

Introduction to Mass Spectrometry-based Proteomics data generation and processing

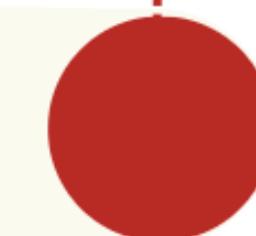
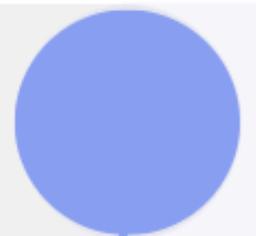
Anna Quaglieri 17/08/2022
Bioinformatics Data Scientist at Mass Dynamics

Mark Robinson laboratory
Department of Molecular Life Sciences, University of Zurich

About me



Degree in
Statistics



PhD Cancer
Genomics
2016



Research assistant
Population Genetics/Epigenetics



Data Science
Consulting
2020



Data Scientist in
Proteomics
Since Mid 2021

Mass Dynamics

**Where life scientists
discover, together.**

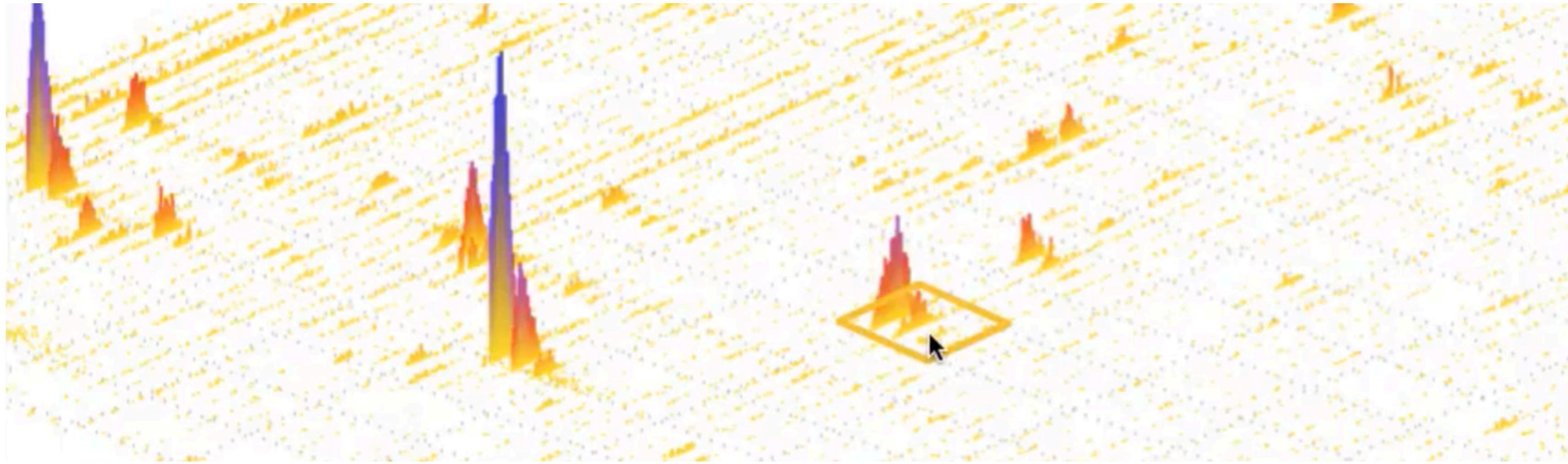


Mass Dynamics enables life scientists to transform complex and quantitative scientific data to knowledge - better, faster, easier and without compromise.

Mass Spectrometry data for proteomics

Outline

1. Mass Spectrometry (MS)-based proteomics data generation
 1. Data generation data
 2. Computational tasks in data processing
 3. Different technologies, different data
2. Common types of analyses
3. Commonly used methods and software
4. Useful links

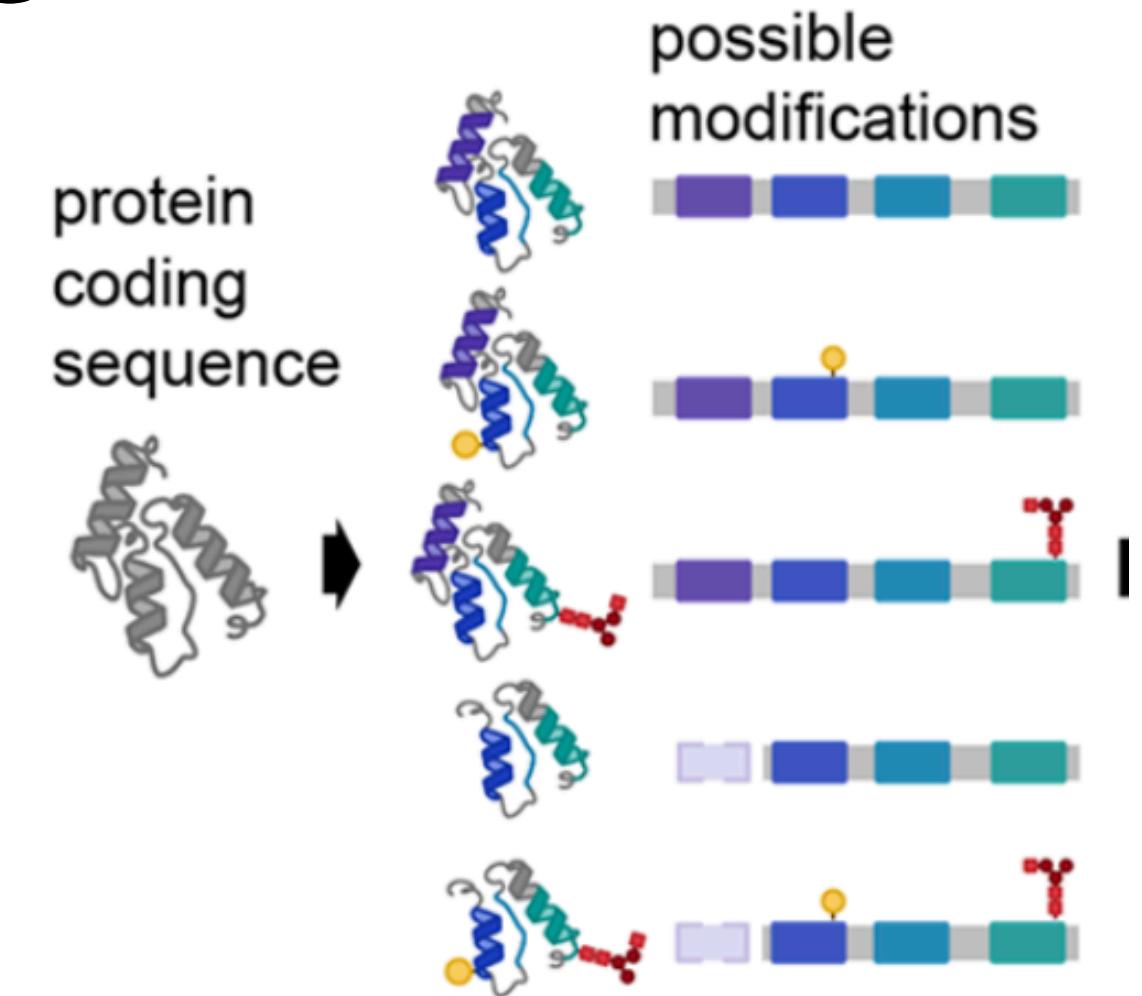


Mass Spectrometry based proteomics data generation

Proteomics vs Genomics

Proteomics

Large scale study of proteins: sequence, structure, abundance, modifications etc..

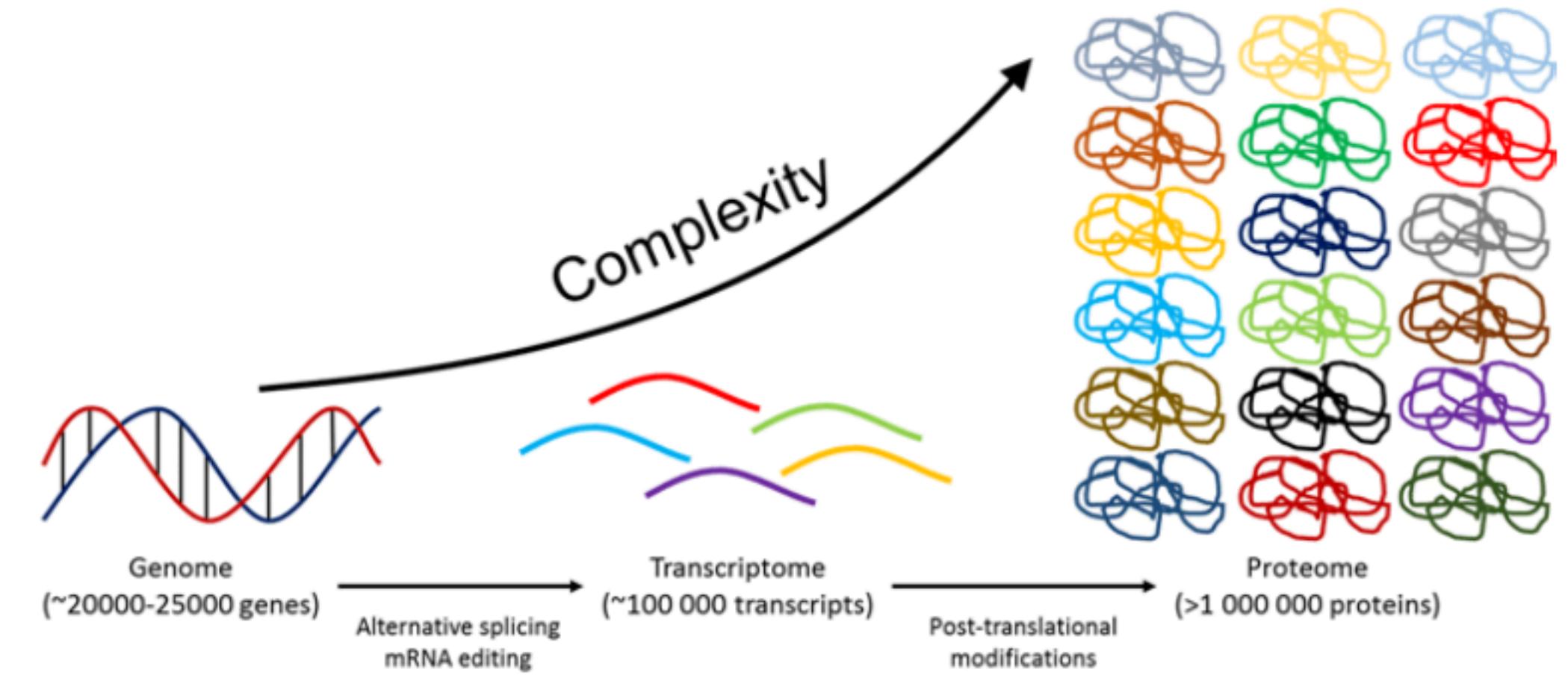


Proteomics is complex:

Human genome: ~20k genes

Human proteome can vary: ~50k - 500k proteins

We cannot amplify proteins



Proteoforms

Mass Spectrometry (MS)

We measure Mass Spectra

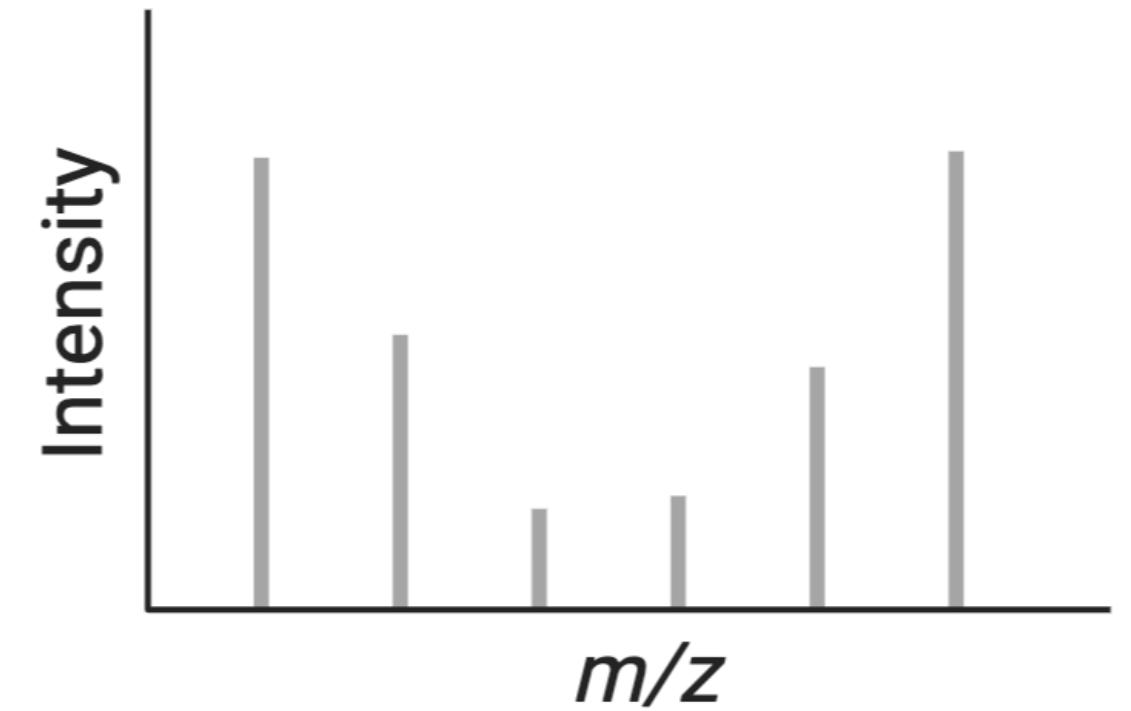
A technique used to measure the **mass-to-charge ratio (m/z) of ions**.

Ions: atoms or groups of atoms carrying one or more + or - electrical charges.

In MS-based proteomics we measure **peptide ions**: peptides that carry a particular electrical charge

MS is one of the most used technique to quantify proteins in a sample.

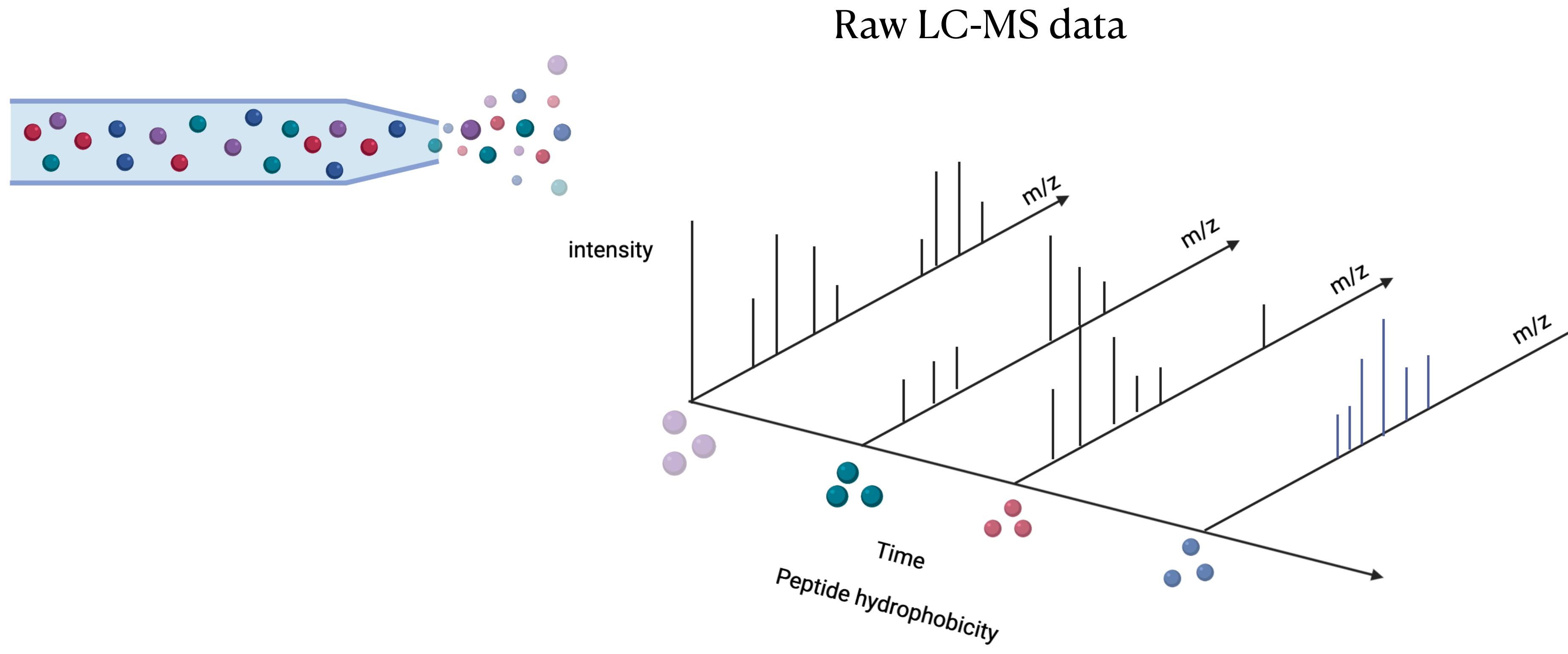
It is also used to quantify lipids and metabolites giving rise to broad fields of research known as **Metabolomics** and **Lipidomics**.



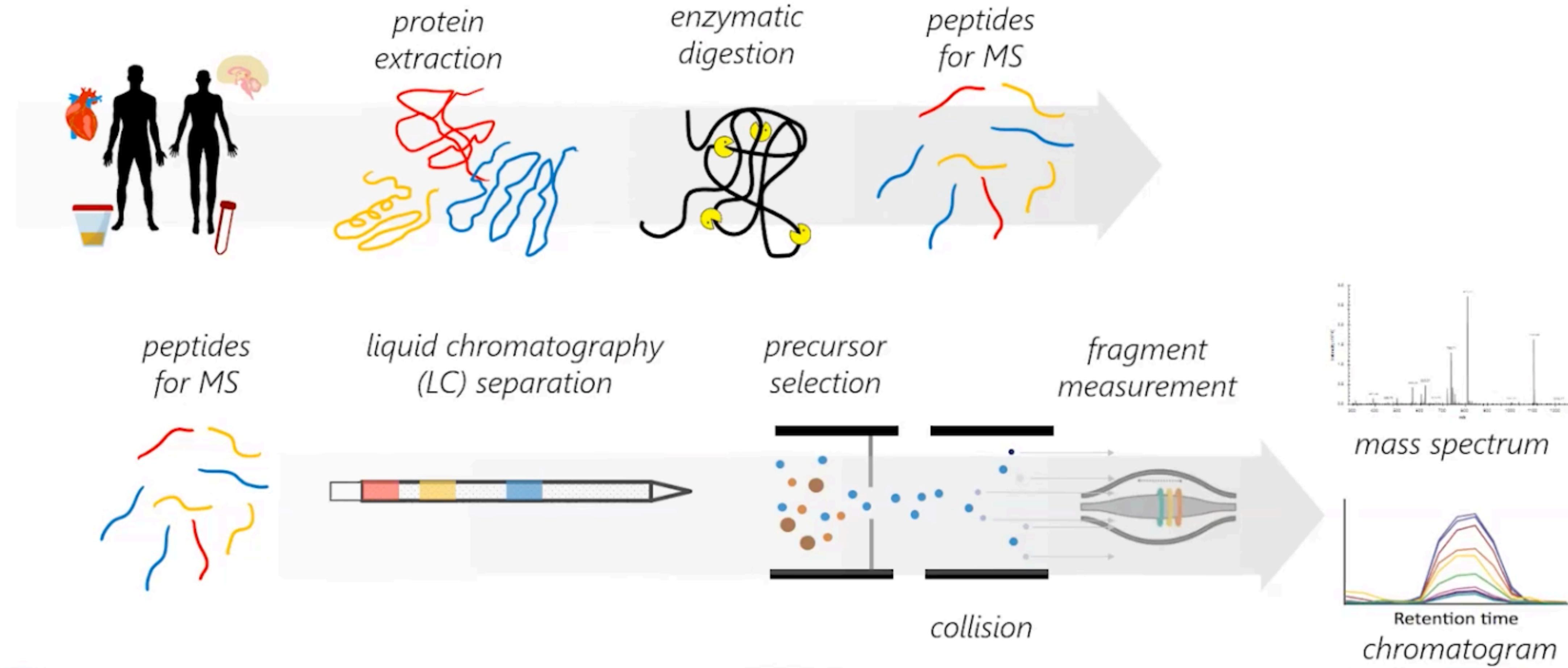
LC -MS/MS

Liquid Chromatography (LC) coupled with MS

MS is coupled with Liquid Chromatography to jointly detect and quantify peptides (and proteins).
LC exploits the hydrophobicity of peptides to separate them.



Overview of the mass spectrometry proteomics workflow



Lindsay K Pino, Talus Bio, USA

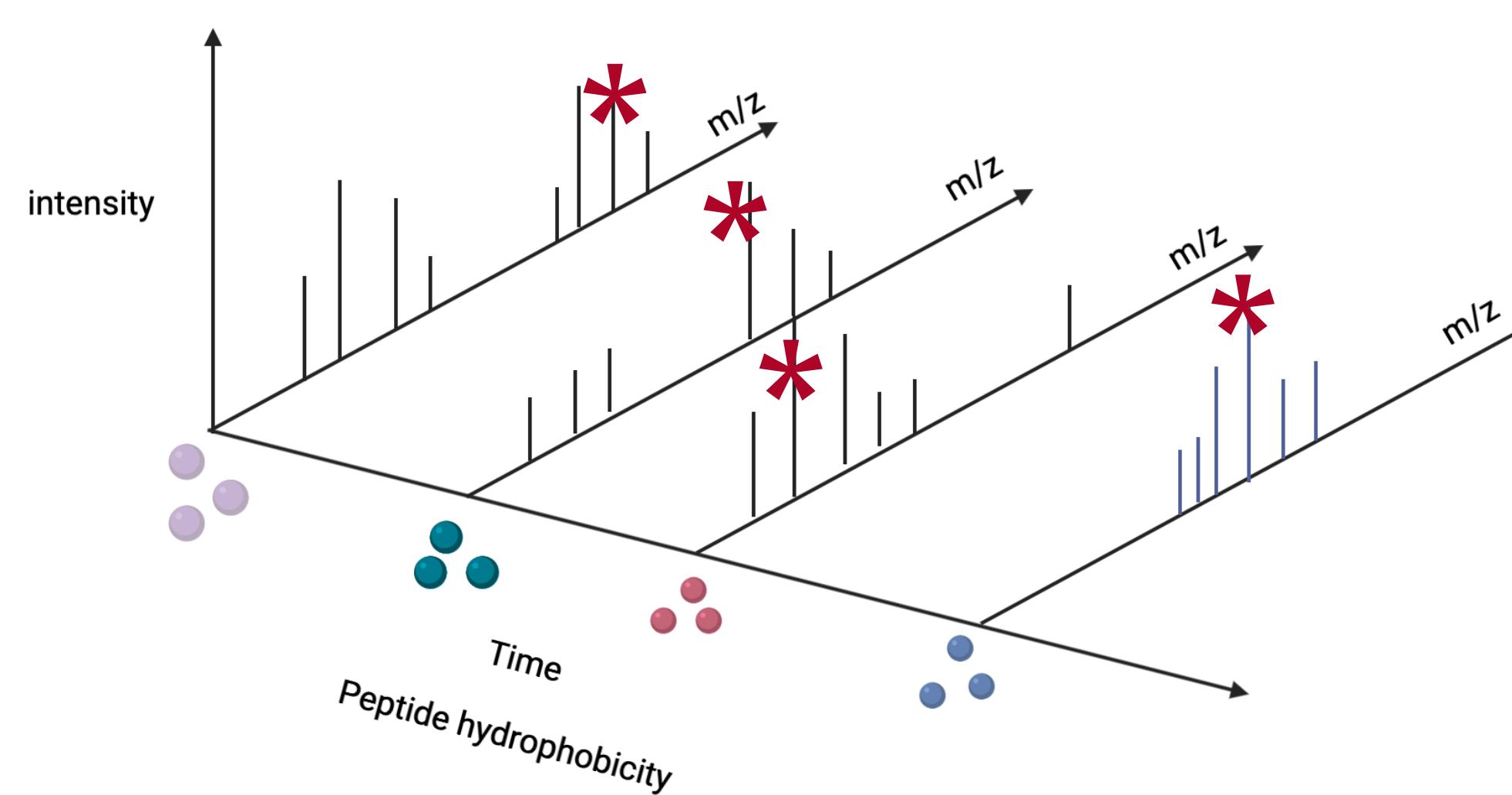
LC-MS/MS data

For each peptide ion going through the mass spec we get:

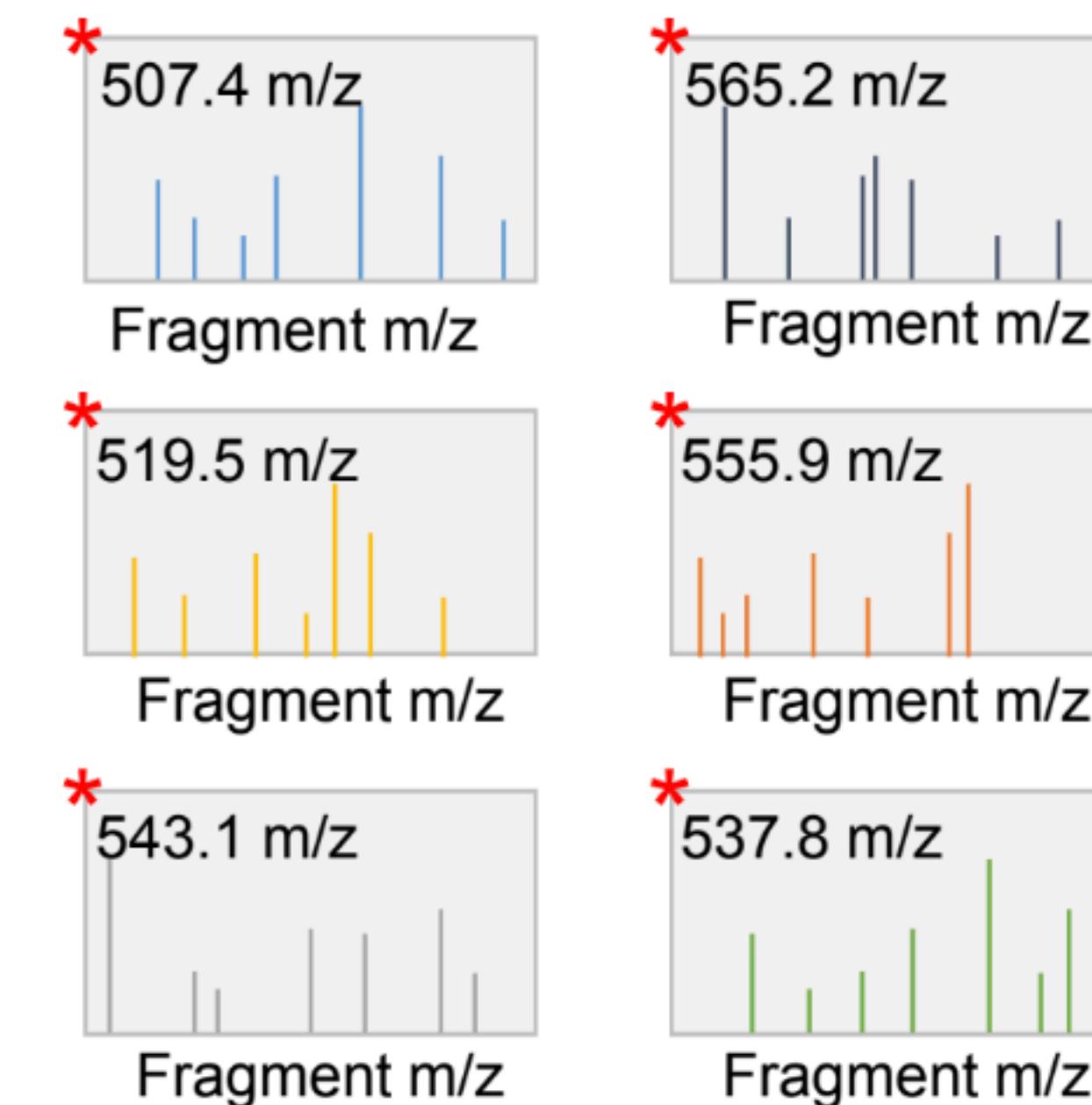
- m/z
- Retention time

The y-axis intensity is a measure of how much biological material we found for that time and m/z

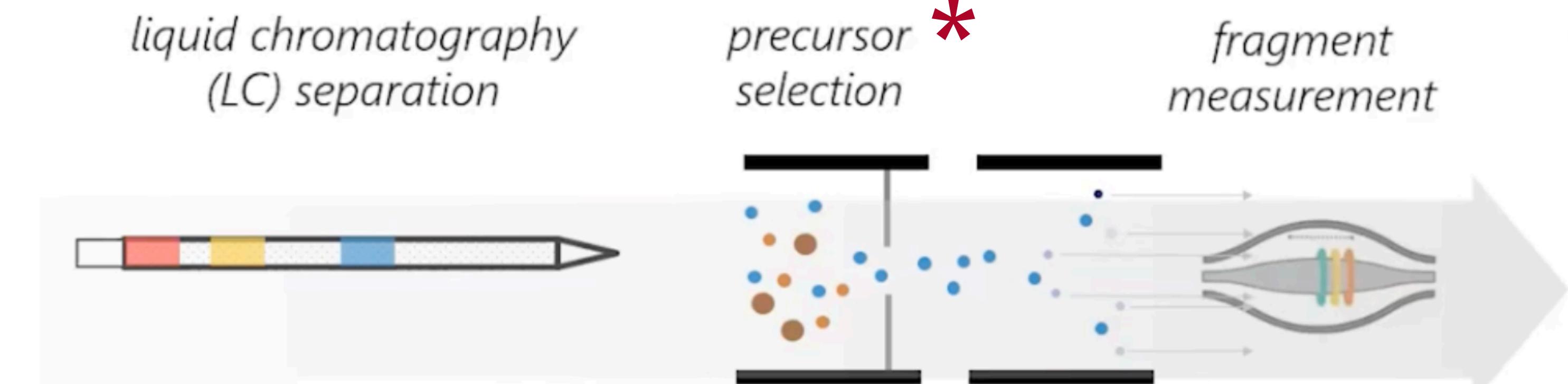
MS1 data



MS2 data



*liquid chromatography
(LC) separation*



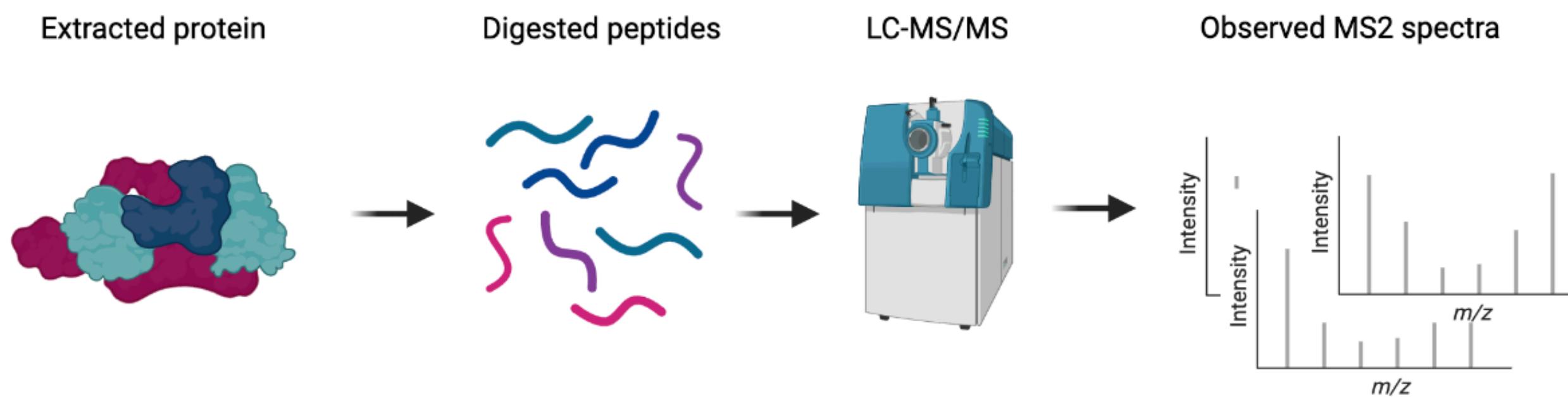
MS/MS (MS1/MS2)

Computational tasks: Identification

Peptide identifications: what peptides are found in a sample?

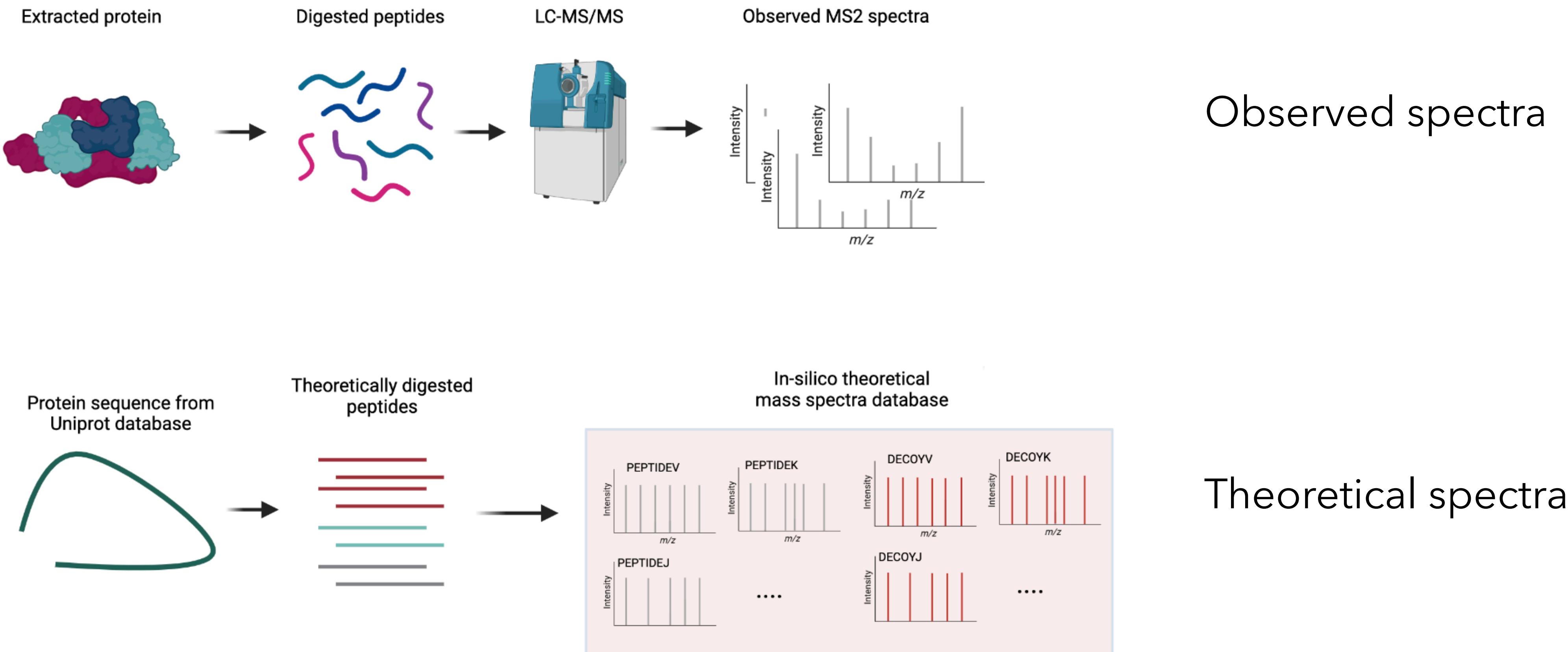
Computational tasks: Identification

Peptide identifications: what peptides are found in a sample?



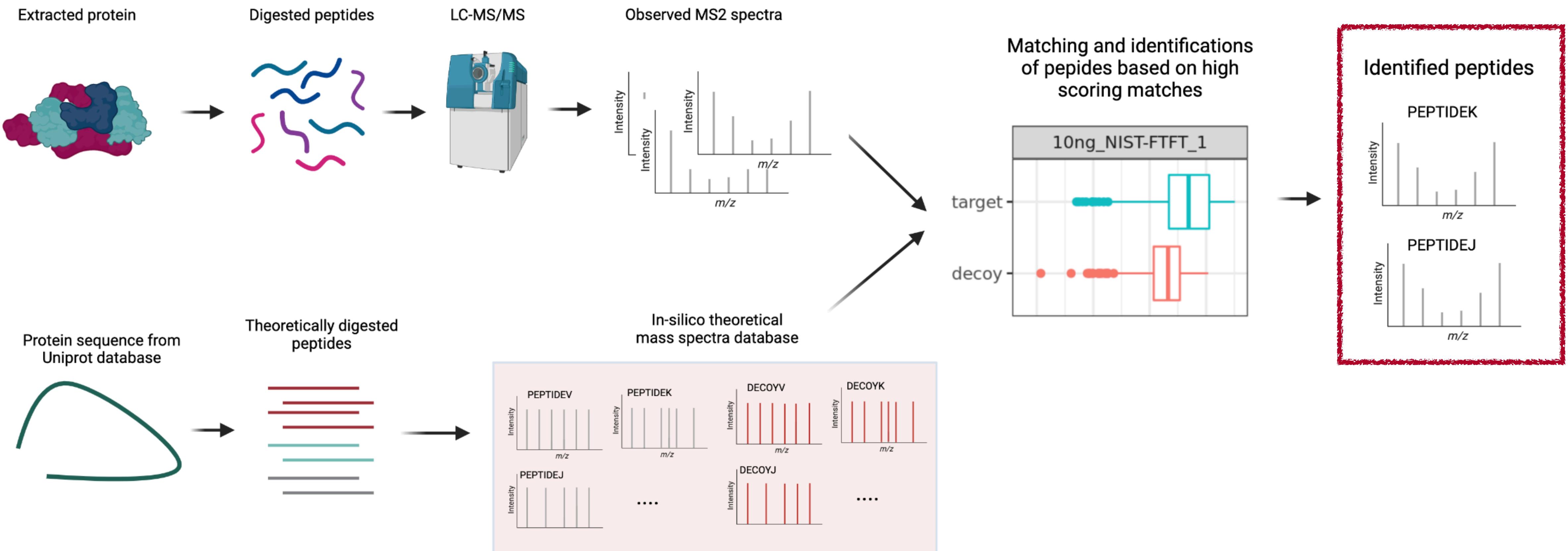
Computational tasks: Identification

Peptide identifications: what peptides are found in a sample?



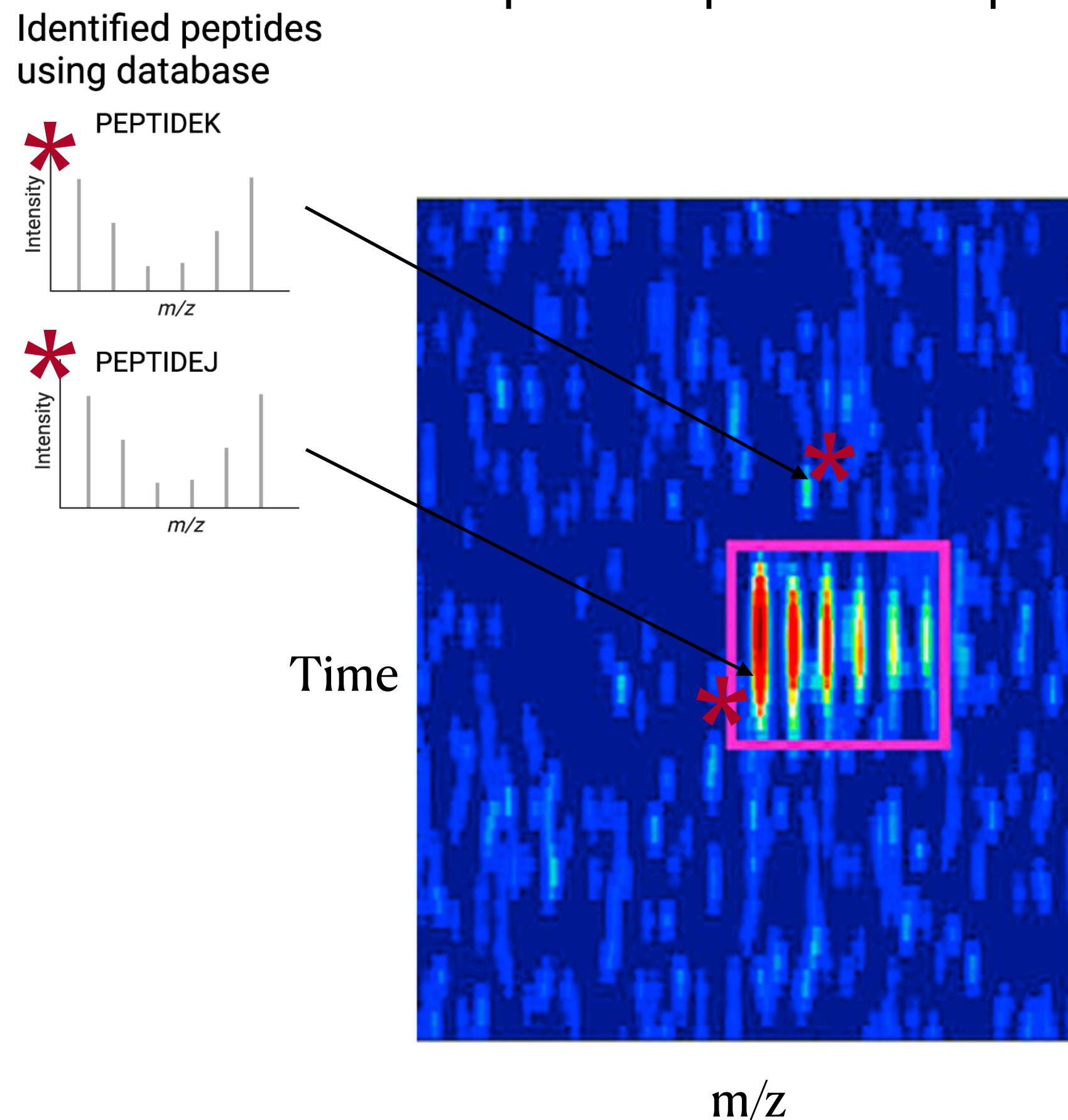
Computational tasks: Identification

Peptide identifications: what peptides are found in a sample?



Computational tasks: Quantification

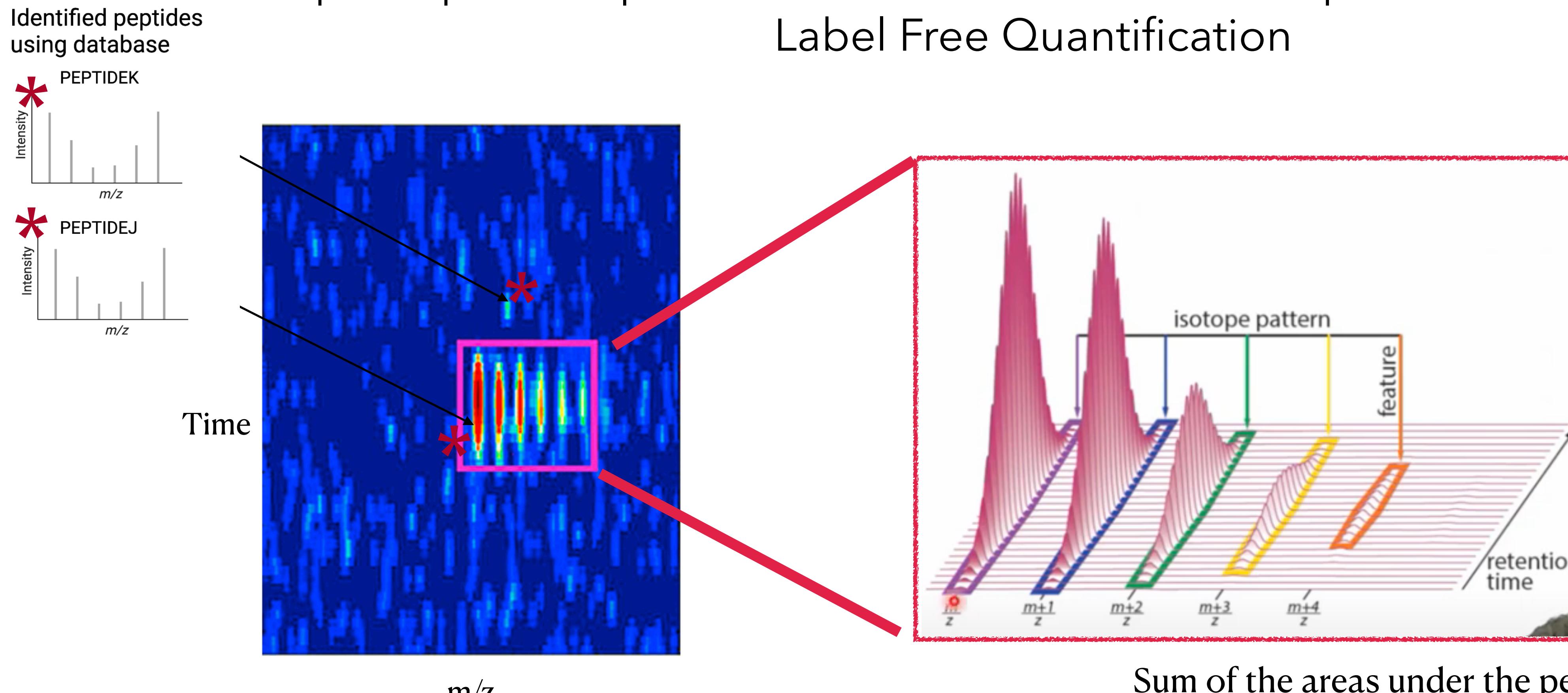
Peptide/protein quantification for downstream comparison across samples
Label Free Quantification



Computational tasks: Quantification

Peptide/protein quantification for downstream comparison across samples

Label Free Quantification



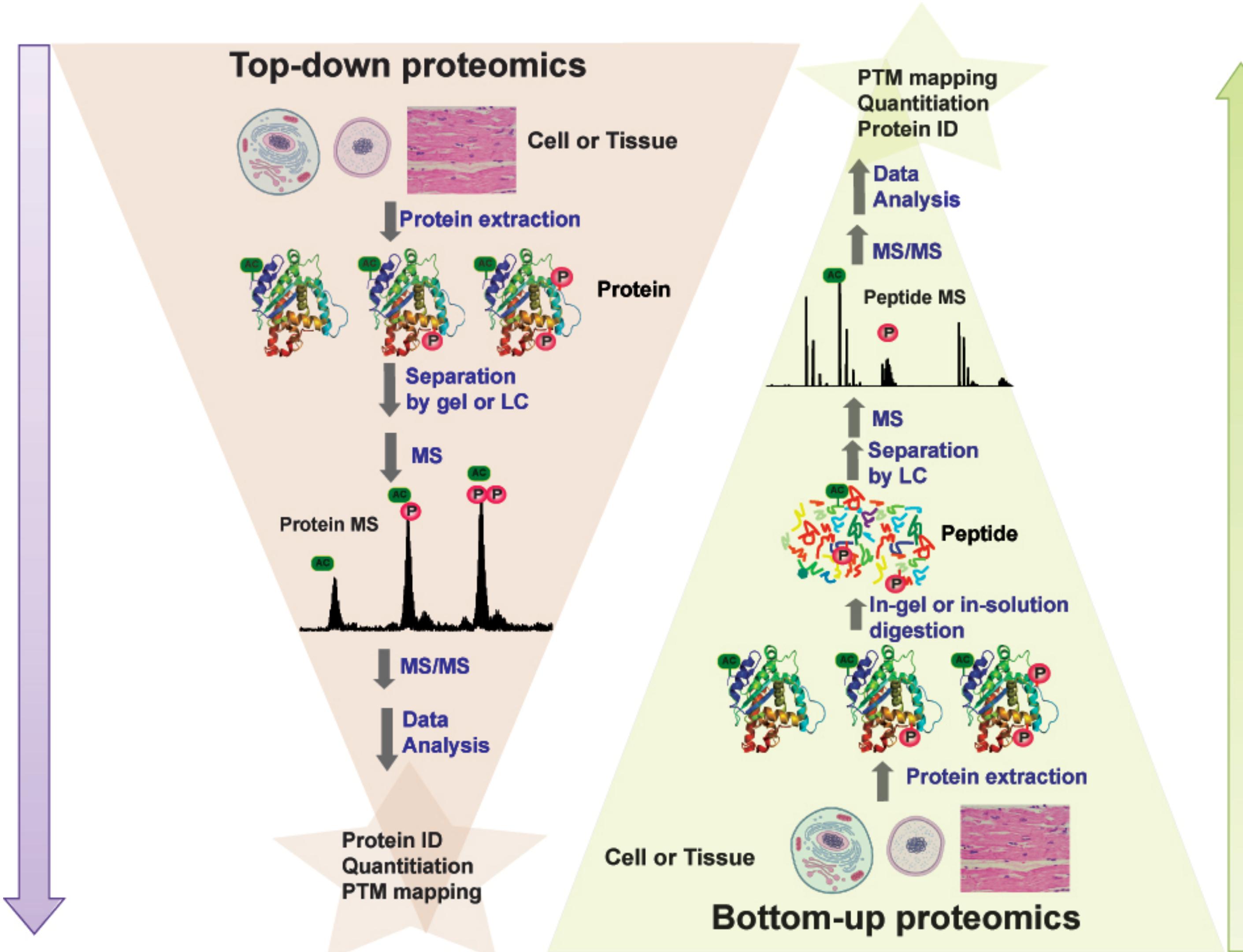
Final table of peptides/proteins intensities

Peptide Sequence	Sample1	Sample2	Sample3	Sample4
AGVVLKQ	12345	100000	0	0
WJKLGGFD	0	4545345	0	1000
TGFSSKGHG	7568564	3523523	45678	1345
JKLKFSAA	0	0	7888	8766
NNFDQKLHG		12578	0	0
FGDCBNJK	0	875645	1322	3432526

So far,
described the workflow for a *shotgun label free DDA*
LC-MS/MS experiment

Different LC-MS/MS technologies generate different data

Top-down vs Bottom-up (shotgun)

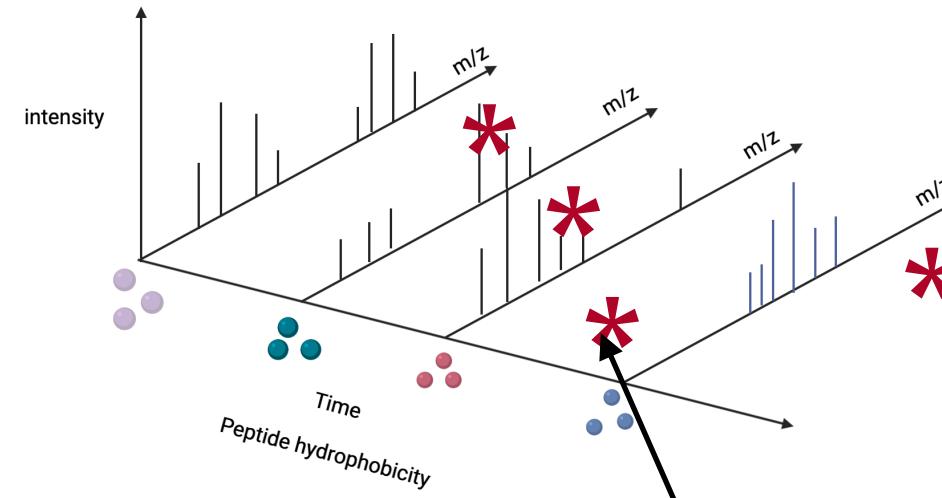


In Bottom-up Proteins are first digested (= divided) into peptides

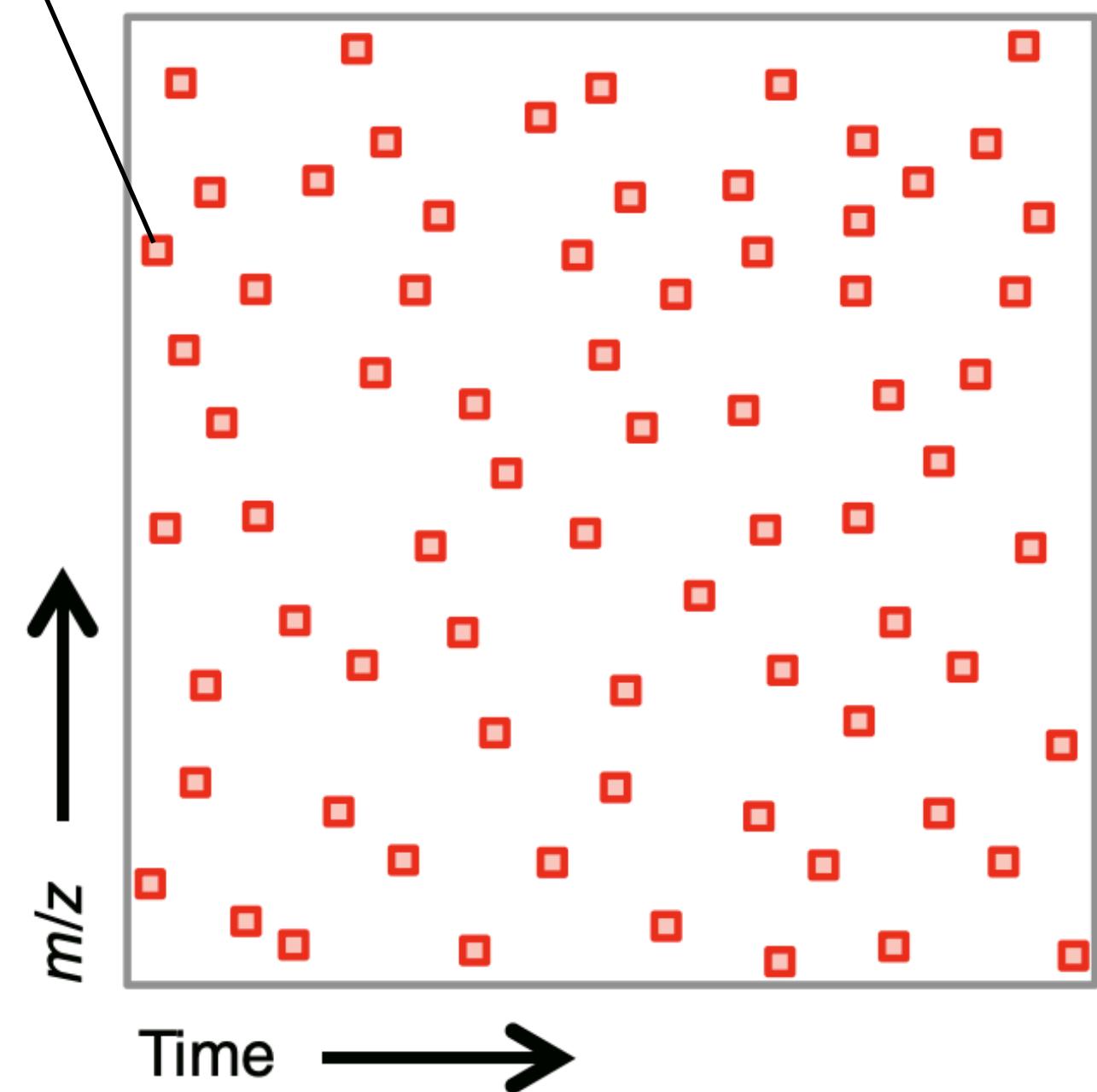
In Top-down Proteins are NOT digested but they are analysed directly:

- Harder
- Need more material
- Proteins can be really big and hard to handle

Data acquisition

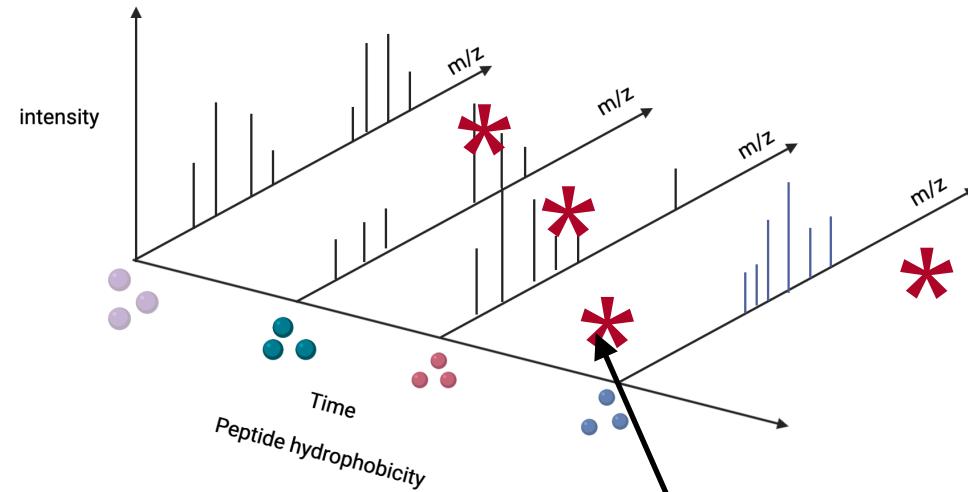


Discovery-driven
Data-dependent

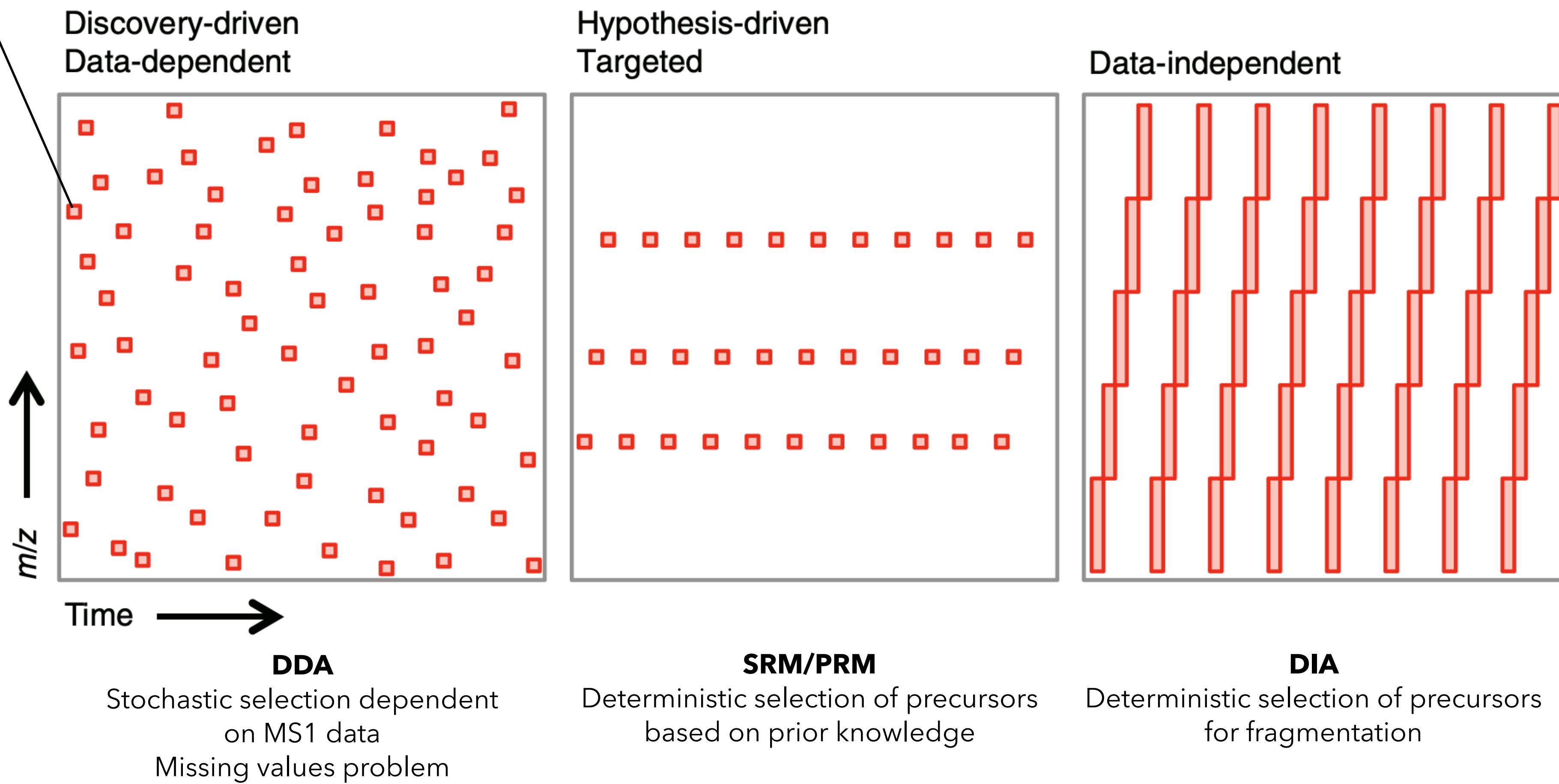


DDA

Stochastic selection dependent
on MS1 data
Missing values problem

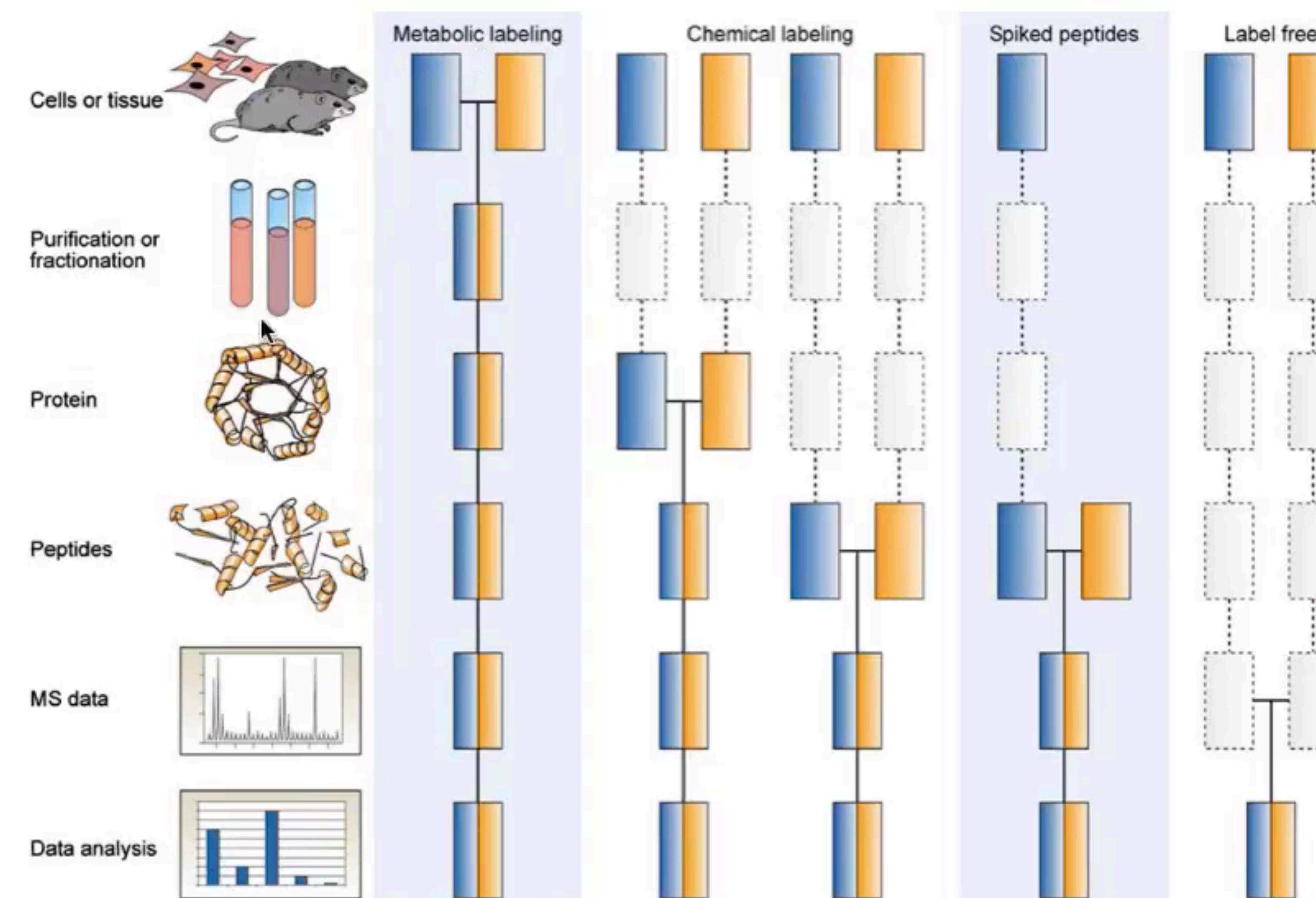


Data acquisition



Labelled vs Label Free Quantitative proteomics

Quantification Strategies

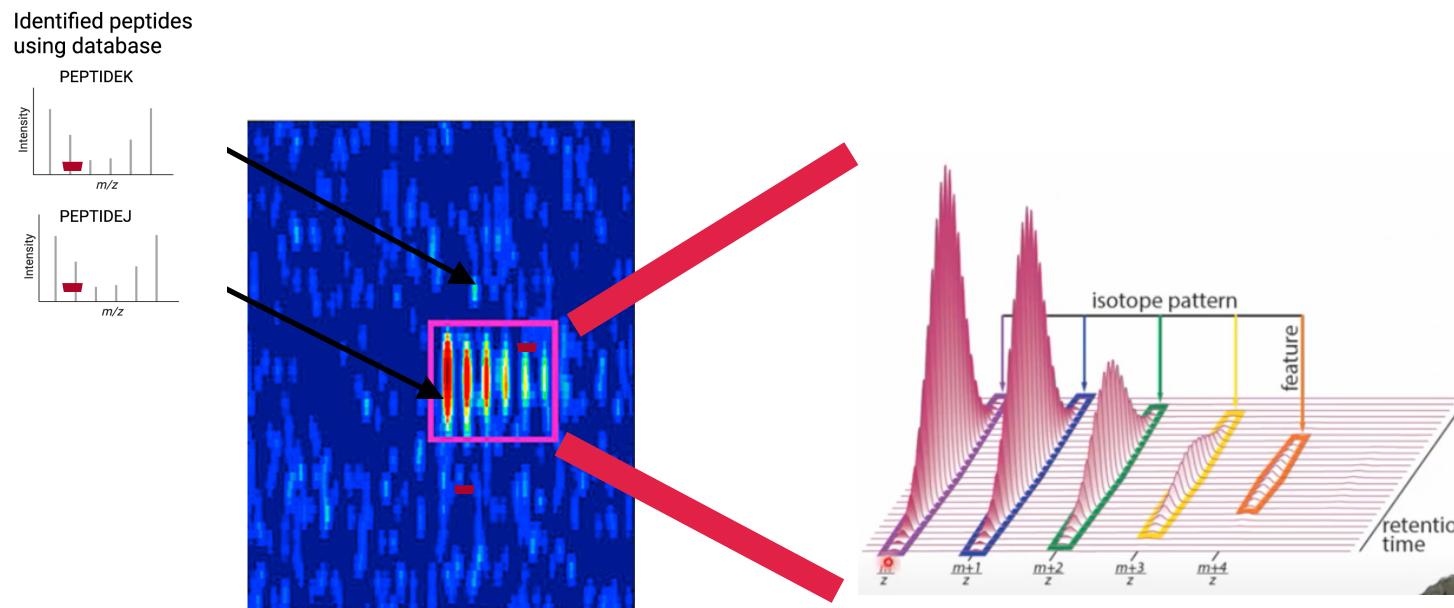


Common types of analyses using MS-based proteomics data

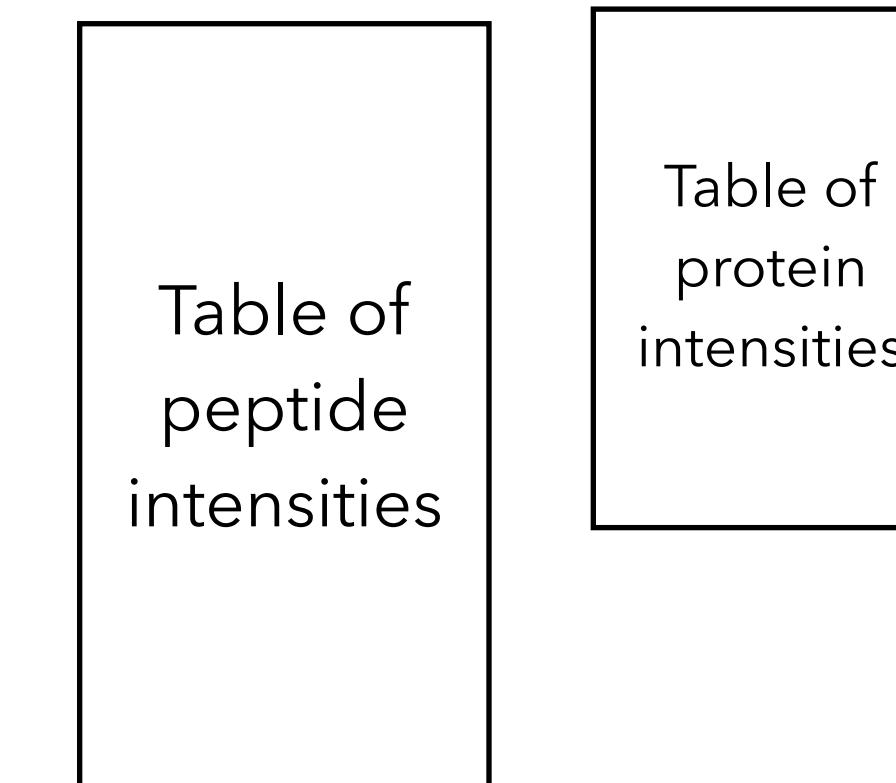
Discovery type analysis

Identify and quantify peptides/proteins proteome wide for differential expression (abundance) analyses/clustering/prediction

Pre-processing

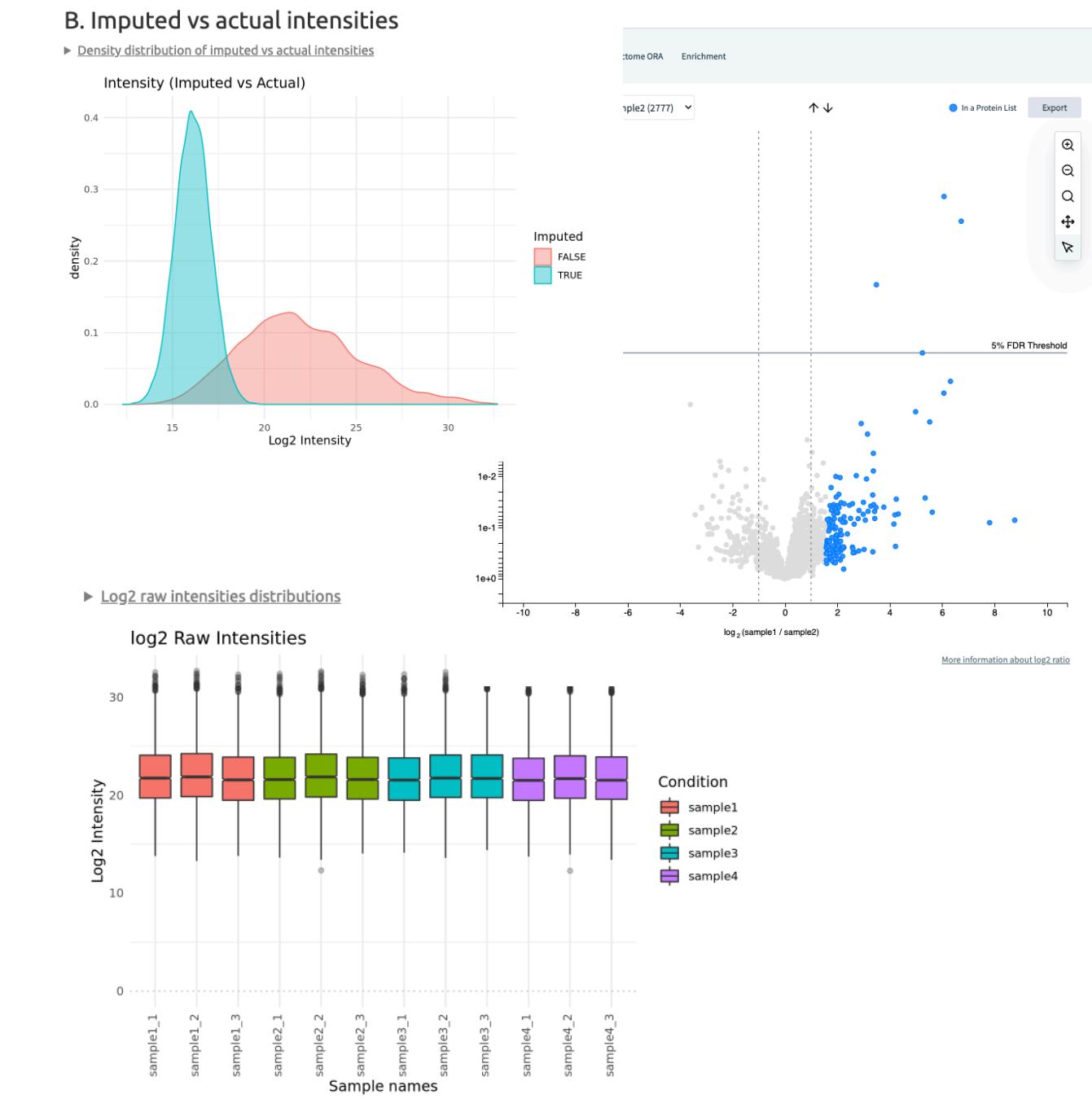


Quantified data



Downstream analysis of intensities

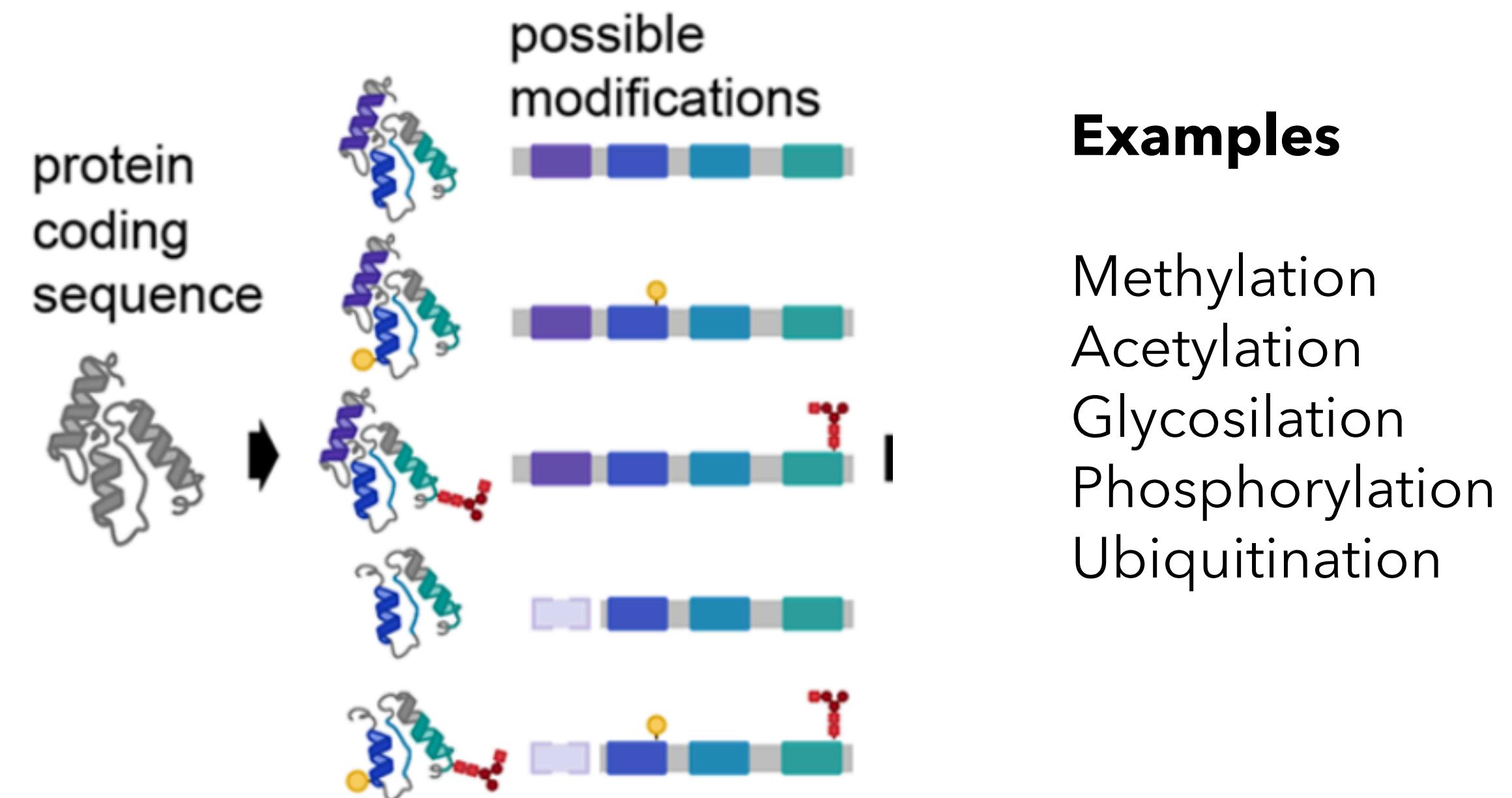
- Normalisation of intensities
- Imputation of missing values
- Differential expression
- Pathway analysis etc...



Post Translational Modifications (PTMs) identification

PTMs are chemical modifications of the proteins that can modify the activity of the protein and how it interacts with other molecules

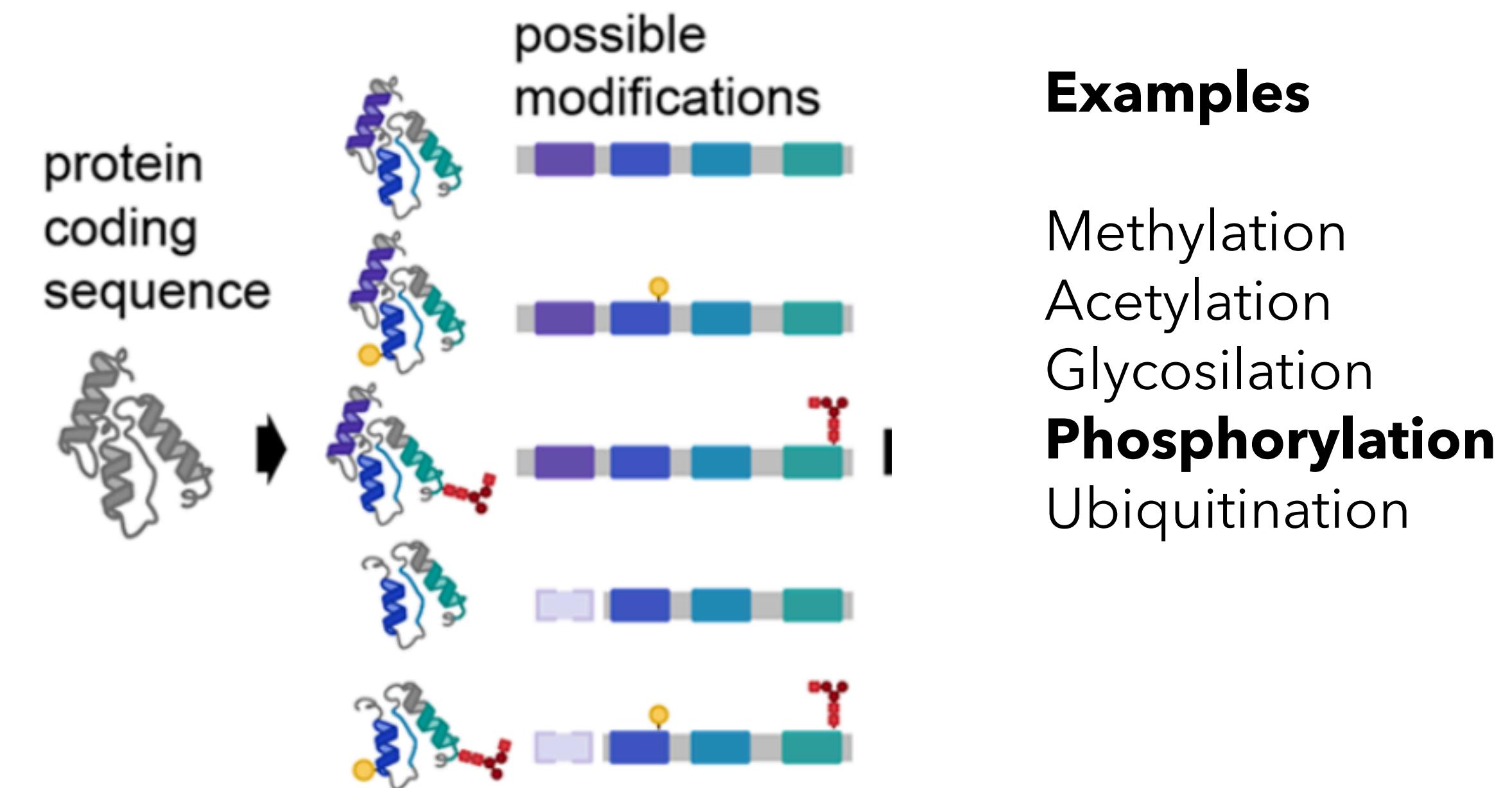
Usually the first step involves enrichment for specific PTMs to avoid missing low abundant ones.



Post Translational Modifications (PTMs) identification

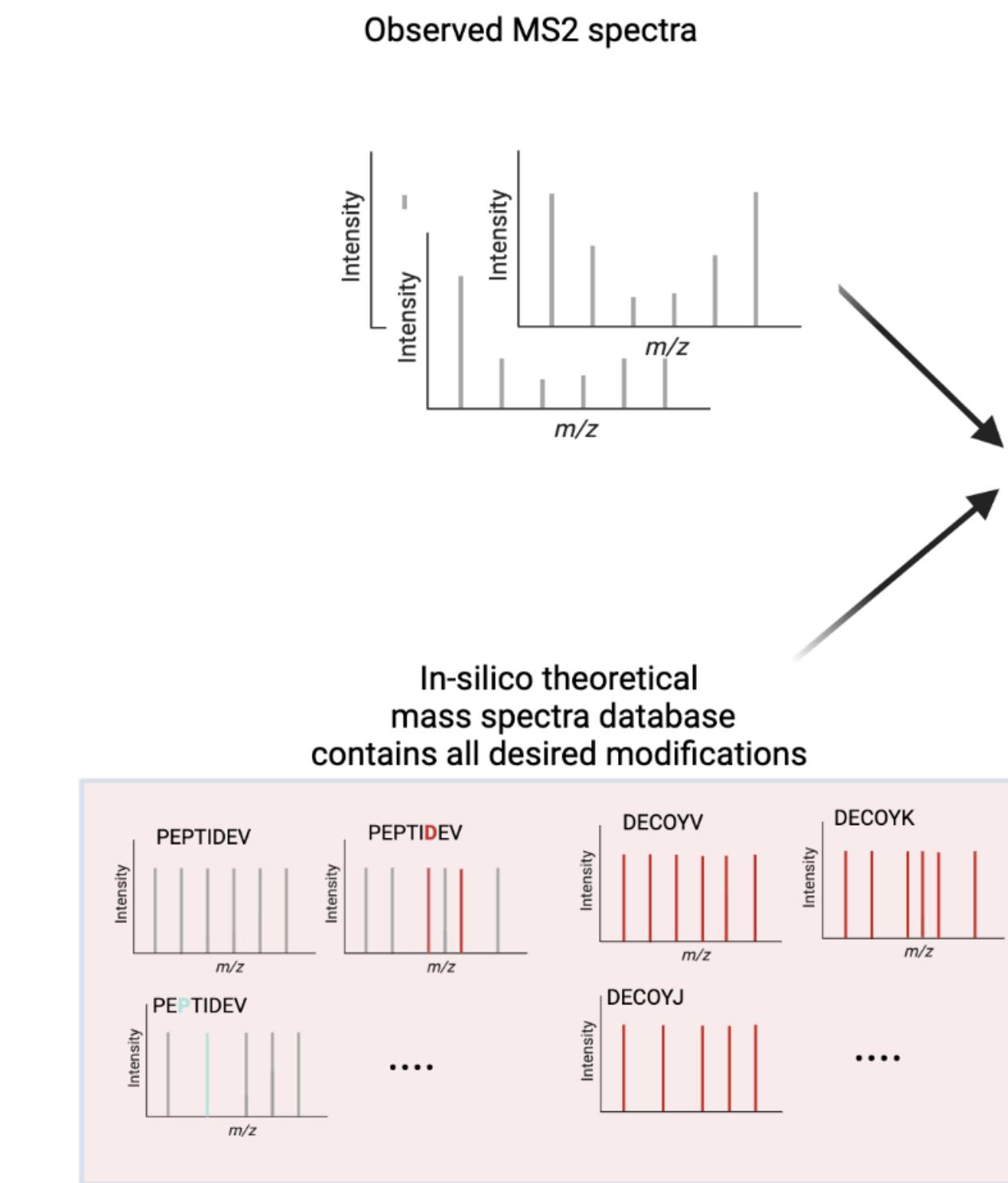
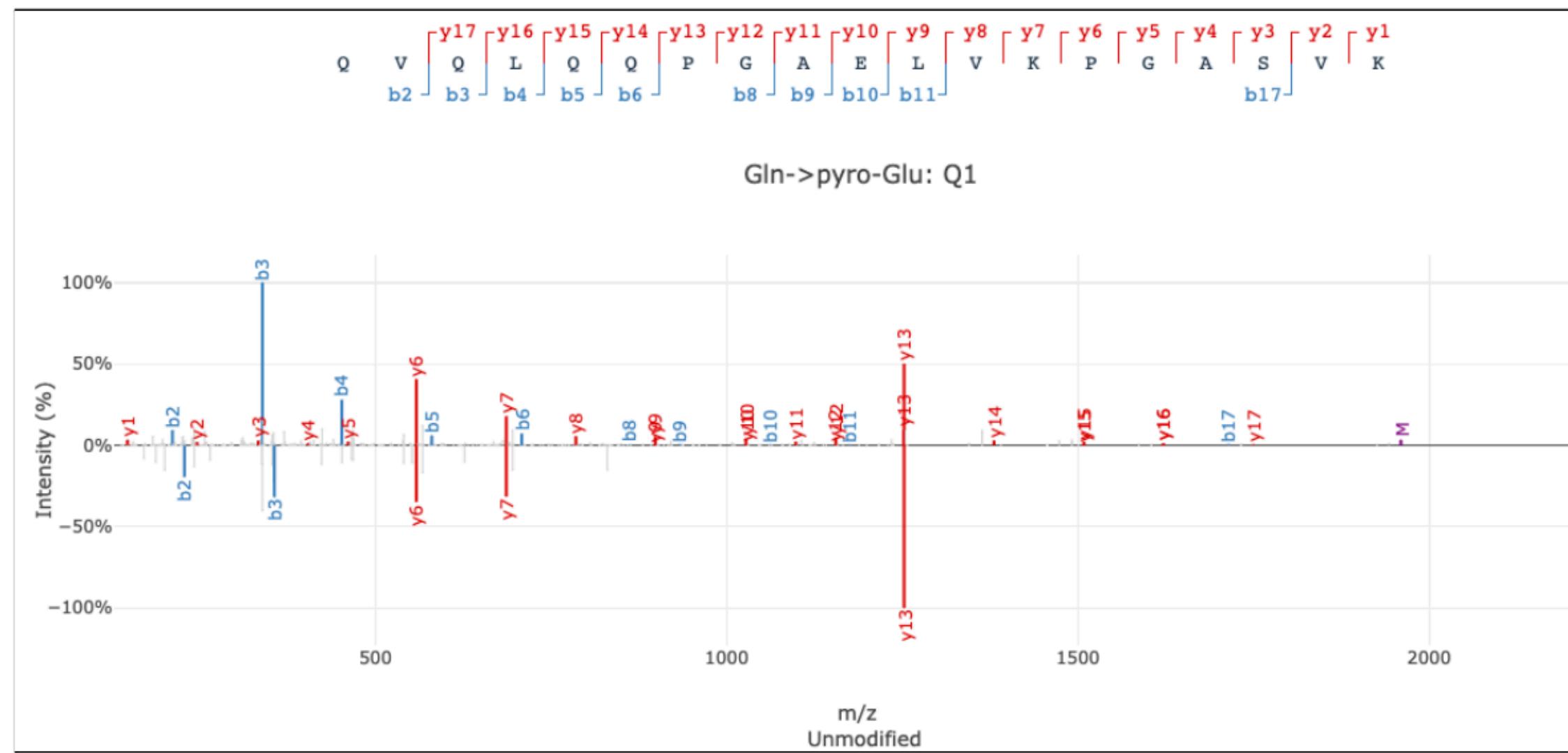
PTMs are chemical modifications of the proteins that can modify the activity of the protein and how it interacts with other molecules

Usually the first step involves enrichment for specific PTMs to avoid missing low abundant ones.

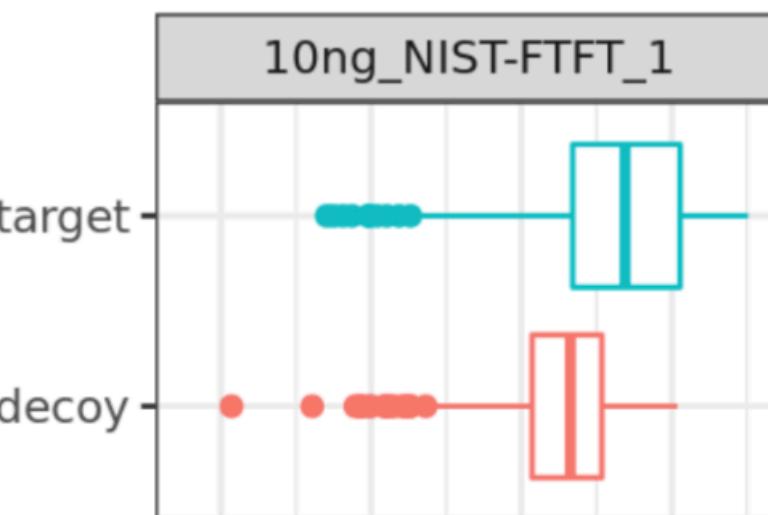


Peptide mapping

Characterise and quantify post-translational modifications (PTMs) in target peptides/proteins.
Widely used for drug compounds characterisation.



Matching and identifications of peptides based on high scoring matches



What tools people use?

The choice of a software depends on

Your computational level expertise:

Is there a Mass Spec facility that can support you?

Do you have statistical/programming knowledge to DIY?

Is there someone in your lab that can help?

The type of data and analysis to be performed

Differential expression vs PTMs

The type of technology

DIA vs DDA vs Targeted

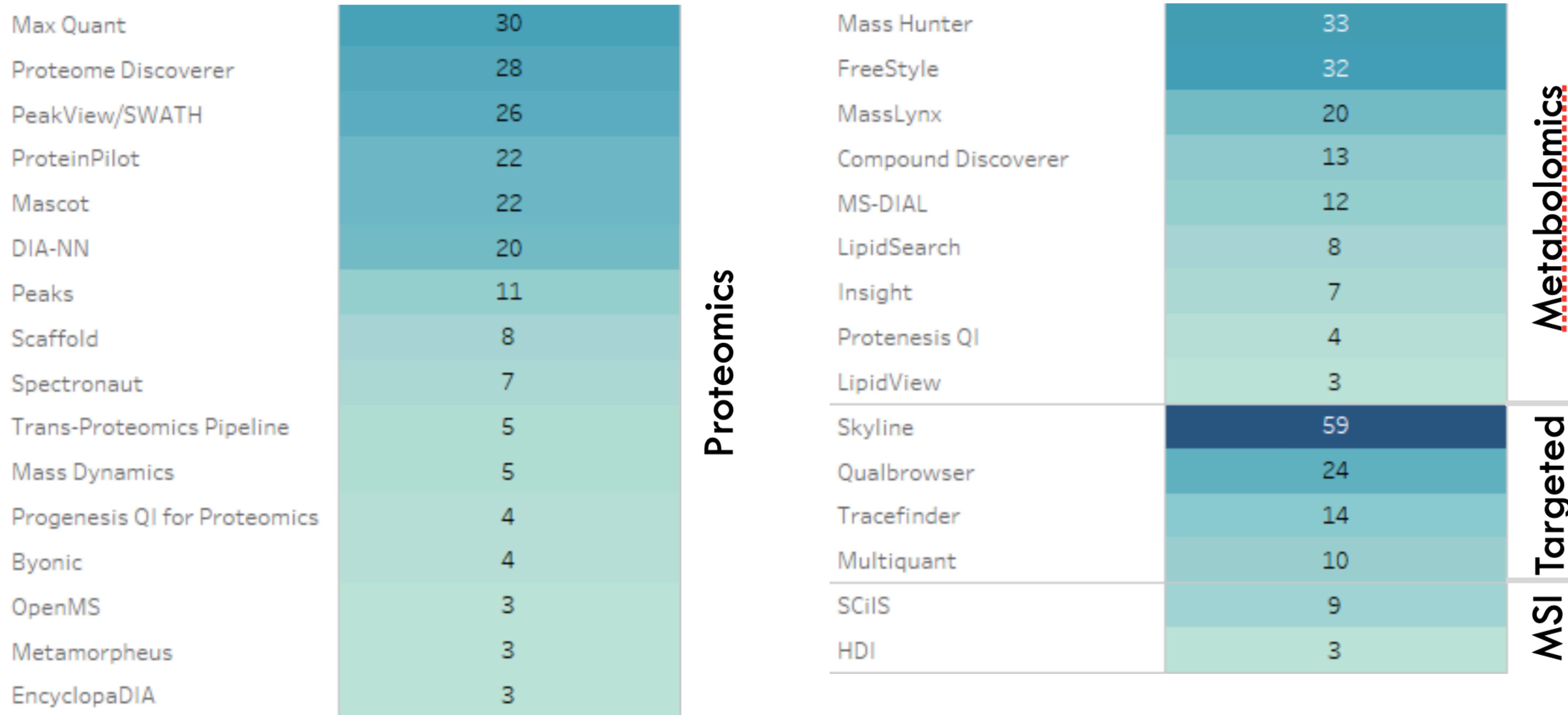
Label free vs Labelled

Australasia Core Facility Survey



- Results from the 3rd survey
- This round was open in July 2021, the majority of questions refer to the 2020 calendar year.
- 48 facilities invited, 37 responses received, inc. 6 from NZ.

Software



My breakdown of methods for data Processing

I divided the methods into 3 classes:

1. Software suite for several types of analyses. Aimed at researchers, .e.g MaxQuant-Perseus/Skyline/Mass Dynamics/LFQ Analyst
2. Do-it-yourself: assemble existing software to make your own workflow
3. Software targeted at BioPharma

Class 1: Software suite

Free, need download and only available for Windows (sometimes Linux)

DDA

- MaxQuant + Perseus: DDA Label Free Quantification data processing and analysis
- LFQ Analyst: downstream analysis of MaxQuant data: <https://bioinformatics.erc.monash.edu/apps/LFQ-Analyst/>

DIA

- Skyline (also open source): DIA/SWATH, Targeted:
- DIA-NN (also open source): DIA data - Free

Needs paid licence

- Spectronaut (from Byognosis Zürich): DIA

Class 1: Software suite - Several workflows

Free: Need download

Progenesis QI (WATERS): several workflows supported

Licence \$: Need download

- ProteomeDiscoverer: wide range of proteomics workflows
- Peaks

The screenshot shows the homepage of the PEAKS Studio Xpro website. At the top, there's a banner for 'PEAKS Studio X PRO' with the tagline 'PUSH THE LIMITS OF DISCOVERY WITH THE BRAND NEW'. Below the banner, there's a navigation bar with links for Home, About Us, Products & Services, Free Trial, Help & Support, and Contact Us. The main content area features a large image of a protein structure and text about the software's capabilities, including peptide/protein identification, de novo sequencing, database search, spectral library search, post-translational modification (PTM) search, sequence variant and mutation search, protein quantification, and detailed visualization. There are also links for 'New Features in PEAKS Xpro', 'Free Trial', and 'Request Quotation'.

The screenshot shows the homepage of the Proteome Discoverer Software website. At the top, there's a navigation bar with links for Home, Popular, Applications & Techniques, Shop All Products, Services, and Support. There's also a search bar. The main content area features a large image of a protein structure and text about the software being 'The intelligent protein informatics platform'. A red 'Contact us' button is visible. There's also a brief description of proteomics and a link to a 'PEAKS Studio Xpro Walkthrough' video.

Class 1: Software suite - Cloud based

\$ Membership or free

- Mass Dynamics (partially open source):
downstream processing or from RAW (DDA data)
DE, pathways analysis + Peptide Characterisation



Create Experiment

My Experiments

Shared with you

Papers from journals

Discovery Proteomics

Use differential expression experiments to understand global protein expression changes in complex mixtures

Start Discovering

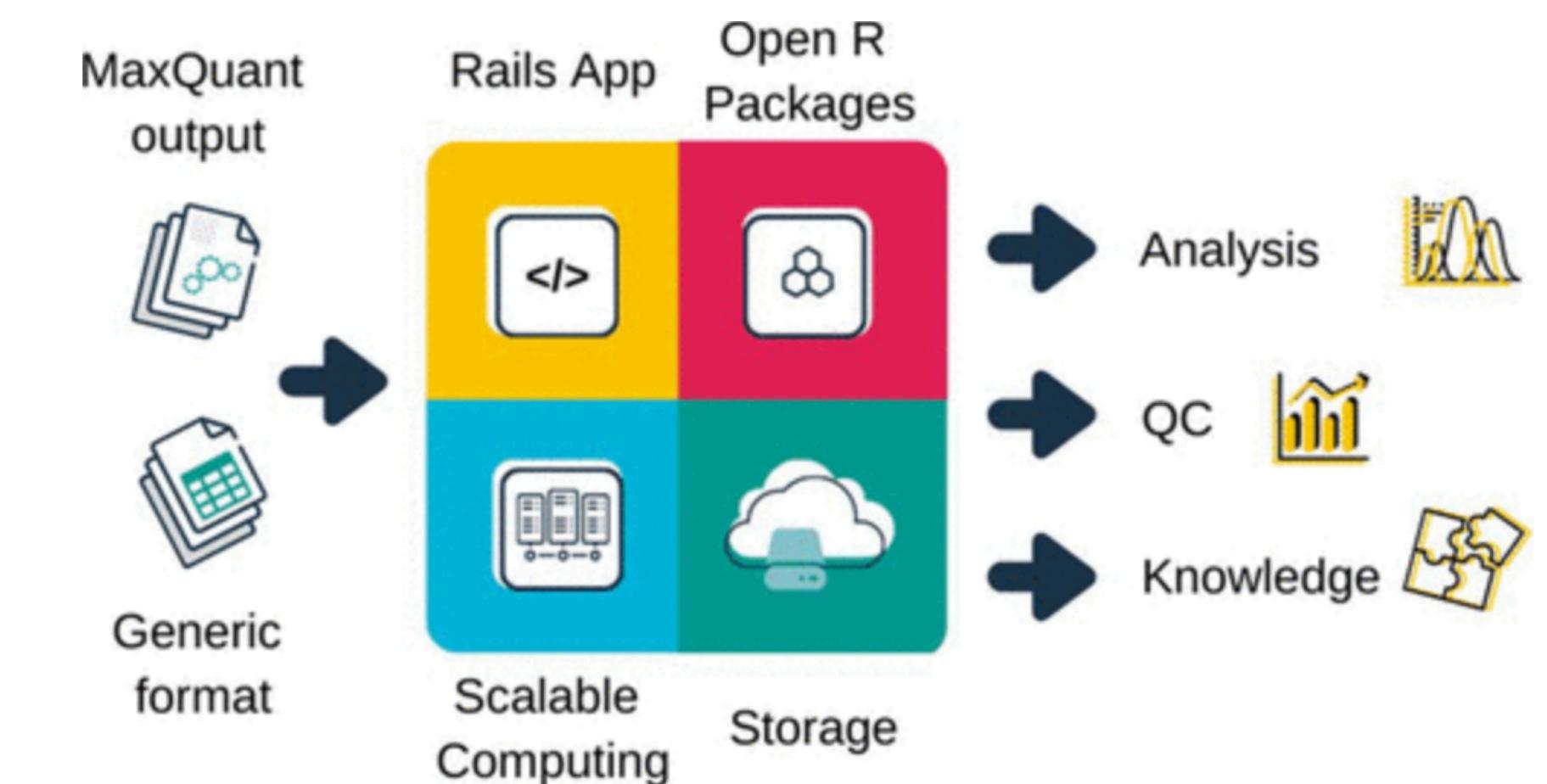
Protein Characterisation

Use peptide mapping (screening and quantitation) or quality attribute monitoring experiments to comprehensively characterise proteins of interest

Start Characterising

Need something else? We listen to the needs of the community and welcome your feedback, thoughts or suggestions

Send Feedback



Class 2: Do-It-Yourself

For biologists/mass spec experts and developers with data analysis and coding skills

Most often open source libraries are used and allow flexibility and extensibility of new features

Class 2: Do-It-Yourself

OpenMS

<https://www.openms.de/>

Several tools for preprocessing of spectra , search, quantification

WHAT IS OPENMS?

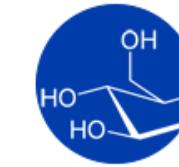
As part of the [deNBI Center for integrative Bioinformatics](#), OpenMS offers an open-source software C++ library (+ python bindings) for LC/MS data management and analyses. It provides an infrastructure for the rapid development of mass spectrometry related software as well as a rich toolset built on top of it. OpenMS is free software available under the three clause BSD license and runs under Windows, macOS and Linux.

 DOWNLOAD OPENMS



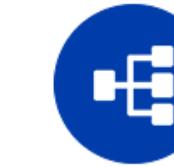
PROTEOMICS

Almost 200 ready-made and customisable tools for analysing your proteomics data. Combined in flexible workflows they provide solutions to a broad range of problems arising in state-of-the-art proteomics labs.



METABOLOMICS

OpenMS provides a wide range of customisable tools, specifically designed for all steps in the analysis of your targeted and untargeted metabolomics data.



WORKFLOWS

Every tool that is developed within OpenMS is available on flexible, scalable and easy-to-use workflow engines (such as KNIME, Galaxy, and nextflow) facilitating reproducible science

[Full list of available tools](#)

Class 2: Do-It-Yourself

MSstats

Olga Vitek Lab (Northeastern University) R package
introduction: https://www.youtube.com/watch?v=IYAqm2ENjkk&ab_channel=FenyoLab

Skyline provides an interface for MSstats



Class 2: Do-It-Yourself

R for Mass Spectrometry (<https://www.rformassspectrometry.org/>) & Bioconductor

Packages

This page lists some core packages from *RforMassSpectrometry*. For a full listing of currently available package see the project's [R universe page](#).

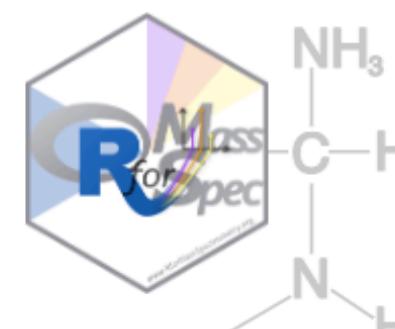
Installation and use

Execute `BiocManager::install("RforMassSpectrometry/RforMassSpectrometry")` to install all the *RforMassSpectrometry* packages.

Load the core packages with `library("RforMassSpectrometry")`.

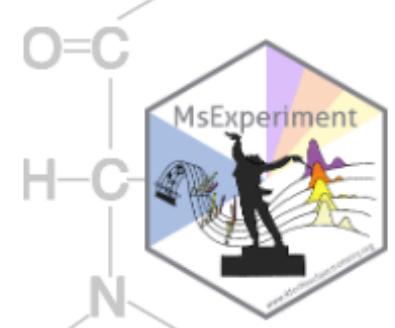
[More informations ...](#)

Packages



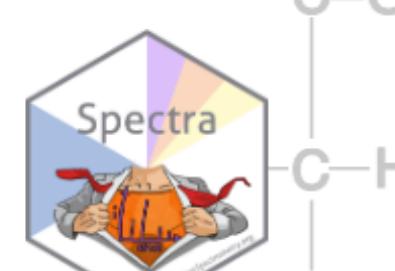
RforMassSpectrometry

RforMassSpectrometry is a meta-package that is used to manage the R for Mass Spectrometry suite and core package versions in a coherent way. Users will rely on this package to install and manage the other packages. [Learn more ...](#)



MsExperiment

The MsExperiment package provides the infrastructure to store and manage all aspects related to a complete proteomics or metabolomics mass spectrometry experiment. It relies on the other RforMassSpectrometry core packages for the data crunching. [Learn more ...](#)



Spectra

The Spectra package provides base classes and processing methods for raw mass spectrometry data. It is designed with efficiency, both in terms of memory footprint and processing time in mind, and can man-

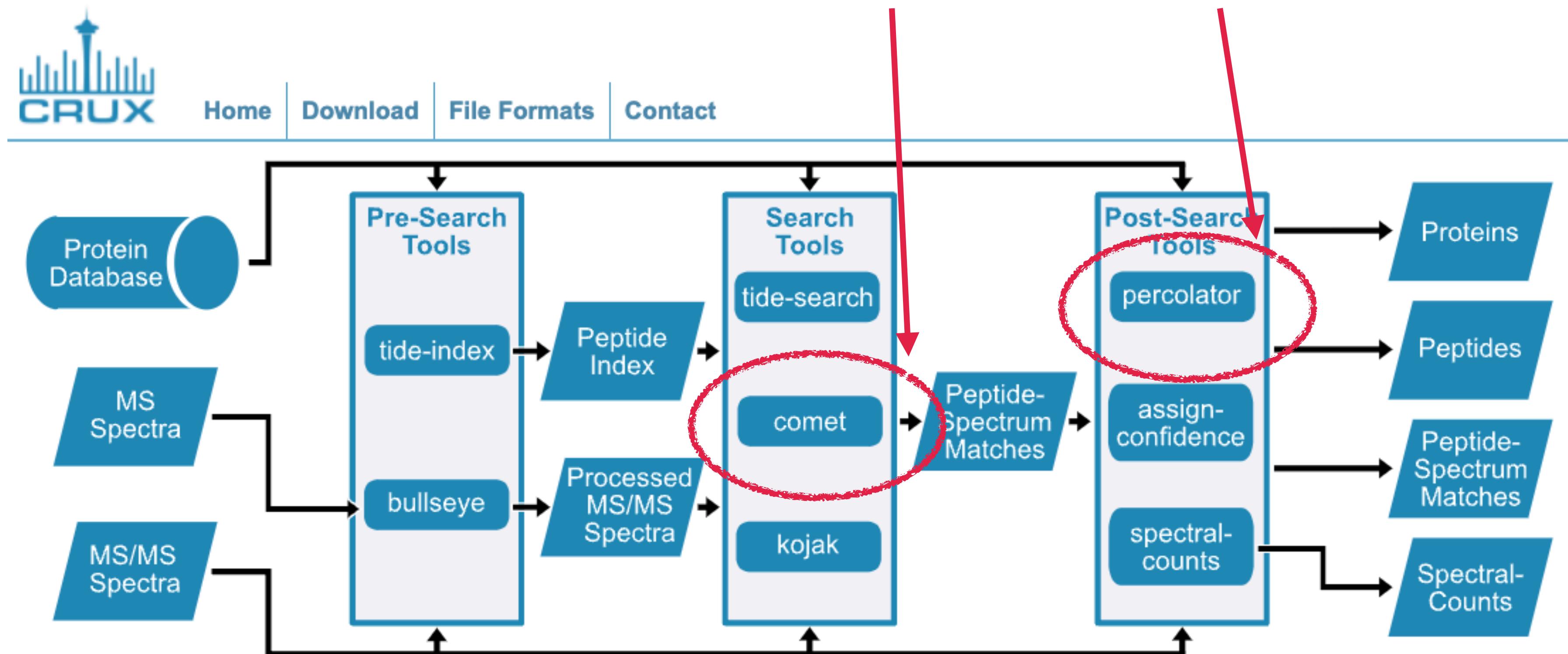
- Mostly useful to store information about experiments

- More useful to perform the downstream analysis once the RAW data have been pre-processed with MaxQuant, Skyline etc..

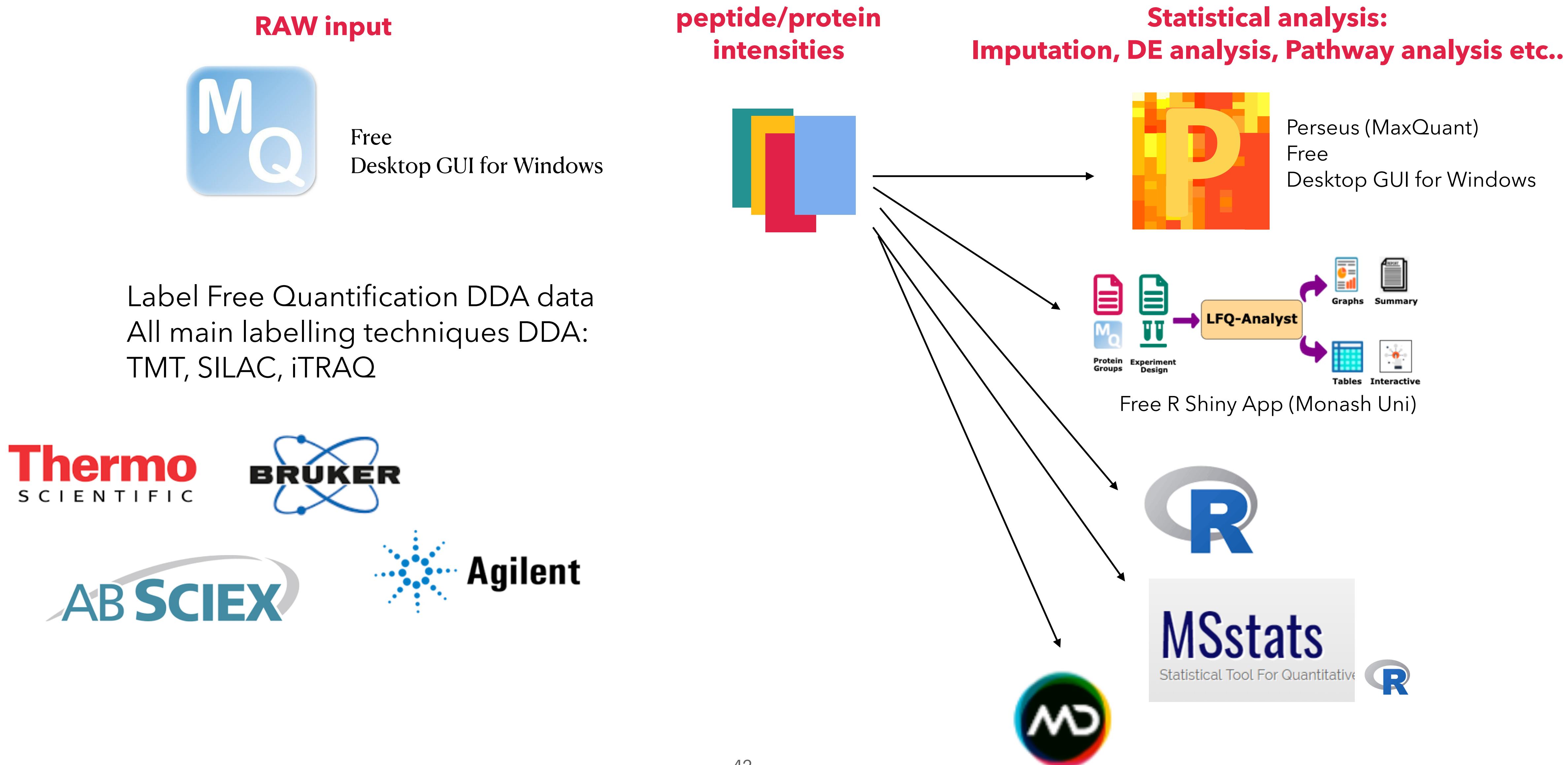
Class 2: Do-It-Yourself

Crux (~GATK for genomics)

Included in Mass Dynamics Peptide Mapping workflow



Example combining class 1 and 2



Example combining class 1 and 2

Statistical analysis:
Imputation, DE analysis, Pathway analysis etc..

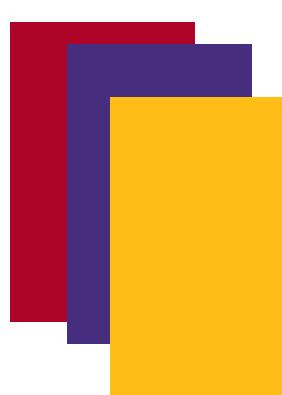
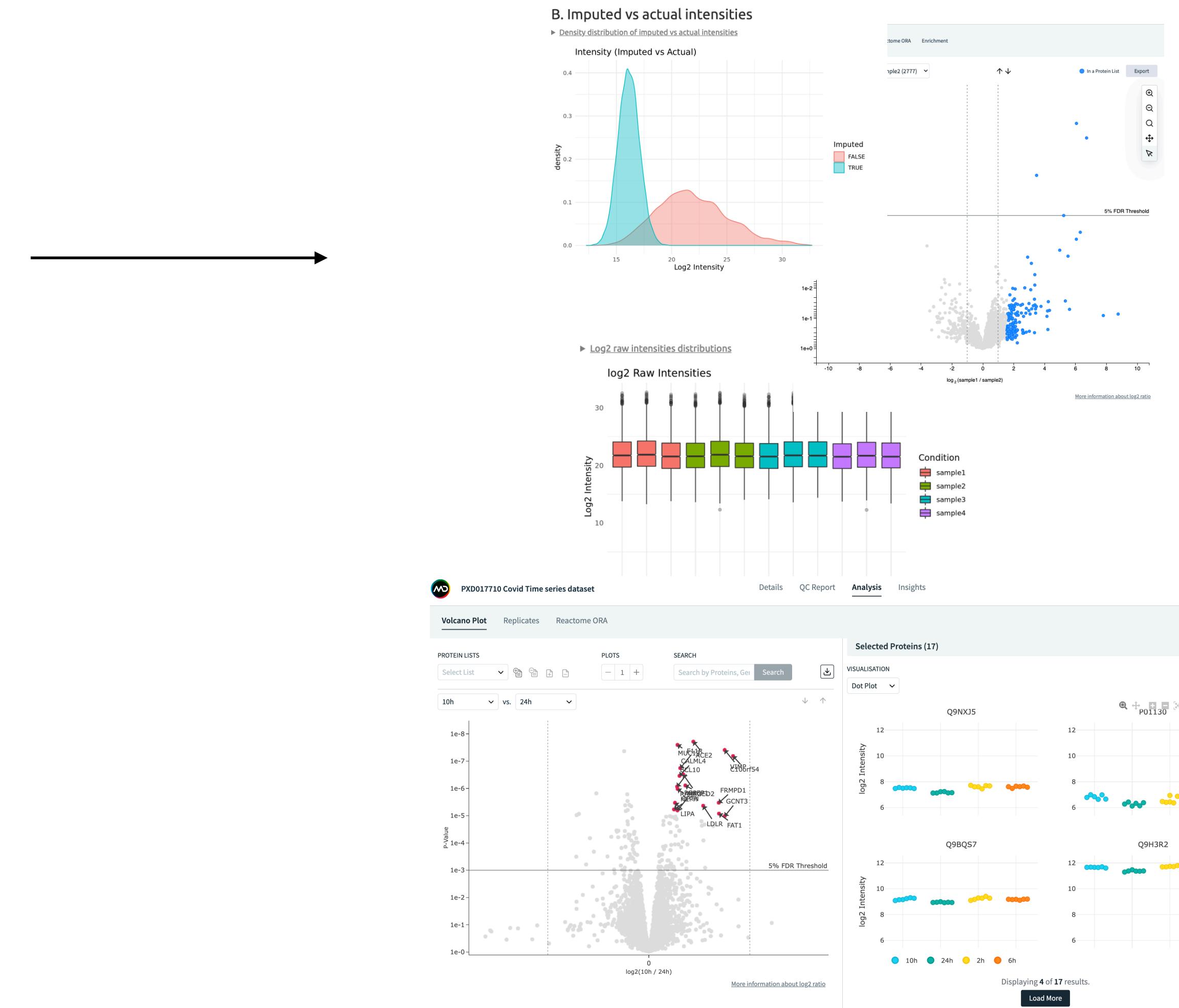
RAW input
(LFQ DDA Thermo RAW)



peptide/protein
Intensities input



ThermoFisher
SCIENTIFIC



Class 3: For bio pharma

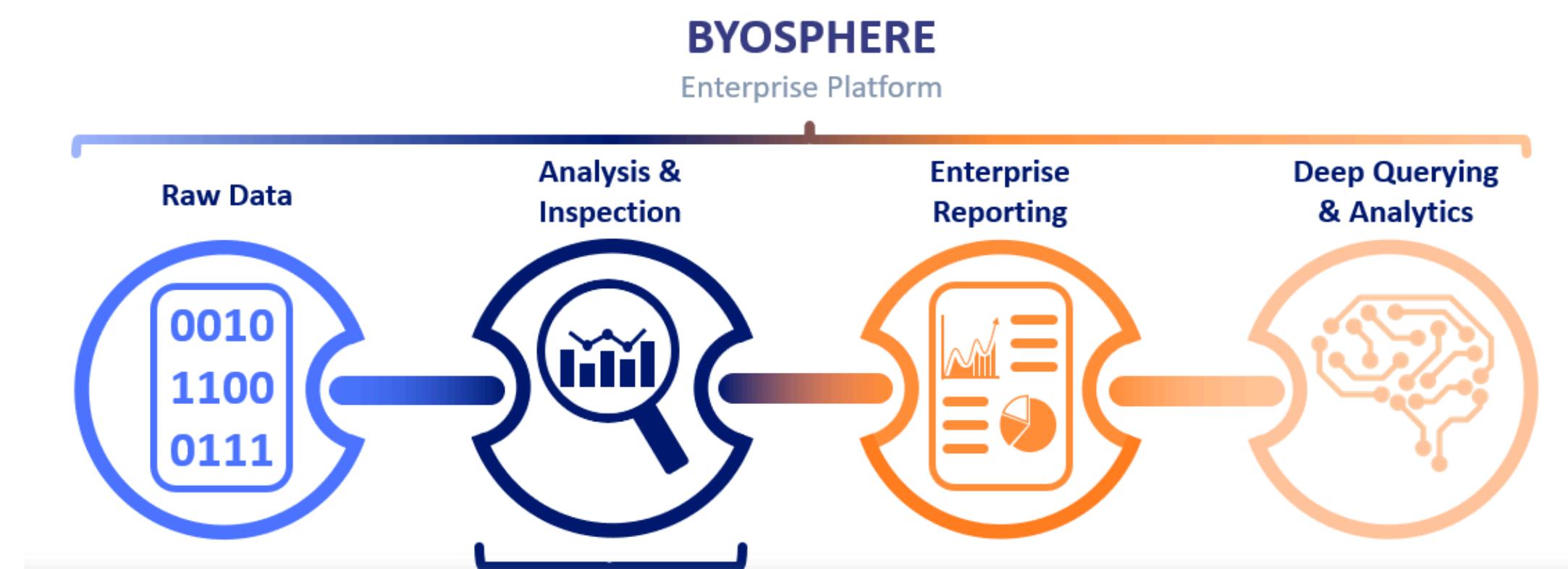
Software to support processing in BioPharmas

ProteinMetrics

Available for private cloud use

GeneData

Bench-to-Enterprise Software for Biotherapeutic Analysis



Digitalizing Biopharma R&D

The Genedata Biopharma Platform is the #1 enterprise software system for achieving operational excellence in biopharma R&D.

[Discover Our Platform](#)



Concluding remarks

The landscape of MS Proteomics data types is varied leading to different computational challenges
I've shared the most common software used to process most types of Proteomics data

When starting your first ever MS Proteomics project as a computational person:

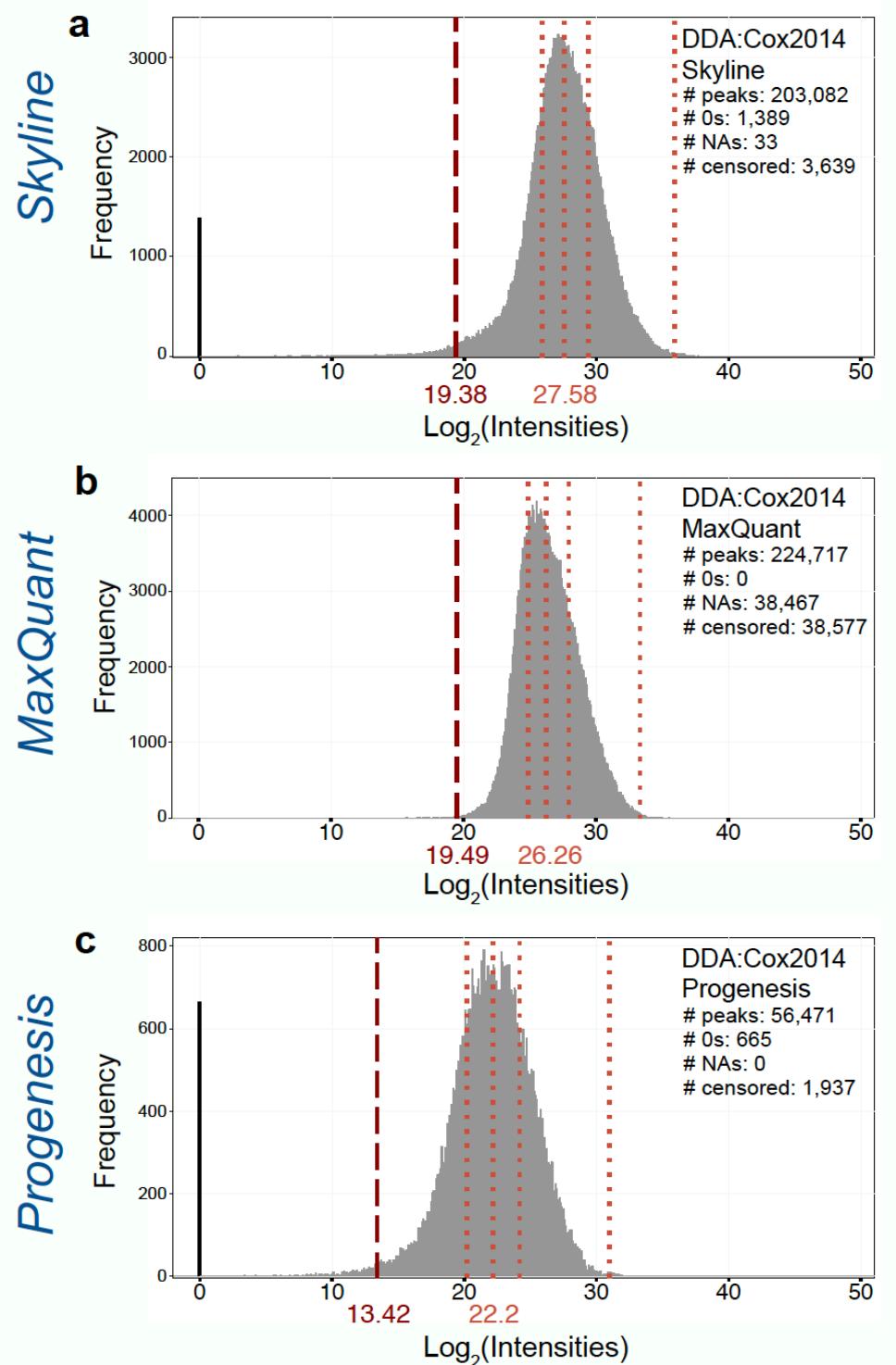
- Understand the type of analysis & experimental design required/used: e.g. DE analysis vs PTM analysis
- Understand the MS experiment design required/used: e.g. Acquisition type, Labelling, Fractionation etc..

I've noticed that it's not common in Proteomics to find big benchmarking comparisons to find what is the best methods

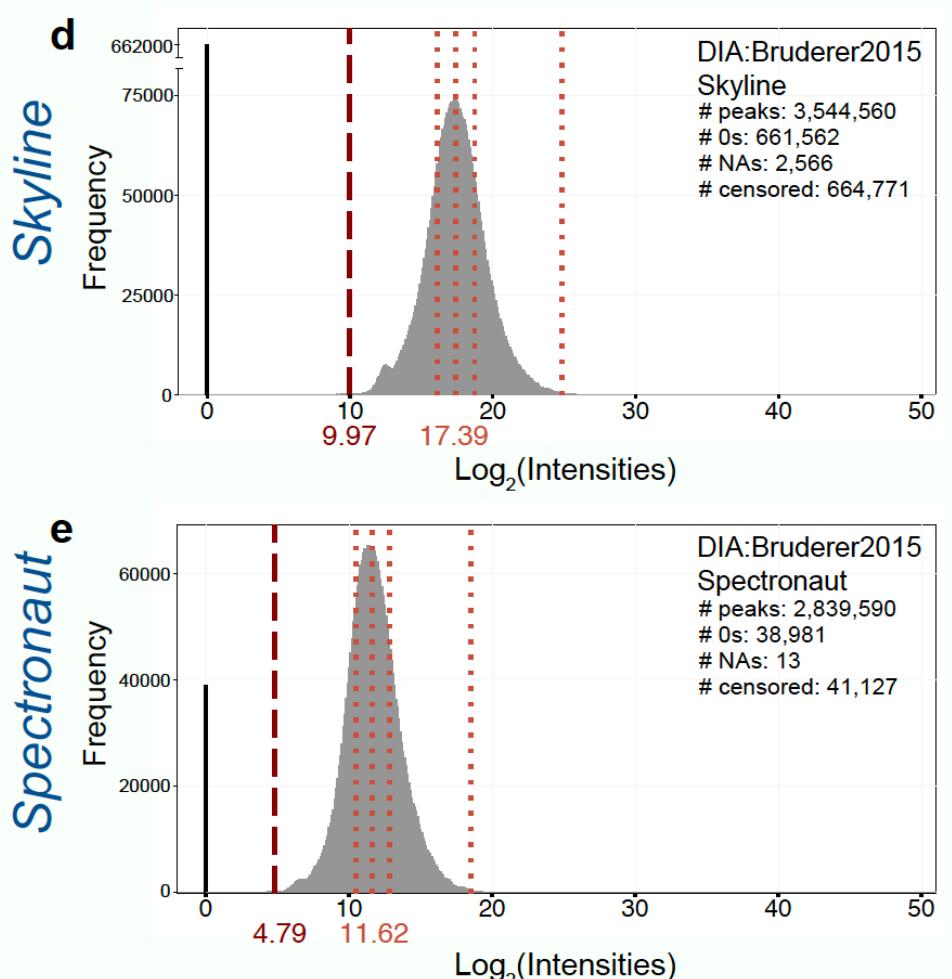
There are many methods, with several applications, all trying to improve something about the status quo
Most choices would be ok, the main complications is in the quality of the MS data

PROPERTIES OF PEAK INTENSITIES VARY BETWEEN DATA PROCESSING TOOLS

DDA: Cox 2014



DIA: Bruderer 2015



- Estimated censoring threshold
- Quantiles of $\log_2(\text{intensity})$
- Frequency of peaks with intensity reported as between 0 and 1

18

Workflow Benchmarking

Mass Dynamics regularly benchmarks all our workflows, to assess quality and measure performance as we continuously improve.

Peptide Mapping - *Probing the Sensitivity of the Lumos Mass Spectrometer using a Standard Reference Protein in a Complex Background**

The workflow is one of two [product characterisation](#) workflows. It is a targeted LFQ-DDA (Label Free Quantification-Data Dependent Acquisition) pipeline involving, feature detection, search, psm-scoring, ms1 recalibration and targeted feature detection (like match-between runs).

Read

* “*Probing the Sensitivity of the Lumos Mass Spectrometer using a Standard Reference Protein in a Complex Background*”

Regular benchmarks against ground truth to be aware of regression in the workflow performance

Olga Vitek, MSstats tutorial, ASMS 2022

Useful links

Intro to MS and LFQ

Introduction to Mass Spectrometry and Proteomics from HUPO 2021: <https://vimeo.com/showcase/8948473/embed>

Olga Vitek: MSstats: an R package for quantitative MS-based proteomic experiments: <https://www.youtube.com/watch?v=IYAqm2ENjkk>

MaxQuant Basics 1, pre conference-course: https://www.youtube.com/watch?v=H_vClGghnNo

Label-Free quantitative proteomics: <https://www.youtube.com/watch?v=p5cVkJTuKWIQ>

DIA-NN

[High throughput proteomics with DIA-NN, Dr. Vadim Demichev](#), Single Cell Proteomics 2021

Single Cell Proteomics

Follow the [Single Cell Proteomics conference YouTube channel](#), 2021/2022 conference introductions from Prof. Nikolai Slavov

Galaxy Training: Proteomics

Galaxy has loads of very useful training materials for how to get started using different software or how to perform different types of analyses.

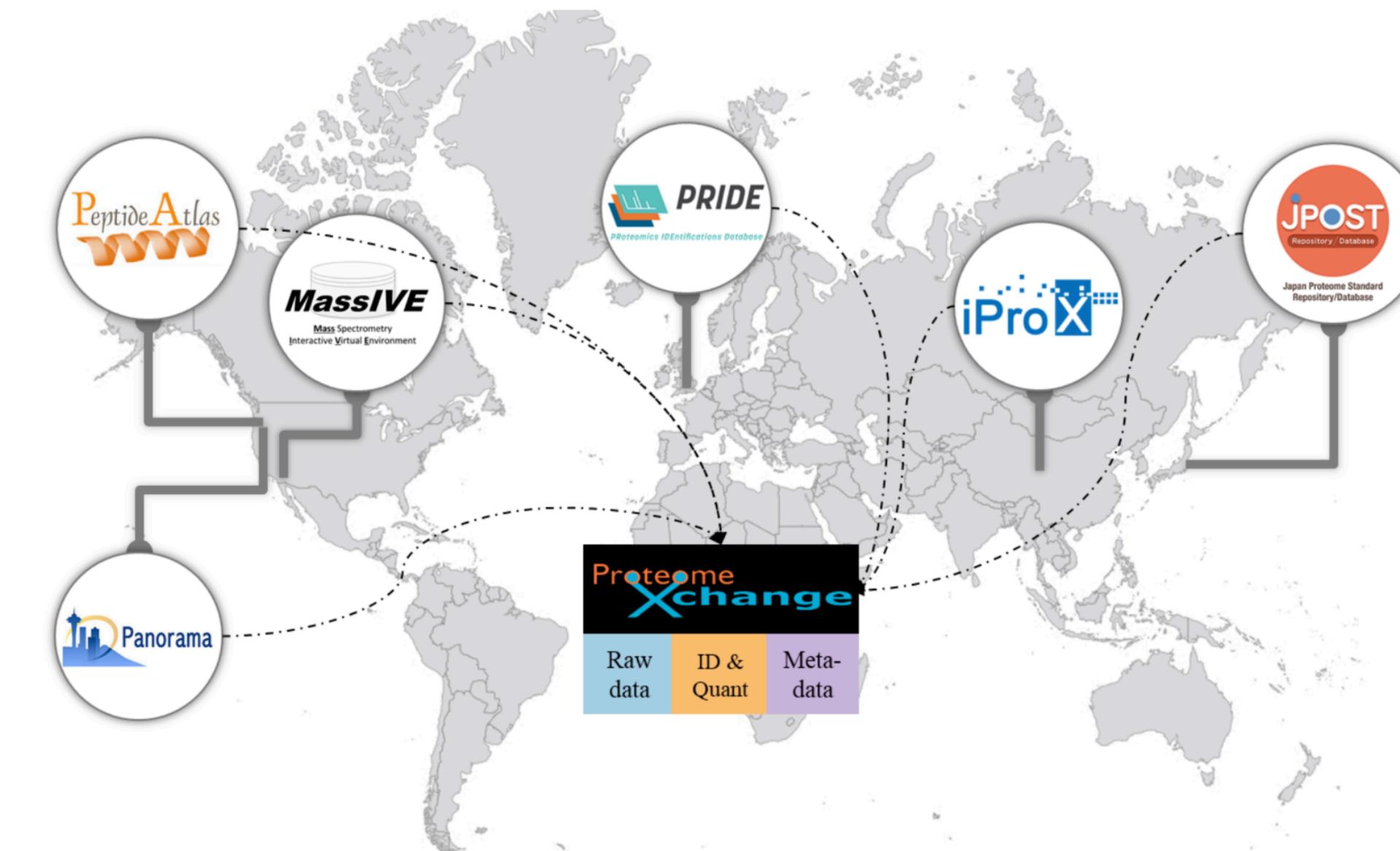
Pioneers in Zurich

Ruedi Aebersold, ETH Zurich, [Biognosys](#) was founded from his lab

Paola Picotti's lab, [pioneer of LiP-MS technology](#)

Proteomics data repositories and guidelines

ProteomeXchange

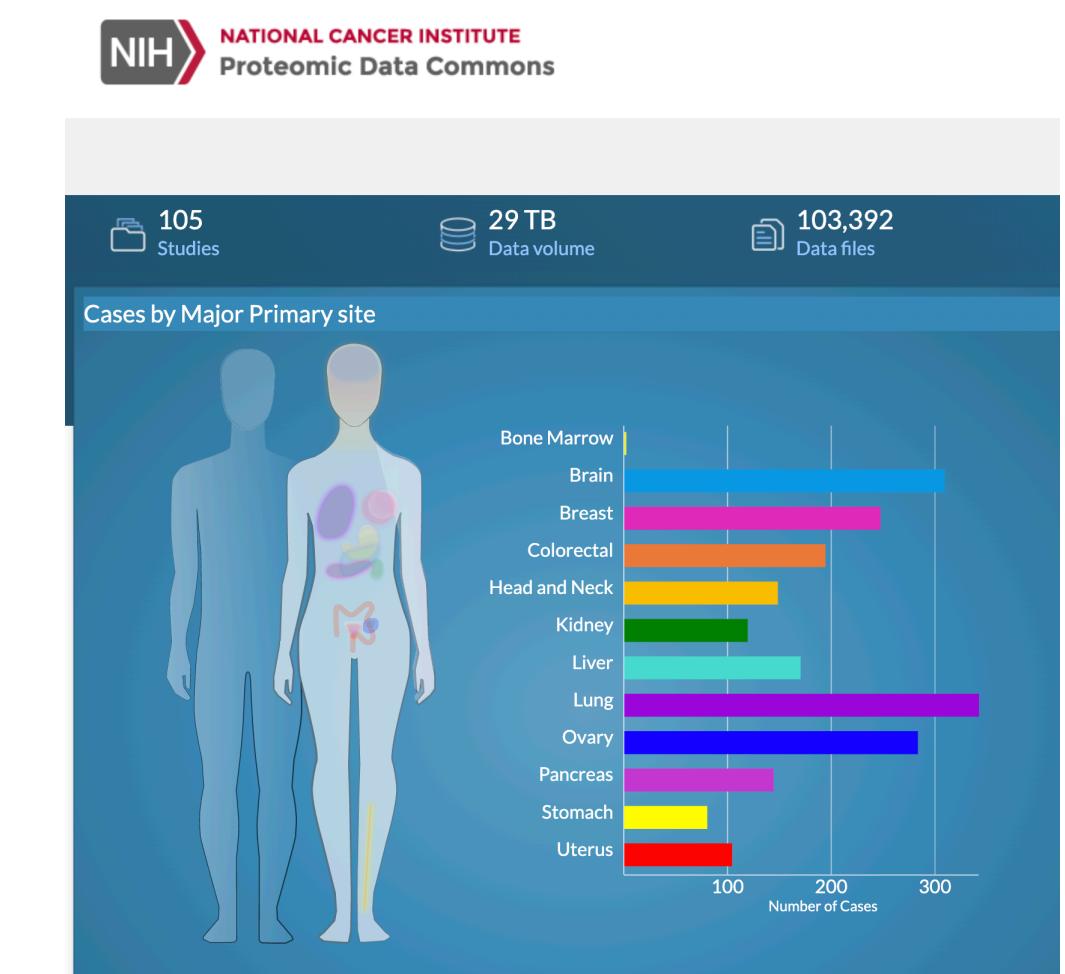


MassIVE.quant

Meena Choi et al, Nature Methods , 2020

Data Portal at CLINICAL PROTEOME TUMOR ANALYSIS CONSORTIUM (CPTAC)

Has genomics and Proteomics data



Acknowledgments

The source of all pictures and resources have been acknowledged in slides

ACF Survey summary acknowledgements

The Australasian Core Facilities Meetings are coordinated by Ben Crossett (USyd) and Ralf Schittenhelm (Monash).

This survey was originally designed with assistance from Paula Burton (Mass Dynamics), Mark Condina (while at UniSA) and the contents has been refined at the last two ACF meetings.

We would like to thank: Matt Padula (UTS), Tara Pukala (Adelaide) and Nick Williamson (UniMelb) for ‘beta testing’; the 37 facilities that complete the survey and Naveed Nadiv (USyd) for expert assistance in compiling this summary.

Any enquiries, please contact: ben.crossett@sydney.edu.au