

# Integrating MS-based proteomics and RNA-Seq data

March 6, 2016

## Use case 1: mapping peptides

**Mapping** of *peptides along protein sequences* (although not explicitly considered a mapping exercise) and *short reads along genome coordinates*.

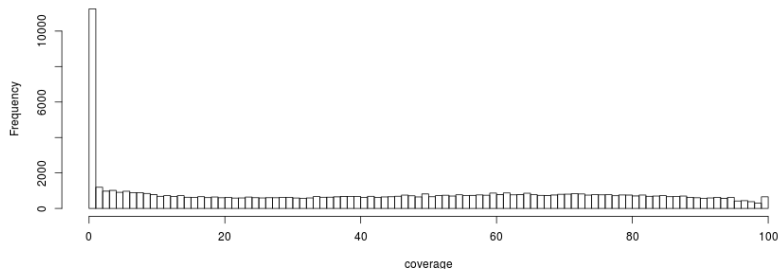
But...

- ▶ coverage
- ▶ protein inference
- ▶ identifier mapping
- ▶ missing values

# Coverage

- ▶ Coverage in proteomics in %
- ▶ Coverage in RNA-Seq in fold X

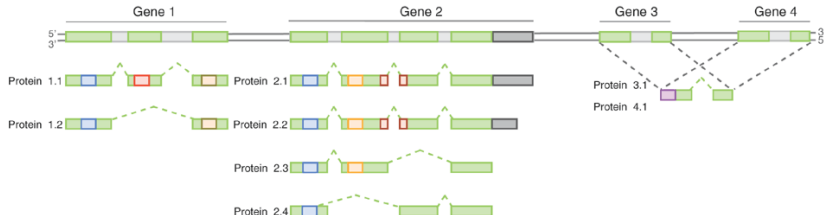
The following values are higher bounds, *without* peptide filtering for about 80000 *gene groups*



# And

- ▶ the majority of peptides map to a minority of proteins
- ▶ different peptides within one protein can be differently detectable in an MS

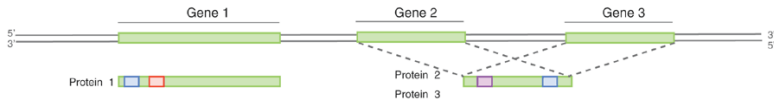
# Protein inference



Eukaryotes

Prokaryotes

Class	Protein sequence(s)	Protein isoform(s)	Gene(s)
1a	Unambiguous	Unambiguous	Unambiguous
1b	Unambiguous	Ambiguous	Unambiguous
2a	Ambiguous	Ambiguous	Unambiguous
2b	Ambiguous	Ambiguous	Unambiguous
3a	Unambiguous	Ambiguous	Ambiguous
3b	Ambiguous	Ambiguous	Ambiguous



From Qeli and Ahrens (2010). See also Nesvizhskii and Aebersold (2005).

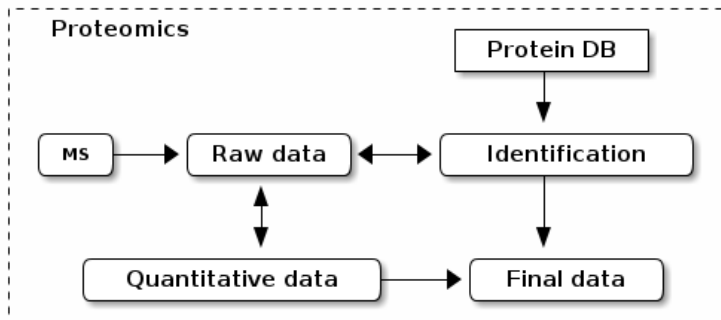
# Protein groups

Often, in proteomics experiments, the features represent single proteins and **groups** of indistinguishable or non-differentiable proteins identified by shared (non-unique) peptides.

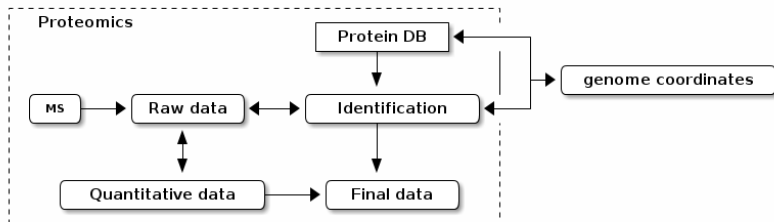
**Caveat:** Mapping between protein groups and unique transcripts?

# Mapping protein and gene identifiers

The protein database and the genome are *independent*, i.e. the proteins do not make explicitly reference to the genome they originate from.



# Mapping protein and gene identifiers



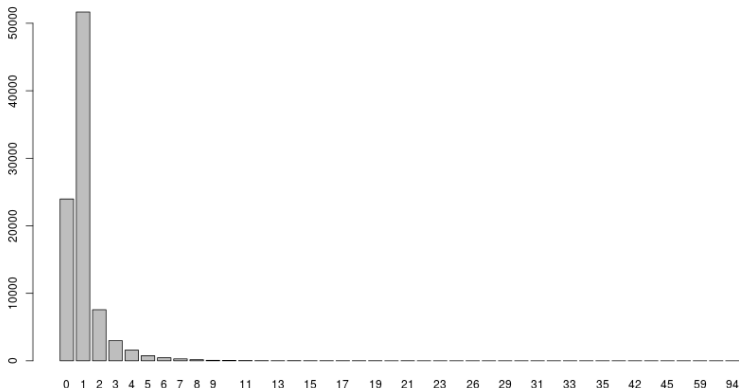
If we want to map UniProt accession numbers to genomic identifiers (Ensembl transcript identifiers):



# Mapping identifiers

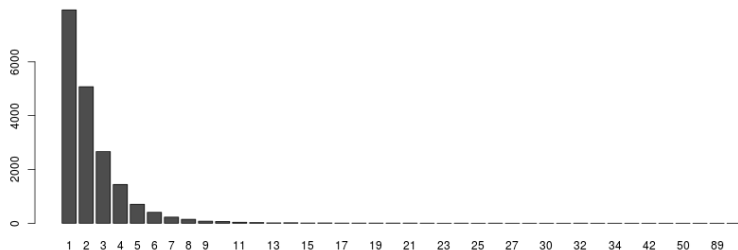
The UniProt human proteome (release 2015\_02) has 89796 entries.  
Using UniProt.ws:

- ▶ 23972 have no transcript identifier
- ▶ 51673 have a unique transcript identifier
- ▶ 14151 have more than one transcript identifier



## Using biomaRt:

Mapping 18911 identifiers, of which



**Caveat:** Mapping between single protein and unique transcripts?

# Missing values

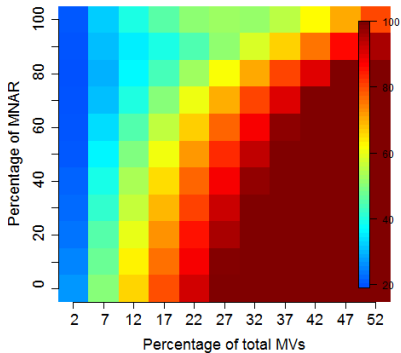
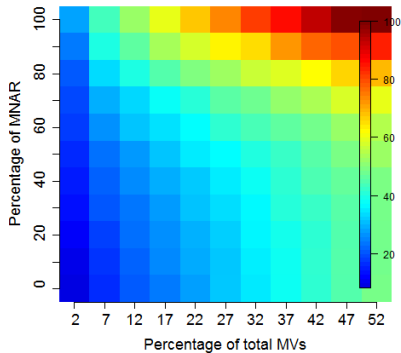
Options are:

- ▶ Filtering: Remove missing values, or at least features or samples with excessive number of missing values:
- ▶ Data imputation: inferring plausible values for missing data.

# Data imputation

There are two types of mechanisms resulting in missing values in LC/MSMS experiments.

- ▶ Missing values resulting from absence of detection of a feature, despite ions being present at detectable concentrations. For example in the case of ion suppression or as a result from the stochastic, data-dependent nature of the MS acquisition method. These missing value are expected to be randomly distributed in the data and are defined as **missing at random** (MAR) or **missing completely at random** (MCAR).
- ▶ Biologically relevant missing values, resulting from the *absence* of the low abundance of ions (below the limit of detection of the instrument). These missing values are not expected to be randomly distributed in the data and are defined as **missing not at random** (MNAR).

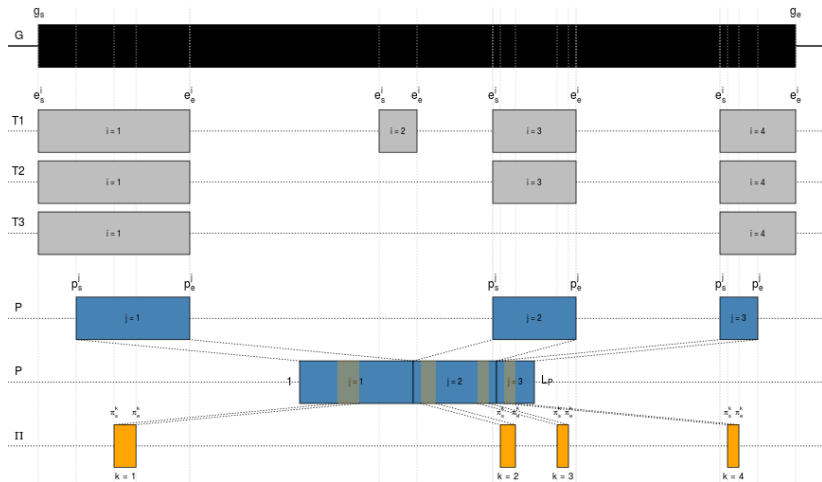


MNAR features should ideally be imputed with a **left-censor** (minimum value, zero, ...) method. Conversely, it is recommended to use **hot deck** methods (nearest neighbour, maximum likelihood, ...) when data are missing at random.



# Mapping peptides to genomic coordinates

The **goal** is to map peptides from protein coordinates (1 to  $L_p$ ) to genomic coordinates.



## Data

Illustration with the Pbase Bioconductor package (devel version).

We have an example data (named `p`) composed of 9 proteins, with UniProt accession numbers and Ensembl transcript identifiers and each protein has a set experimentally observed peptides (see table below). This `p` object is generated from the protein database (fasta file) and the MS identification results (`mzIdentML` file) against this very same protein database.

Acc	ENST	npep
A4UGR9	ENST00000409195	36
A6H8Y1	ENST00000358731	23
O43707	ENST00000252699	6
O75369	ENST00000295956	13
P00558	ENST00000373316	5
P02545	ENST00000368300	12



# Multiple transcripts per protein

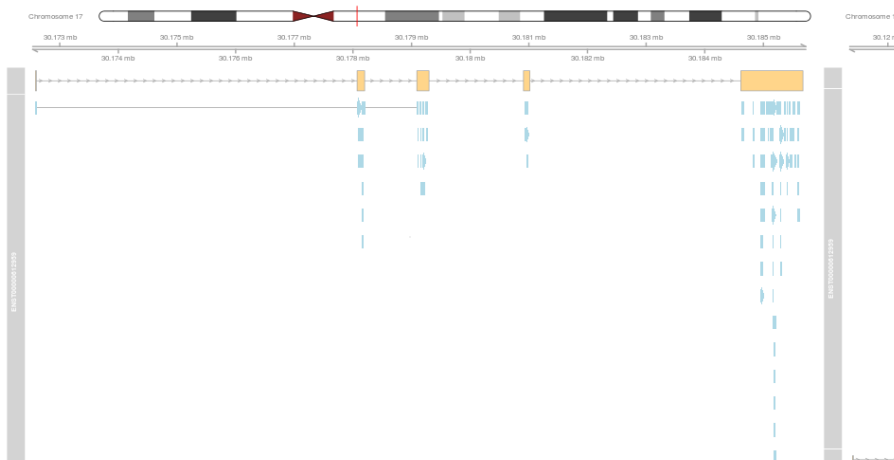
## Data

We use our example data, named `p`, composed of 9 proteins, with UniProt accession numbers and Ensembl transcript identifiers and each protein has a set experimentally observed peptides:

Acc	ENST	npep
A4UGR9	ENST00000409195	36
A6H8Y1	ENST00000358731	23
O43707	ENST00000252699	6
O75369	ENST00000295956	13
P00558	ENST00000373316	5
P02545	ENST00000368300	12
P04075	ENST00000338110	21
P04075-2	ENST00000395248	20

# Mapping MS peptides and RNA-Seq short reads

The last step of the mapping process is to combine the newly mapped peptides and reads from RNA-Seq experiments. The figures below illustrate this with data from Sheynkman et al. (2013, 2014) from the Jurkat cell line (TIB-152). The mass spectrometry and RNA-Seq (SRR791580) were processed with standard pipelines.



## References

Laurent Gatto and Sebastian Gibb (2016). Pbase: Manipulating and exploring protein and proteomics data. R package version 0.11.3. <https://github.com/ComputationalProteomicsUnit/Pbase>

Bernd Fischer, Steffen Neumann, Laurent Gatto and Qiang Kou. mzR: parser for netCDF, mzXML, mzData and mzML and mzIdentML files (mass spectrometry data). R package version 2.5.3.

Lazar C, Gatto L, Ferro M, Bruley C, and Burger T. Accounting for the Multiple Natures of Missing Values in Label-Free Quantitative Proteomics Data Sets to Compare Imputation Strategies.

Publication Date: February 23, 2016

DOI: { []10.1021/acs.jproteome.5b00981[] }(http://pubs.acs.org/doi/abs/10.1021/acs.jproteome.5b00981).

Pang et al. Tools to covisualize and coanalyze proteomic data with genomes and transcriptomes: validation of genes and alternative mRNA splicing. J Proteome Res. 2014 Jan 3;13(1):84-98. doi: 10.1021/pr400820p. Epub 2013 Nov 12. PubMed PMID: 24152167.

Shevinkman GM, Johnson JE, Jagtap PD, Shortreed MR, Onsongo

## Session information

```
sessionInfo()
```

```
## R Under development (unstable) (2016-03-03 r70270)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 14.04.3 LTS
##
## locale:
##   [1] LC_CTYPE=en_GB.UTF-8      LC_NUMERIC=C
##   [3] LC_TIME=en_GB.UTF-8      LC_COLLATE=en_GB.UTF-8
##   [5] LC_MONETARY=en_GB.UTF-8  LC_MESSAGES=en_GB.UTF-8
##   [7] LC_PAPER=en_GB.UTF-8     LC_NAME=C
##   [9] LC_ADDRESS=C             LC_TELEPHONE=C
##  [11] LC_MEASUREMENT=en_GB.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
##   [1] grid      stats4    parallel  stats      graphics g
##   [8] datasets  methods  base
##
```