



# Bioinformatics Master Thesis



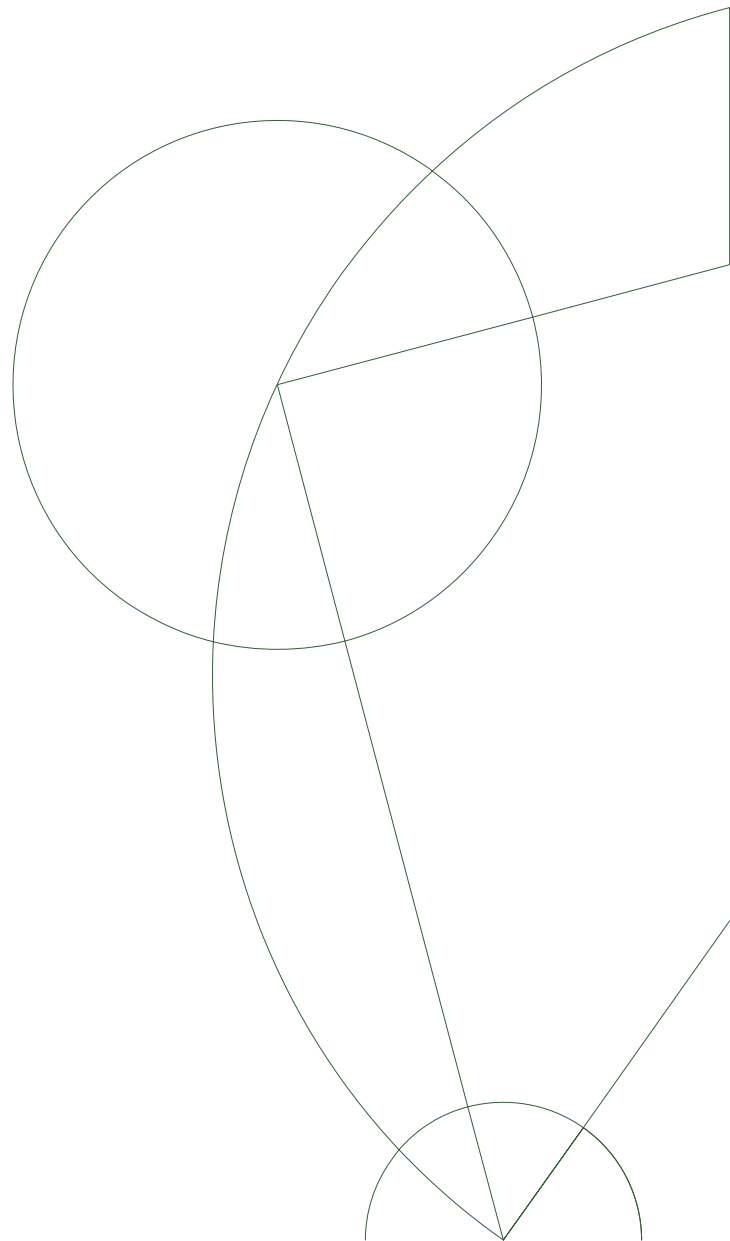
## Topics in Mass Spectrometry and Proteomics

Antonio Ortega Jiménez      <ntoniohu@gmail.com>

### Supervisors

Thomas Hamelryck      <thamelry@gmail.com>

Mathias F. Gruber      <mafg@novozymes.com>



# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Aminoacids and proteins . . . . .	2
1.2	The protein-focused biotechnology industry . . . . .	3
1.3	Mass spectrometry and shotgun proteomics . . . . .	4
1.4	The mass spectrometer . . . . .	5
1.4.1	The ion source . . . . .	6
1.4.2	The mass analyzer . . . . .	6
1.4.3	The detector . . . . .	6
1.4.4	Sample preprocessing . . . . .	7
1.4.5	Peptide Mass Fingerprint (PMF) . . . . .	8
1.4.6	Tandem MS (MS/MS) . . . . .	8
1.5	Spectra processing: search engines . . . . .	9
1.5.1	X!Tandem . . . . .	10
1.6	Results validation . . . . .	11
1.7	Protein quantification . . . . .	12
<b>2</b>	<b>Materials and Methods</b>	<b>13</b>
<b>3</b>	<b>Results and Discussion</b>	<b>14</b>
<b>4</b>	<b>Conclusion</b>	<b>15</b>
4.1	Overview of proteomics analysis software . . . . .	15

4.1.1	MaxQuant . . . . .	15
4.1.2	Compomics group . . . . .	17
4.1.3	Other suites . . . . .	18
4.2	Label-free quantification . . . . .	18
4.2.1	MaxLFQ . . . . .	18
4.2.2	Bayesian model . . . . .	22
4.2.3	Data . . . . .	23
4.3	Extension of pipelines . . . . .	24
<b>5</b>	<b>Bibliography</b>	<b>26</b>

## Abbreviations

**MS1** first tandem MS analyzer

**MS2** second tandem MS analyzer

## **Preface**

# Chapter 1

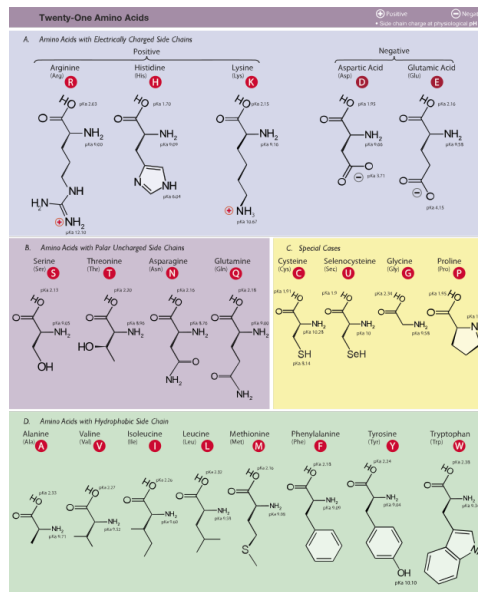
## Introduction

### 1.1 Aminoacids and proteins

Proteins represent the last link in the central dogma of biology, where information encoded in DNA, is transcribed to RNA for posterior translation into proteins at the ribosome.

Proteins are made up of 20 basic units, called aminoacids. All aminoacids share a common chemical structure, where a carbon atom ( $C_{\alpha}$ ) is covalently bonded to a hydrogen atom, a carboxyl group, an amino group, and last but not least, a radical, also called side chain of the aminoacid (see figure ??). The side chain differs between aminaocids and generates them from each other. A slight deviation from this pattern exists in proline, where the radical is bound to the nitrogen atom, making it an iminoacid. Even though the side chains are all different, they can be classified into four different groups: aliphatic, polar, positively charged and negatively charged (see figure 1.1).

Two aminoacids are joined together through the formation of a peptidic (covalent) bond between them. Such a linkage is formed by removal of the



**Figure 1.1**

elements of water (dehydration) from the  $\alpha$ -carboxyl group of one amino acid and the  $\alpha$ -amino group of another. CITE LEHNINGER PAG 82. The remaining  $\alpha$ -amino and  $\alpha$ -carboxyl groups are available for linkage to other aminoacids, and in this way peptidic chains or peptides can be created.

While there are 20 basic units that constitute the majority of naturally observable proteins, their side chains can be modified both by physiological processes and by experimental procedures CITE. One frequent instance of such modifications is the oxidation of methionine

## 1.2 The protein-focused biotechnology industry

Proteins carry out most of the cell's molecular functions, they work as molecular agents that can perform an extremely wide range of tasks. The advent of biotechnology has sought to take advantage of this power, either by using proteins as present in natural conditions (wild type) or engineered by

humans. This potential economic activity is carried out by several biotech companies, including Novozymes. Novozymes A/S is a company whose line of business consists of the development of enzymatic products performing chemical transformations in different industrial processes. The application of these products, instead of conventional chemical-based solutions, has the advantage that they require less chemical substances, potentially simplifying industrial processes, reducing their costs and their environmental impact. Notorious examples of such applications include wastewater treatment, household care and the baking industry.

In order to place Novozymes ahead of its competitors, the company has utmost interest in developing protein research and several departments in the organization approach the study of proteins and their translation to the market through among other things the application of mass spectrometry.

### **1.3 Mass spectrometry and shotgun proteomics**

Life systems consist of complex systems, meaning their behaviour cannot be easily explained by analyzing the individual elements alone. Moreover, they present multiple layers of complexity, given by the nature of the elements that make it up. The layer provided by proteins is one of them, and its study is called proteomics. It is a complex layer because thousands of different proteins can be present in a single cell at any time, and their exact composition and quantities constantly change, responding to the stimuli of the surrounding environment. The study of the protein-specific complexity is called proteomics. With proteomics, one endeavors to infer the protein composition of a sample, and eventually quantify its protein amounts.

It may be useful to divide the existing approaches into two types of paradigms: top-down and bottom-up. In the top-down paradigm, intact proteins are

directly used for the analysis. In the bottom-up paradigm, the proteins are first cleaved into smaller parts, and these parts are then used for identification, characterization, and quantification. These smaller parts are called peptides. CITE 1.6.1 COMPUTATIONAL METHODS. Such peptides acquire physicochemical properties fitting the requirements of the downstream analytical methods, mainly mass spectrometry (MS), which performs the data acquisition. The bottom-up paradigm is most often used because peptides are much more suitable to analysis by mass spectrometry. The interested reader is directed to CITE 1.6.1 COMPUTATIONAL METHODS to learn why. The top-down paradigm will be ignored in the rest of the manuscript. A summary of the bottom-up approach MS analytical pipeline follows.

MS is performed by means of a mass spectrometer, a piece of equipment that can acquire plentiful mass measurements for enough sample components, generating a dataset that, with adequate computational analysis tools, is enough to perform the inference steps required to gather knowledge about the original protein sample. These inference steps can be condensed to the peptide and protein inference problems. A third computational problem needs to be solved if quantitative, and not just qualitative information, is to be gained from the experiment. This is the quantification problem CITE.

## **1.4 The mass spectrometer**

The mass spectrometer consists of three main parts: an ion source, a mass analyzer, and a detector SEE FIGURE.



### **1.4.1 The ion source**

All mass spectrometers exploit the physical properties of mass and electric charge exhibited by the analyzed components. Thus, ionization of these is required prior to any measurement. This is achieved in the ionization source CITE 5.1 COMPUTATIONAL METHODS. The most frequent ionization methods in proteomics are Matrix-Assisted Laser Desorption-Ionization (MALDI) and Electro Spray Ionization (ESI) CITE. As a result of ionization, the components acquire positive charge, usually by the addition of 1, 2, or 3 protons, that correspondingly give +1, +2 or +3 charge. The acquired charge yields a mass/charge ( $m/z$ ) ratio, a property that can be applied in the following component separation and mass measurement of the separated components steps in the analyzer.

### **1.4.2 The mass analyzer**

The plethora of separation methods is reflected upon the range of different analyzers available, mainly quadrupole (Q), time of flight (TOF), Ion trap (IT) and Fourier Transform (FT). These apply different principles to perform the same task: handling of the individual molecular components of the ionized sample for posterior sequential measurements of the components one by one.

### **1.4.3 The detector**

the intensity of a given  $m/z$  ratio within an observable window. This analysis returns the MS spectrum. It is up to the MS technician to decide the best protocol, according to the particularities of the equipment, and the dataset.

#### 1.4.4 Sample preprocessing

An MS experiment starts with the generation of protein mixtures samples. They are first separated in order to sort the proteins via physicochemical criteria. This is most frequently carried out via SDS-PAGE based on mass or isoelectric point. Once the electrophoresis is completed, protein bands can be excised from the gel. Each band will contain a subset, or even only one of the proteins originally available, thus making the downstream analysis simpler

REFERENCE COMPUTATIONAL METHODS CHAPTER 2

The reduced complexity protein mix is extracted in a denaturalized state and subjected to enzymatic digestion with specific enzymes, that can cut the chain of aminoacids following a predictable pattern, key for the downstream analysis. The enzyme most frequently sued for this is Trypsin, which cuts peptidic bonds whenever a positively charged residue, either Lysine or Arginine, lies on the carboxyl side of the peptidic bond. Even though enzymes are very specific, they can miss some of their targets, due to steric inaccessability or the presence of specific aminoacids that can weaken their function. This is the case of Trypsin whenever the residue on the other side of the peptidic bond is Proline. This variability, though limited, needs to be taken care of in downstream analysis.

The result of this process is a mix of peptides following a length distribution given by the cleavage sites frequency and each protein's aminoacidic composition. As explained above, this length distribution is fitted to the resolution of the MS analyzer, thus optimizing the throughput of the method. The downstream workflow now diverges based on the simplicity of the original protein sample. When it consisted of a single protein, Peptide Mass Fingerprinting (PMF) is used, otherwise tandem MS (MS/MS) shall be performed.

#### 1.4.5 Peptide Mass Fingerprint (PMF)

If the original sample was known to contain a single protein, PMF, or *protein-centric* proteomics, is conducted. In PMF, the mixture of peptides can be already transferred to the spectrometer, where a spectrum containing a peak for every  $m/z$  ratio present in the ionized peptide mix will be recorded. Thus, spectra generated this way can be considered a pattern, or fingerprint, of the peptides making up the original protein. Therefore, the spectra are intended to contain enough information to identify a protein.

#### 1.4.6 Tandem MS (MS/MS)

When presented with the problem of analyzing a mixture of proteins, the capacities of mass spectrometers are easily overcome by a too complex mixture, resulting in the analysis of only a minor part of the total protein complement of the sample. This can be overcome by dividing the initial sample into fractions, and using a series, or tandem, of spectrometers in the analysis. The spectrometers are used to analyze each obtained fraction separately, using different schemes. Fractionation is usually achieved by different methods of separation CITE 1.7, most commonly via HPLC methods.

HPLC methods work by loading the peptide mix in a column containing a stationary and a solid phase. These phases create an environment where peptides interact differently based on their physico-chemical properties, set by the nature of the phases. The output of the column, called elute, will consist of subsets or fractions of peptides leaving the column at different retention times. Therefore, the input to the machine will consist of a few peptides at a time. The two most common methods in proteomics are reverse phase chromatography (separating on hydrophobicity) and strong cation exchange chromatography (separating on charge) CITE 4.2 compu-

tational methods.

The tandem MS analysis starts when a peptide accesses the analyzer. It is ionized in the ion source and accesses the first mass spectrometer. While different ways of handling the peptides are available, we will focus on the product ion scan. In this protocol, the first analyzer is used to select ionized peptides within a narrow  $m/z$  window. The peptides passing this first scan are then subjected to fragmentation, most often via collision-induced dissociation (CID). Briefly, the peptides enter a collision cell containing an inert gas. Given enough kinetic energy, hits of ionized peptides and the gas will trigger the fragmentation of the peptide into smaller units. PAG 123 COMPUTATIONAL METHODS. The most frequently occurring fragment types by far for CID are the b and y ions and partly a ions. pag 134 computational methods. The produced fragments then enter the second analyzer, where a  $m/z$  spectrum of the fragments is recorded. Thus, unlike in PMF, where the spectrum recorded reflects the  $m/z$  ratios acquired by the protein peptides cleaved by the enzyme, tandem MS spectra on product ion scan mode record the  $m/z$  ratio of the fragments produced by an ionized peptide with a given  $m/z$  ratio. The  $m/z$  ratio of this precursor ion is changed during the run, thus, multiple spectra are obtained where PMF would create only one.

## 1.5 Spectra processing: search engines

A search engine can be used to map the recorded pattern of  $m/z$  ratios to a protein entry in a database. This task is performed by *in silico* cleavage of each sequence entry in the database based on the specific cleavage pattern of the enzyme used, coupled with the simulation of the expected spectrum based on the expected peptides.

Given the stochastic nature of the cleavage and spectra recording process, the resulting spectra exhibit variability manifested in missing peaks or spurious ones. Furthermore, the peptide to spectrum matching (PSM) process against a sufficiently big database can, at random, return wrong matches. This translates to the obtention of multiple matches, of which one, if any, will be correct. Therefore, the lists of matches need to be somehow ranked. The issue is addressed by search engines through the usage of scoring systems that measure the goodness of the match. Assuming the correct protein is present in the database, a good system should give the protein the best score. This way, proteins can be identified.

Two steps in protein identification can be distinguished:

1. Peptide inference: infer the peptides present in the sample.
2. Protein inference: based on the inferred peptides, infer what proteins generated them.

Both are taken care of by the search engine. Multiple search engines exist that implement different matching and scoring algorithms.

### **1.5.1 X!Tandem**

X!Tandem was one of the first open source search engines in mass spectrometry. It produces a score based on the dot product between the theoretical and the experimental tandem mass spectra CITE Overview of Tandem Mass Spectrometry (MS/MS) Database Search Algorithms. The scores assigned to wrong matches are assumed to follow a hypergeometric distribution, allowing the program to extrapolate  $E(s)$  (expected number of wrong matches at a given score) for any score.

## 1.6 Results validation

The scoring system implemented by the search engines can be used to filter and validate the results. A basic common filter is the false discovery rate (FDR), usually set to 1%, indicating that 1 out of a hundred filtered matches are expected to be false positives. The most commonly used method to compute the FDR is the target-decoy search. The search engine tries to match the same spectra against a decoy database, usually generated by reversing the sequences present in the original database (target). All matches to the decoy will be regarded as wrong. Thanks to the fact that the basic properties of the decoy (size, composition, etc) are identical to the target, whenever a mistake is made, it is as likely to happen in both databases CITE COMPOMICS TUTORIAL 1.5, thus the number of matches in the decoy provides an accurate estimate of the number of random matches, or false positives, against the target database ( $n_{fp}$ ). Together with the number of PSMs passing a threshold score ( $n_{tp} + n_{fp}$ ), the FDR can be computed using the formula below.

$$FDR = \frac{n_{fp}}{n_{tp} + n_{fp}}$$

The minimal FDR at which a given PSM is considered a positive match constitutes the PSM's q-value

$$q(PSM_i) = \min FDR \forall PSM_i \in Positives$$

Finally, a posterior error probability (PEP) can be defined as the probability of the match being random ( $P(s_i|T_i = 0)$ ). This can only be done if a statistical model describing the distribution of scores for correct and wrong matches is fit to the dataset.

## **1.7 Protein quantification**

CONTINUE maybe add reason for the project.

## **Chapter 2**

# **Materials and Methods**



## **Chapter 3**

# **Results and Discussion**

## Chapter 4

# Conclusion

### 4.1 Overview of proteomics analysis software

#### 4.1.1 MaxQuant

One of the most used program for the quantification of label free proteins is the MaxQuant suite [1]. It can be downloaded for free and is developed by the Max Planck institute in Germany. It consists of a user friendly GUI that provides the most needed steps in a proteomics pipeline.

MaxQuant has been successfully adopted by the scientific community due to its ease of use and a comprehensive pipeline. A Google group and a tutorial are available <sup>1 2</sup>.

The main processing steps incorporated in MaxQuant consist of:

- Read MS spectra files in the .RAW format, the closed format produced by ThermoScientific MS analyzers.

---

<sup>1</sup><http://www.coxdocs.org/doku.php?id=maxquant:viewer:tutorial>

<sup>2</sup>[https://www.youtube.com/watch?v=\\_AJHFHi5CxM](https://www.youtube.com/watch?v=_AJHFHi5CxM)

- Contains the Andromeda search engine [2] for matching of MS2 spectra to a proteome database.
- Andromeda supports configuring aminoacid modifications and decoy search for FDR (false discovery rate) estimation.
- Increase peptide identifications by performing match between runs (MBR), which aims at transferring peptide to spectrum matches (PSMs) across replicate runs, based on precursor mass and retention times.
- Filters Andromeda results to provide a list of inferred peptides and proteins with a cutoff FDR value.
- Supports label-based and label-free quantification.
- Its output can be passed to the Perseus software for data visualization.

However, MaxQuant also suffers from some problems:

1. MaxQuant only runs on Microsoft Windows, impeding its integration in automated pipelines on cluster environments [3]. The processing steps cannot easily be fine tuned or exchanged with other pipelines. For example, only the Andromeda search engine is supported. This is a serious drawback, as the integration of results from several search engines, like MSGF+ [4] or MS-Amanda [5] has been shown to further improve results [6].
2. It does not provide a command line interface (CLI), thus all analyses must be performed through the GUI. This hinders reproducibility and scalability.
3. While it is supported in Linux through Mono, the user experience is best in Windows. As Linux is by far the most used OS in Bioinformatics, this implies that an additional OS is required to get the best

experience with MaxQuant.

4. Lacks a well written official documentation, as most information is made available only through talks published on Youtube or third party tutorials.

#### 4.1.2 Compomics group

The Compomics group at Ghent University <sup>3</sup> provides high quality analysis software for the identification and integration of results through the SearchGUI, PeptideShaker and moFF (modest Feature Finder) programs [7] [8] [3].

- SearchGUI provides a common interface to a range of search engines so that multiple engines can be used in a straightforward manner.
- PeptideShaker reads the SearchGUI output and performs quality control, gene set enrichment analysis (GSEA), and implements multiple validation filters to provide the best results.
- moFF reads from PeptideShaker output and implements a match between runs (MBR) model analogous to that in MaxQuant, to increase the number of spectra that are matched to a peptide. Furthermore, it summarises the MS1 peaks into refined features, making downstream analyses more sensitive.

SearchGUI and PeptideShaker are very well integrated and documented <sup>4</sup>. Both provide not just beautiful GUIs, but also a comprehensive and extremely well documented CLI <sup>5</sup> <sup>6</sup>. They easily enable the development of analytical pipelines finely fitted to each use case.

---

<sup>3</sup><https://compomics.com/>

<sup>4</sup><https://compomics.com/bioinformatics-for-proteomics/>

<sup>5</sup><https://github.com/compomics/searchgui/wiki/SearchCLI>

<sup>6</sup><https://github.com/compomics/peptide-shaker/wiki/PeptideShakerCLI>

moFF is less documented and integrated in the workflow. Furthermore, it does not offer a GUI yet. The group is currently working to better integrate moFF and PeptideShaker<sup>7</sup>.

The main issue in the Compomics suite of tools is the lack of a robust label-free quantification step, as they currently recommend the export of results to a third party tool. Only the spectral count based methods emPAi (Exponentially Modified Protein Abundance Index) and NSAF (Normalized Spectral Abundance Factor), which are not robust [9], are supported.

#### **4.1.3 Other suites**

Other software suites like ProteomeDiscoverer, Progenesis QI, GeneData Expressionist, etc, are available. They are not open, requiring a license to function. Furthermore, they are GUI oriented and make use of different data formats thus, making exchange across pipelines extremely difficult. The OpenMS [10] and Trans Proteomic Pipeline (TPP) [11] suites are open, but they also suffer from data exchangeability and lack comprehensive documentations.

## **4.2 Label-free quantification**

### **4.2.1 MaxLFQ**

The label-free quantification engine implemented in MaxQuant is MaxLFQ [12]. It performs 2 minimisation steps to infer protein quantities from extracted ion currents (XIC), as stored in the spectra files.

---

<sup>7</sup><https://groups.google.com/forum/#!topic/peptide-shaker/Lqe7lYKLcHI>

### Fractionated XIC aggregation

First, due to the prefractionation of samples during the upstream mass spectrometry (MS) analysis, XIC signals for each peptide are distributed along multiple fractions. MaxLFQ summarises the XIC from several fractions for each peptide and sample into a single peptide intensity by minimising the sum of the square of the logarithm of the intensity between samples via Levenberg–Marquardt optimization (see figure 4.1 for a graphical explanation).

The intensity of peptide  $P$  in sample  $A$  is defined as a weighted average of  $P$ 's XIC signals across fractions  $\{1..k\}$ , where the weights are the normalization factors that the minimisation algorithm seeks to find.

$$I_{P,A}(N) = \sum_{j=1}^k N_{P,A,j} \times \text{XIC}_{P,A,j} \quad (4.1)$$

The minimisation assumes that the best normalization factors minimize the sum of squared logarithm pairwise-ratios  $H$ .

$$H_P(N) = \sum_{a,b} \left| \frac{I_{P,a}(N)}{I_{P,b}(N)} \right|^2 \quad (4.2)$$

$$H(N) = \sum_p H_p(N) \quad (4.3)$$

where  $a, b$  iterates over all possible pair-wise combinations of samples where  $P$  was detected.

The peptide intensities can be computed from the estimated normalization factors and the observed XIC signals.

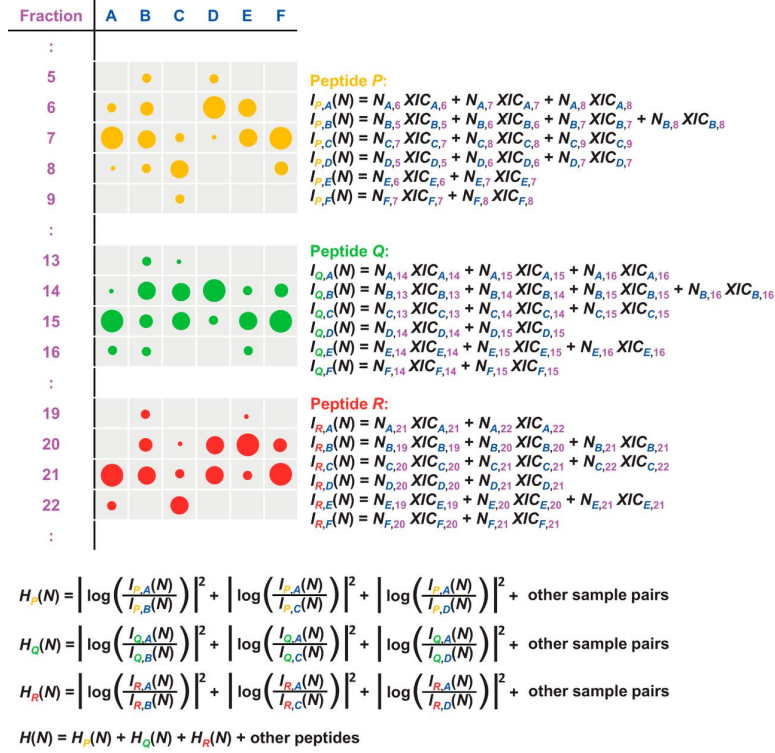


Figure 4.1 Taken from [12]

## Protein intensity inference

Second, once the peptide intensities have been summarised, the intensity of the proteins that the peptides originated from are inferred. This is achieved in three steps (see figure 4.2 for a graphical explanation).

1. Peptide intensity ratios are computed for all possible pairwise combinations by dividing the normalized intensities obtained in the previous step.
2. For every inferred protein A, its protein ratio  $r_{A,a,b}$  for the pair of samples a,b is set to the median intensity of its children peptide ratio in samples a,b. The median is selected as a summarising statistic to protect from outliers. A minimal number of non-zero intensity peptides are required for the median to be valid, usually 2. Otherwise,

the protein ratio for the pair of samples is set to 0.

3. The correct protein intensities are assumed to minimise the sum of the squared difference of the logarithm of the protein ratios and logarithm of the protein intensities  $I$ .

$$\sum_{a,b} (\log(r_{A,a,b}) - \log(I_{A,a}) + \log(I_{A,b}))^2 \quad (4.4)$$

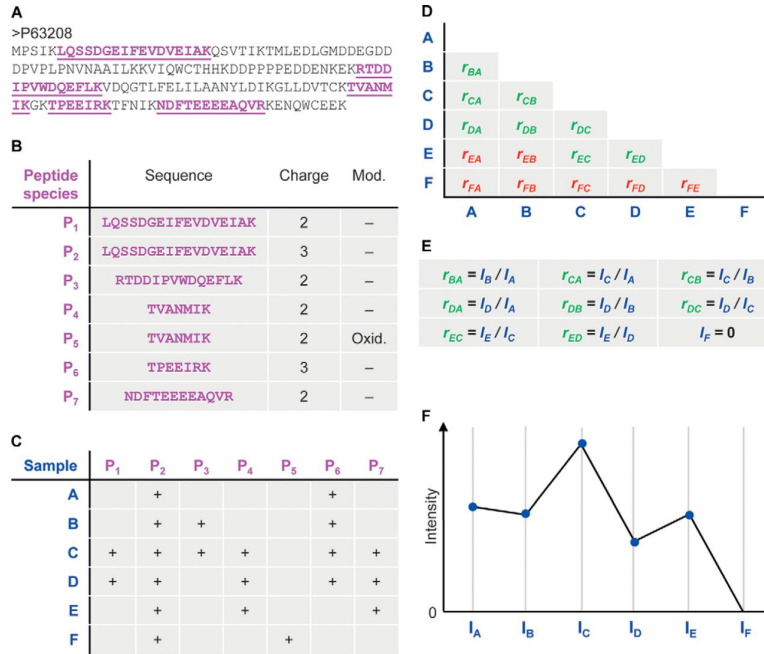


Figure 4.2 Taken from [12]

For the three most intense tryptic peptides, the signal per mole of protein was shown to be constant within a coefficient of variation of  $\pm 10\%$  [13]. Thus, protein intensities are linearly proportional to protein quantities and protein quantity ratios can be computed from protein intensity ratios. These values are reported in the MaxQuant output in the `proteinGroups.txt` file under columns named *LFQ Intensity*.



## Drawbacks of the model

The LFQ model provides a powerful method to infer protein quantities from a label-free experiment. However, it makes use of point estimates in different steps of the process, namely (1) during peptide intensity aggregation, where the Levenberg-Marquardt minimum is selected (2) during protein ratio estimation, where the peptide median is selected and (3) during protein intensity estimation, where the least squares minimum is selected. These point estimation bottlenecks, specially (2), discard increasing amounts of data that could otherwise be used for more accurate results. Moreover, not just quantities, but also the uncertainty behind them, could also be provided. The only way of MaxQuant to provide uncertainty measurements is the number of peptides supporting the quantification measurement.

### 4.2.2 Bayesian model

Bayesian statistics provides probability measurements for observed data assuming underlying mathematical models. A model for protein quantification based on the MS1 intensities or XICs could be developed to assess both quantities and uncertainties behind the estimated quantities.

While MS1 intensity measurements are dependent upon the underlying protein quantities, other factors influence the final measurement, thus adding noise and distorting the results. Besides the actual quantity, two more factors can be distinguished:

- **Sequence derived factors:** some peptides are easier to cleave for the cleaving enzyme (Trypsin, etc) than others. Moreover, the sequence of the peptide could influence the final measurement detected in the analyzer, by having different elution dynamics in the column or the

analyzer.

- **Random noise:** the stochastic processes intrinsic to MS1 measurements could also have an impact.

The model would take as input:

1. **Features extracted from the precursor sequence**, including the surrounding aminoacids in the original protein. They can be used to model the sequenced derived factors.
2. **Observed MS1 precursor intensities/XICs.**

and return as output a protein quantification value for each protein, together with a measurement of the uncertainty behind it.

#### 4.2.3 Data

The MaxLFQ paper dataset is available <sup>8</sup>. The authors submitted the RAW files (>50 GB) and the search files containing the final results. They can be used to benchmark any new quantification tool or method, for example, it has been used to benchmark third party quantification tools like StPeter from the TPP [14].

The dataset contains the results of the quantification of proteins in 2 different proteomes (*E. coli* and *Homo sapiens*) from 6 samples organized in 2 conditions with 3 replicates each. In the first condition (H), both proteomes were mixed in a 1:1 ratio, while in the second condition (L), the *E. coli* proteome had 3 times more contribution to the mix (3:1). Thus, the fold change of protein quantities across conditions should be 1 for human proteins and 3 for bacterial proteins.

---

<sup>8</sup><https://www.ebi.ac.uk/pride/archive/projects/PXD000279>

A list of the protein groups identified that could be mapped unambiguously to each species was made available <sup>9</sup>. Thus, a dataset of thousands of proteins of known ratios between conditions, producing tens of thousands of peptides with known sequence, MS1 intensity/XIC and parent protein can be built from the publication, and used to train new quantification methods (figure 4.3).

	Taxonomy	Peptides	Proteins
1	Escherichia coli (strain K12)	14483	1556
2	Homo sapiens	32647	3444
	Total	47130	5000

**Table 4.1** Summary table of the dataset

[illegible]

**Figure 4.3** Extract of the available dataset. Total of 47130 entries, one per peptide, for a global of 5000 proteins.

This compiled dataset can be downloaded here [https://mega.nz/#!/AgcTgJYa!w9DoAKYRc6u-SaRy\\_UIMz3aileUHXoaWgrxf-UycqiQ](https://mega.nz/#!/AgcTgJYa!w9DoAKYRc6u-SaRy_UIMz3aileUHXoaWgrxf-UycqiQ)

### 4.3 Extension of pipelines

The Compomics suite of programmes provides the best documented and accessible set of tools for proteomics analysis, but still lacks proper quantification tools. The implementation of a MS1 intensity-based downstream quantification tool compatible with the output of the Compomics tools

<sup>9</sup><http://www.mcponline.org/content/13/9/2513/suppl/DC1>

would yield a free Linux compatible, CLI supported, complete proteomics pipeline featuring a robust quantification method.

## Chapter 5

# Bibliography

- [1] Jürgen Cox and Matthias Mann. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nature Biotechnology*, 26(12):1367–1372, 2008.
- [2] Jürgen Cox, Nadin Neuhauser, Annette Michalski, et al. Andromeda: A peptide search engine integrated into the MaxQuant environment. *Journal of Proteome Research*, 10(4):1794–1805, 2011.
- [3] Andrea Argentini, Ludger J.E. Goeminne, Kenneth Verheggen, et al. MoFF: A robust and automated approach to extract peptide ion intensities. *Nature Methods*, 13(12):964–966, 2016.
- [4] Sangtae Kim and Pavel A. Pevzner. MS-GF+ makes progress towards a universal database search tool for proteomics. *Nature Communications*, 5:1–10, 2014.
- [5] Viktoria Dorfer, Peter Pichler, Thomas Stranzl, et al. MS Amanda, a Universal Identification Algorithm Optimized for High Accuracy Tandem Mass Spectra. *J. Proteome Res.*, 13(8):3679–3684, 2014.
- [6] David Shteynberg, Alexey I. Nesvizhskii, Robert L. Moritz, and Eric W. Deutsch. Combining Results of Multiple Search Engines in Proteomics. *Molecular & Cellular Proteomics*, 12(9):2383–2393, 2013.
- [7] Harald Barsnes and Marc Vaudel. SearchGUI: a highly adaptable common interface for proteomics search and de novo engines. *Journal of Proteome Research*, page acs.jproteome.8b00175, 2018.

- [8] Marc Vaudel, Julia M. Burkhardt, René P. Zahedi, et al. PeptideShaker enables reanalysis of MS-derived proteomics data sets: To the editor. *Nature Biotechnology*, 33(1):22–24, 2015.
- [9] Noelle M. Griffin, Jingyi Yu, Fred Long, et al. Label-free, normalized quantification of complex mass spectrometry data for proteomic analysis. *Nature Biotechnology*, 28(1):83–89, 2010.
- [10] Marc Sturm, Andreas Bertsch, Clemens Gröpl, et al. Openms – an open-source software framework for mass spectrometry. *BMC Bioinformatics*, 9(1):163, Mar 2008.
- [11] Eric W Deutsch, Luis Mendoza, David Shteynberg, et al. A Guided Tour of the Trans-Proteomic Pipeline. *Proteomics*, 10(6):1150–1159, 2011.
- [12] Jürgen Cox, Marco Y. Hein, Christian A. Luber, et al. Accurate Proteome-wide Label-free Quantification by Delayed Normalization and Maximal Peptide Ratio Extraction, Termed MaxLFQ. *Molecular & Cellular Proteomics*, 13(9):2513–2526, 2014.
- [13] Ingvar Eidhammer, Kristian Flikka, Lennart Martens, Svein-Ole Mikalsen. *Computational Methods for Mass Spectrometry Proteomics*. Wiley, 2007.
- [14] Michael R. Hoopmann, Jason M. Winget, Luis Mendoza, and Robert L. Moritz. StPeter: Seamless Label-Free Quantification with the Trans-Proteomic Pipeline. *Journal of Proteome Research*, 17(3):1314–1320, 2018.