# P·MART

# Interactive Informatics Resource for Research-driven Proteomics

CONTACT:
BOBBIE-JO WEBB-ROBERTSON
BJ@PNNL.GOV

# Exploring Proteomic Data - The Need

▶ Wealth of proteomic cancer data available in the public domain

▶ Peptide level data difficult to access

▶ Protein level data

- ■ Often not reproducible
- ■ Doesn't offer in depth view of peptide evidence
- ■ Isoforms are usually lost

OFFICE OF CANCER CLINICAL
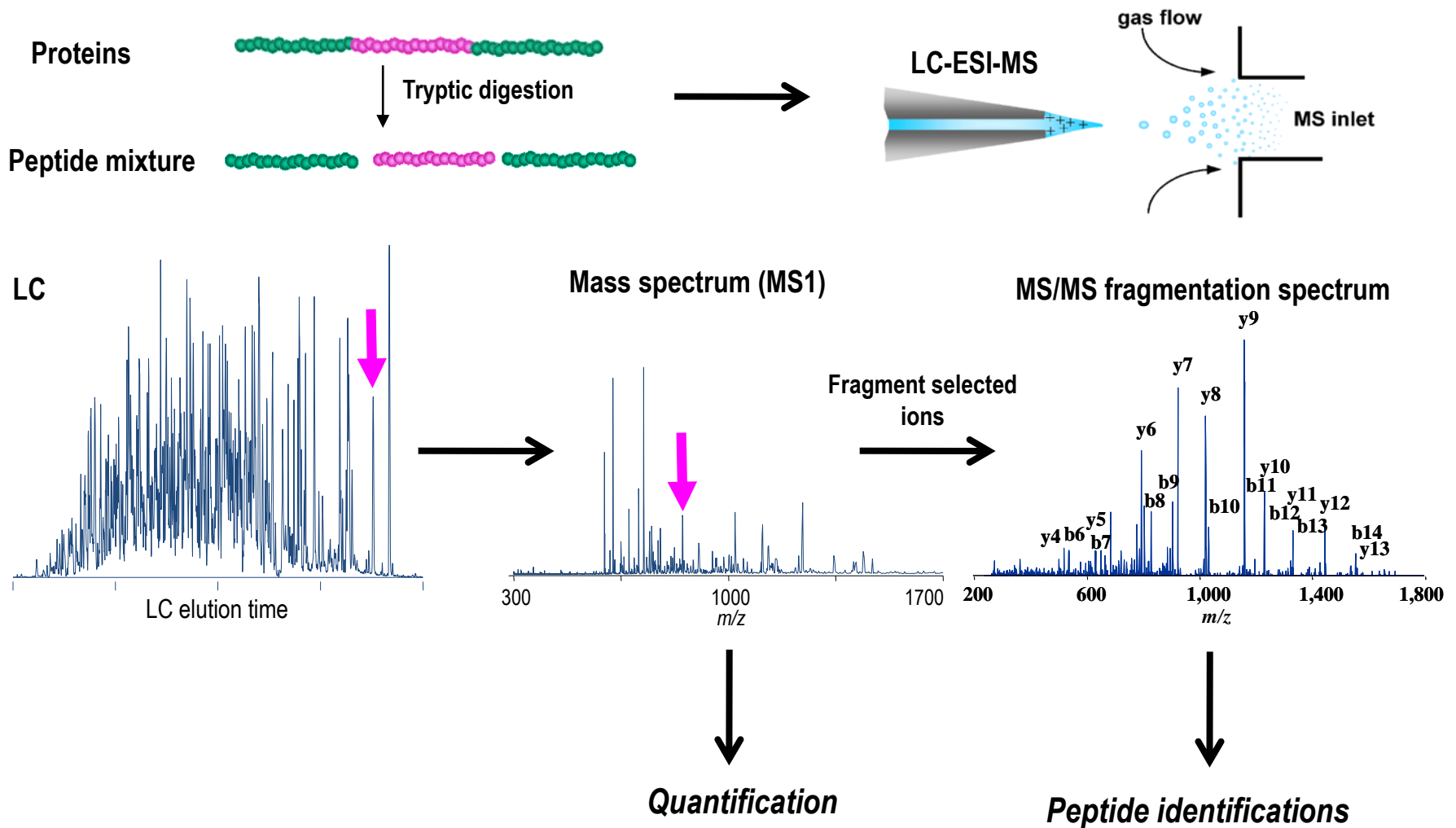**PROTEOMICS** RESEARCH
CPTAC DATA PORTAL

**Early Detection Research Network**
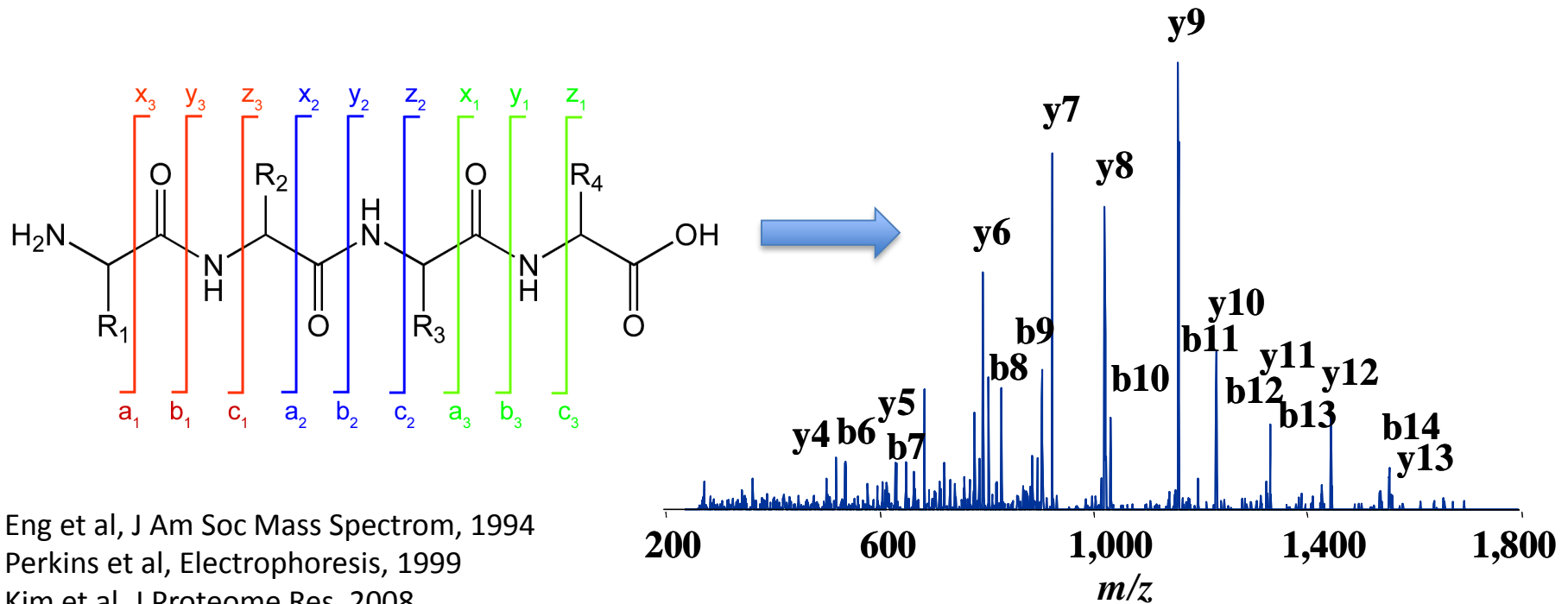*Biomarkers: the key to early detection*

**Anyone who wants to explore peptide and/or protein level data associated with cancer.**

# Typical Bottom-up Global MS-based Proteomics Workflow

# Peptide Identification

In theory the peptide should fragment into a MS/MS spectra based on measurable mass shifts
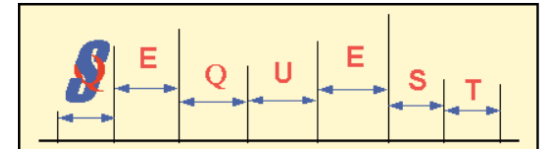


Eng et al, J Am Soc Mass Spectrom, 1994
Perkins et al, Electrophoresis, 1999
Kim et al, J Proteome Res, 2008

# Peptides are "Initially" Identified via a Database Search



## Protein Database

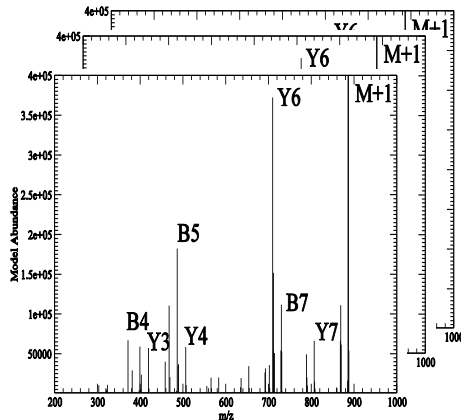>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]
LCLYTHIGRNIYYGSYLYSETWNKGIMLLLITMATAFMGYVLPWKQMSFWGATVITNLFSAIPYIGTNLV
EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFHPYYTIKDFLG
LLILILLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGGVLALFLSIVIL
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFLPIAGX
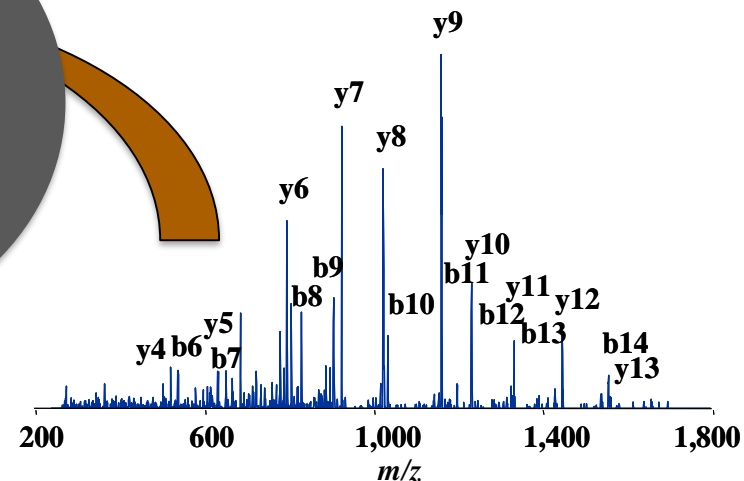IENY

LCLYTHIGR•N
R •NIYYGSYLYSETWNK•G
.......

**Does the Experimental Match a Theoretical Spectra?**

Peptide Match (scoring threshold)
K •IRYQVTSVSNK•G

MS/MS Experimental fragmentation spectrum

# Peptides are "Initially" Identified via a Database Search

# The challenges for proteomics compared to genomics?

**http://proteomics.cancer.gov/whatisproteomics**

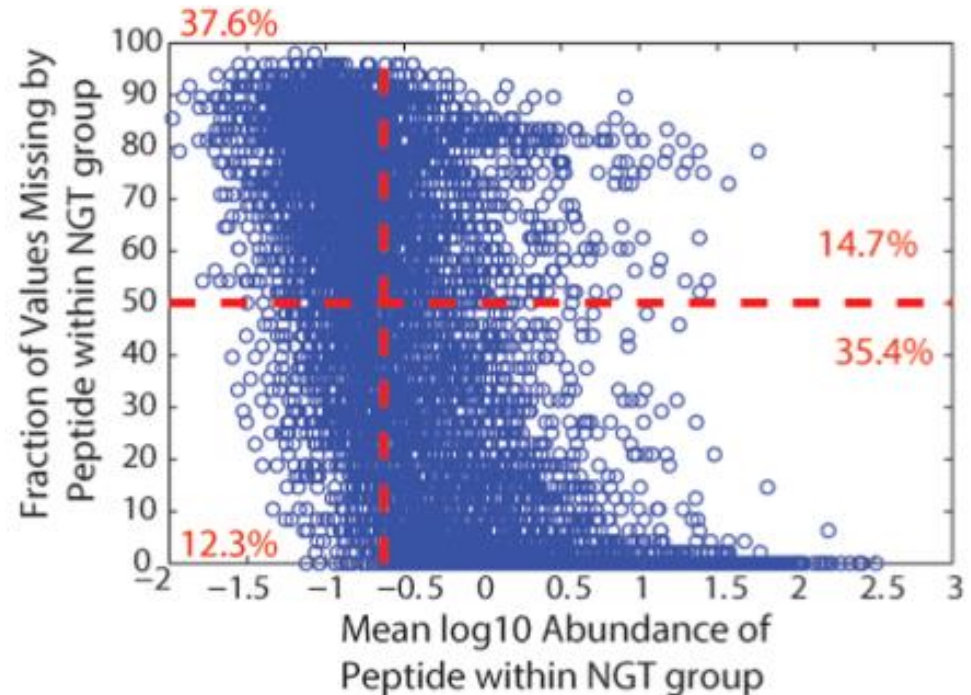*Exact quote from the cited website*

▶ One gene can encode more than one protein (even up to 1,000). The human genome contains about 21,000 protein-encoding genes, but the total number of proteins in human cells is estimated to be between 250,000 to one million.

▶ Proteins are dynamic. Proteins are continually undergoing changes, e.g., binding to the cell membrane, partnering with other proteins to form complexes, or undergoing synthesis and degradation. The genome, on the other hand, is relatively static.

▶ Proteins are co- and post-translationally modified. As a result, the types of proteins measured can vary considerably from one person to another under different environmental conditions, or even within the same person at different ages or states of health.

Additionally, certain modifications can regulate the dynamics of proteins.

▶ Proteins exist in a wide range of concentrations in the body. For example, the concentration of the protein albumin in blood is more than a billion times greater than that of interleukin-6, making it extremely difficult to detect the low abundance proteins in a complex biological matrix such as blood. Scientists believe that the most important proteins for cancer may be those found in the lowest concentrations.

# The challenges of statistical analysis of proteomics data

► Missing data

■ There is a lot of it
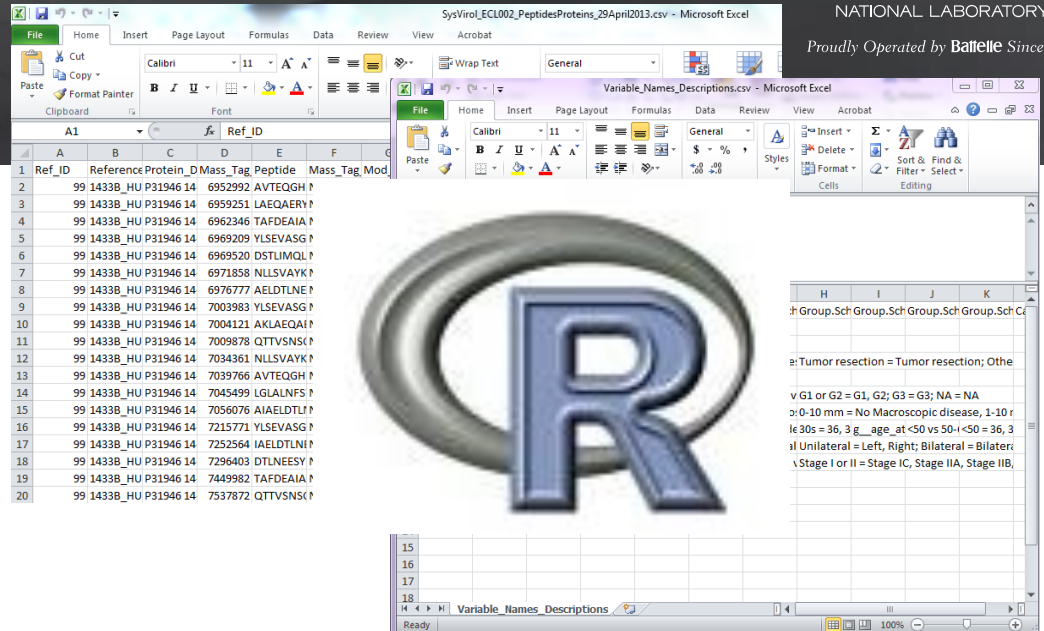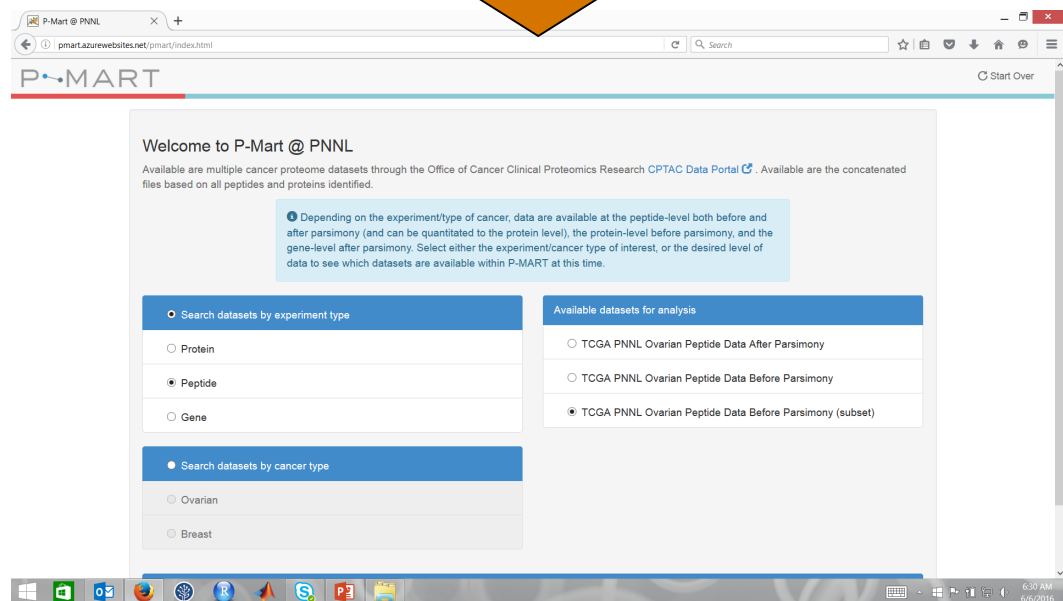
■ The mechanisms by which it is missing is complex

*E.g., In CPTAC ovarian cancer data there is not a single peptide that has 100% coverage across all biological samples.*
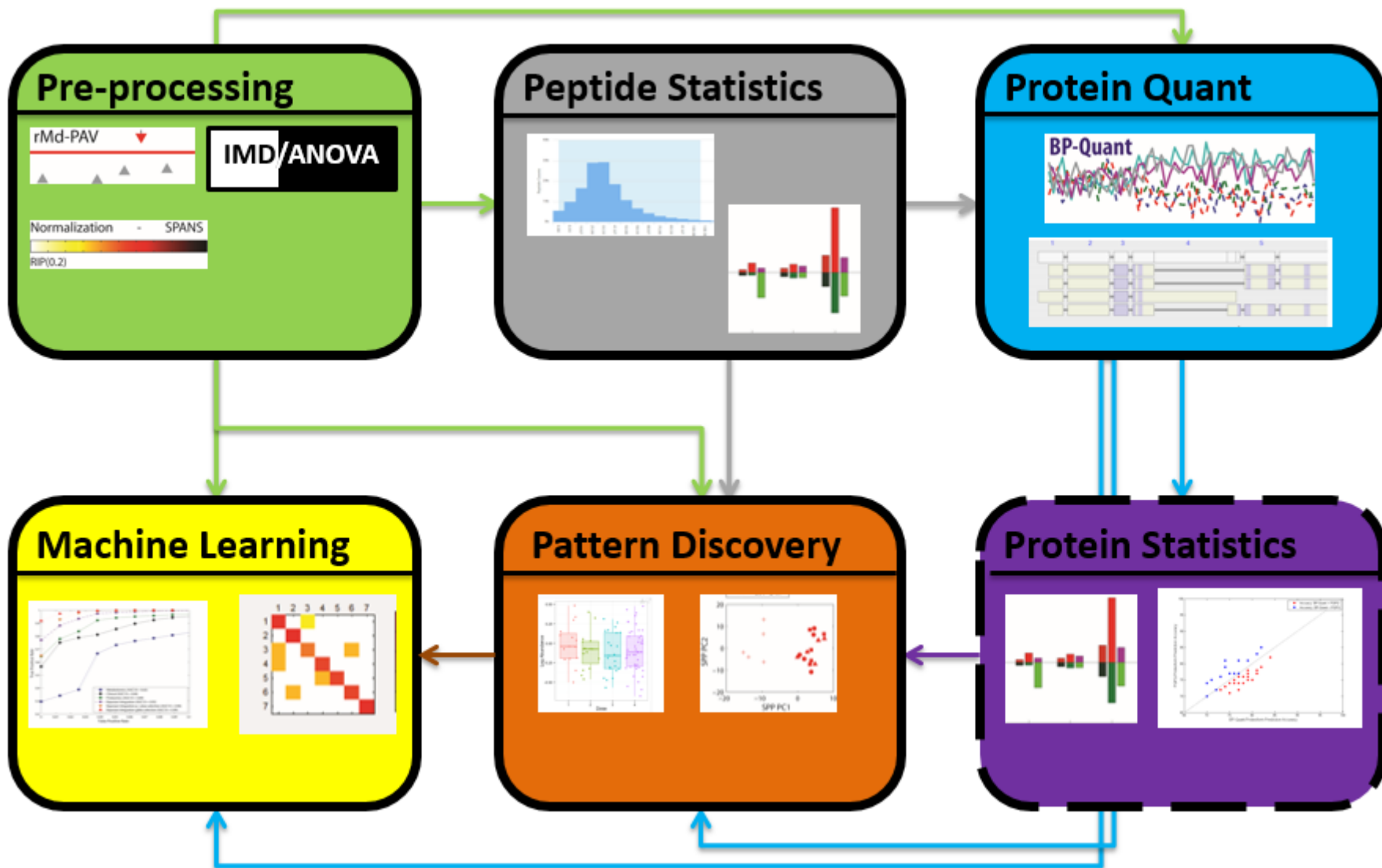


► Mis-identifications

► Variability (technical and biological

# P-Mart Goal

Enable Biomarker Discovery from large complex proteomics datasets by facilitating reproducible statistical processing and complex exploration of high-throughput peptide and protein level data generated through MS robust to missing values.

# P-Mart Capabilities

# P-Mart - Data

▶ Access to peptide, protein or gene data

- Peptide data is extracted from the PSM files
- Protein and Gene data is as provided by CPTAC DCC

# P-Mart - Data

▶ Summaries of the data available

▶ All clinical outcomes reported by CPTAC can be used as a primary factor of interest or as a covariate

# P-Mart – Customized Workflows

▶ User may create a workflow (under constraints)

▶ P-Mart will make a suggestion based on data type

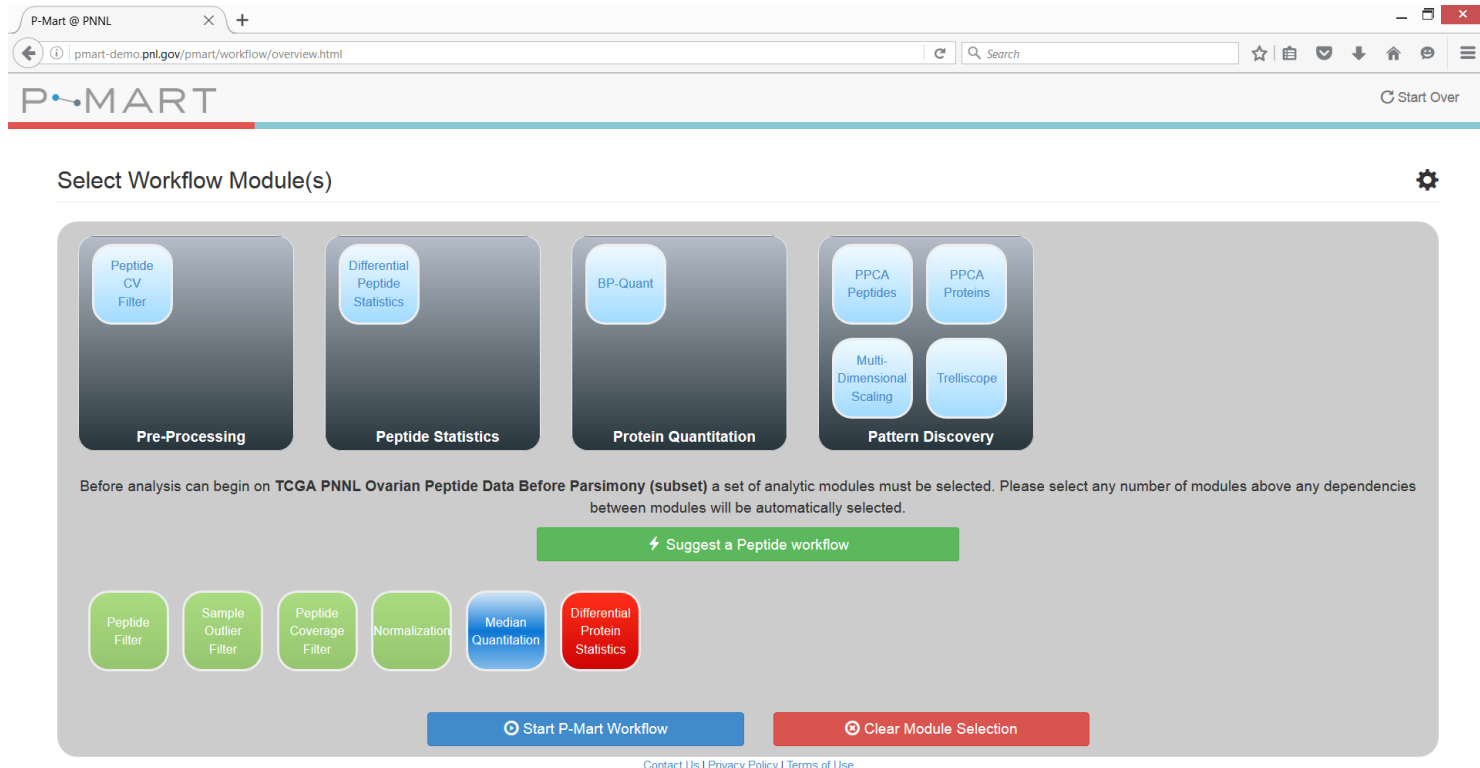# P-Mart – Customized Workflows

▶ User may create a workflow (under constraints)

▶ P-Mart will make a suggestion based on data type
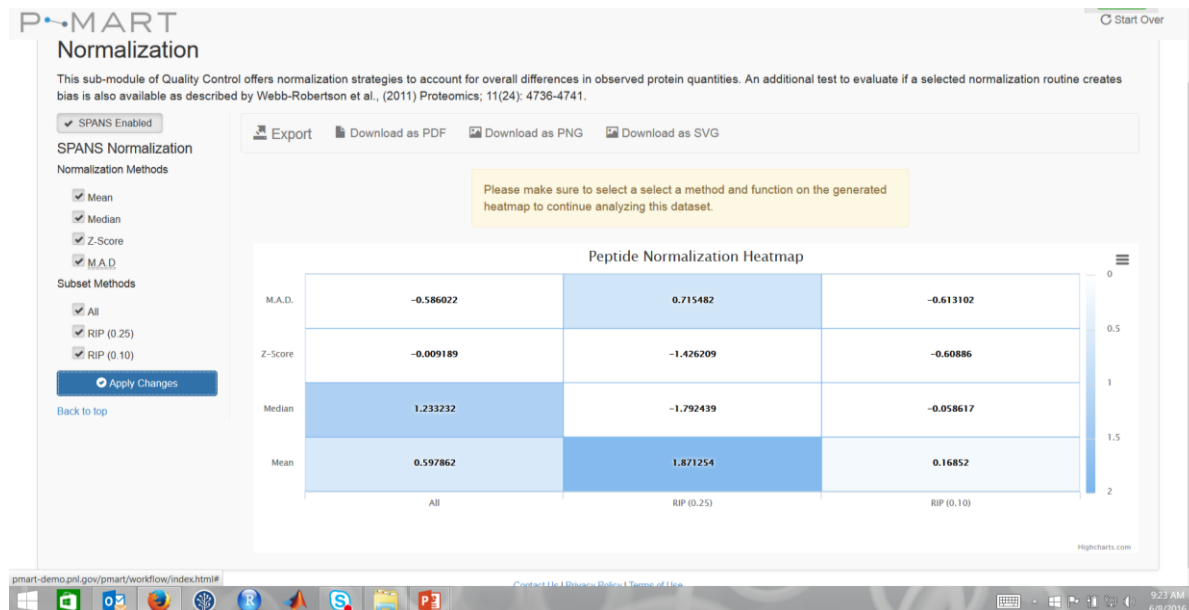
# P-Mart – Pre-Processing

▶ Samples

■ Looks at statistical metrics of the distribution of peptide abundances for anomalies

▶ Peptides/Protein

■ Statistical evaluation of peptide/protein coverage coefficient of variation within peptides

▶ Samples/Peptides

■ Data normalization approaches evaluated to identify the most robust approaches

# P-Mart – Protein Quantification

► Offers multiple standard protein quantification methods

► Uses unique approaches to identify proteoforms



BP-Quant

Proteome Complexity

Genome
~20-25,000 genes

Alternative promoters
Alternative splicing
mRNA editing

Transcriptome
~100,000 transcripts

Post-translational
modifications

Proteome
>1,000,000 proteins

# P-Mart – Exploratory Data Analysis

▶ Perform PCA in a manner that is robust to missing data

▶ Allow exploration of potential biomarkers through smart queries defined by the user

# P-Mart – Exploratory Data Analysis

▶ Perform PCA in a manner that is robust to missing data

▶ Allow exploration of potential biomarkers through smart queries defined by the user

# P-Mart – Exploratory Data Analysis

► Perform PCA in a manner that is robust to missing data

► Allow exploration of potential biomarkers through smart queries defined by the user

# P-Mart – Exploratory Data Analysis

► Perform PCA in a manner that is robust to missing data

► Allow exploration of potential biomarkers through smart queries defined by the user

# P-Mart – Exploratory Data Analysis

► Perform PCA in a manner that is robust to missing data

► Allow exploration of potential biomarkers through smart queries defined by the user
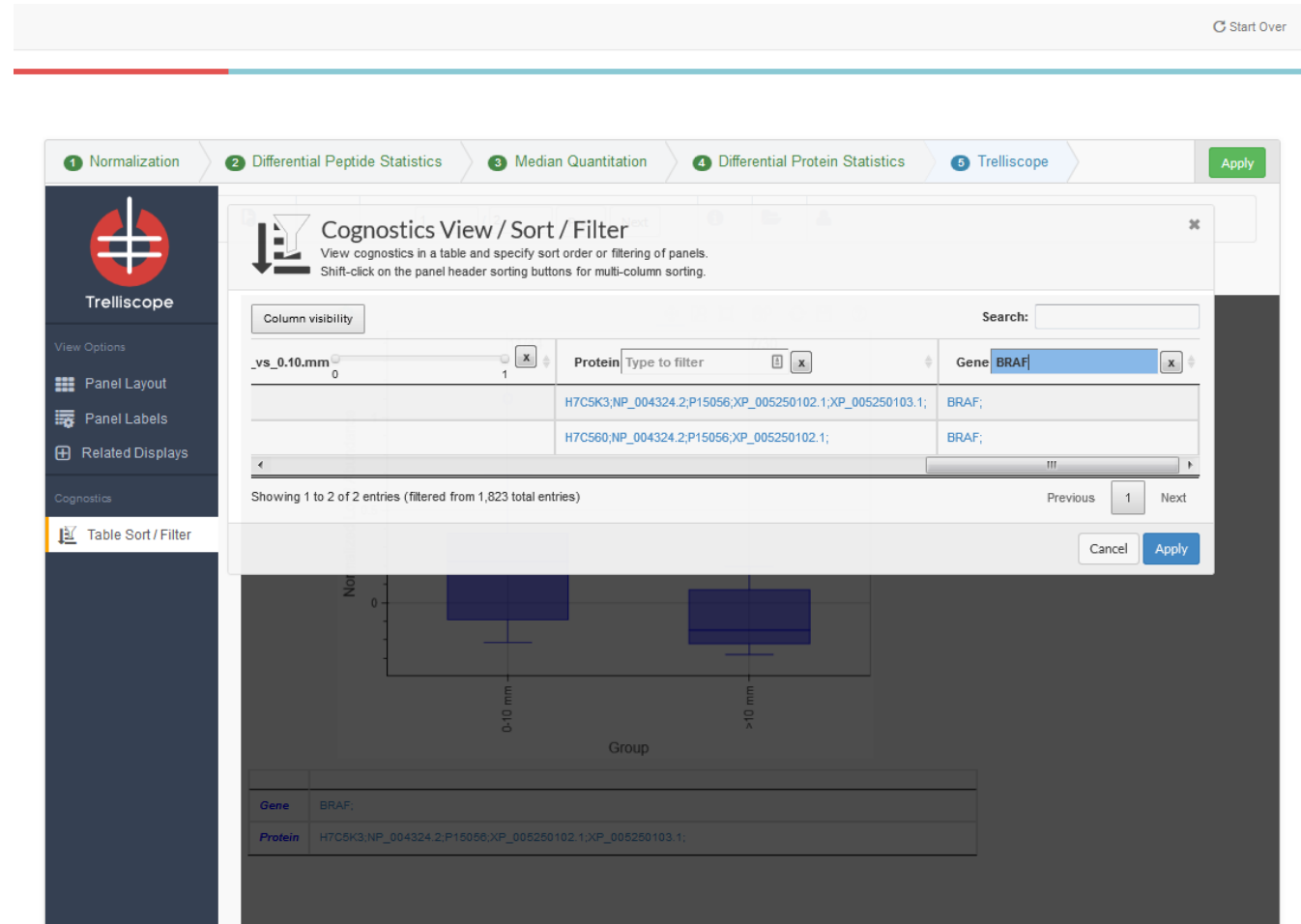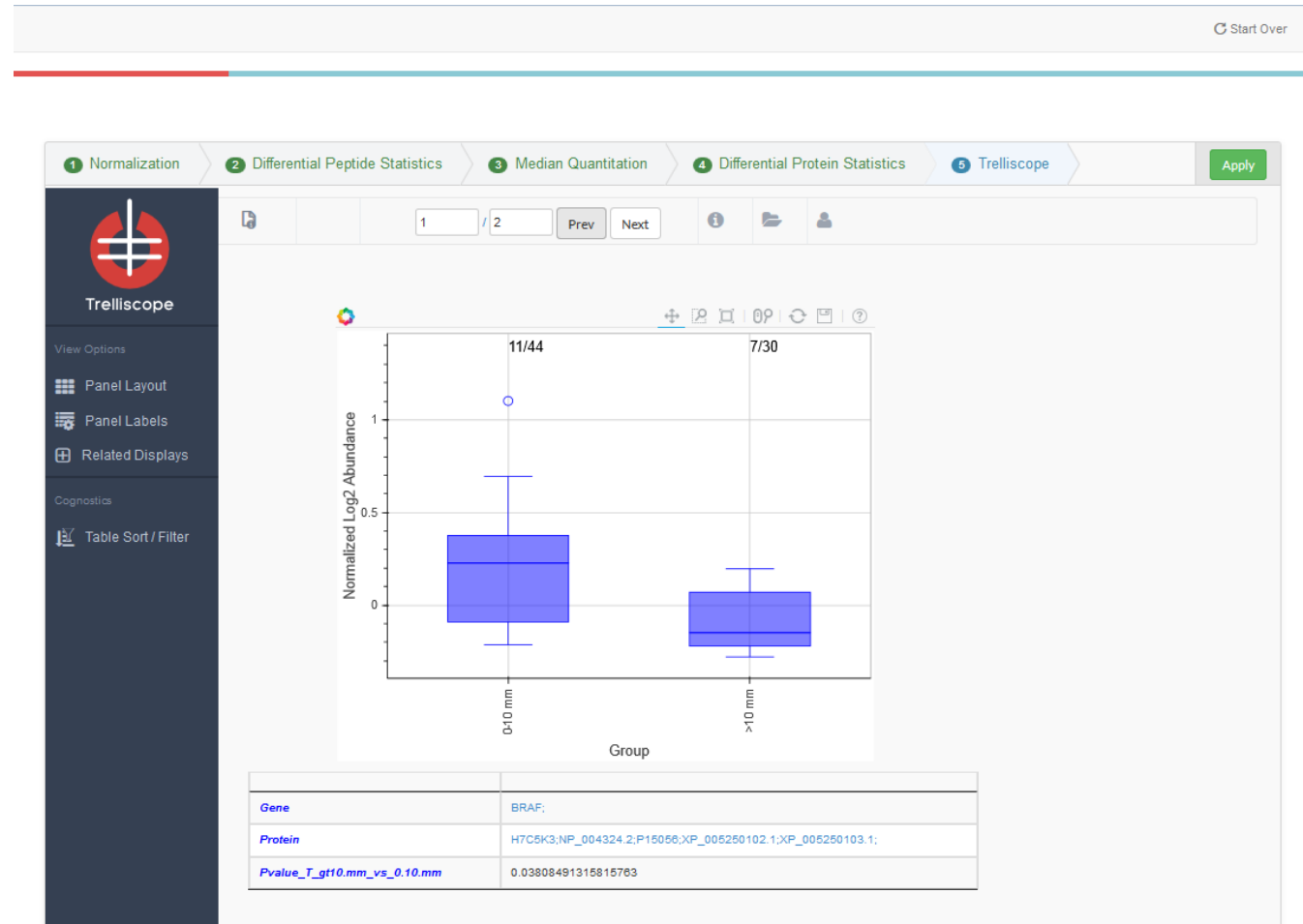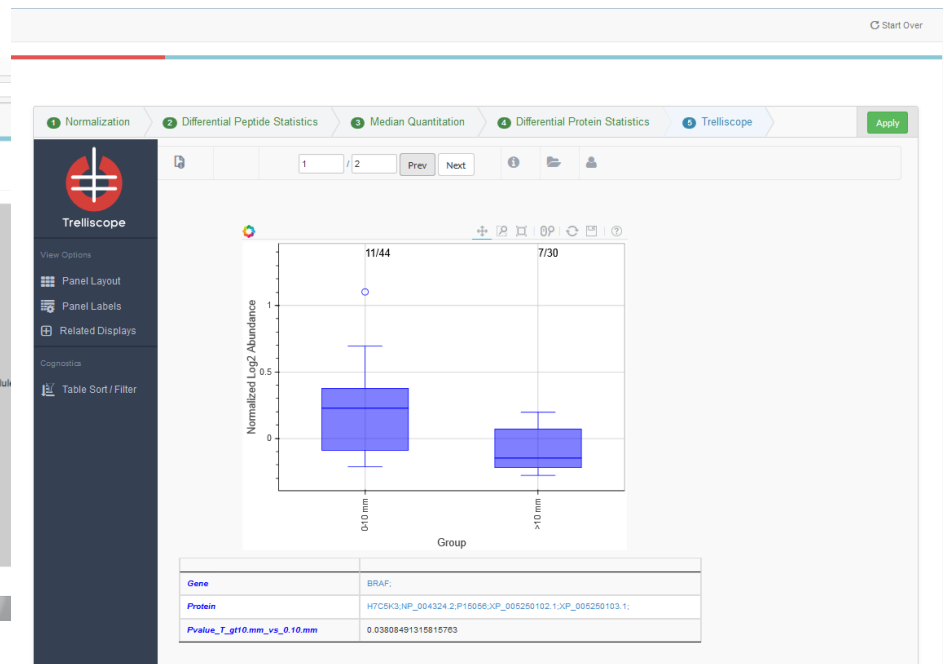
# P-Mart Outcomes

▶ Easy to use online software driven by visual interactions

▶ Customizable workflows documented for reproducibility

▶ Open-source code of statistical methods (R)

▶ Easy exports of data at multiple levels of processing

# P-Mart Status

► Completed

- ■ Website deployed with Azure cloud technology
  [http://pmart.azurewebsites.net/pmart/index.html](http://pmart.azurewebsites.net/pmart/index.html)

- ■ Statistical methods robust to missing data released on GitHub
  [https://github.com/MSomics-StatTools/MSomicsQC](https://github.com/MSomics-StatTools/MSomicsQC)

- ■ Peptide and Protein level statistics and associated meta-data can be queried to rapidly identify biomarkers of interest

► In Progress

- ■ Documentation allows for reproducible analyses
- ■ Exports allow the user access to data at various levels of quality control and processing to enable downstream analyses (e.g., pathway analysis)
- ■ Machine learning aides in validation and feature selection tasks to select candidates for validation

# Acknowledgements 1U01CA184783-01

Pacific Northwest
NATIONAL LABORATORY
*Proudly Operated by Battelle Since 1965*

Clinical Proteomic Tumor Analysis Consortium

VANDERBILT UNIVERSITY

BROAD INSTITUTE

ESAC
Enterprise Science And Computing

Pacific Northwest
NATIONAL LABORATORY

Early Detection Research Network
*Biomarkers: the key to early detection*

NASA Jet Propulsion Laboratory
California Institute of Technology