Lecturers:

- Marco Reverenna PhD student: marcor@dtu.dk
- Henry Emanuel Webel Senior Data Scientist: <u>heweb@dtu.dk</u>
- Alberto Pallejà Caro Team Lead Data Science Platform: apca@biosustain.dtu.dk
- Alberto Santos Director, Informatics Platform : <u>albsad@biosustain.dtu.dk</u>

Date: 14-05-2025

Groups: Multi-omics Networks Analytics & Data Science Platform

Novo Nordisk Foundation Center for Biosustainability

Time	Торіс	Lecturer
8.30 - 10.00	 Proteomics Peptide and protein identification Protein quantification 	Marco Reverenna
10.00 -10.30	Little break	
10.30 - 12.00	 Steps in data processing (using <u>QuantMS</u>) Sample Data Relationship Format (SDRF) FASTA file to define search space Spectrum files from Mass-spectrometer Running QuantMS to process spectra to identified and quantified peptide sequences 	Henry Emanuel Webel
12.00 -13.00	Lunch (sandwiches are provided)	
13.00 - 14.30	Steps in statistical analysis - Brief reference to output formats overview, but focus on using QuantMS	Alberto Santos
15.00 - 16.30	 Basic statistical analysis of a two-group experiment with one timepoint (option1) or four timepoints (option2) Peptide to protein (group) aggregation Downstream data analysis of proteins (using <u>analytical core library</u> developed at DTU Biosustain and other Python libraries) Building a report with <u>vuegen</u> reports (developed at DTU Biosustain) 	Henry Emanuel Webel

Our trip into the proteomics world!



0^{pportunities} Applications spectrometer portuger of the proteomics of the proteomi

Nass spectrum Precursor mass protein patters people intering protein interence protein protein interence protein protein interence protein protein interence protein interence

Proteomics overview



Protein quantification



Proteomics

Forbes

Jun 23, 2021, 04:38pm EDT | 8,993 views

Proteomics: The Next Truly Massive Investing Opportunity



Stephen McBride Former Contributor ① Markets *The editor of RiskHedge Report*

Follow

Proteomics is like a superpower. It lets doctors peek inside your body to see what's happening in real time. Proteins signal when something is wrong inside your body. They can confirm when an illness is underway, long before we even feel sick.

As Joshua LaBaer, founder of Harvard's Proteomics Institute, said: *"It's the proteins we can measure before anything else*." Illnesses like the flu, COVID, and HIV are already diagnosed through protein tests.

DESEASE DETECTION

Jan van Oostrum from Novartis reported that the company's proteomics program has already resulted in two new drug candidates in clinical development. Additionally seven drug candidates are in the pipeline in preclinical development. Van Oostrum also revealed that in addition to the drug candidates, two biomarkers are now being evaluated in clinical studies.





Global Proteomics Market Poised for Significant Growth, Projected to Reach USD 134.82 Billion by 2035 | Future

Market Insights, Inc.

The USA's proteomics market grows steadily, driven by NIH funding in personalized medicine, while Canada is expected to grow at a 12.1% CAGR during the forecast period. Increasing adoption of precision medicine and sophisticated diagnostics in specific target disease treatment, as well as the increased concentration of target diseases, are significant drivers of growth prospects.

January 09, 2025 10:30 ET | Source: Future Market Insights Global and Consulting Pvt. Ltd.

OPPORTUNITIES

INVESTMENT



Introducing proteomics





"Proteomics is the study of interactions, function, composition, and structure of proteins and their cellular activities."

Proteins drive function and phenotype



Proteins exist in many modified forms, called proteoforms. These result from combinations of **posttranslational modifications** (PTMs).

Protein kinase A has 45 potential phosphorylation sites \rightarrow With just 5 events: 1.2 million proteoforms (C(45, 5)) \rightarrow With 9 events: 890 million proteoforms (C(45,9))

Even small system get complex fast \rightarrow 11 sites + 5 events = 462 potential proteoforms



Different post translational modifications



2 most common PTMs

- Phosphorylation: Addition of a phosphate
 group to proteins, often regulating cell
 signalling and protein function.
- Acetylation: Attachment of an acetyl group,
 commonly influencing gene expression and
 protein stability.

Ölzscha, H. Biological Chemistry, vol. 400, no. 7, 2019, pp. 895–915 (https://doi.org/10.1515/hsz-2018-0458) Shahin Ramazi, Javad Zahiri, Post-translational modifications in proteins: resources, tools and prediction methods, Database, Volume 2021, 2021, baab012, https://doi.org/10.1093/database/baab012

biomarker discovery single-cell proteomics medical microbiology cancer proteomics evolution archaeology microbiome studies crop development protein identification host parasite interactions immunoproteomics functional annotations cell signalling clinical proteomics drug design

Decoding peptides with tandem mass spectrometry





From proteins to peptides: the role of proteases

Proteases are enzymes that cleave proteins into smaller peptides by cutting specific peptide bonds. Each protease has its own cleavage specificity, depending on the amino acid sequence. **Trypsin** is the most used protease in proteomics.

- It cleaves after lysine (K) and arginine (R)
- Generates mostly small (0.5-3kDa) di-charged peptides (N-term and R/K)
- Can't get 100% of coverage by using only trypsin



NRRPCHSHTKECESAWKNRPCHSHTKKPCHSHTKKNRKVWKIPPFFW



Bottom-up proteomics: digest first, identify later



Most of the experiment are done by following the **bottom-up approach**

This method involves the fragmentation of proteins into smaller peptides, followed by the reconstruction of the protein sequences.

Protein identification

What does a MS1 spectrum tell us?

Y axis How abundant that specific ion was in the sample at the time of detection; higher peak indicate ions that are more abundant and lower peaks indicate ions that are less abundant



X axis

Mass-to-charge ratio is calculated by dividing the ion's mass by its charge. Each peak position along the axis tells you the size of the detected ion



Precursor masses provide the overall mass of the peptide, but not the sequence!

The sequence define which peptide it is, and so which protein belong to, that 's why it is so import to know!

In MS/MS, this precursor is isolated and broken into fragment ions for identification

Residue masses: the basis for peptide sequencing



Residue mass = Amino acid mass – H2O

Amino Acid	Abbreviation		Residue Mass (Da)
Glycine	G		57.05
Alanine	А		71.08
Serine	S		87.08
Valine	V		99.13
Leucine	L	1	113.16
Isoleucine	I	Isomers	113.16
Threonine	Т		101.11
Cysteine	С		103.15
Proline	Р		97.12
Phenylalanine	F		147.18
Tyrosine	Y		163.18
Tryptophan	W		186.21
Aspartic Acid	D		115.09
Glutamic Acid	Е		129.12
Asparagine	Ν		114.10
Glutamine	Q		128.13
Lysine	K		128.17
Arginine	R		156.19
Histidine	Н		137.14
Methionine	М		131.19





Break point around the peptide bond depends on the fragmentation method

Peptide fragmentation can be controlled by instrument parameters (pressure, collision energy ...)

How fragment ions reveal the peptide sequence



Bear in mind

1. Practical to start with y1: Lysine (147) or

Arginine (175)

- 2. y1 may not be observed
- 3. b1 almost never observed
- 4. Leu/lle are isobaric
- 5. Read the sequence from the end





Massive amount of mass spectra data



1h LC-MS/MS produces 5000 MS1 spectra, the equivalent of 100000 MS2 spectra







Protein database can help to identify our peptides!



sequenced genomes

Peptide spectra matching (PSM) can be still wrong





Pros

- Easily automated for high throughput applications (millions of spectra)
- Can get matches from spectra that are difficult to interpret

Cons

- Can produce matches from marginal data
- It is slow if no enzyme specificity is used
- Can be slow if many variable modifications are allowed
- Can be slow if large data sets are searched
- You can detect only peptides which are in the peptide/protein database

In large scale analysis we will always produce matches, even if they might be wrong.

To identify proteins, we need high quality PSMs obtained by mass spectrometry instruments with high accuracy

We need to create a scoring metric to find out which of the matches are plausible

Scoring peptides help to select the best candidates



Rank	Peptide	Score
1	TADKLQEFLQTLR	225
2	TADKNQKFLQTLR	152
3	TANELQEFLQTLR	89
4		

PSM with highest score is chosen and used for protein identifications

PEP1	#	Y ions	PE
Т	13		Г
А	12	1461.8	A
D	11	1390.8	
K	10	1275.7	ľ
N	9	1147.6	
Q	8	1033.6	
K	7	905.5	k
F	6	777.46	F
L	5	630.39	
Q	4	517.3	
Т	3	389.3	Г
L	2	288.2	
R	1	175.1	F

PEP2	#	Y ions
Т	13	
А	12	1461.8
D	11	1390.8
K	10	1275.7
Ν	9	1147.6
Q	8	1033.6
К	7	905.5
K F	7 6	905.5 777.46
K F L	7 6 5	905.5 777.46 630.39
K F L Q	7 6 5 4	905.5 777.46 630.39 517.3
K F L Q T	7 6 5 4 3	905.5 777.46 630.39 517.3 389.3
K F L Q T L	7 6 5 4 3 2	905.5 777.46 630.39 517.3 389.3 288.2

DTU

score: 152

score: 225

Mascot Ion Score

Usually, we deal with more than 100'000 PSMs, not all of them are good quality, resulting in a distribution of high and low scoring PSMs. How do we decide which peptides is wrong or right? FDR permits to control the number of incorrect matches and so minimise mistakes!



Bottom-up proteomics identifies peptides not proteins!

Many proteins share identical sequence parts, so proteins can only be identified if a unique peptide is present



No info about presence or absence of protein 3



Manual spectrum interpretation

Spectrum matching

Protein quantification

Protein quantification refers to measuring the **relative** or **absolute abundance** of peptides or proteins in a biological sample.

- Enables comparison between different conditions (e.g. treated vs control)
- Critical for understanding cellular changes, disease states or drug response
- Can be performed label-free or with stable isotopic labelling
- Used in biomarker discovery, systems biology and clinical proteomics





Protein quantification can be performed either relatively, comparing expression between conditions, or absolutely, determining exact molecule counts using internal standards.

Relative quantification

Compares same protein across conditions Reports fold changes No info on absolute amount Good for global, untargeted profiling Used in discovery studies Hypothesis generating

Absolute quantification Determines number of molecules Reports molecular counts Requires internal standards Ideal for targeted, hypothesis-driven analysis Enables stoichiometry comparisons Hypothesis testing

There are two main acquisition methods are used in proteomics:

- 1) Data Dependent Acquisition (DDA) selects ions based on intensity
- 2) Data Independent Acquisition (DIA) fragments all ions in a defined m/z window.

Fragmentation trigger

Selection

Identification

Quantification

Use case

Drawback

DDA (Data Dependent)	DIA (Data Independent)
Most intense precursor ions	All precursors in predefined windows
Stochastic (intensity-based)	Systematic and unbiased
MS2 spectra matched to libraries	Requires deconvolution and libraries
Label-free or labeled (TMT, SILAC)	Label-free or plexDIA
Discovery workflows	Large-scale, high reproducibility
Missing values	Complex data analysis

In label-free quantification, peptide intensity is measured across multiple runs. Matching ions by both m/z and retention time ensures accurate quantification.

- 1) m/z \rightarrow identifies peptides
- 2) Intensity \rightarrow reflects amount of peptide
- 3) Retention time \rightarrow increase match reliability
- 4) XIC \rightarrow tack specific ions across runs



Label-free quantification (LFQ)

In label-free quantification (LFQ), protein amounts are inferred by comparing ion intensities across multiple LC-MS runs and **no labelling is needed.**

- Each sample is analysed separately by LC-MS.
- Peak intensities of same ions are compared
- Identification is typically done via MS2
- No isotope labelling is involved

Pros: cost-effective, simple workflow, no chemical labelling steps.

Cons: run-to-run variability, missing values, requires normalization.



	SILAC & SILAM (in vivo)	iTRAQ & TMT (in vitro)
Labelling method	Isotopic amino acids fed to cells	Chemical tagging with isobaric labels
When labeled	During protein synthesis (in cell culture)	After protein digestion
Samples processed	Separately, then mixed before MS	Pooled and analyzed together
MS comparison	Mass shift seen in MS1	Reporter ions detected in MS2 or MS3
Multiplexing	Usually 2–3 samples	Up to 16–35 samples (TMTpro)
Pros	High precision, direct incorporation	High throughput, compatible with any sample
Cons	Limited to cultured cells	Ratio compression (fixed with SPS- MS3)

iBAQ (Intensity Based Absolute Quantification) is a method used to estimate absolute protein abundance from mass spectrometry data.

iBAQ= Total intensity of all detected peptides/ Number of theoretically observable peptides

It accounts for protein length and tryptic behaviour, making protein intensities comparable across the proteome.

Why use it?

- 1. Allows comparison between proteins, not just across samples
- 2. Compatible with DDA and DIA workflows
- 3. With spike-in standards, gives absolute protein copy numbers

QuantMS to analyse our data

nature methods

Explore content V About the journal V Publish with us V

<u>nature</u> > <u>nature methods</u> > <u>brief communications</u> > article

Brief Communication Open access Published: 04 July 2024

quantms: a cloud-based pipeline for quantitative proteomics enables the reanalysis of public proteomics data

<u>Chengxin Dai, Julianus Pfeuffer, Hong Wang, Ping Zheng, Lukas Käll, Timo Sachsenberg, Vadim</u> <u>Demichev, Mingze Bai, Oliver Kohlbacher & Yasset Perez-Riverol</u>

Nature Methods 21, 1603–1607 (2024) Cite this article

11k Accesses | 17 Citations | 22 Altmetric | Metrics

Abstract

The volume of public proteomics data is rapidly increasing, causing a computational challenge for large-scale reanalysis. Here, we introduce quantms (https://quant.ms.org/), an open-source cloud-based pipeline for massively parallel proteomics data analysis. We used quantms to reanalyze 83 public ProteomeXchange datasets, comprising 29,354 instrument files from 13,132 human samples, to quantify 16,599 proteins based on 1.03 million unique peptides. quantms is based on standard file formats improving the reproducibility, submission and dissemination of the data to ProteomeXchange.



Let's have a break!



Acknowledgments





Multi-omics Networks Analytics



The Novo Nordisk Foundation Center for Biosustainability