# Algorithms for automatic spectra interpretation

Ralf Gabriels

VIB-UGENT CENTER FOR MEDICAL BIOTECHNOLOGY

UNIVERSITEIT GENT

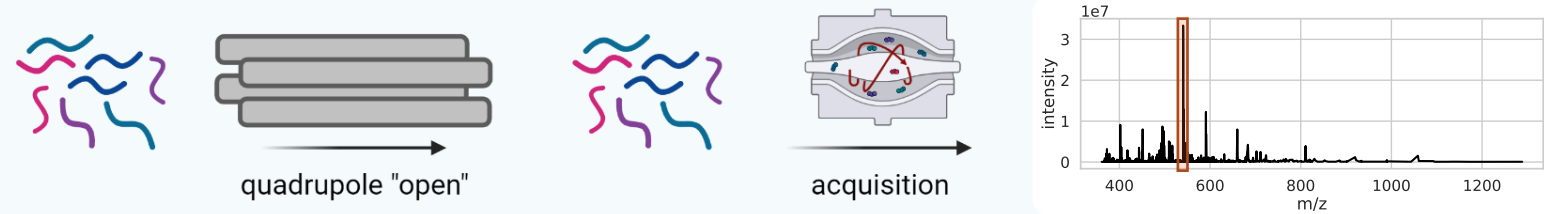Comp omics

# Summary

- LC-MS/MS recap
- Database search engines: Replicating the LC-MS/MS workflow *in silico*
- Specialized methods
  - Open-modification searching
  - *De novo* and sequence tag-based searching
  - Spectral library searching
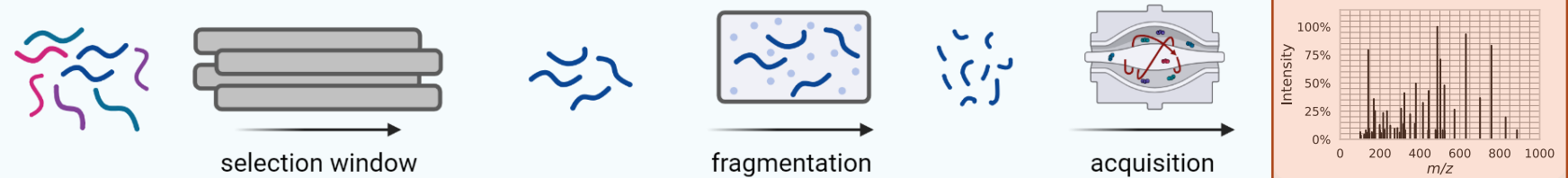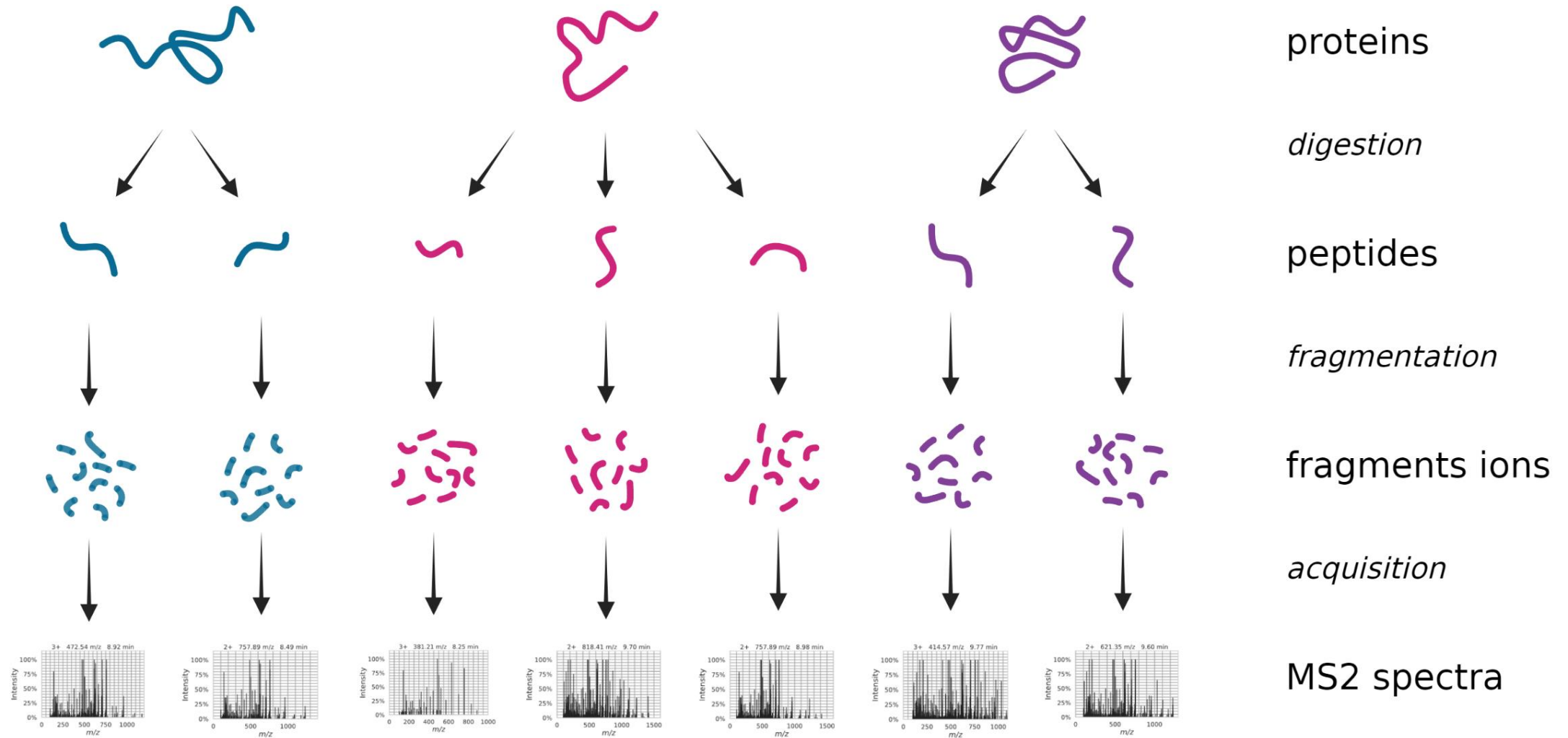
# LC-MS/MS recap

# LC-MS/MS recap

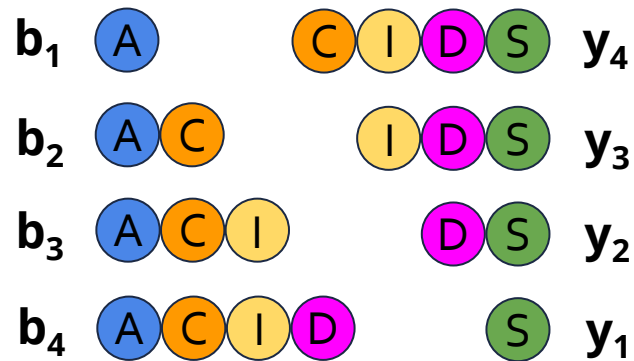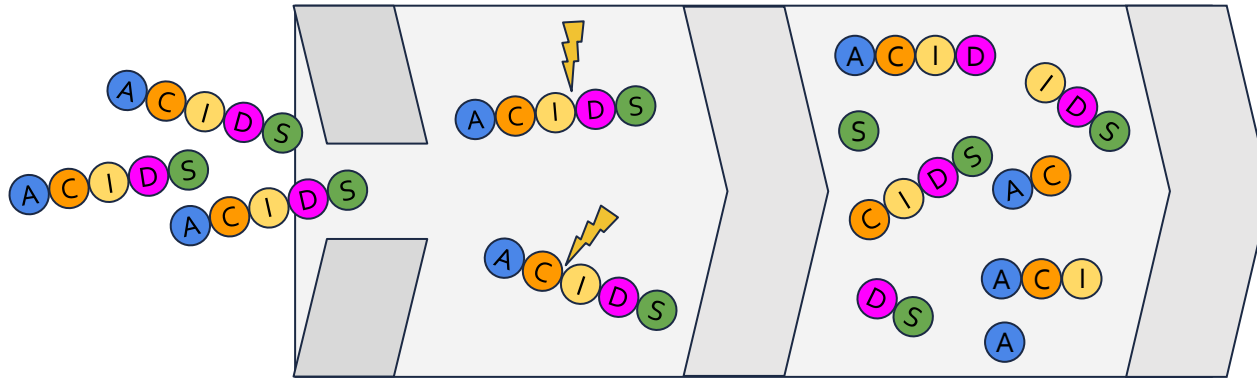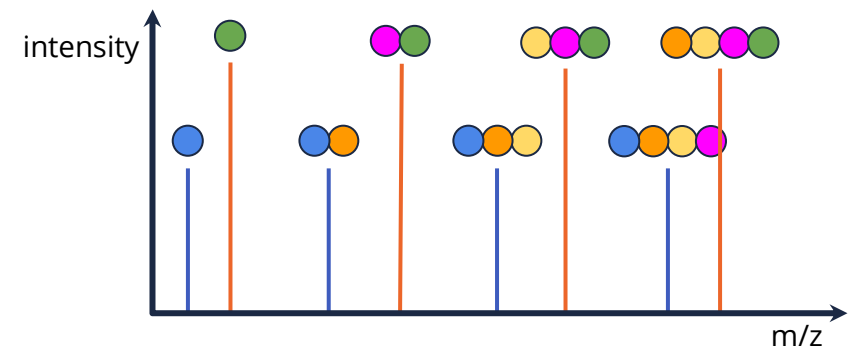# How do we link MS2 spectra back to the original proteins?

# Peptides can be identified by their fragmentation spectra

# Peptides can be identified by their fragmentation spectra

# *De novo* peptide spectrum identification is not straightforward



theory

reality

# Database search engines replicate the LC-MS steps *in silico*



Proteins → **Digestion** → Peptides → **Fragmentation & detection** → MS2 spectra → **Identification** → H$_2$N-LPVWSQRNFK-COOH

Protein sequence DB → *In silico* digestion → Peptide search space → Fragment *m/z* calculation → Theoretical MS2 spectra

>LPVGTSSGDSPK
PDYKGATTVRDLT
PENEWYVKEITGS
TKTKVIKLNWQRS
TVK

>LPVGTSSGDSPK
>PDYKGATTVR
>DLTPENEWYVK
>EITGSTKTK
>VIKLNWQRSTVK

# General proteomics search workflow

1. Define search space
2. *In silico* digestion of proteins to peptides
3. For each MS2 spectrum
   1. Select candidate peptide-spectrum matches (PSMs)
   2. Calculate theoretical peptide spectra
   3. Match candidate theoretical spectra to observed spectrum
   4. Select best match
4. PSM score post-processing

# The search space should include ALL proteins that are expected to be present in the sample



UniProt + gpm

Reference proteome

Common contaminants

common **R**epository
of **A**dventitious **P**roteins

⭐ Swiss-Prot

📄 Swiss-Prot + TrEMBL

**Search parameters**
- Database selection

# A mismatched search space can lead to false identifications



SGVSRKPAPG/2

ATASPPRQK/2

# *In silico* digestion replicates the cleavage rules of the enzyme

```
>protein_sequence
LPVGTSSGDSPKPDYKGATTVRDLTPENEWYVKEITGSTKTKVI
KLNWQRSTVK
```

**0 missed cleavages**

LPVGTSSGDSPKPDYK
GATTVR
DLTPENEWYVK
EITGSTK
TK
VIK
LNWQR
STVK

**1 missed cleavage**

LPVGTSSGDSPKPDYKGATTVR
GATTVRDLTPENEWYVK
DLTPENEWYVKEITGSTK
EITGSTKTK
TKVIK
VIKLNWQR
LNWQRSTVK

**2 missed cleavages**

LPVGTSSGDSPKPDYKGATTVRDLTPENEWYVK
GATTVRDLTPENEWYVKEITGSTK
DLTPENEWYVKEITGSTKTK
EITGSTKTKVIK
TKVIKLNWQR
VIKLNWQRSTVK

**Search parameters**
- Allowed missed cleavages

# *In silico* digestion replicates the cleavage rules of the enzyme

```
>protein
LPVGTSSGDSPKPDYKGATTVRDLTPENEWYVKEITGSTKTKVI
KLNWQRSTVK
```

**0 missed cleavages**

LPVGTSSGDSPKPDYK
GATTVR
DLTPENEWYVK
EITGSTK
TK
VIK
LNWQR
STVK

**1 missed cleavage**

LPVGTSSGDSPKPDYKGATTVR
GATTVRDLTPENEWYVK
DLTPENEWYVKEITGSTK
EITGSTKTK
TKVIK
VIKLNWQR
LNWQRSTVK

**2 missed cleavages**

LPVGTSSGDSPKPDYKGATTVRDLTPENEWYVK
GATTVRDLTPENEWYVKEITGSTK
DLTPENEWYVKEITGSTKTK
EITGSTKTKVIK
TKVIKLNWQR
VIKLNWQRSTVK

> **Search parameters**
> - Allowed missed cleavages
> - Minimum peptide length / mass
> - Maximum peptide length / mass

# Very large peptide search space needs to be filtered to candidate peptide-to-spectrum matches



allowed missed cleavages

⭐ Swiss-Prot         📄 Swiss-Prot + TrEMBL         protein database

Computationally expensive
High probability of false positives          ⟶          Filter to a limited set of PSMs

# Candidate PSMs are usually selected by their precursor mass

| Peptide | Mass | +1 *m/z* | +2 *m/z* | +3 *m/z* |
|---|---|---|---|---|
| LPVGTSSGDSPKPDYK | 1646.81 | 1647.82 | 824.41 | 549.49 |
| GATTVR | 603.33 | 604.34 | 302.67 | 202.11 |
| DLTPENEWYVK | 1392.65 | 1393.66 | 697.33 | 465.22 |
| EITGSTK | 734.38 | 735.38 | 368.19 | 245.80 |
| ... | | | | |

**MS1**

**Search parameters**
- Precursor mass window

# For each candidate PSM, the theoretical spectrum is matched against the observed spectrum



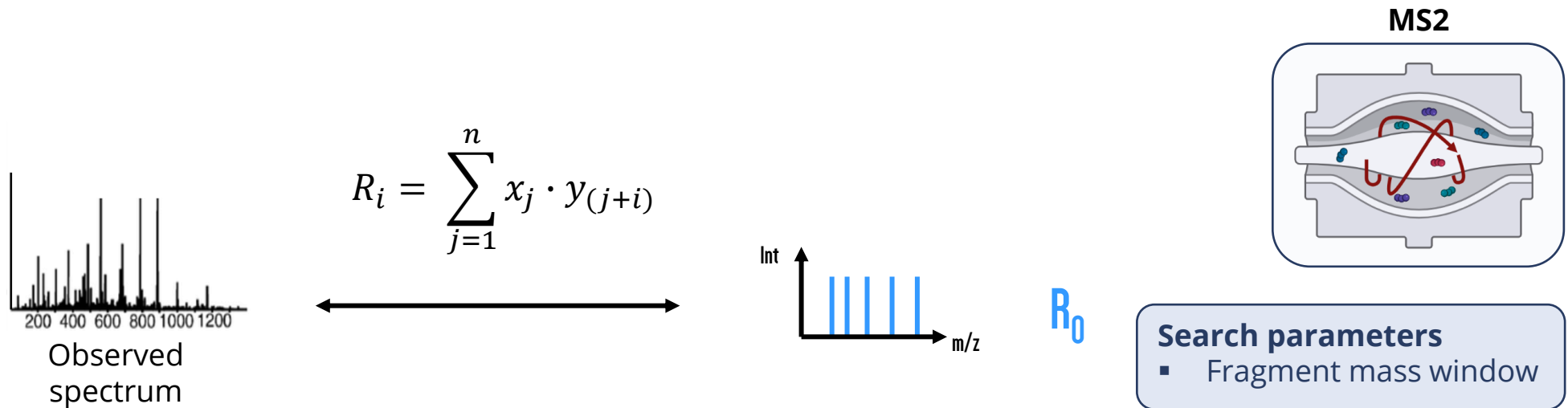Observed spectrum

Theoretical spectra for candidate PSMs

**Scoring functions**
- Explained intensity (e.g. SEQUEST)
- Peak counting (e.g. Mascot, Andromeda)
- Combinations (e.g. X!Tandem)

PSM scores for each candidate PSM

Select best-scoring PSM

# The SEQUEST scoring function calculates the total intensity that can be explained by the theoretical spectrum

**MS2**

$$R_i = \sum_{j=1}^{n} x_j \cdot y_{(j+i)}$$

Int

m/z

$R_0$

Observed spectrum

200 400 600 800 1000 1200

**Search parameters**
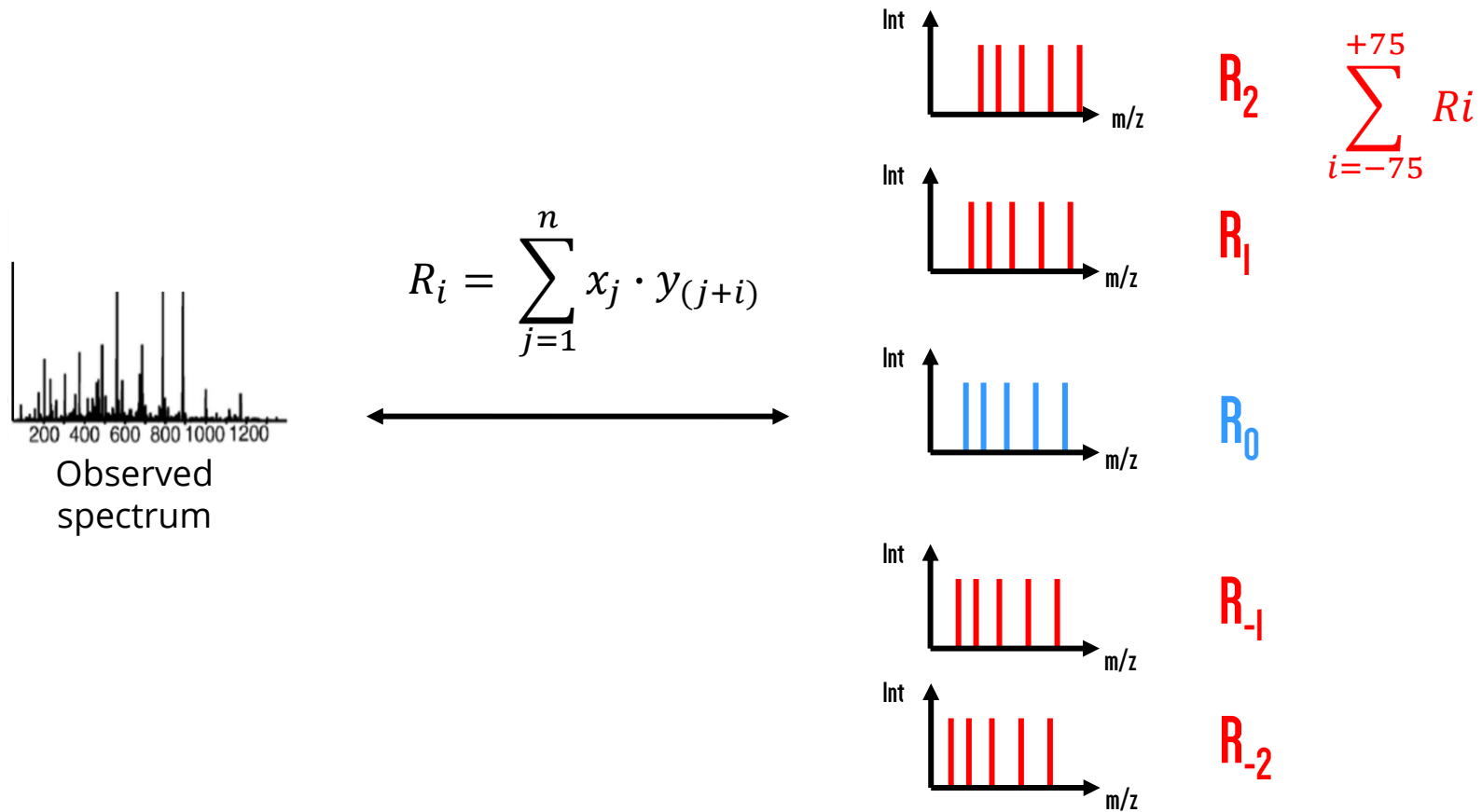- Fragment mass window

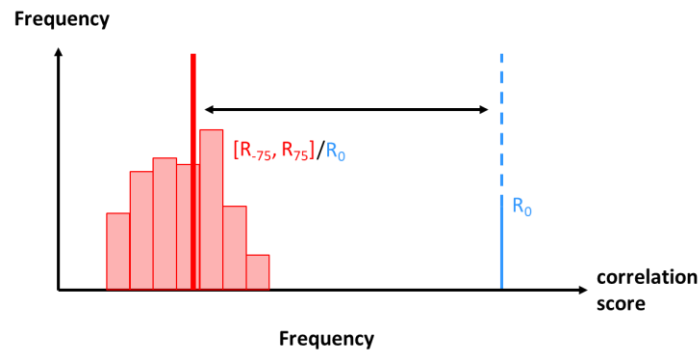Resolution and accuracy: Fragment peak error tolerance

ion trap    0.5 Da

orbitrap   0.02 Da

# The SEQUEST scoring function calculates the total intensity that can be explained by the theoretical spectrum

$$R_i = \sum_{j=1}^{n} x_j \cdot y_{(j+i)}$$

Observed spectrum

$R_2$

$R_1$

$R_0$

$R_{-1}$
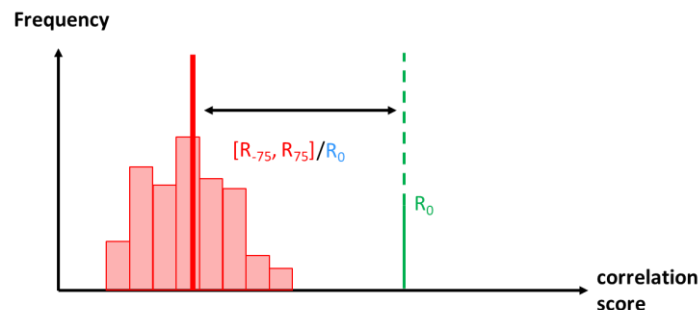
$R_{-2}$

$$\sum_{i=-75}^{+75} Ri$$

# This score is then calibrated to the distribution of random matches, then the difference between the best and second-best score is calculated

XCorr$_1$

XCorr$_2$



deltaCn = $\dfrac{XCorr_1 - XCorr_2}{XCorr_1}$

Resulting PSM score for best PSM

# The raw output of a search engine is a list of scores for the best-scoring PSM for each spectrum
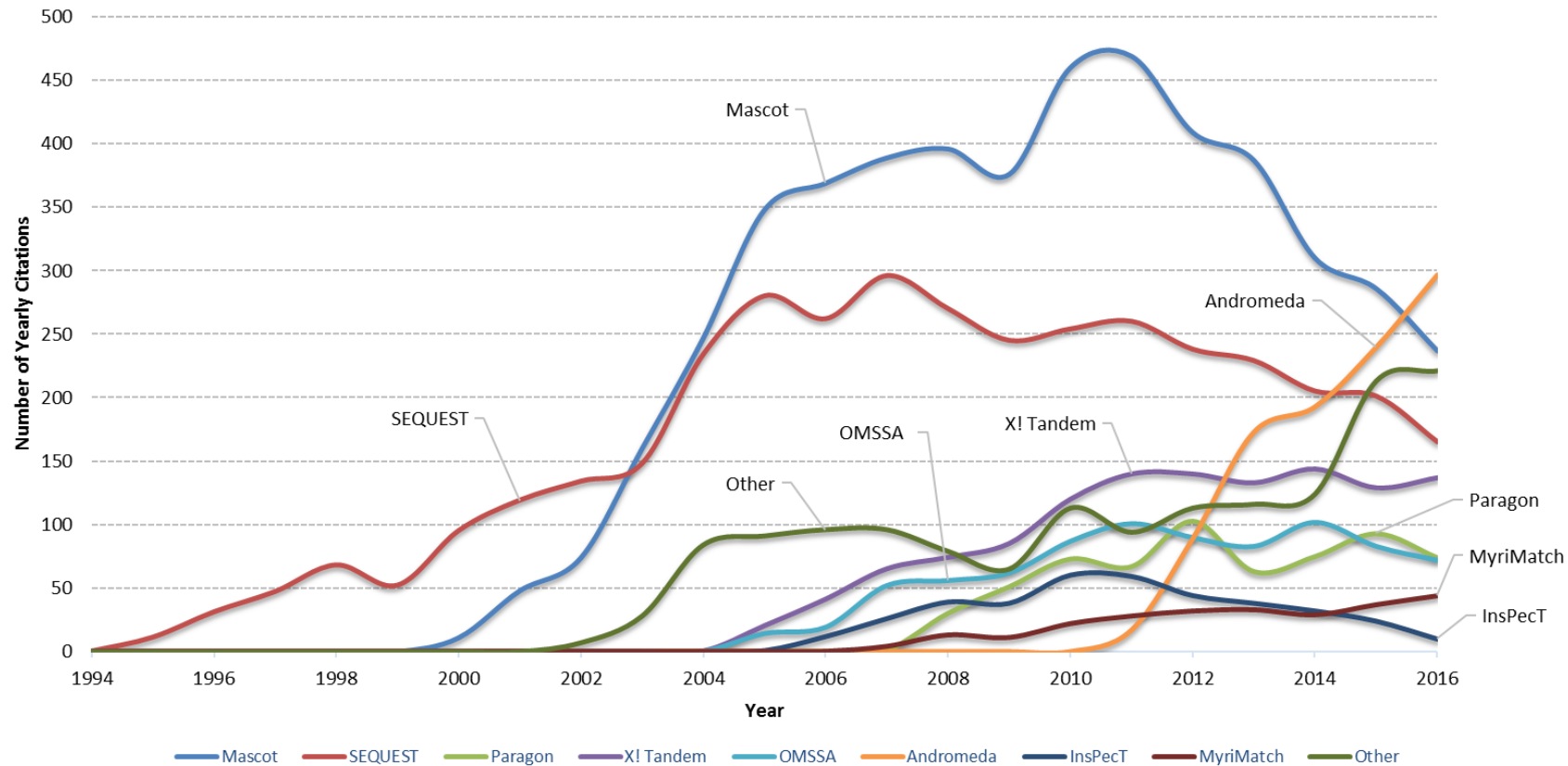
| Spectrum | Best PSM | Score |
|---|---|---|
| 1 | LPVGTSSGDSPKPDYK | 34.12 |
| 2 | GATTVR | 95.43 |
| 3 | DLTPENEWYVK | 134.87 |
| 4 | EITGSTK | 12.9 |
| ... | ... | 245.67 |

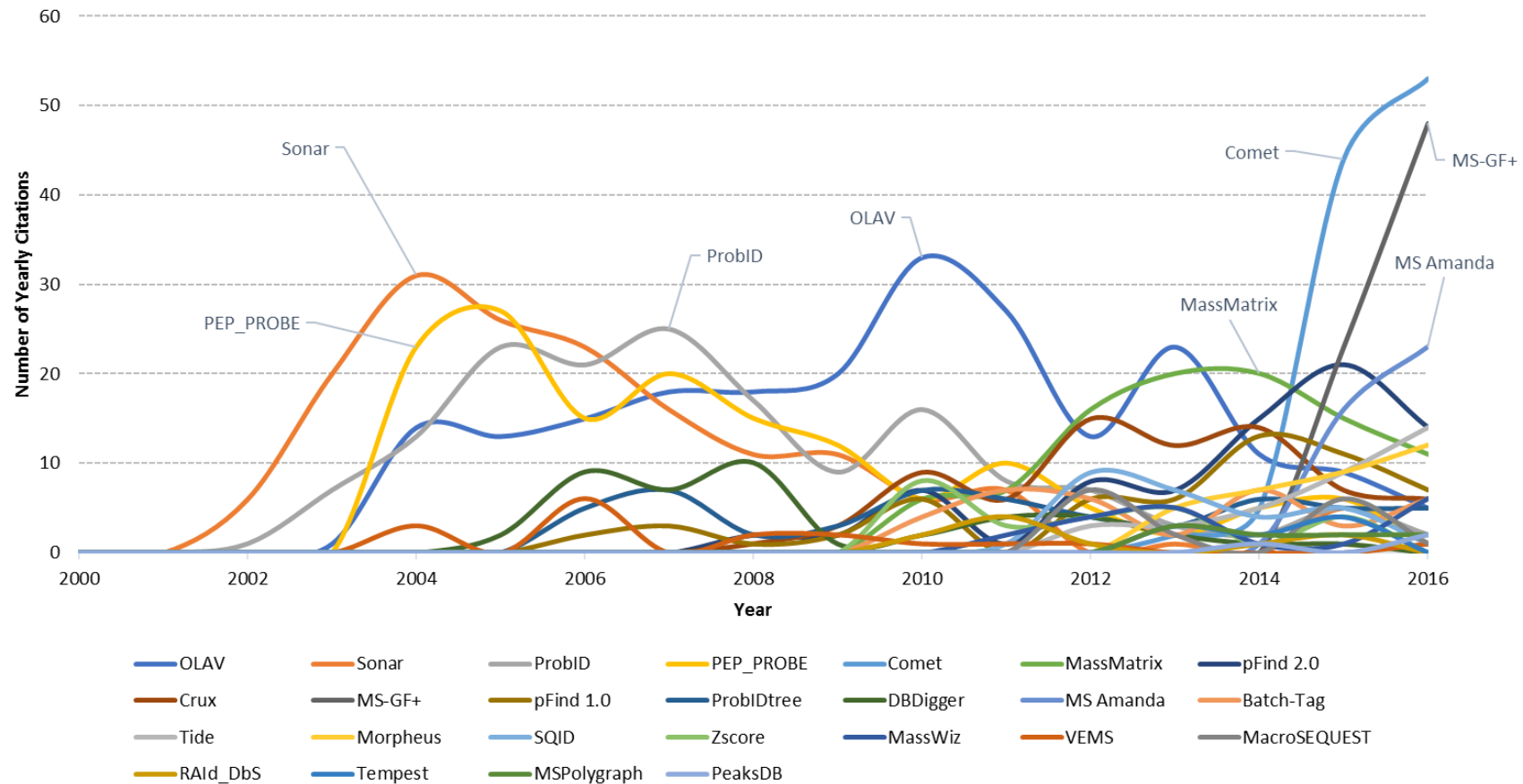Now what? Which scores are good? Which are bad?

# General proteomics search workflow

1. Define search space

2. *In silico* digestion of proteins to peptides

3. For each spectrum
   1. Select candidate peptide-spectrum matches (PSMs)
   2. Calculate theoretical peptide spectra
   3. Match candidate theoretical spectra to observed spectrum
   4. Select best match

4. PSM score post-processing: Next lecture

# Historically, Mascot and SEQUEST were heavily used; now Andromeda, Comet, MS-GF+, and MS-Amanda are more favored.
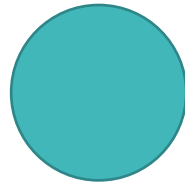
# Historically, Mascot and SEQUEST were heavily used; now Andromeda, Comet, MS-GF+, and MS-Amanda are more favored.

# Summary

- LC-MS/MS recap

- Database search engines: Replicating the LC-MS/MS workflow *in silico*

- Specialized methods

  ▸ Open-modification searching

  ▸ *De novo* and sequence tag-based searching
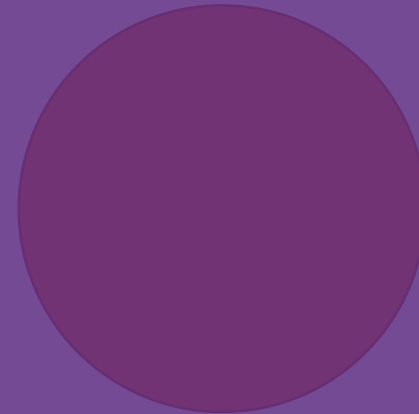
  ▸ Spectral library searching

# A few modifications can be added to the search space. To identify more PTMS, open-modification search engines are required.

**No variable modifications**

**Common modifications**

**All known modifications**

*MSFragger*
*open-pFind*
*TagGraph*
*ionbot*
*...*

*Combinatorial explosion*

# Sequence tag-based and *de novo* searching

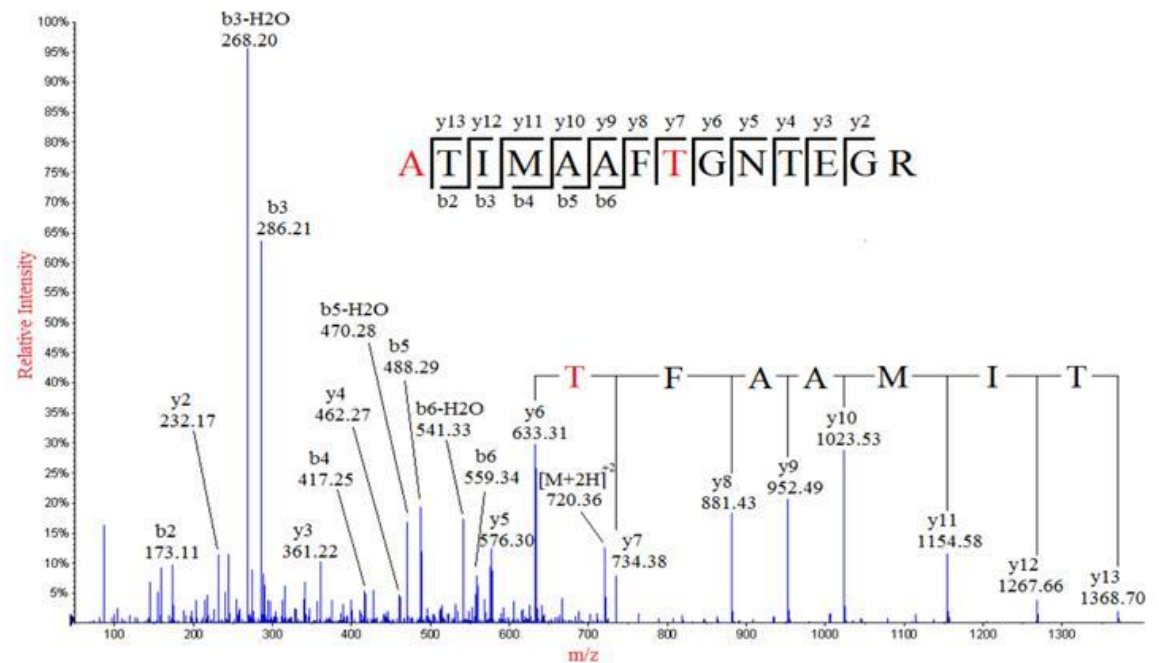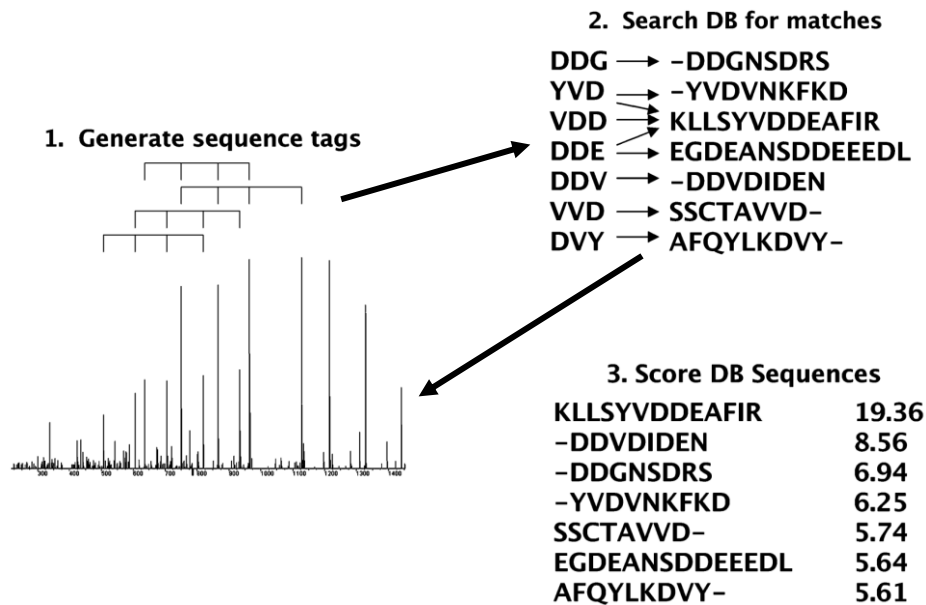**Tag-based**
*DirecTag*
*TagGraph*
*ionbot*
*...*

Find partial sequences in spectra
Usually to limit search space
Open modification searching!
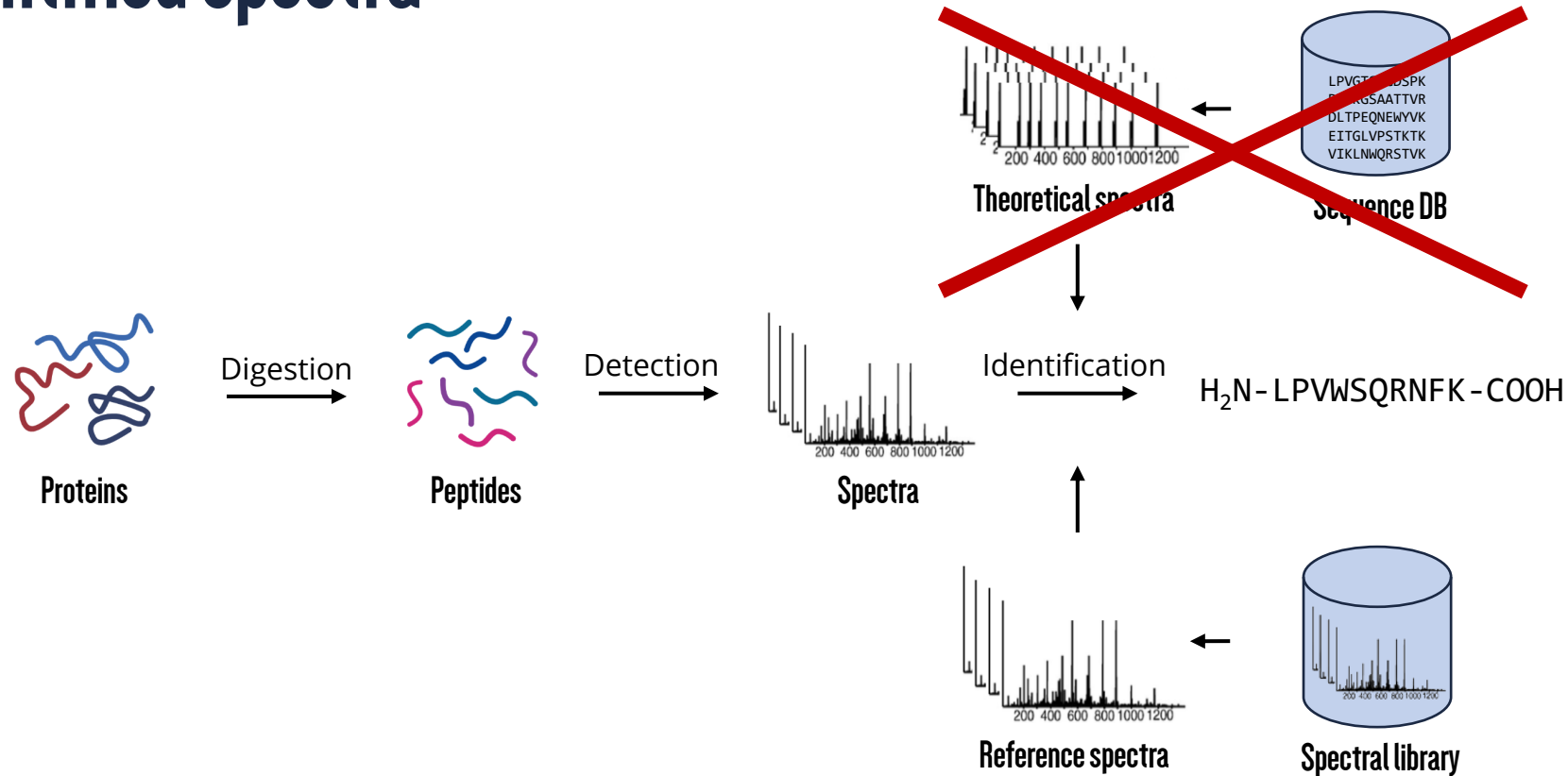
**De novo**
*Novor*
*PepNovo*
*PEAKS*
*...*

No database at all
Try to elude full sequence
Rare / unknown species

# Spectral library search engines match spectra against previously identified spectra



Theoretical spectra    Sequence DB
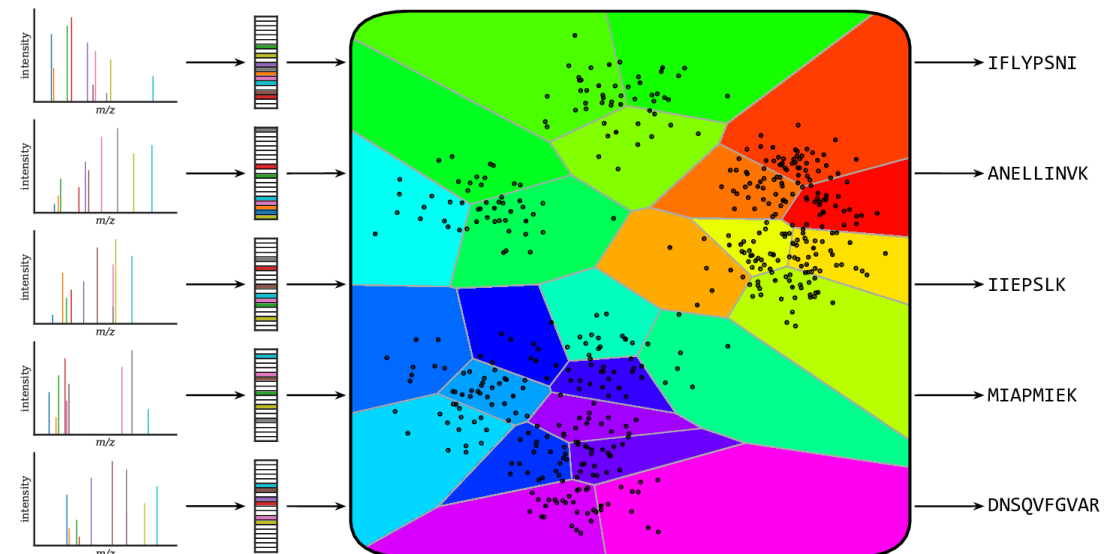
Proteins → Digestion → Peptides → Detection → Spectra → Identification → H₂N-LPVWSQRNFK-COOH

Reference spectra    Spectral library

**+ more sensitive scoring**

**- only identifies what has been seen before**

# Spectral library search engines match spectra against previously identified spectra

## ANN-SoLo



**SpectraST**
**Trans-proteomic pipeline (TPP)**

**COSS**
**CompOmics Spectral Library Searching**

Feature hashing          GPU-powered ANN searching

# Summary

- LC-MS/MS recap
- Database search engines: Replicating the LC-MS/MS workflow *in silico*
- Specialized methods
  - Open-modification searching
  - *De novo* and sequence tag-based searching
  - Spectral library searching

Ralf Gabriels

Ralf.Gabriels@UGent.be

@RalfGabriels | @CompOmics

www.compomics.com