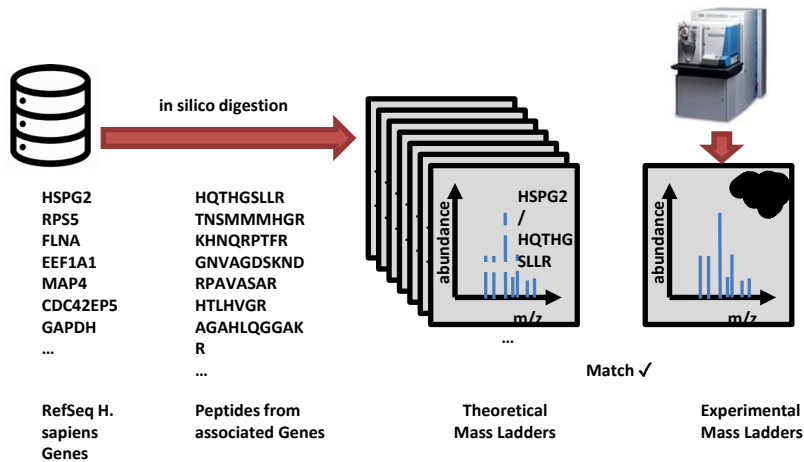
- Proteins extracted from tumor biospecimens from **matched TCGA** samples and proteins were tryptically digested into peptides, small segments of 7-30 amino acids
- Multi-stage high performance liquid chromatography produced homogenous fractions
- High resolution tandem mass spectrometry measures individual peptides at a time

Experimental Quantification



- Processing raw data produces a mass ladder for each peptide or group of few peptides
- Each peak corresponds to a sub-peptide
- These mass ladders are compared to theoretical mass ladders to identify peptides & proteins using the [RefSeq](#) database

Bioinformatics: Peptide Matching



- Before using abundance measurements for genomic analysis mass ladders must be linked to associated peptides and genes
- Mass ladders are identified using the RefSeq Database
 - Sequences for individual genes were fragmented in silico into peptides
 - Composition of peptides used to generate theoretical mass ladders
 - Mass ladders from experimental samples were matched by searching against theoretical mass ladders

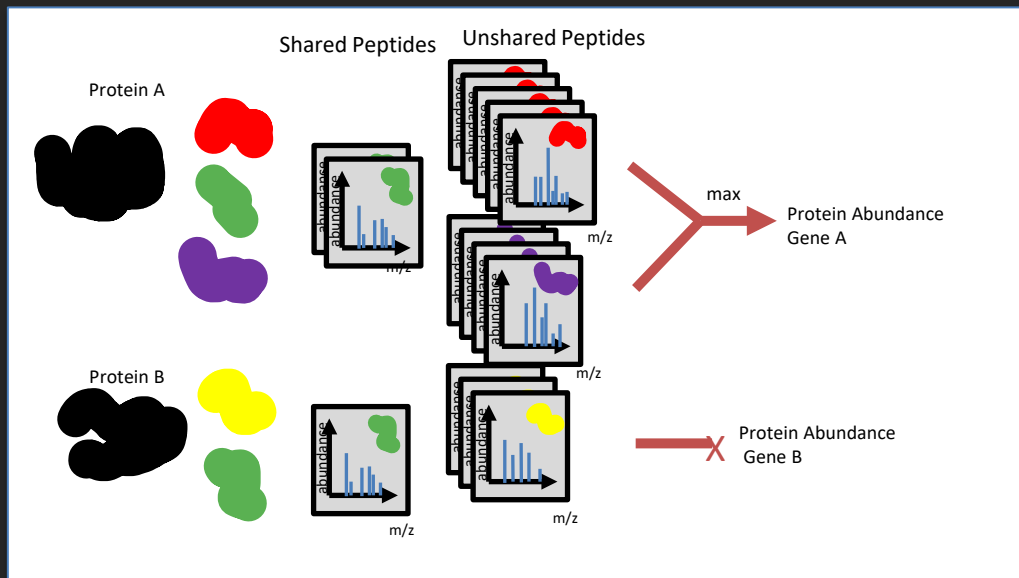
CDAP: Common Data Analysis Pipeline

	CDAP	Broad	JHU	PNNL	Vanderbilt
FASTA	RefSeq-Human-v37-Trypsin.fasta (32,800 entries)	Same as CDAP	Same as CDAP	Same as CDAP	humanRefSeq_v54_trypsin.fasta (34,589 entries)
Search Engine(s)	MS-GF+ (v9733)	SpectrumMill 4.0 (Beta)	MS-GF+ (v9146)	1. MS-GF+ v9324 (2/27/2013) v9358 (3/05/2013) v9593 (05/06/2013) v9699 (07/26/2013) v9736 (09/16/2013)	1. Pepitome 1.0.42 (library) 2. MyriMatch 2.1.87 3. MS-GF+ (v9176)
Ambiguous matches flagged?	Yes	No	No	Yes	Yes
Variable Protein Mods searched	MetOx(+16) Deamidation(+1)	MetOx(+16) Glu->pyro-Glu(-18) Gln->pyro-Glu(-17) Deamidation (N)(+1)	MetOx(+16)	MetOx(+16)	MetOx(+16) Glu->pyro-Glu(-18) Gln->pyro-Glu(-17) Acetylation (+42)
Semi-tryptic searched	Yes	No	Yes	Yes	Yes
Precursor tolerance	20 ppm	20 ppm	10 ppm	10 ppm (post-DTARefinery)	20 ppm
Missed Cleavages	No limit	<5	<2 post search	No limit	By search engine
False Discovery Rate	1% PSM	1% PSM	1 % peptide	1% Peptide	1% PSM

1

- **CDAP** is a standard for analyzing proteomic data proposed by CPTAC
- Some historical data from previous phases deviates from **CDAP**
 - Individual research institutions were free to select different analysis methods for quantifications based on specific needs
- However the Common Data Analysis Pipeline (CDAP) methods are the standard moving forward

Bioinformatics: Gene Level Assembly



- Peptide identifications used to map abundance values to genes; multiple peptides per gene must be considered
- Peptides sequences which are not unique to single gene, **shared peptides**, can be excluded depending on data analysis choices
- Remaining peptides (min 2) aggregated at gene level using max abundance per peptide

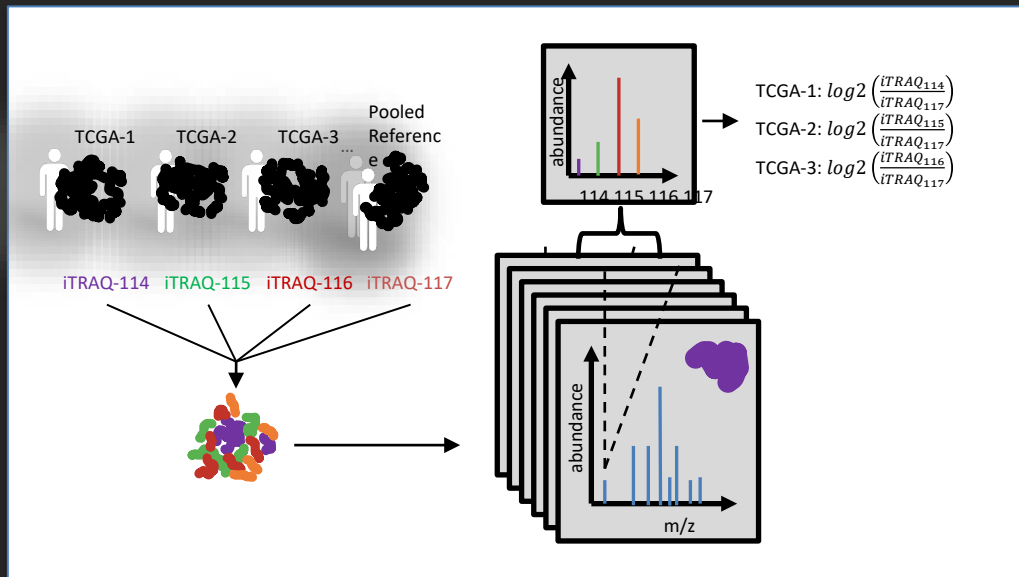
CDAP more info

- <https://pdc.esacinc.com/data-dictionary/harmonization.html>
- <https://www.ncbi.nlm.nih.gov/pubmed/26860878>

Proteomic technologies

- Two versions:
 - Label free
 - iTRAQ

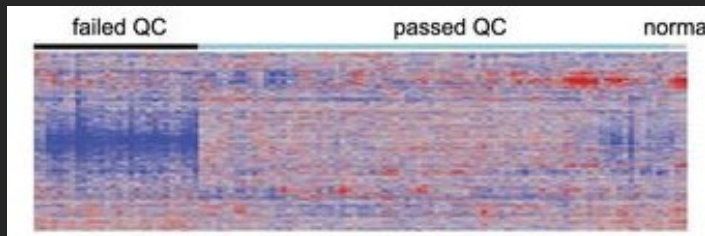
iTRAQ Data Format (BRCA & OV)



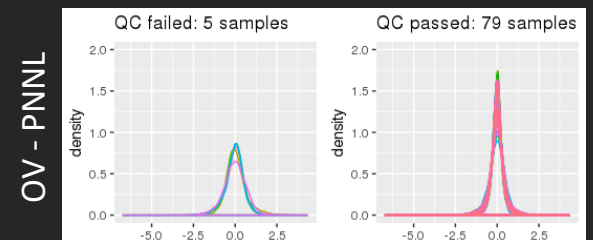
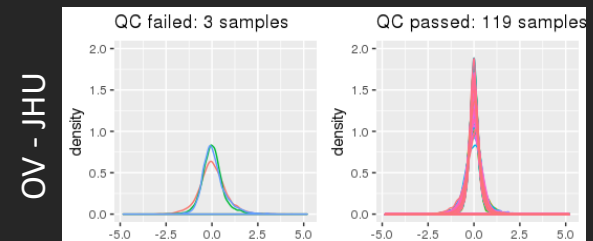
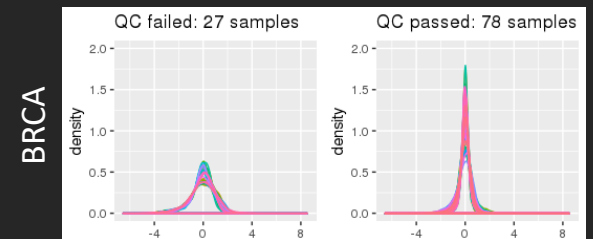
- iTRAQ: Isobaric tags for relative and absolute quantitation – low molecular weight ions are used to tag peptides from each sample
→ tighter quantification
- 4-plex measurements are made: where 3 samples are compared against pooled reference from 40 tumors
- Protein abundance reported as relative log₂ ratio between sample versus reference

iTRAQ Processing: Quality Control issue

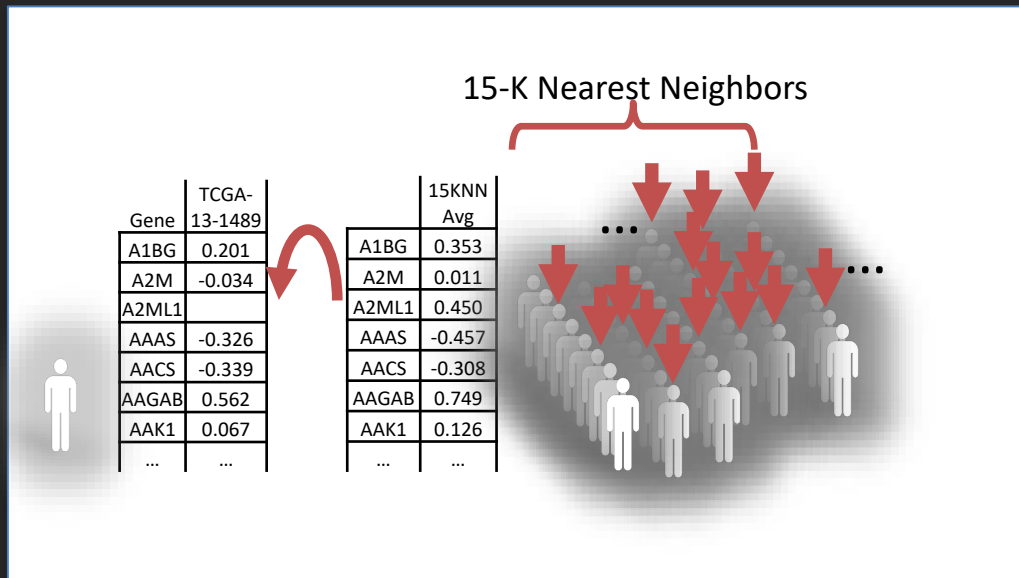
- Mertins et al. identified compromised samples with **excessive low abundance proteins due to protein degradation**



- Bimodal or skewed protein abundance distribution
- Std Dev is a natural QC statistic
- Two component Gaussian mixture model using Std Dev per sample to identify compromised samples



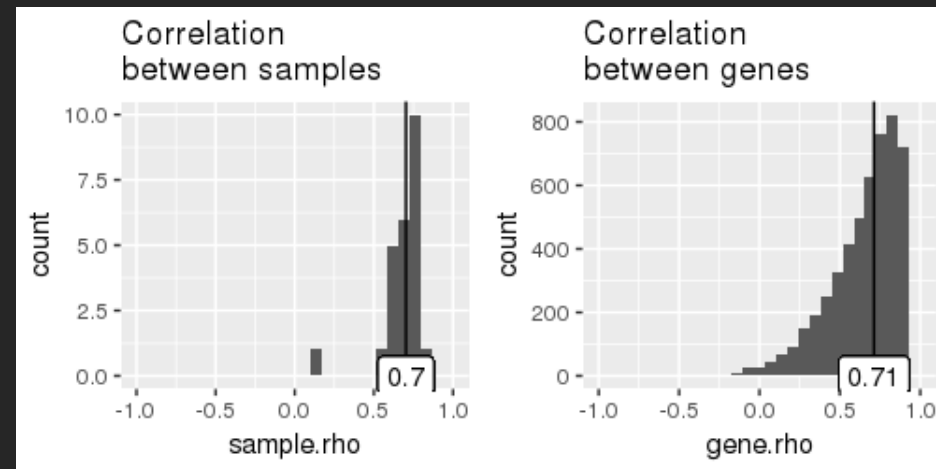
iTRAQ Processing: Missing Data



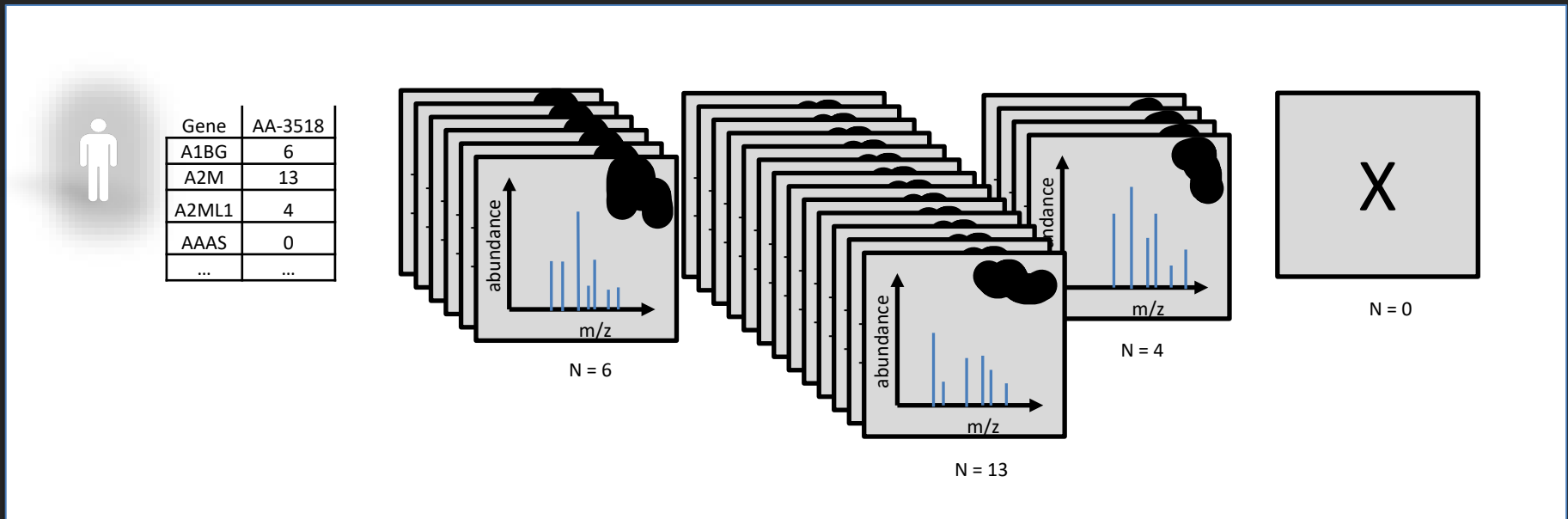
- iTRAQ measurement produces quality measures to **remove low confidence measurements** → missing values
- Filtering: samples and genes with >25% missing data were removed from analysis
- Imputation
 - To impute missing values each sample is matched to 15 patients with the most similar protein abundances across all measured genes (Euclidian Distance)
 - Use the **average abundance among neighbors** to impute missing value

iTRAQ Processing: Merging Multi-Site Datasets

- Ovarian Cancer samples measured at two sites Johns Hopkins University and Pacific Northwest National Laboratory
- 32 samples measured at both sites show high correspondence :
 - median Spearman Rho of 0.70 for intra-tumor variation
 - median Spearman Rho of 0.71 for inter-tumor variation
- To unify datasets
 - Remove duplicate samples from JHU
 - Combine datasets and remove remaining bias using **ComBat** batch correction.



Label-free Data Format (COADREAD)

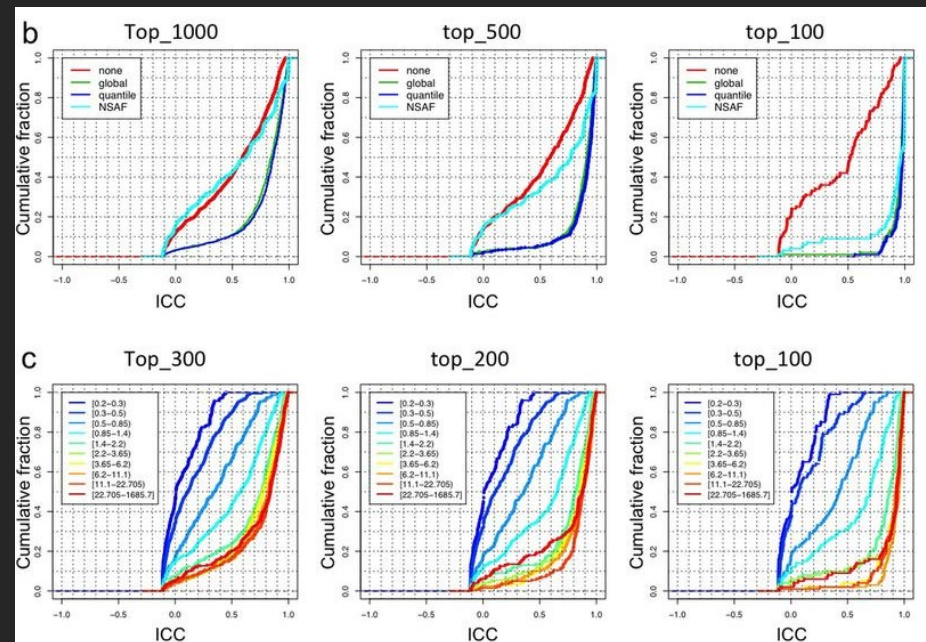


- Unlike iTRAQ, label-free quantification only provides absolute values per sample
- Protein abundance is quantified using the number of spectra measured, spectral counts

Label-free Processing: Zhang et al.

Normalization

- To guide analysis Zhang et al. (COADREAD) quality control data set and examined groups of peptides which map to same protein and assessed quality using intraclass consensus (ICC)
- Avg value of 1.4 across samples set as **min threshold for low abundance peptides**
- **Quantile normalization** used to make distribution of protein abundance in each sample comparable
- Last values are **log₂ transformed**



Processed Data Overview

- CDAP preprocessing & additional QC
 - Mapping mass ladders to peptides
 - Normalization
 - QC
 - Log2 transform
 - Missing value estimation
 - Batch correction
- Resulting data set can be processed similarly to other genomic data using similar statistical tools

Dataset	Number Samples	Number Proteins
BRCA – BROAD Institute	78	8662
OV – JHU & PNNL	150	5233
COADREAD – Vanderbilt University	95	2889



CTPAC data analysis

Linking with other omics data



Stanford
MEDICINE

CPTAC multi-omics data fusion

- Same samples also have (phase 2 & 3)
 - RNA sequencing
 - DNA methylation
 - DNA copy number
 - Etc.
- Example from our work
 - Linking DNA methylation with Proteomic data

CPTAC data: overlapping samples & genes/proteins

	Nr Genes	Samples with mRNA & protein expression		
BRCA	2514	78		
COADREAD	2848	85		
OV	1896	168		

- Only focus on genes with **both** mRNA & protein expression

CPTAC data

	Nr Genes	Samples with mRNA & protein expression	Samples with methylation data	Sample with normal methylation data
BRCA	2514	78	972	123
COADREAD	2848	85	614	78
OV	1896	168	582	8

- All samples with DNA methylation data used for methylation states
- Varying # normal samples

ProteoMix: a statistical model

For each CpG site

Raw methylation data

Identify
mixture
components

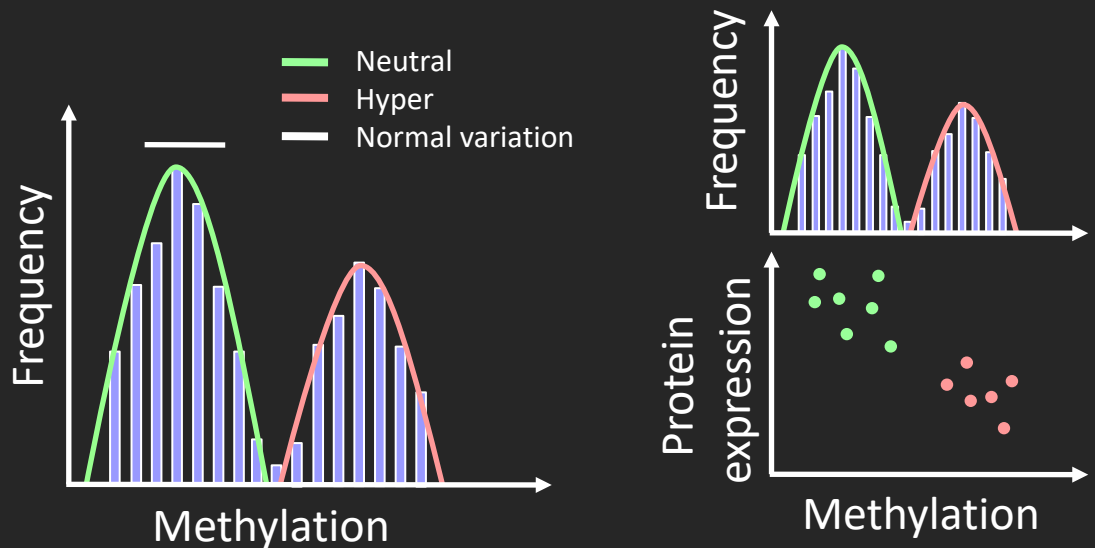
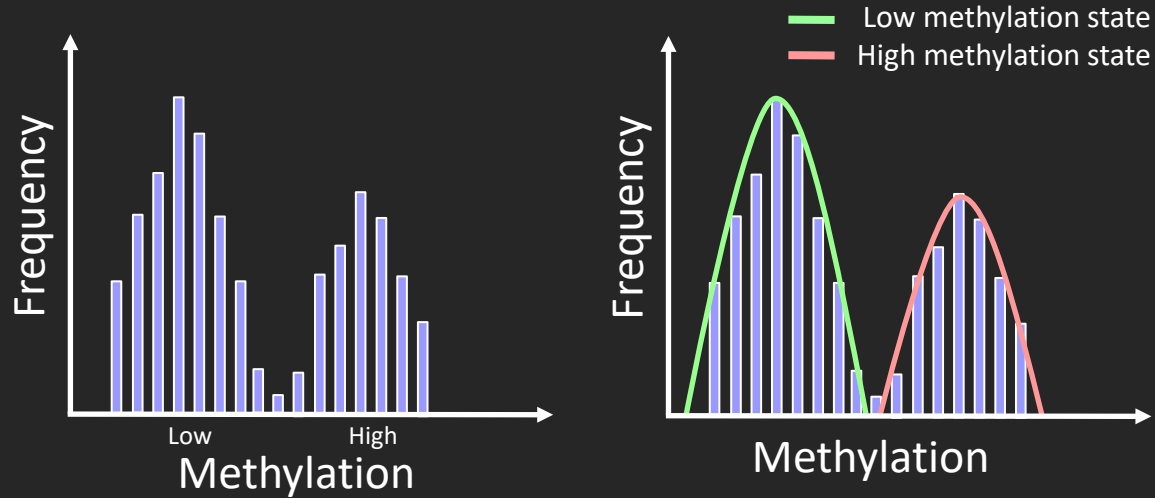
Identified methylation states

Map normal DNA
methylation
variation

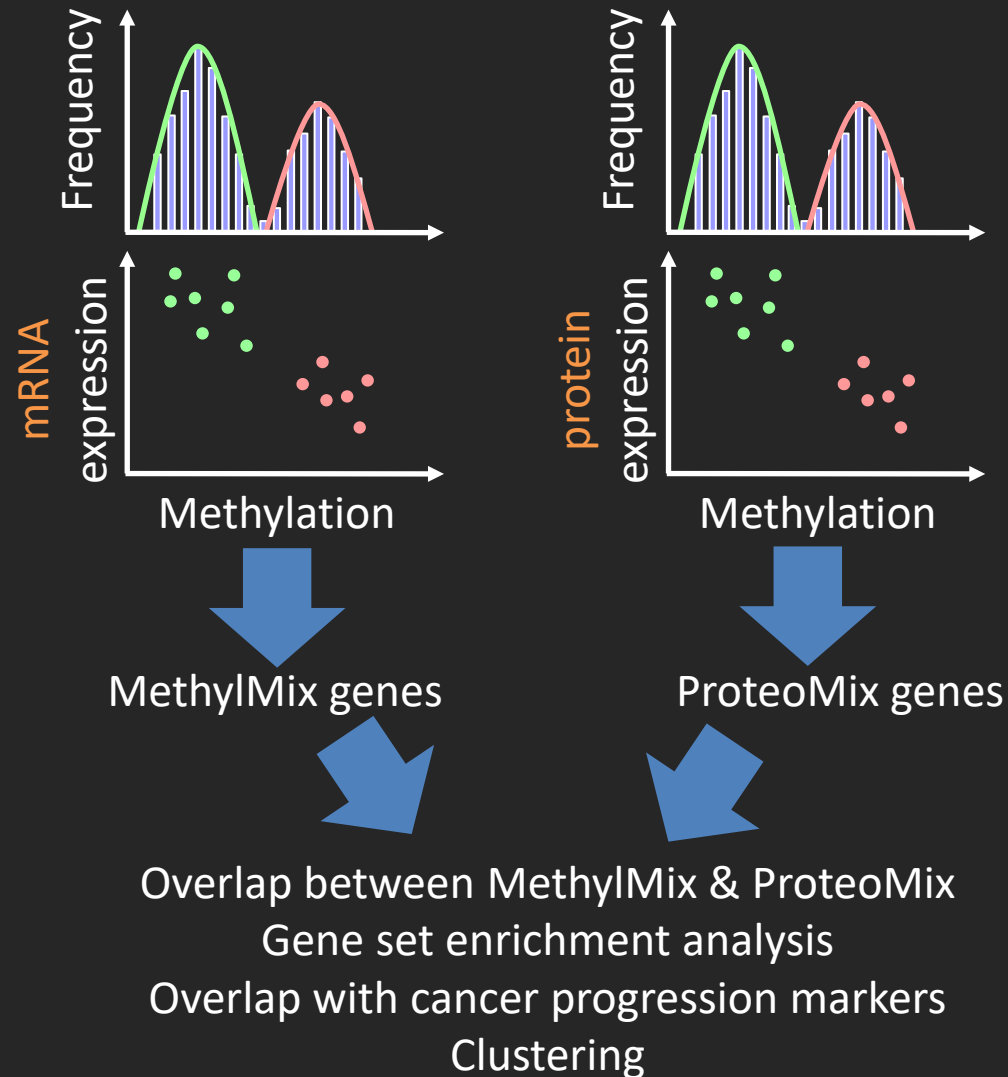
Differential methylation

Incorporate
protein expression
data

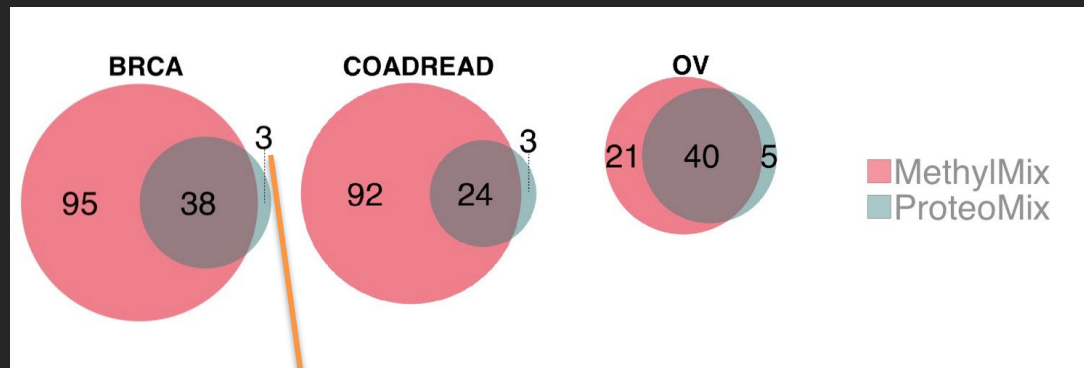
*Differential & Functional
methylation states*



MethylMix vs. ProteoMix

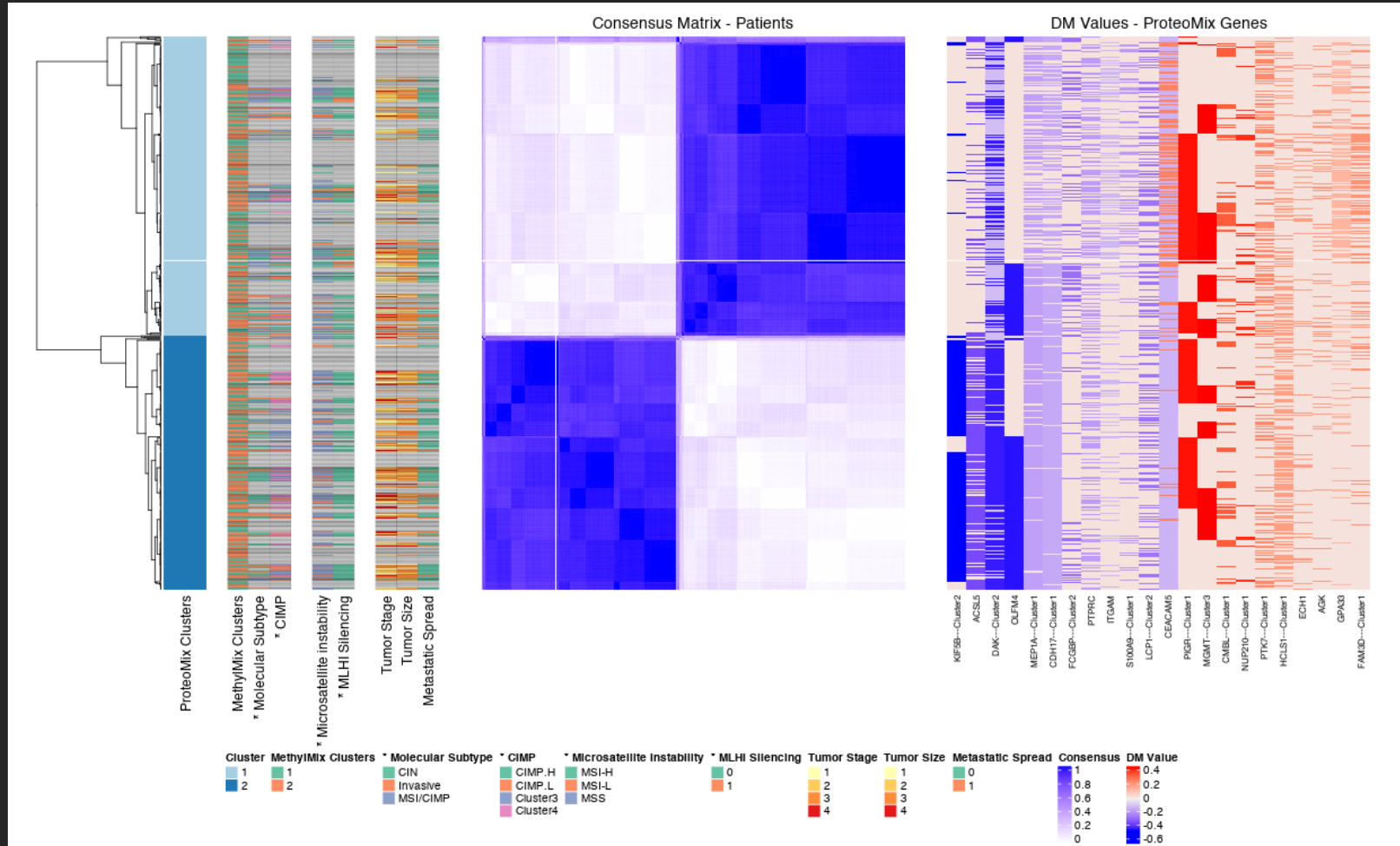


Results



- Breast cancer:
 - hypo-methylation in the UTR of **EHF** well-studied transcription factor involved in HER2 mediated epithelial differentiation
 - knockdown of **EHF** has been shown to inhibit tumor invasion and proliferation

Colorectal cancer



Significant association with CIMP vs non-CIMP



CPTAC data

Radioproteomic maps

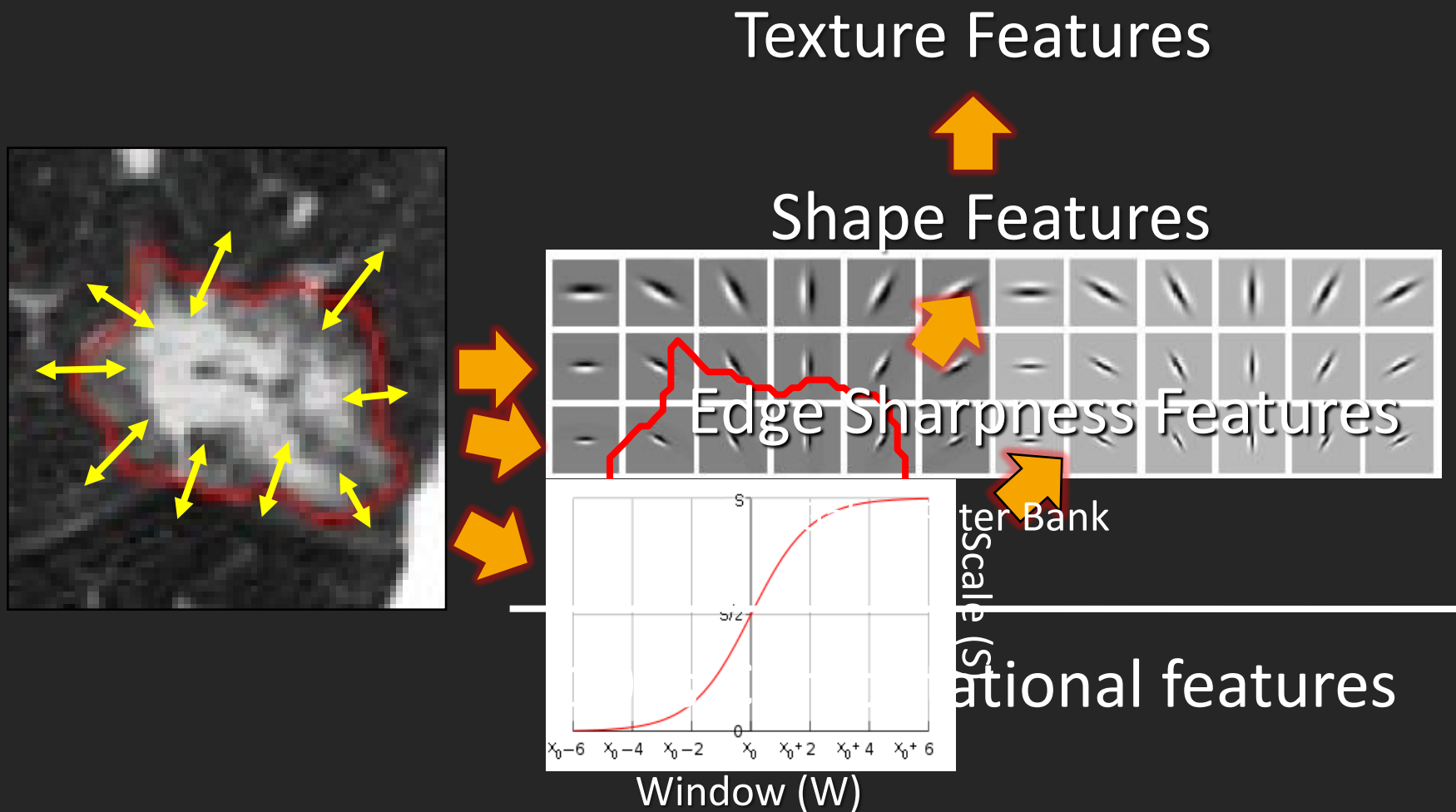
Ongoing work



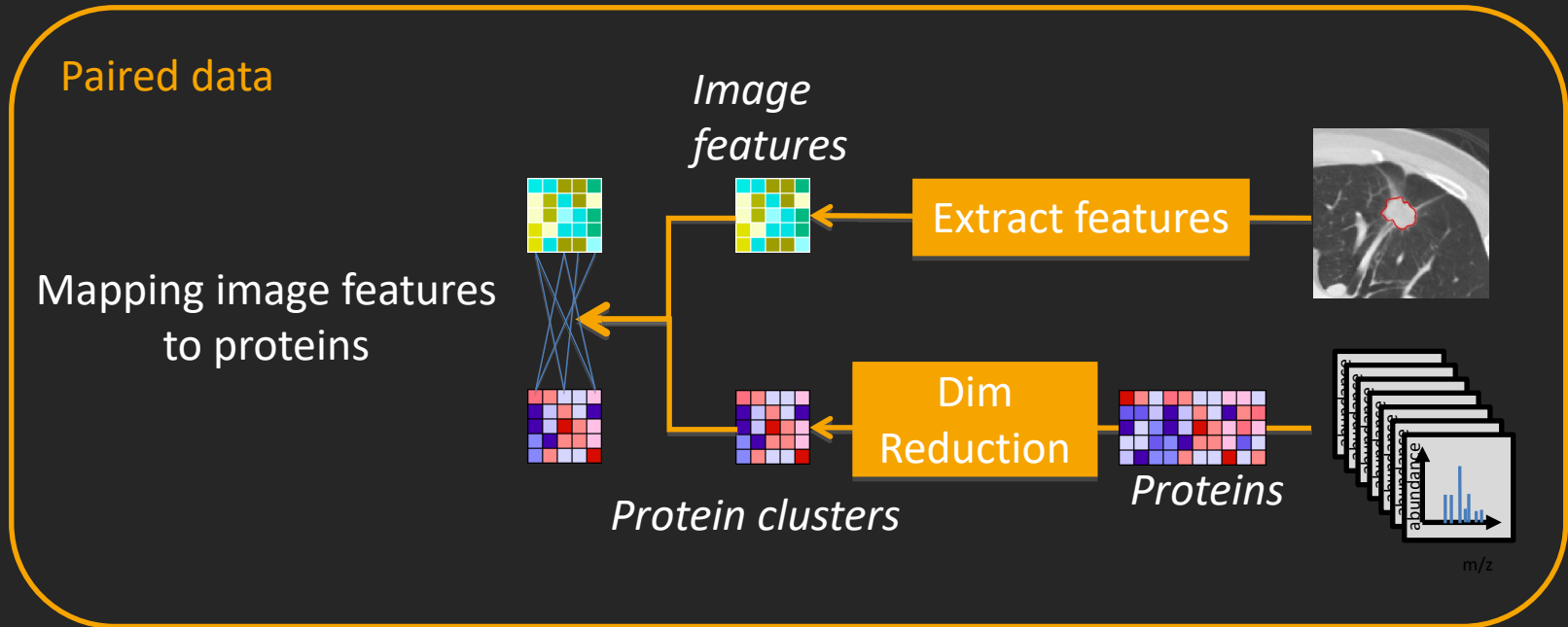
Stanford
MEDICINE

Radiomics/Quantitative imaging

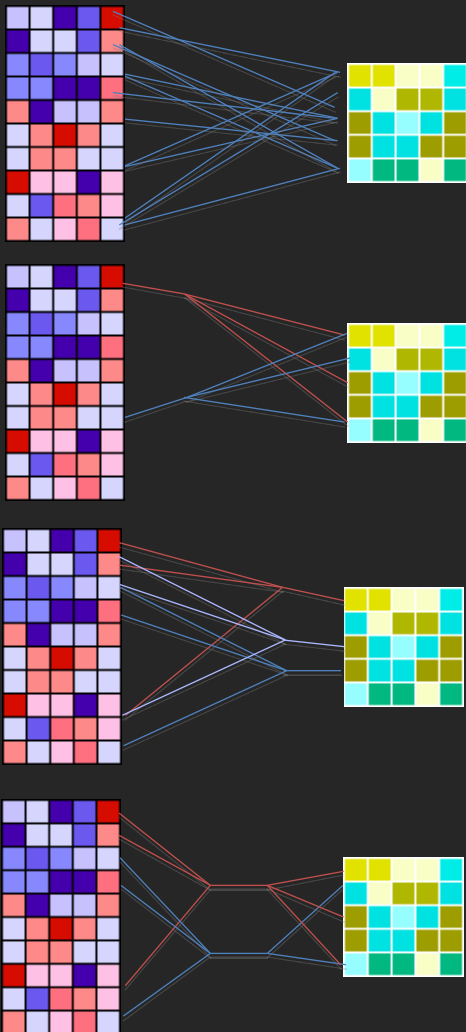
Extraction of computational features



Radiogenomics mapping



Radioproteomic mapping methods



- Integration methods
 - Two-way univariate
- Univariate – multivariate
 - Multivariate model of genes
 - Multivariate model of image features
- Two-way multivariate

Applications

- Predict protein expression clusters from imaging
 - Non-invasive biomarkers
- Predict imaging phenotype from protein data
 - Study how pathways lead to imaging phenotypes
 - Annotate protein function
- Compare with and validate in “traditional” **radiogenomics** maps built on RNA expression data



Acknowledgements

Gevaert Lab

Majed Magzoub

Kevin Brennan

Hong Zheng

Jay Shinde

Pritam Mukherjee

Lucas Patel

Radiology

Parag Mallick

Biomedical Data Science

Robert Tibshirani

Trevor Hastie

Funding



NIH/NIBIB R01 EB020527

NIH/NCI U01 DE025188

NIH/NCI R01 CA184968

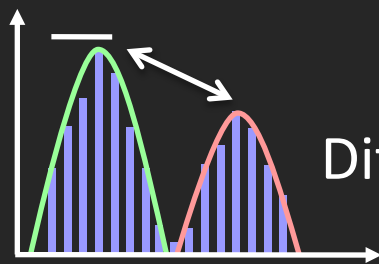
NIH/NCI R01 CA176299



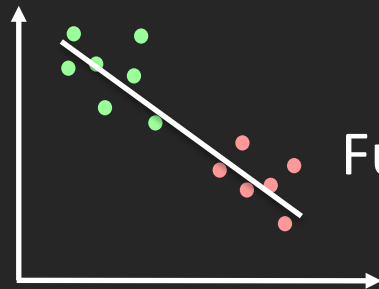
Stanford
MEDICINE



MethylMix & ProteoMix: identifying DNA methylation-driven genes in cancer



Differential



Functional

MethylMix & ProteoMix
R package

Available on Bioconductor & github:
<https://github.com/gevaertlab>

GenePattern module
in development

BioRxiv:



Stanford
MEDICINE