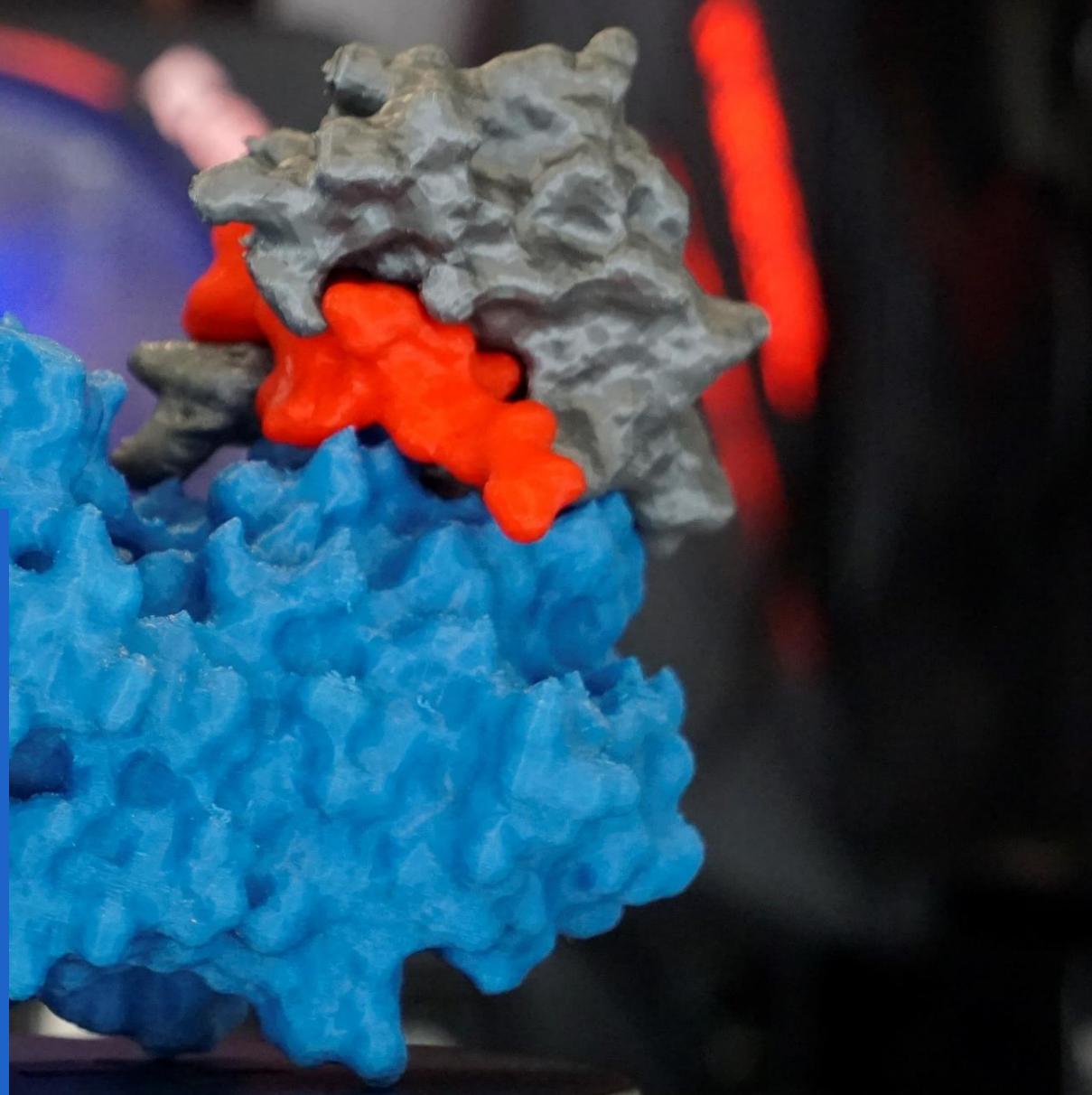


MACHINE LEARNING TO THE RESCUE: ENABLING NOVEL PROTEOMICS WORKFLOWS WITH DATA-DRIVEN BIOINFORMATICS METHODS

Ralf Gabriels

Dissertation submitted to obtain the degree Doctor in Health Sciences
Academic year 2021-2022



Introduction (Dutch)
proteins and proteomics
mass spectrometry
machine learning

Results (English)

MS²PIP: Peptide spectrum prediction for multiple fragmentation methods,
instruments, and labeling techniques

Removing the hidden data dependency of DIA with predicted spectral libraries

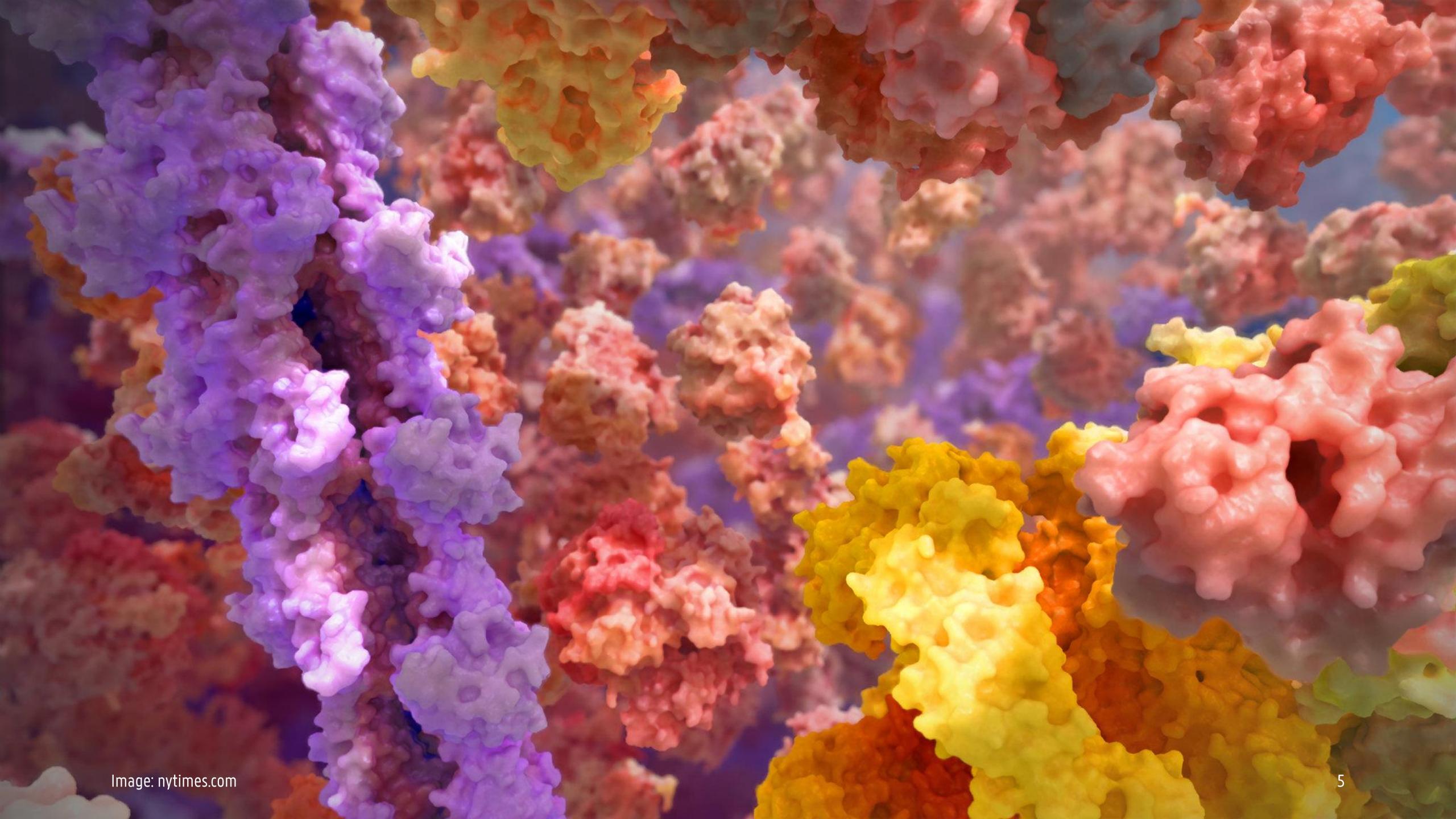
MS²Rescore: Leveraging spectrum predictions to enable novel proteomics workflows

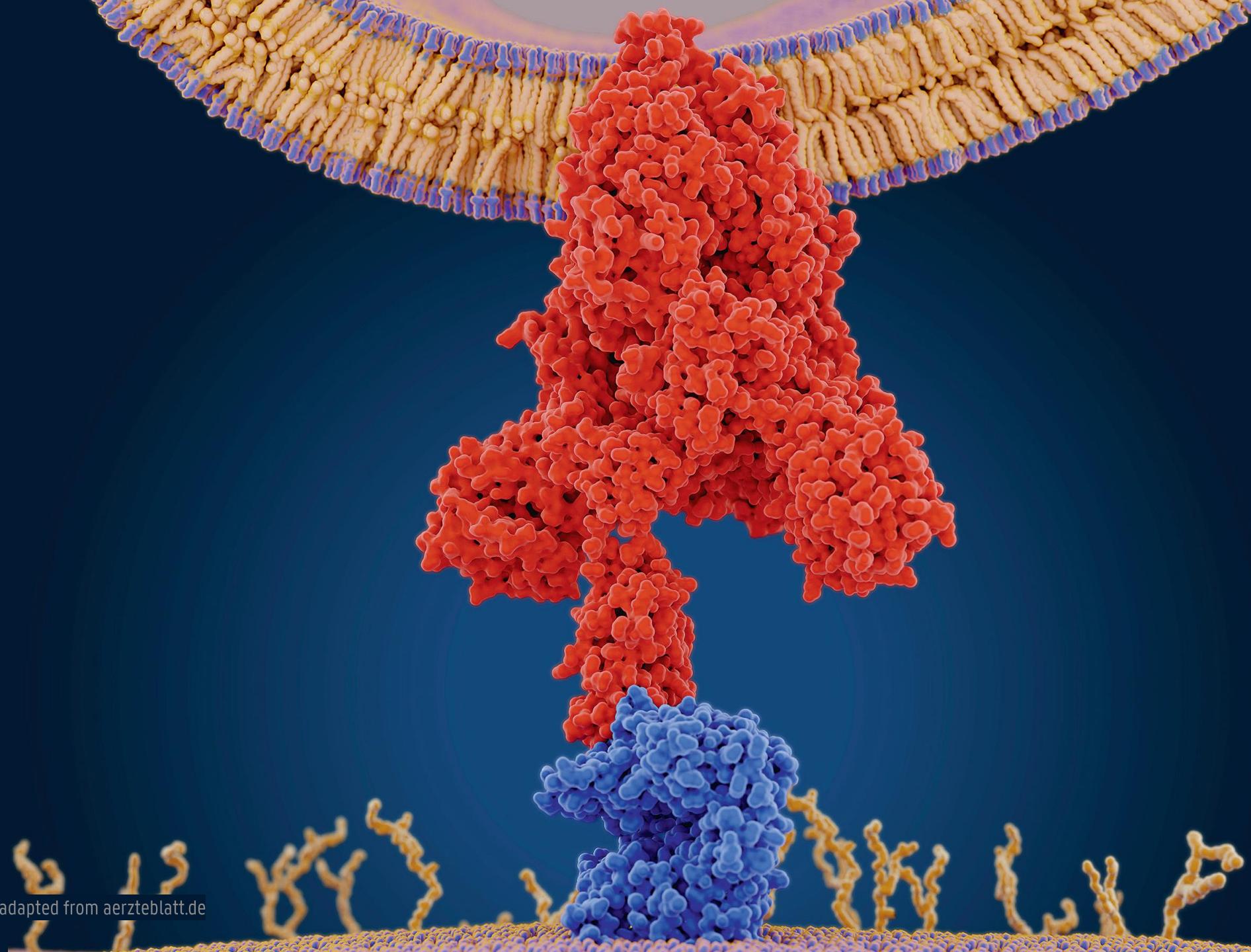
MS²DIP: spectrum prediction for modified peptides

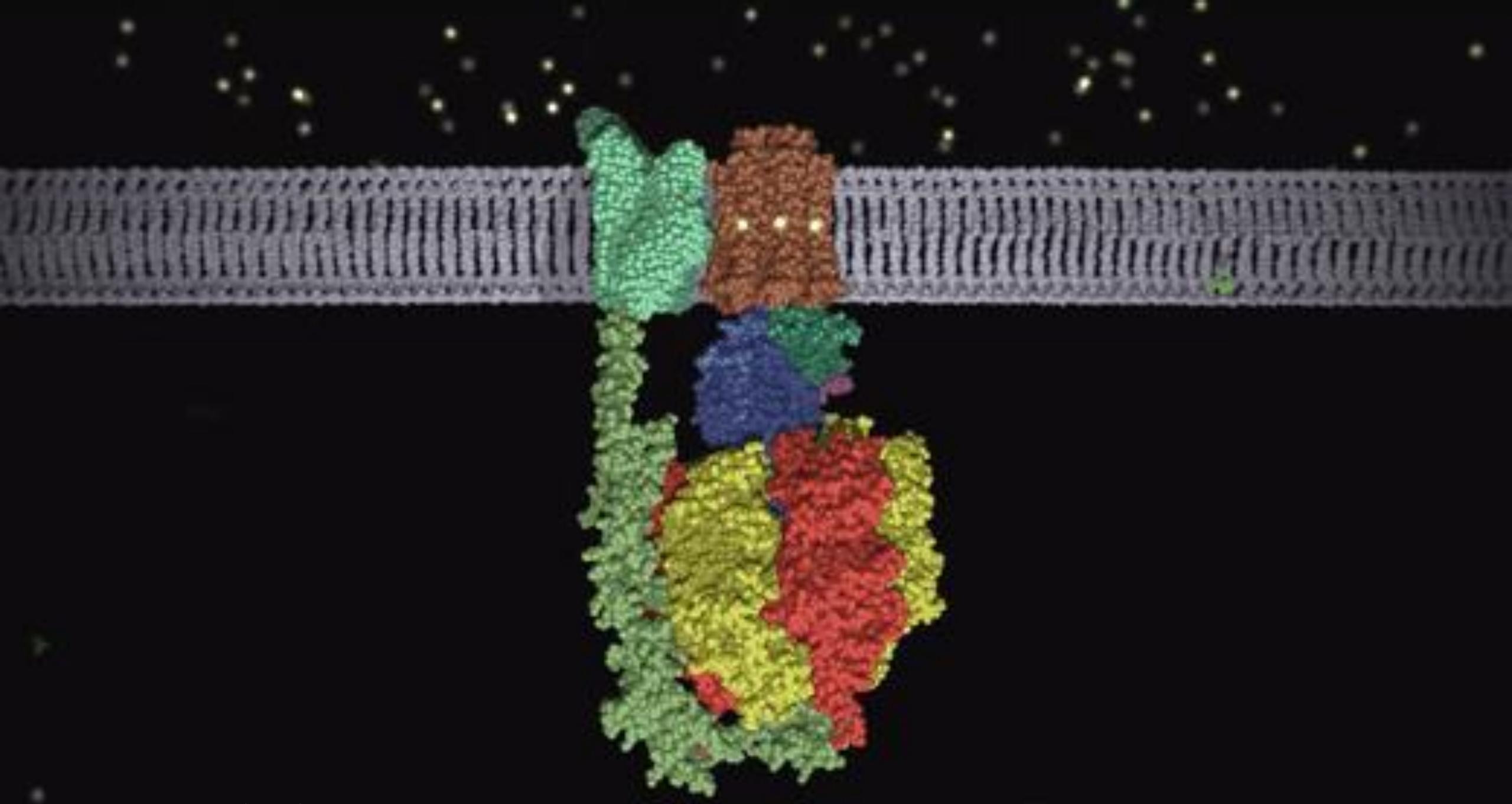
Discussion and future perspectives (English)

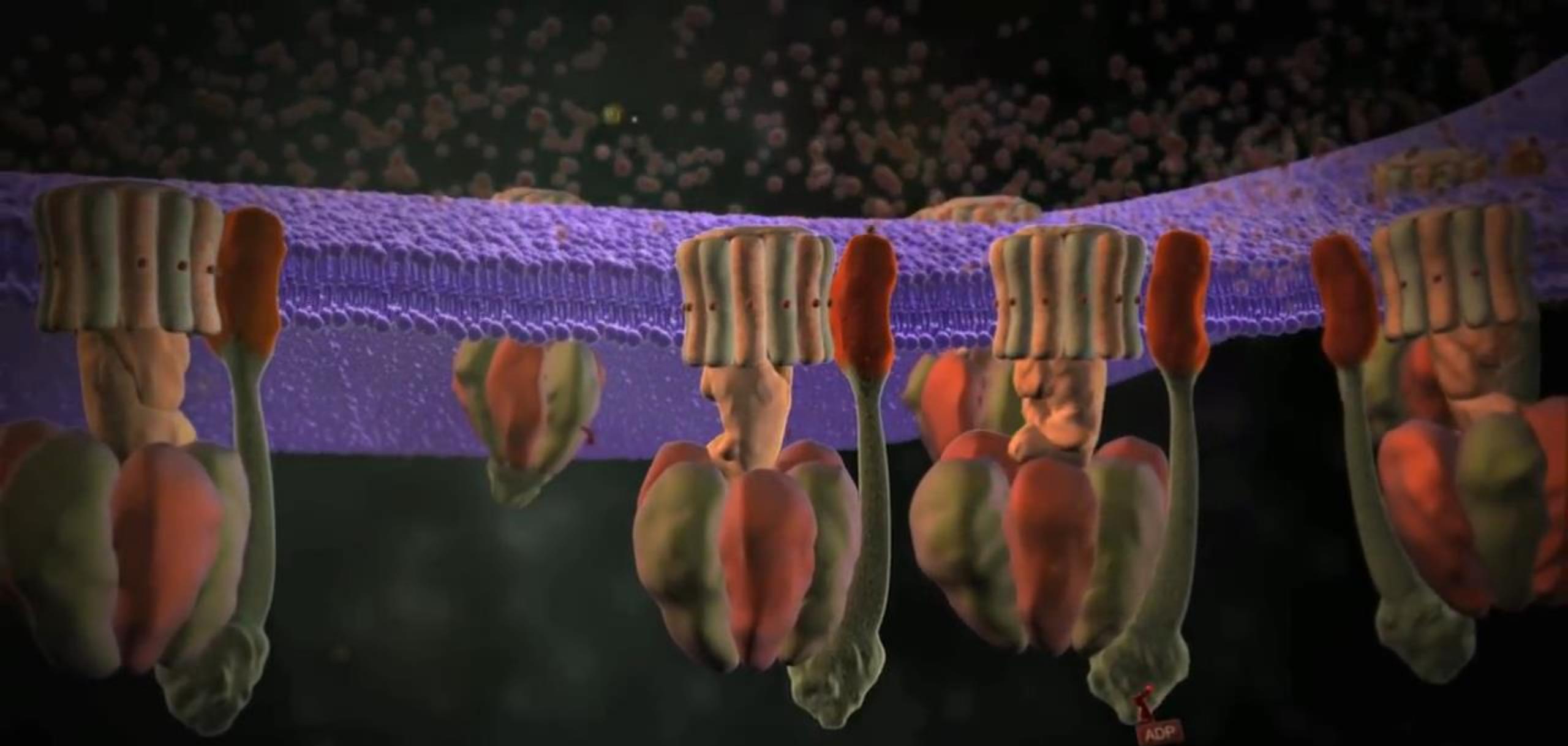
proteins and proteomics

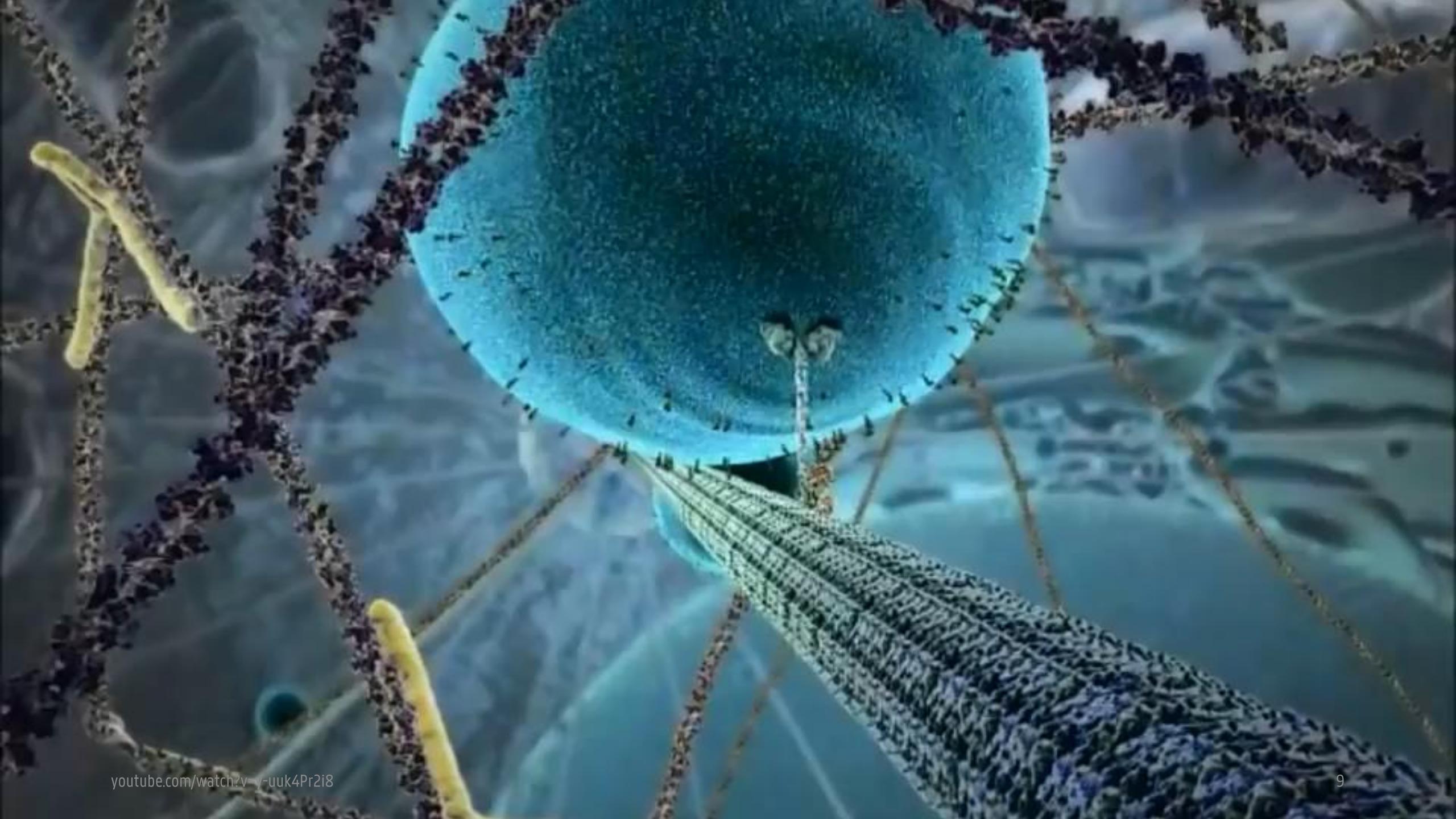










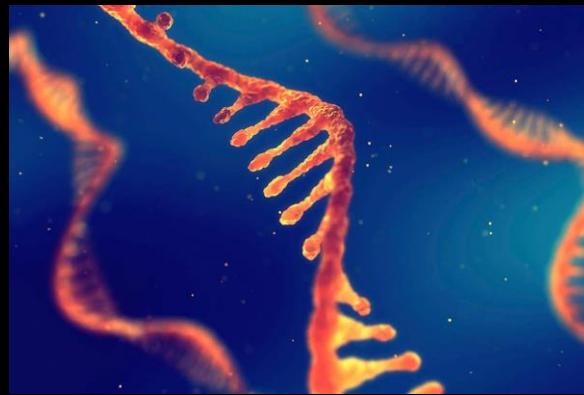


Central dogma of molecular biology



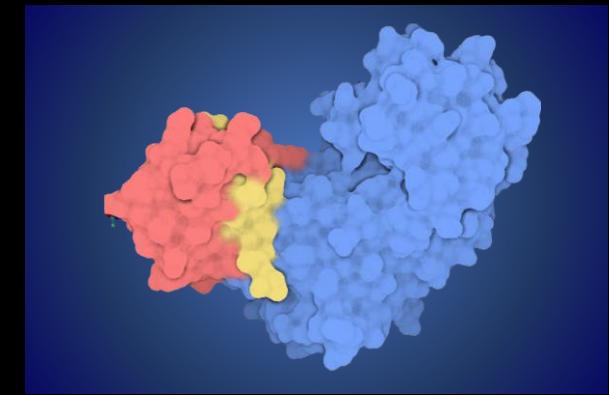
DNA / GENE

ATGGCGGCCGGCCGGAGATAG
TGGGGCCTACGATTGGG [...]



RNA / TRANSCRIPT

AUGGCGGCGGCCGGAGAUAG
UGGGGCCUACGAUUGGG [...]



PROTEIN

MGVPSGLILCVLIGAFFISM
AAAGDSGAYDWVMPARS [...]



source



recipe



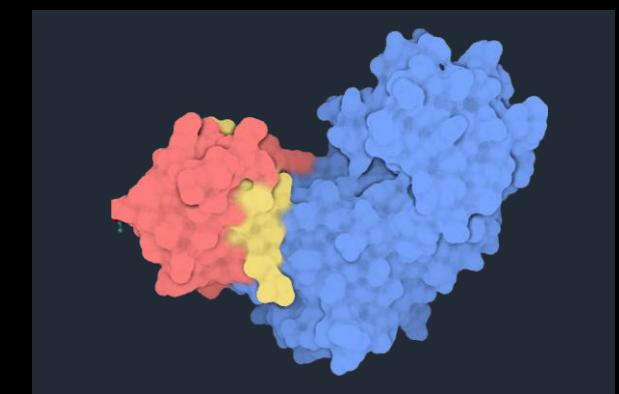
delicious waffles



DNA / GENE



RNA / TRANSCRIPT

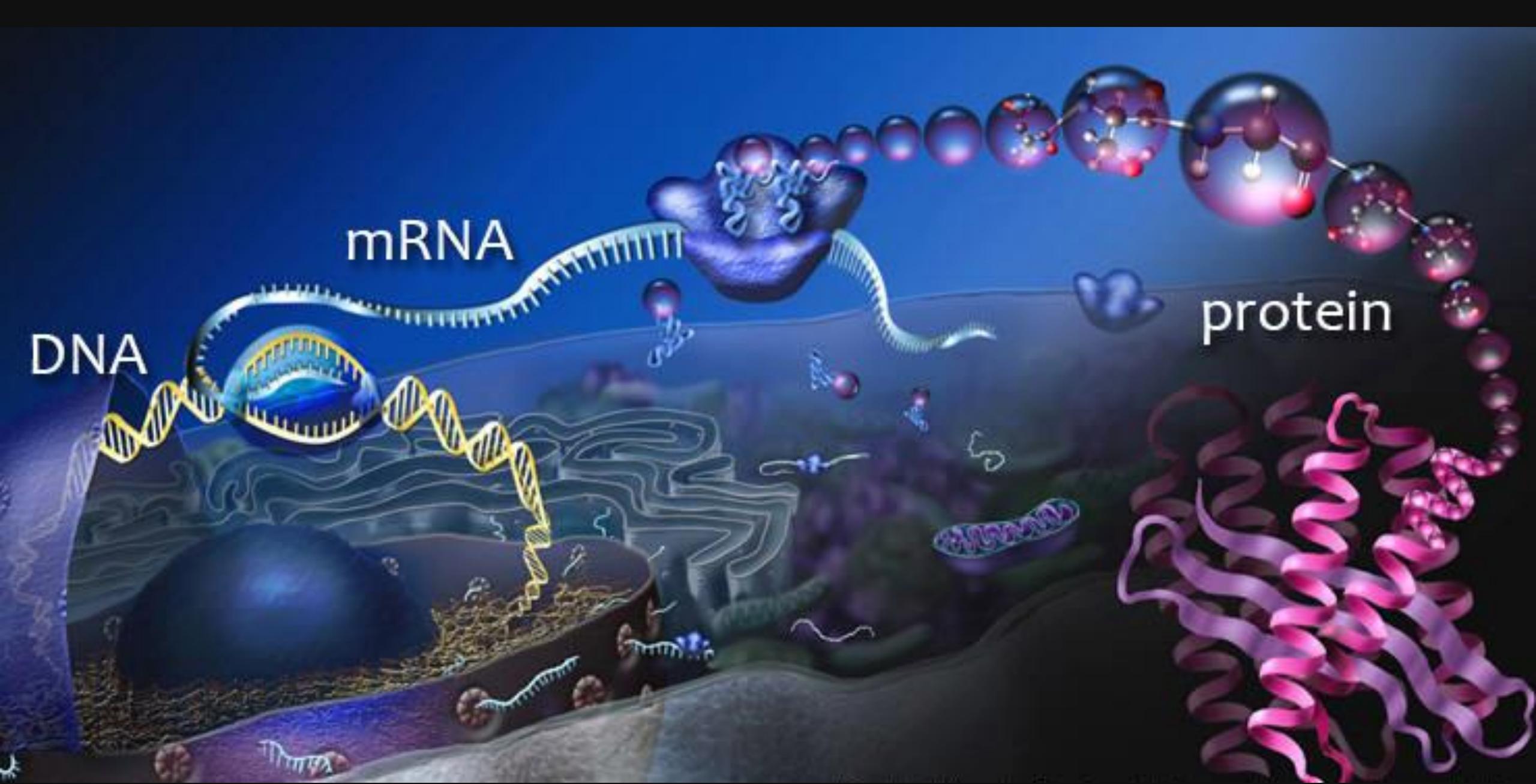


PROTEIN

TTGCCCATTAGCTGTTGTCT
TAGTTCAAGATTTGTT [...]

AUGGCAGGCGGCCGGAGAUAG
UGGGGCCUACGAUUGGG [...]

MGVPSGLILCVLIGAFFISM
AAAGDSGAYDWVMPARS [...]



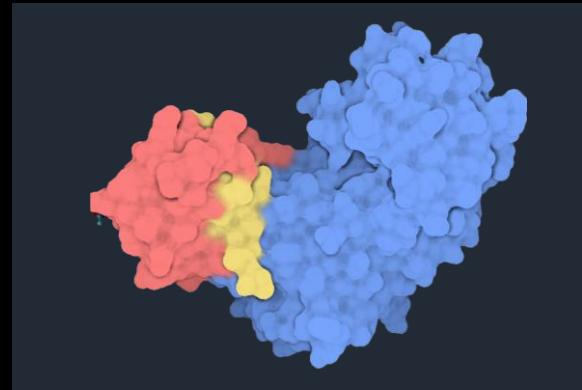


DNA / GENE

TTGCCCATTAGCTGTTGTCT
TAGTTCAAGATTGGTT [...]

4 characters
linear

1 or 2 copies per cell



PROTEIN

MGVPSGLILCVLIGAFFISM
AAAGDSGAYDWVMPARS [...]

20 characters
linear folded into 3D structure
1 – 10^7 copies per cell
post-translational modifications

-omics

= study of biomolecules in full context (holistic)

molecule

DNA

RNA

protein

-omics

= study of biomolecules in full context (holistic)

molecule	entity name
DNA	gene
RNA	transcript
protein	protein

-omics

= study of biomolecules in full context (holistic)

molecule	entity name	collection
DNA	gene	genome
RNA	transcript	transcriptome
protein	protein	proteome

-omics

= study of biomolecules in full context (holistic)

molecule	entity name	collection	field of study
DNA	gene	genome	genomics
RNA	transcript	transcriptome	transcriptomics
protein	protein	proteome	proteomics

proteomics

= study of proteins in full context (holistic)

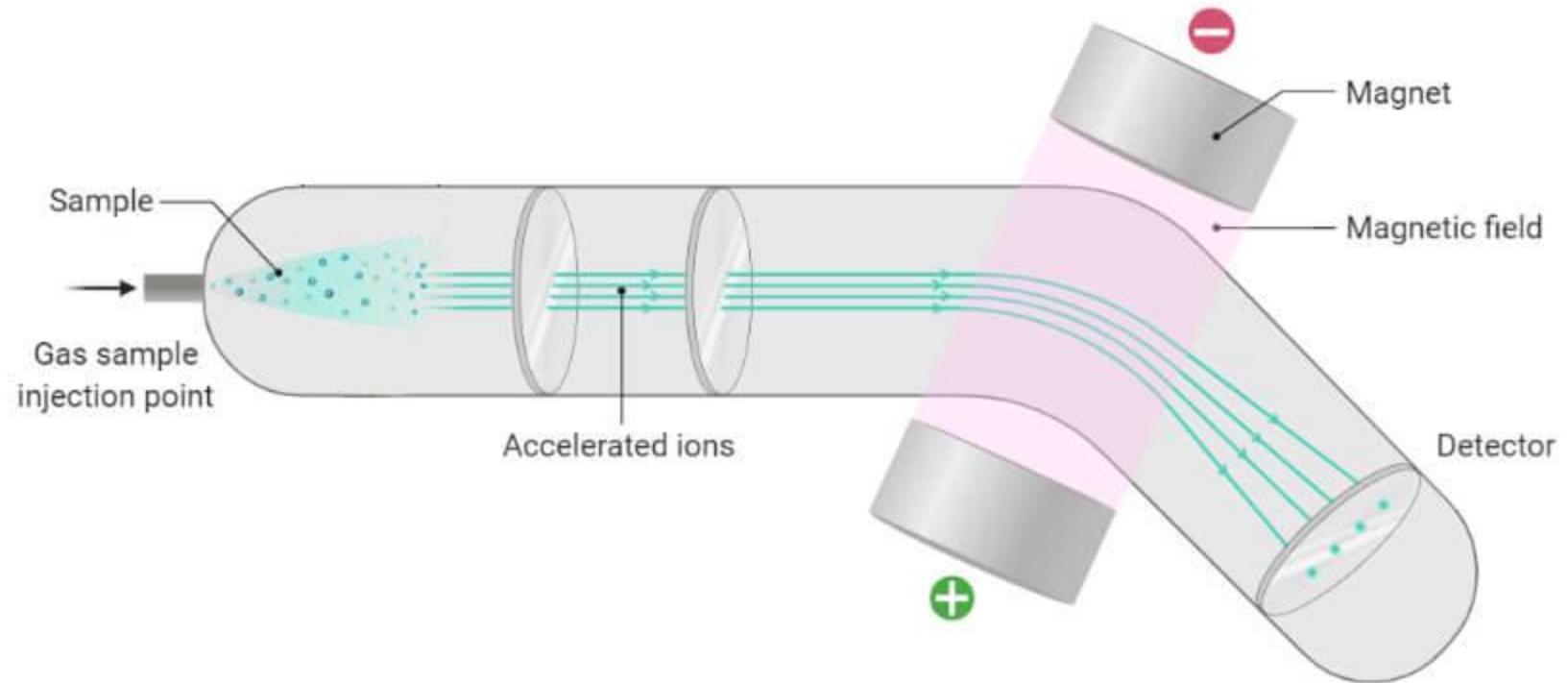


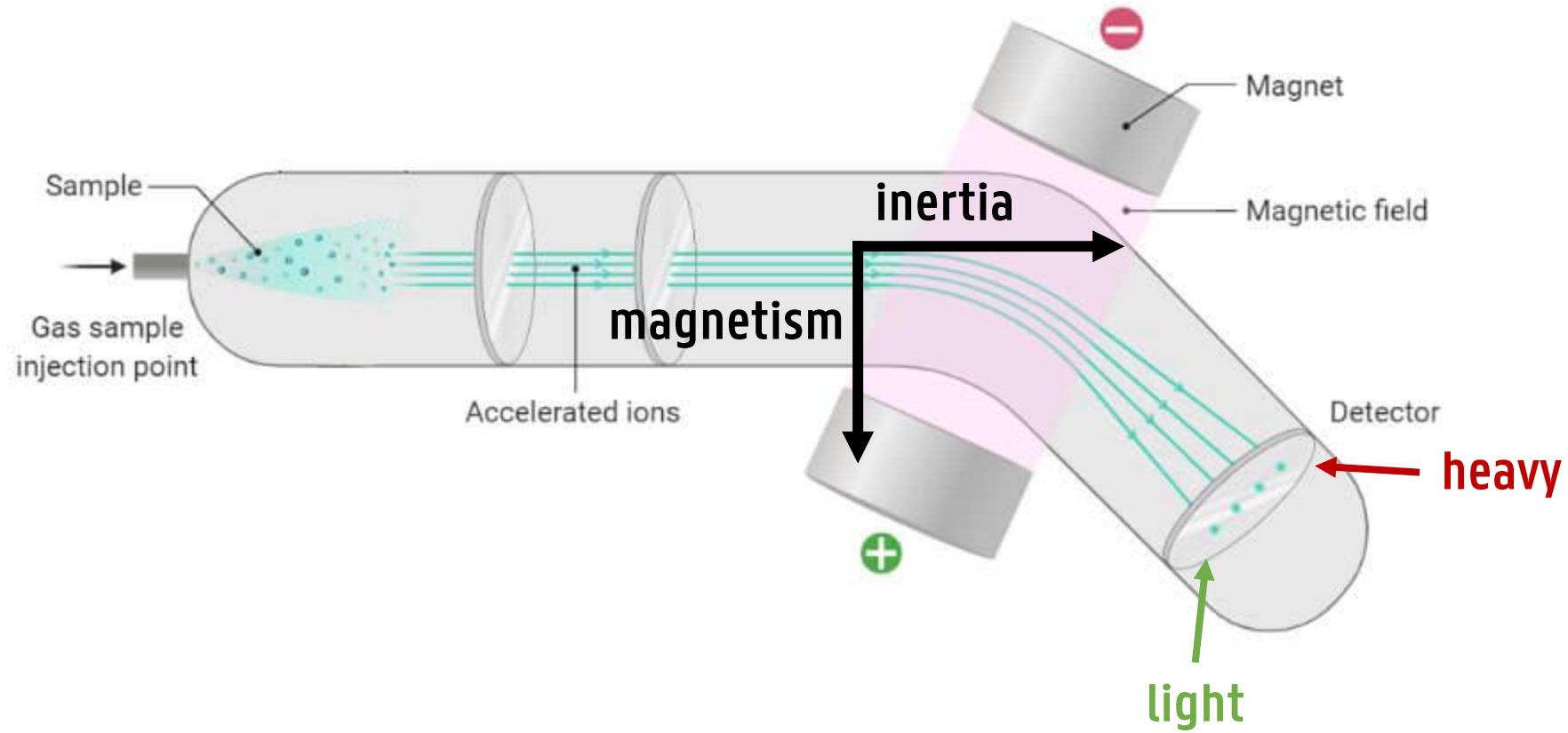
mass spectrometry

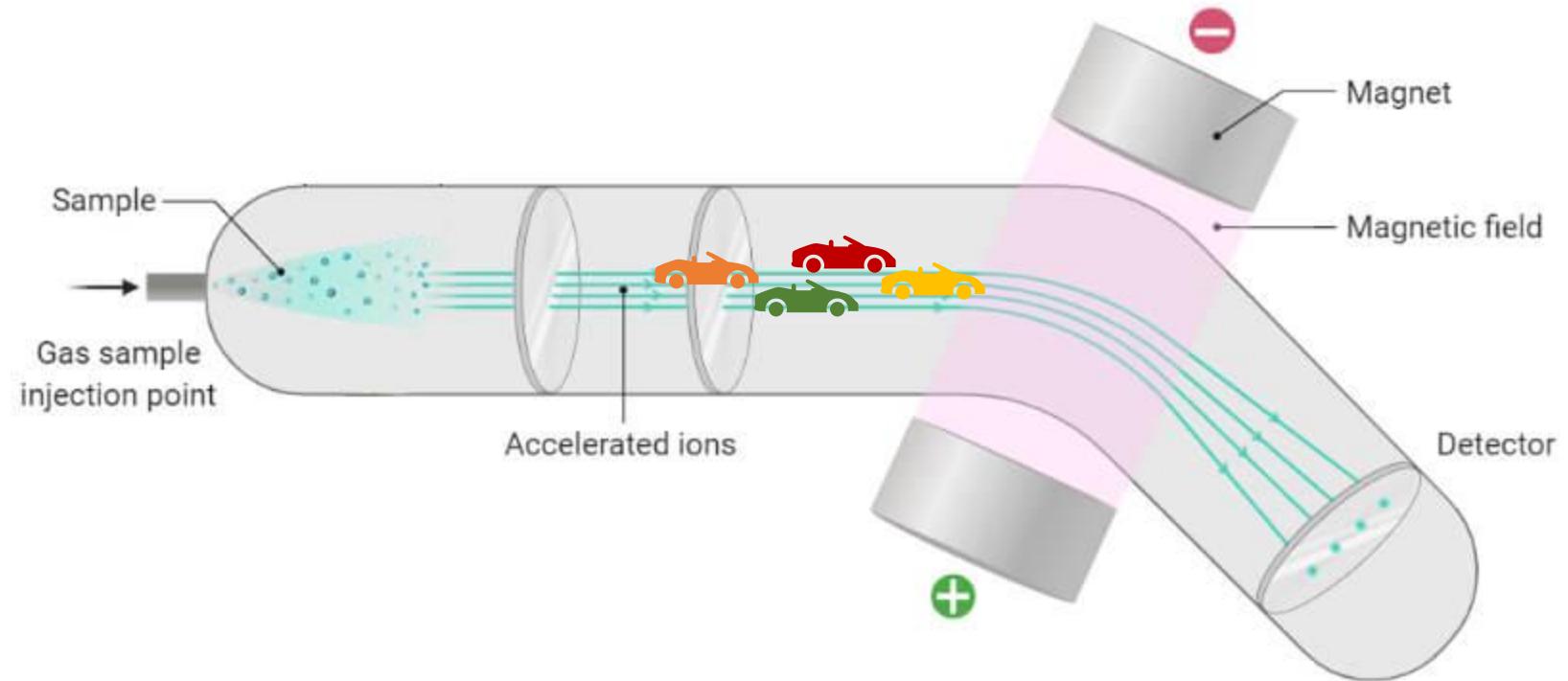
mass spectrometer

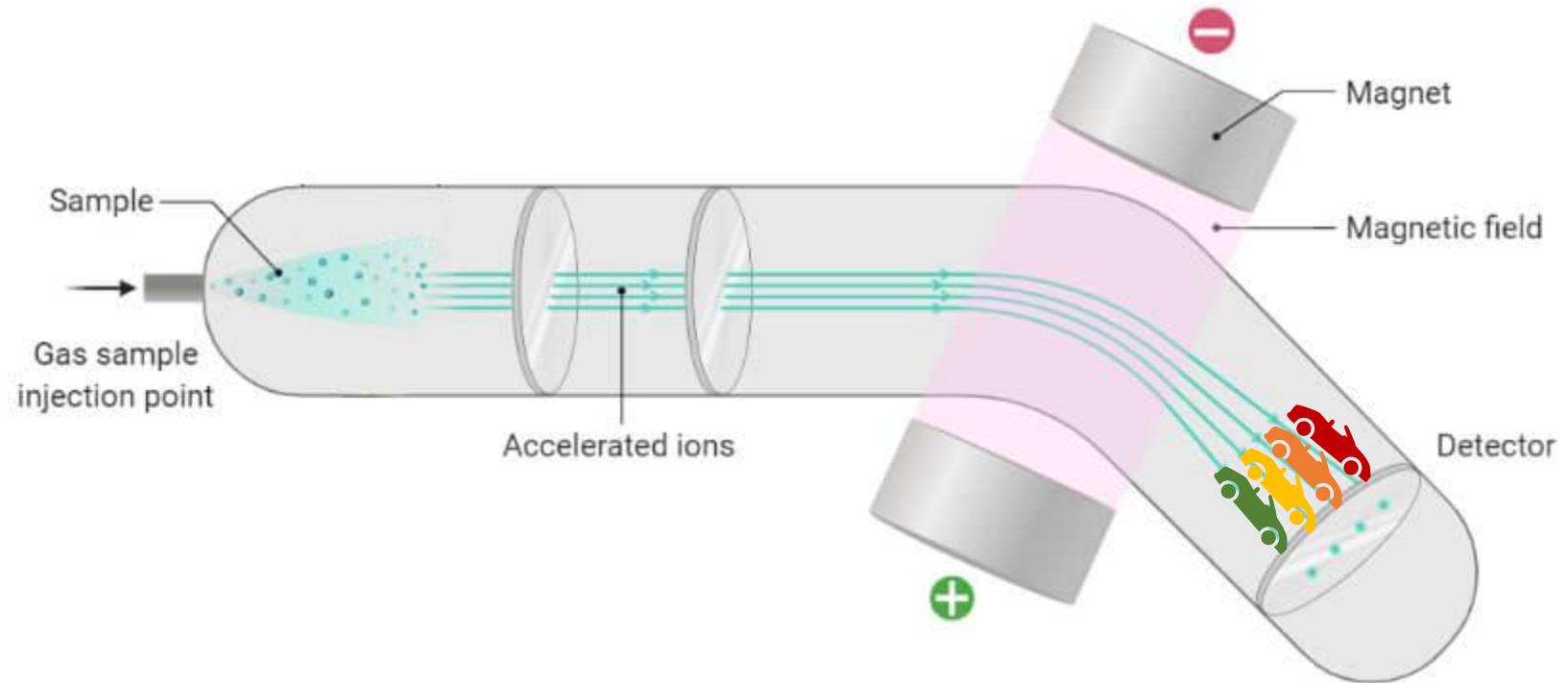
very sensitive molecular “scale”

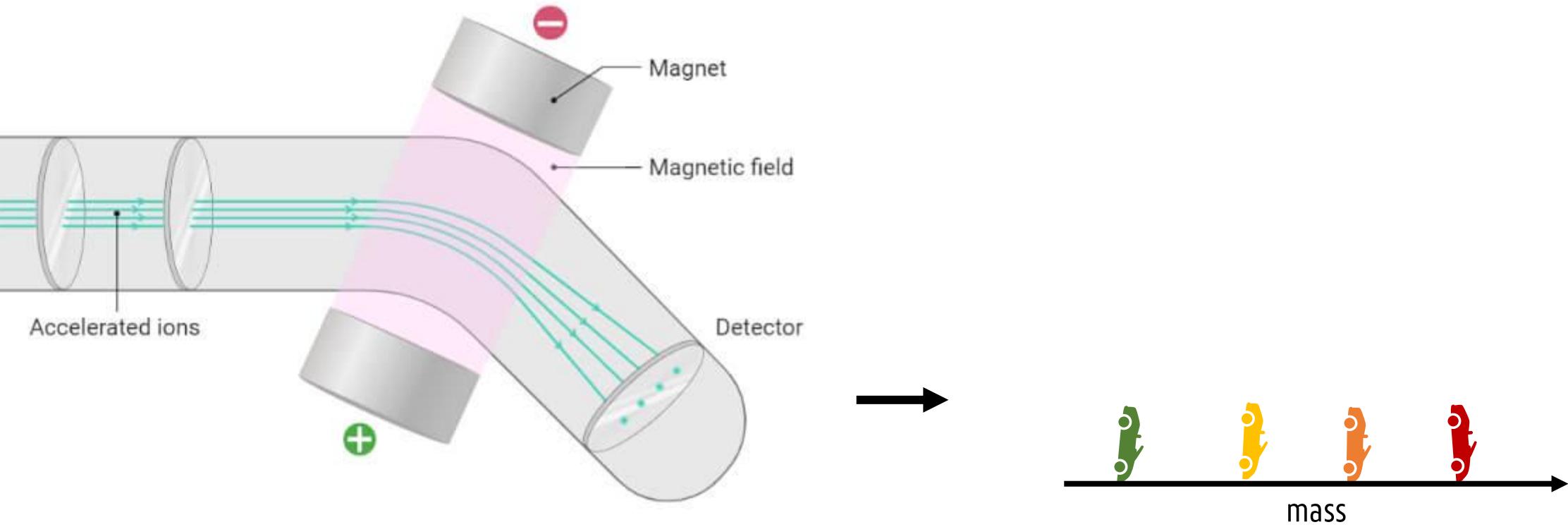


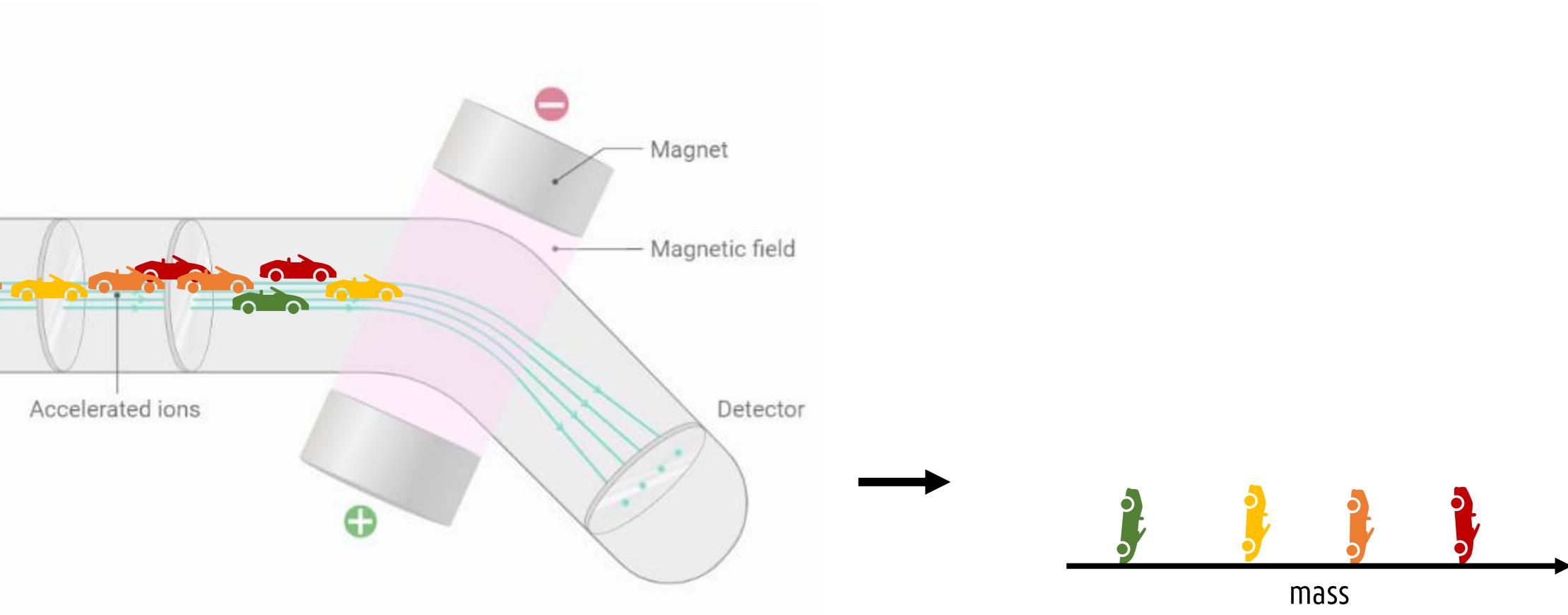


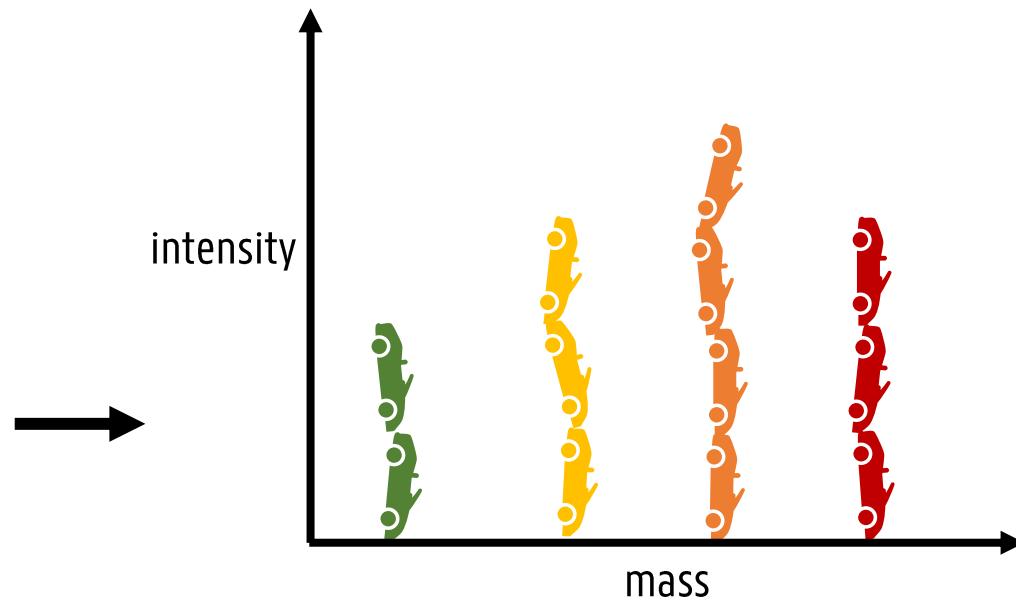
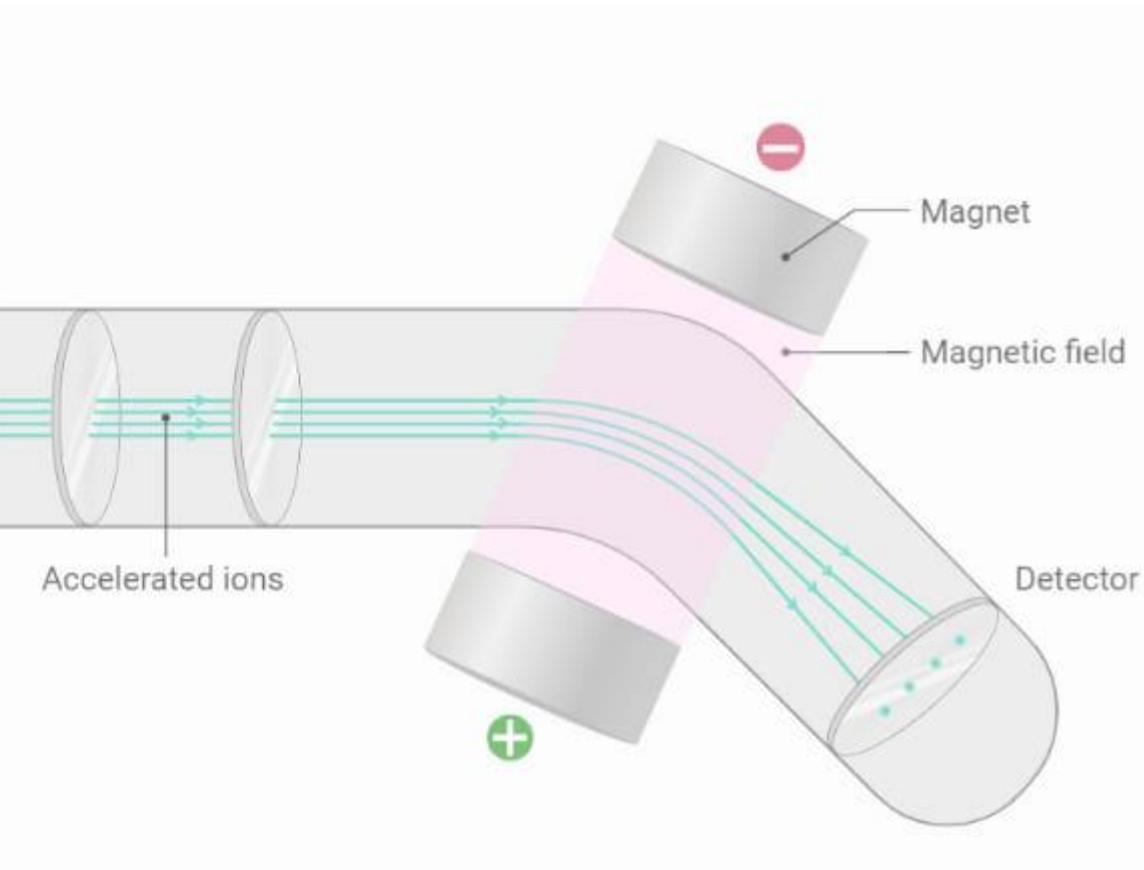








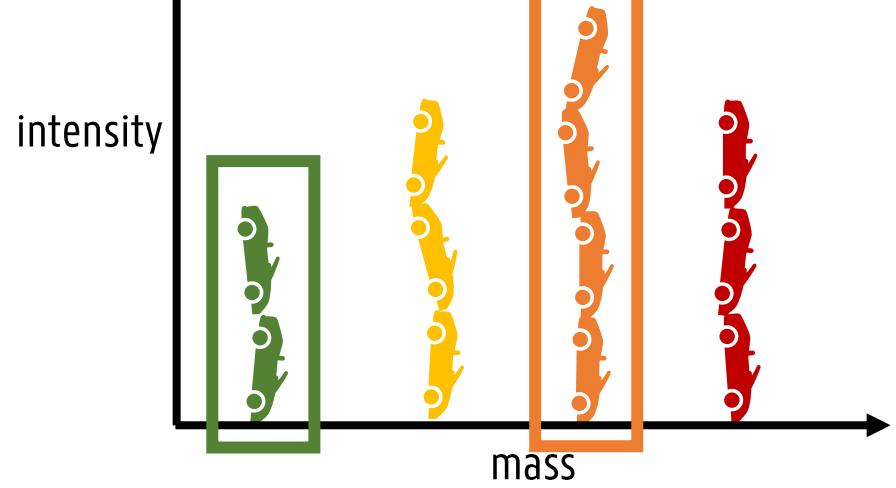




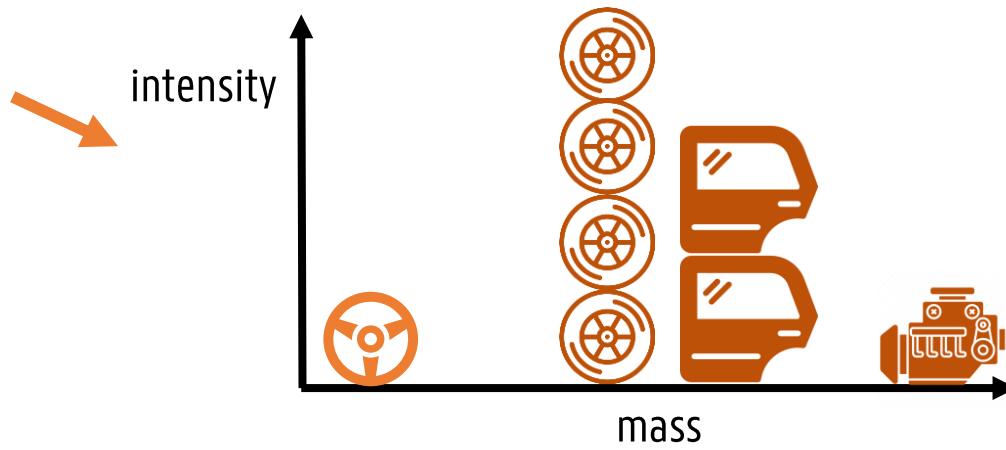
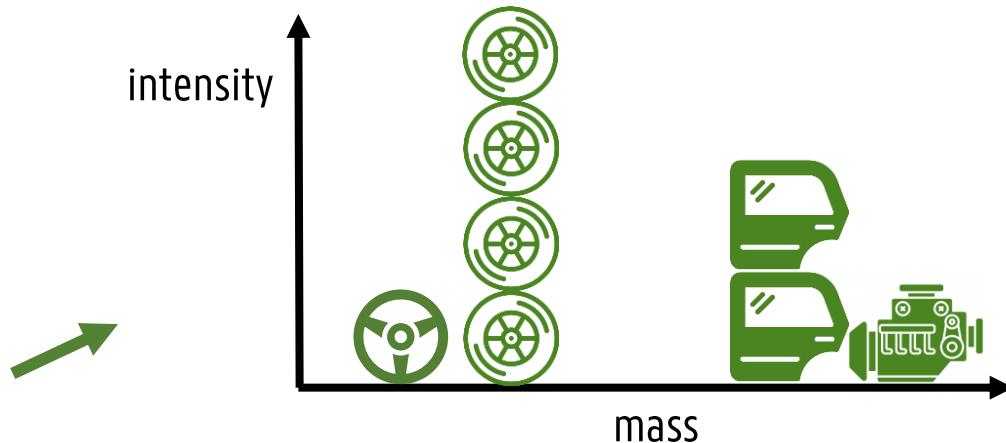


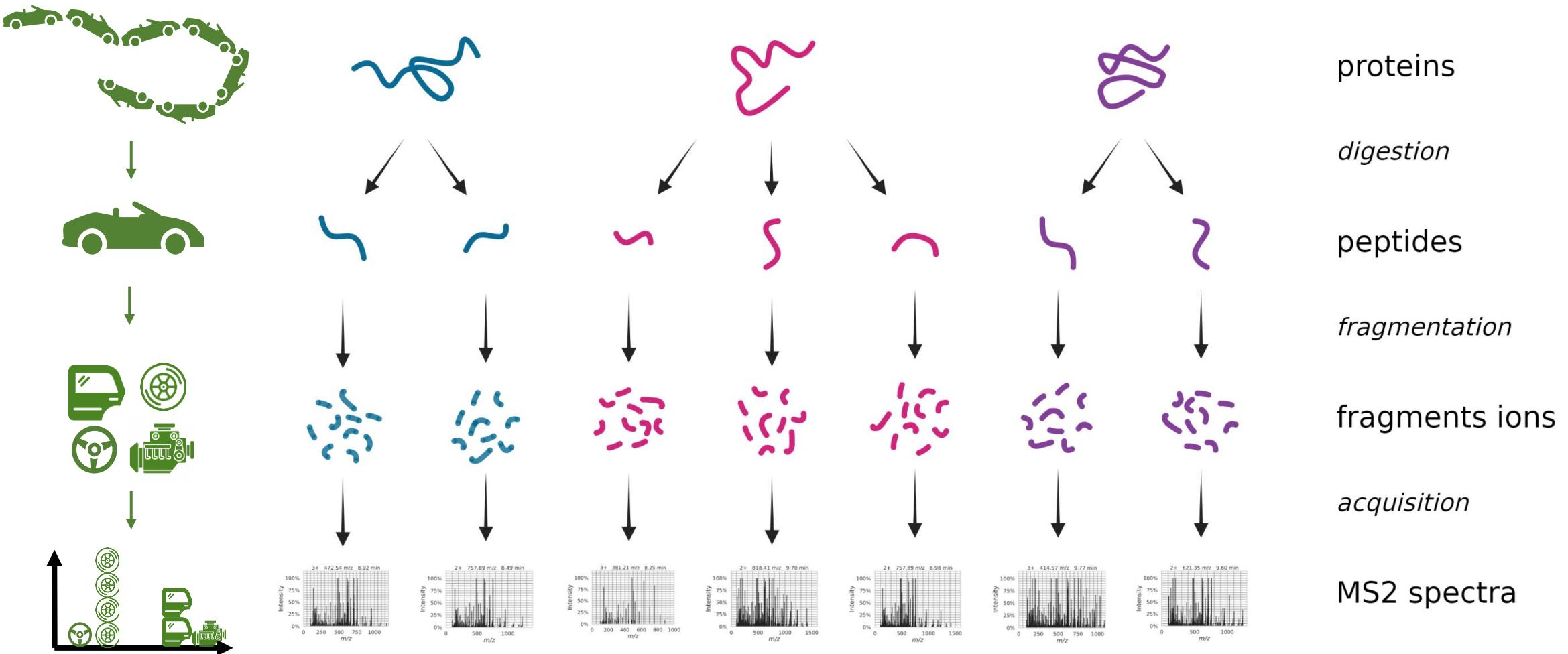


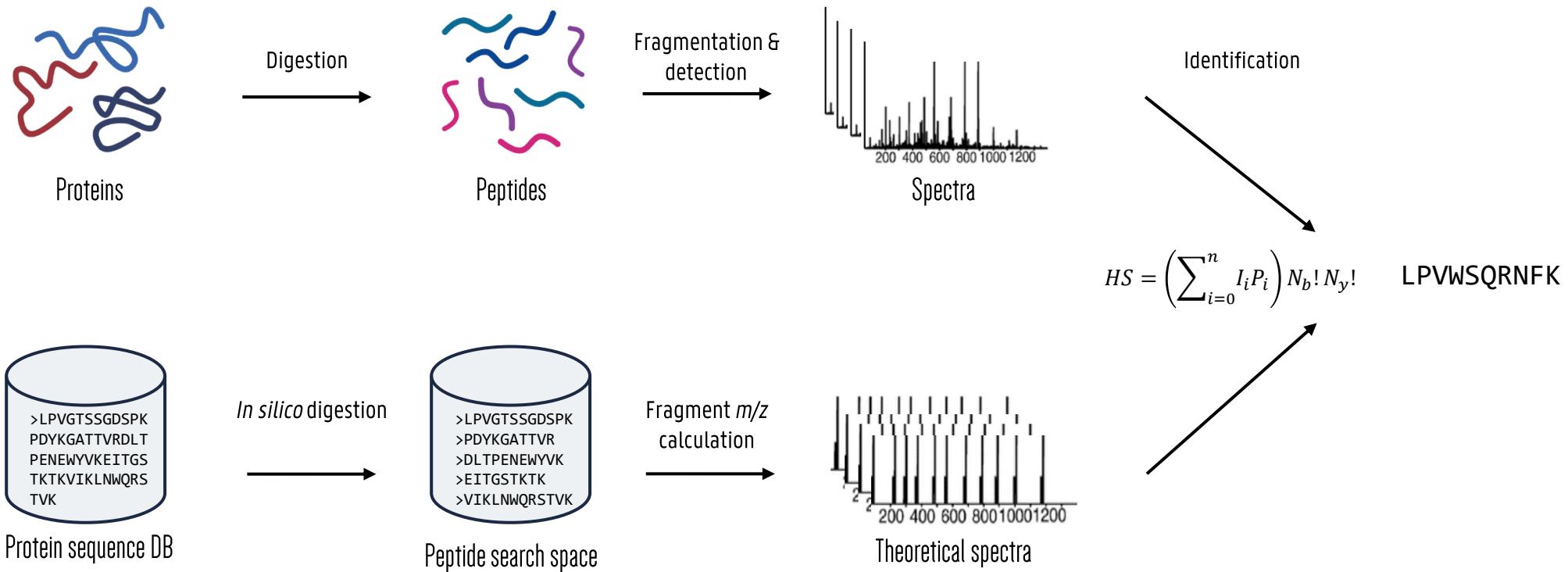
MS1

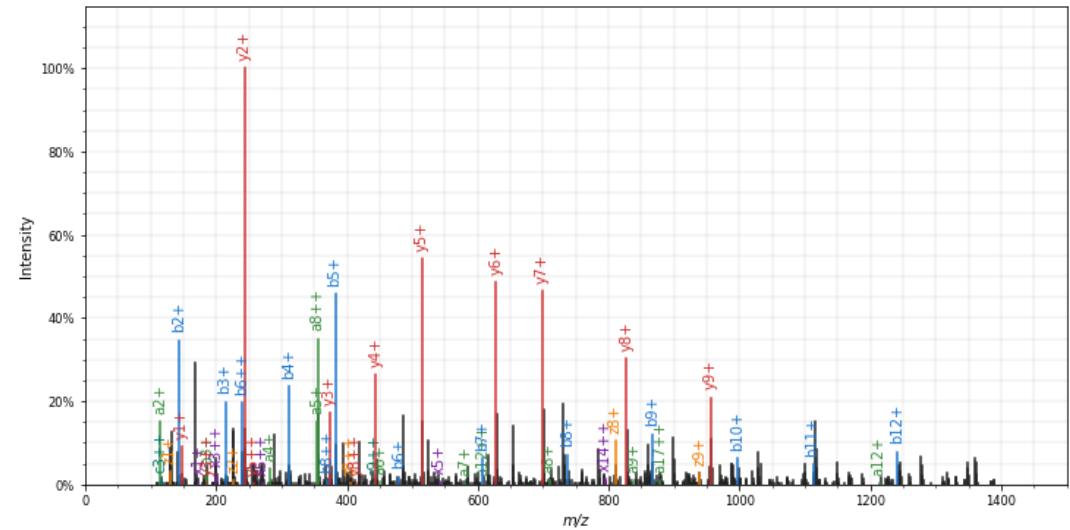
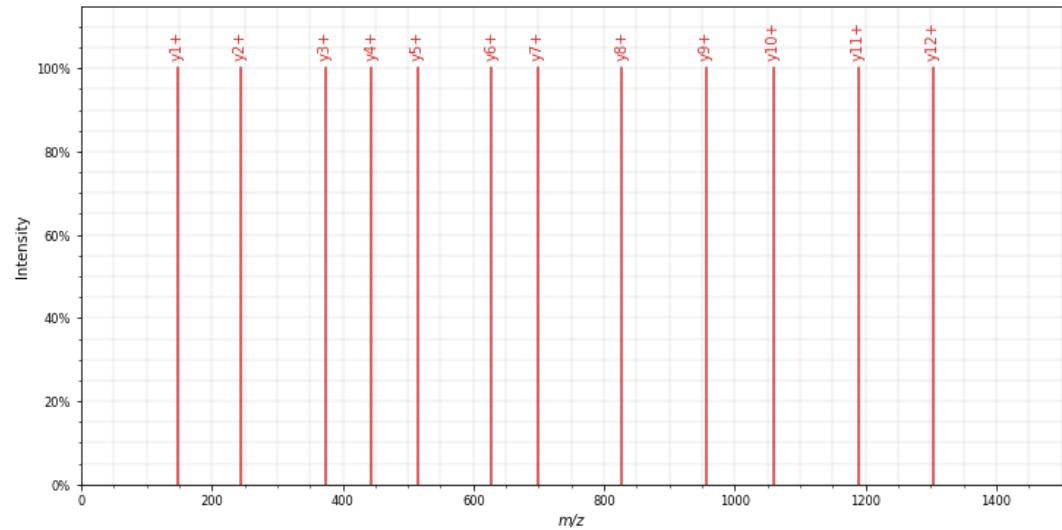


MS2







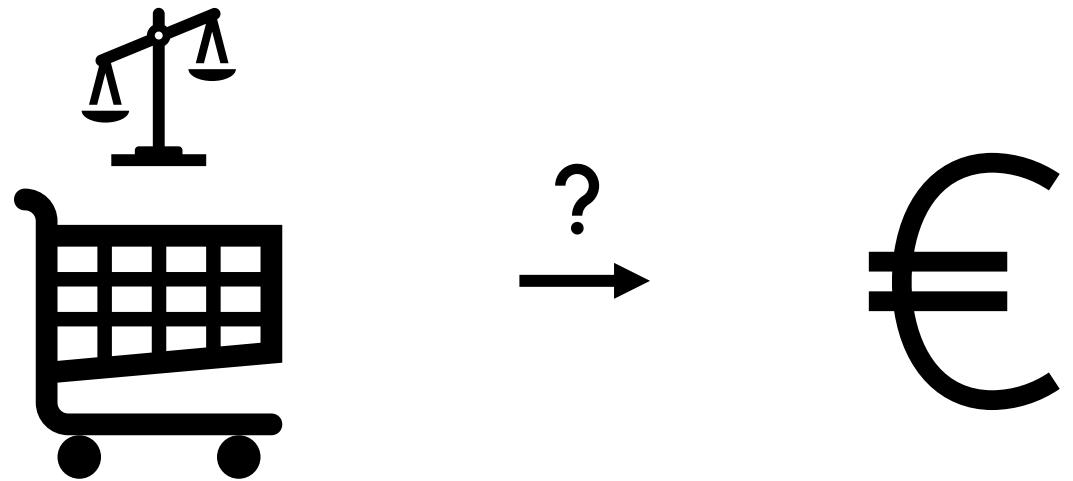


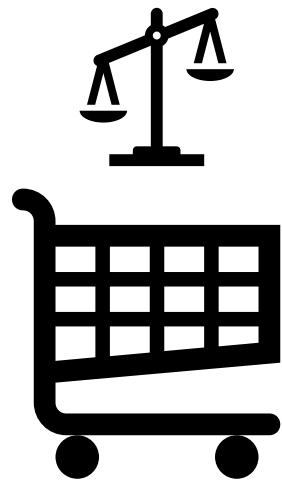
machine learning

artificial intelligence

machine learning

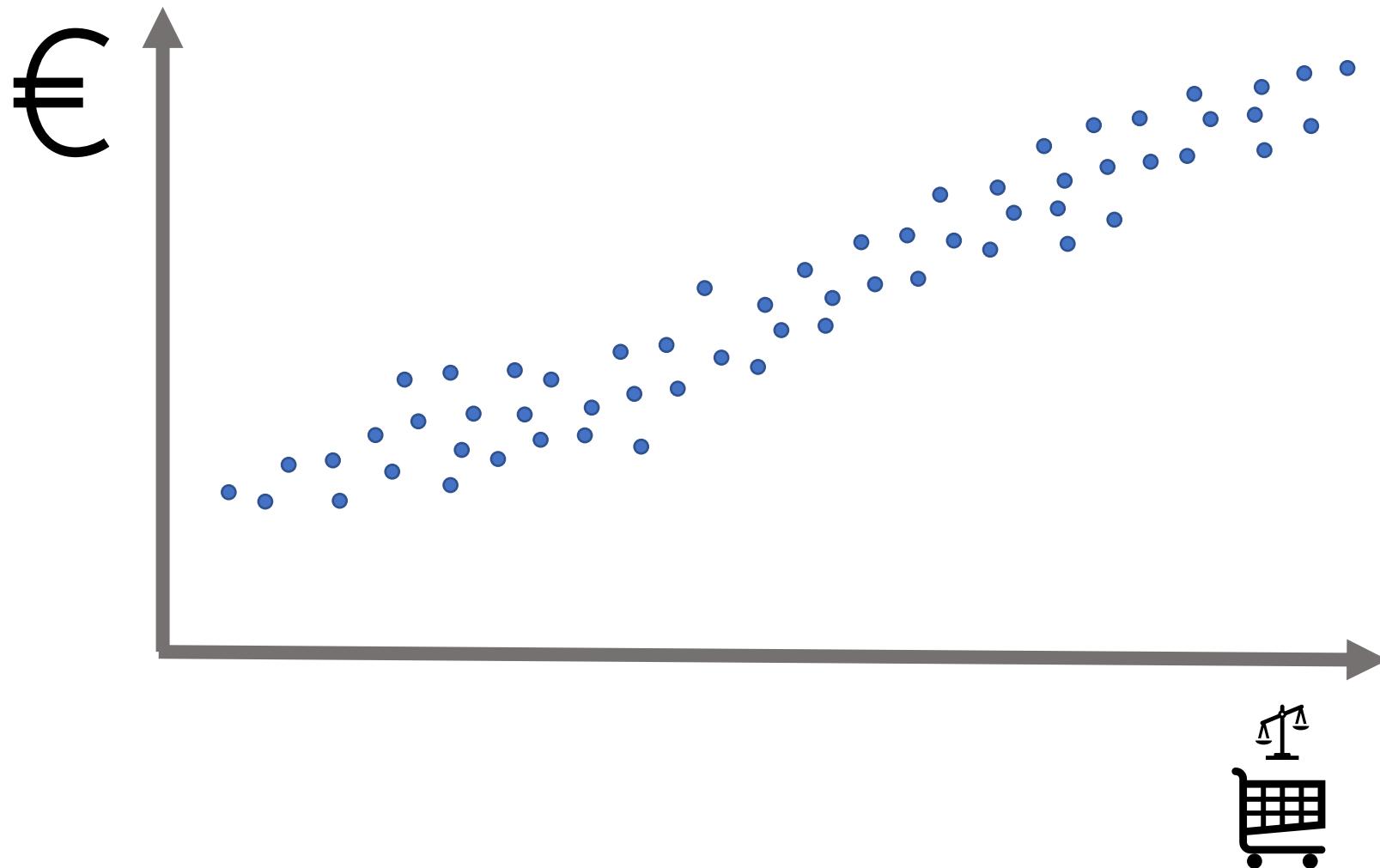
deep learning

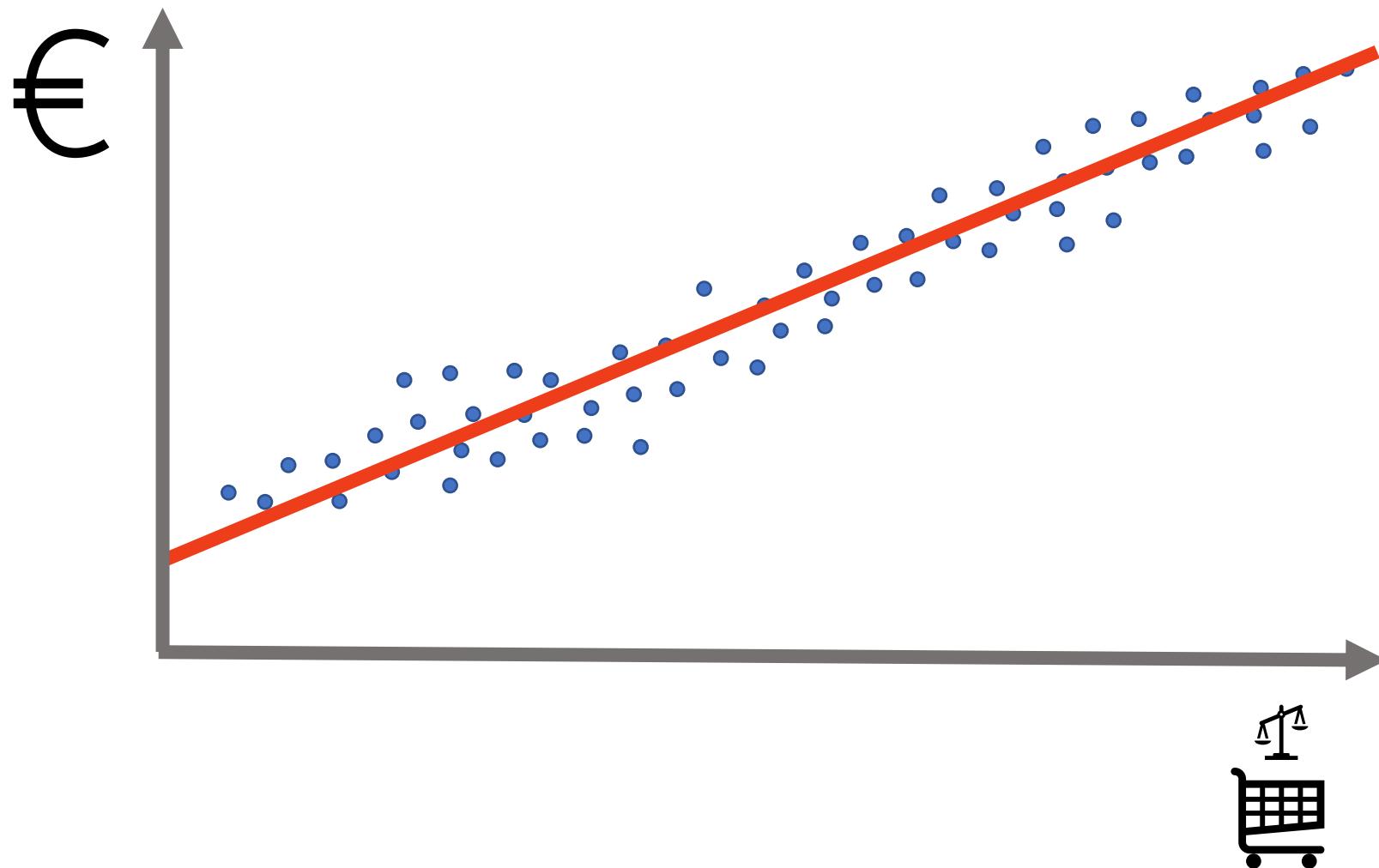


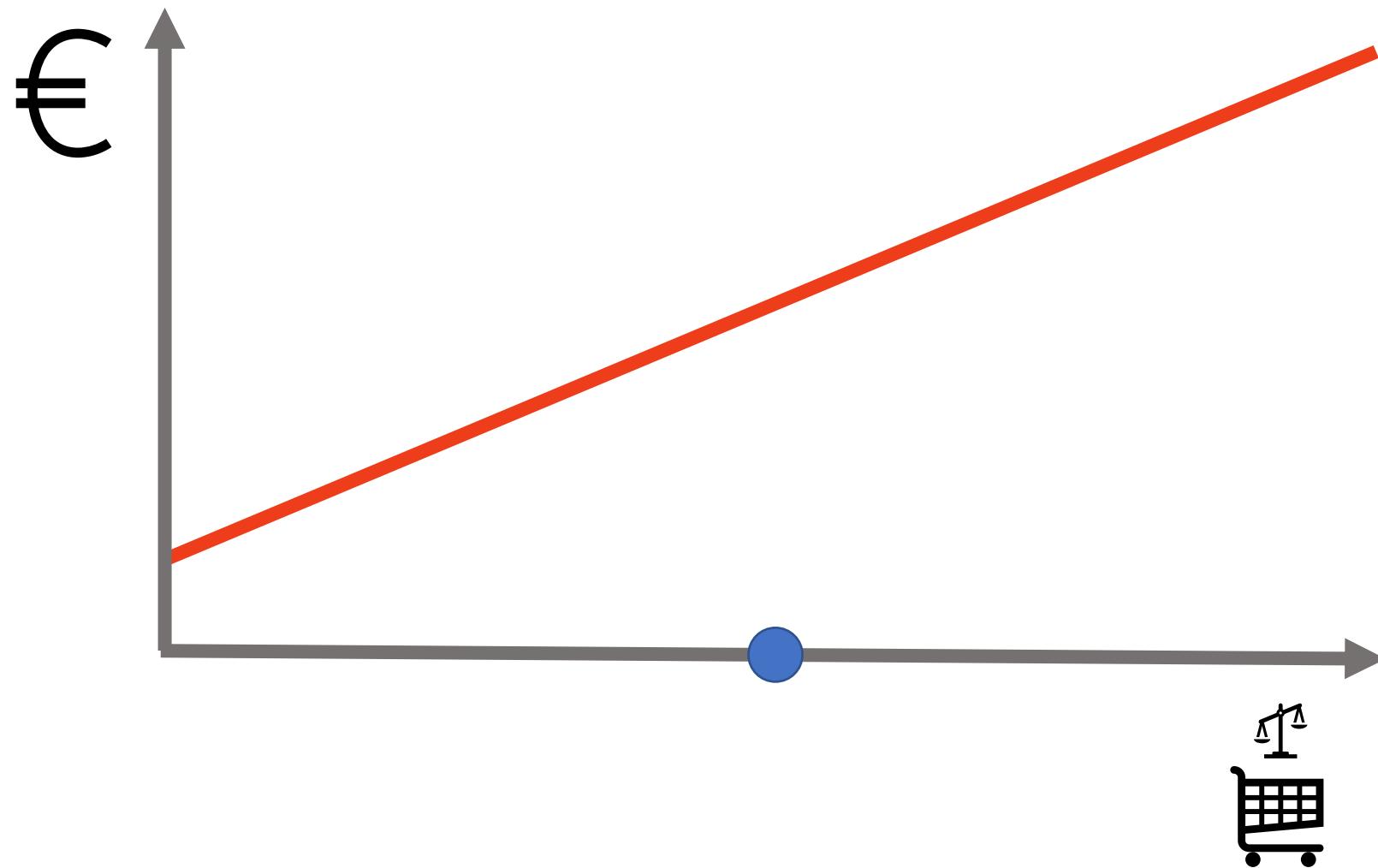


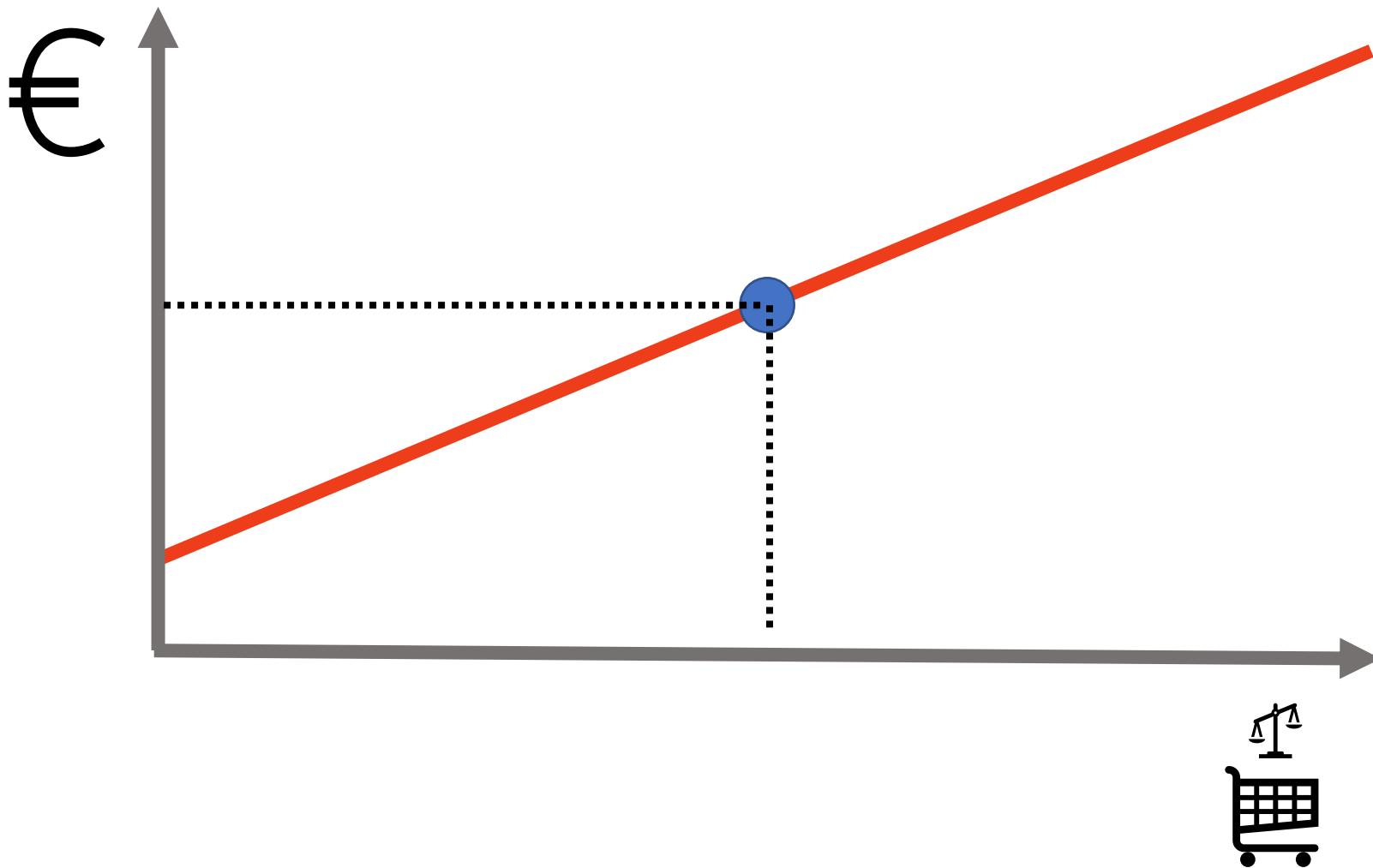
A large black euro symbol (€), representing the currency unit for the cost.

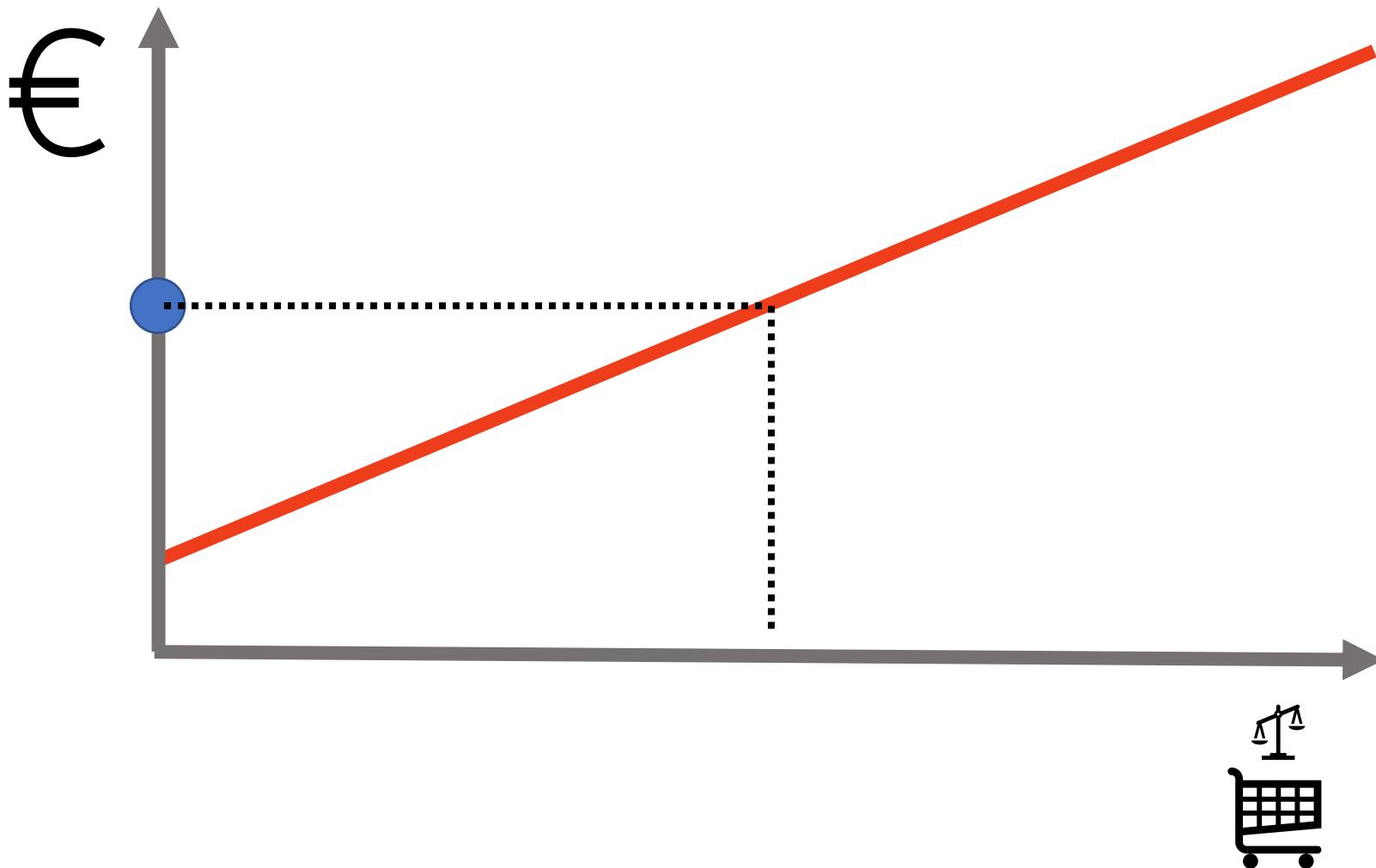
weight	cost
34 kg	€54
12 kg	€21
64 kg	€108
...	...

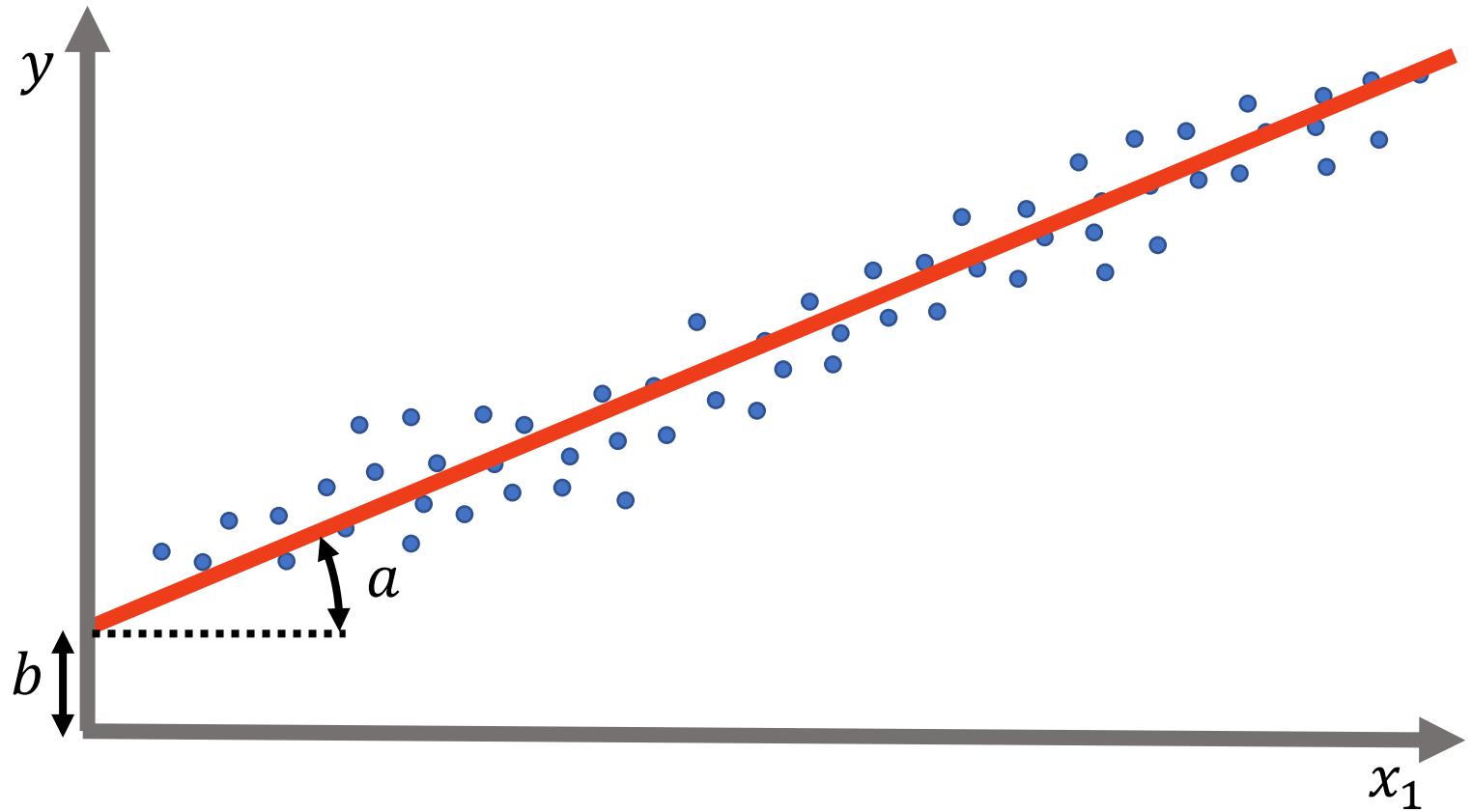




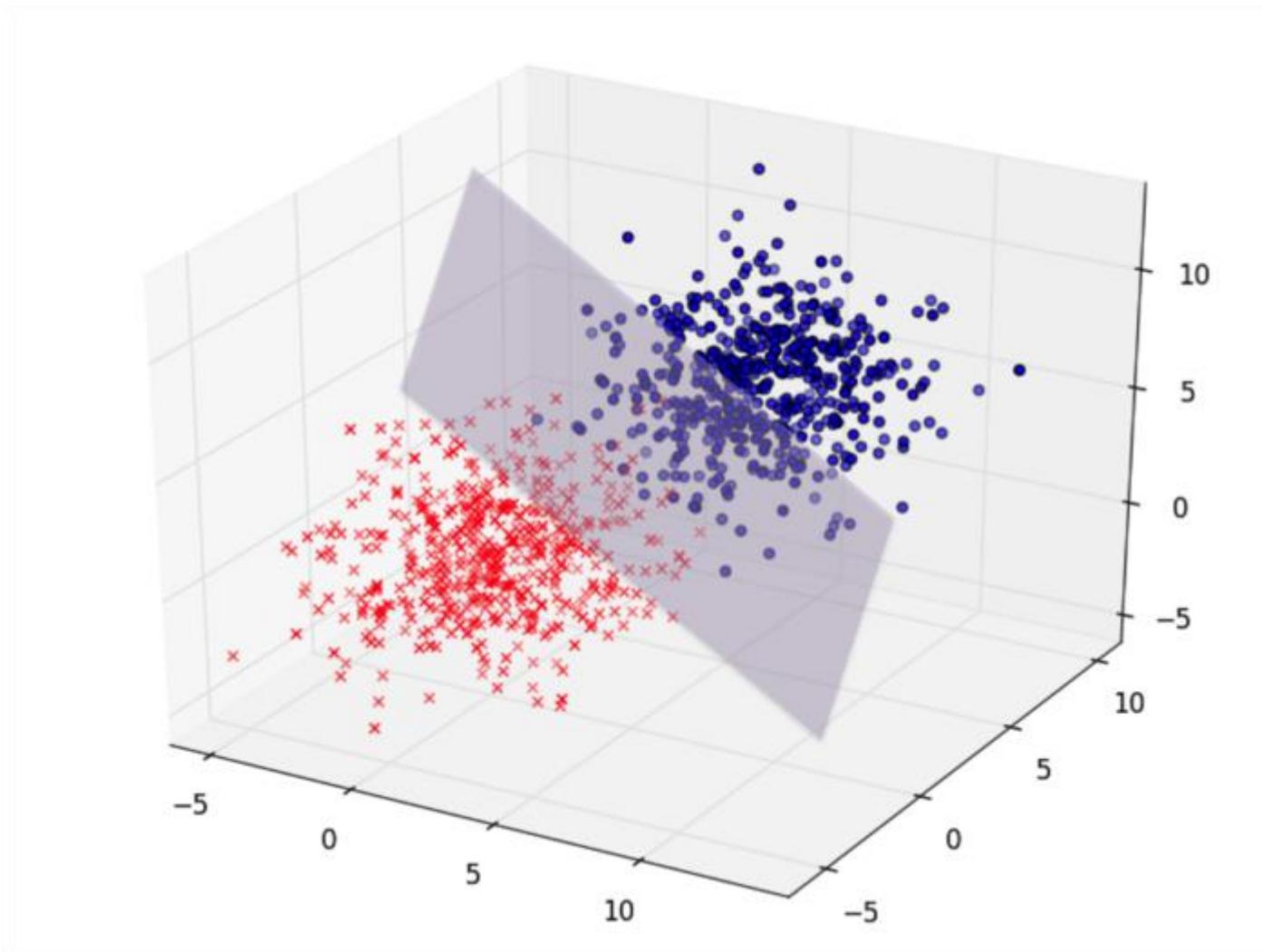






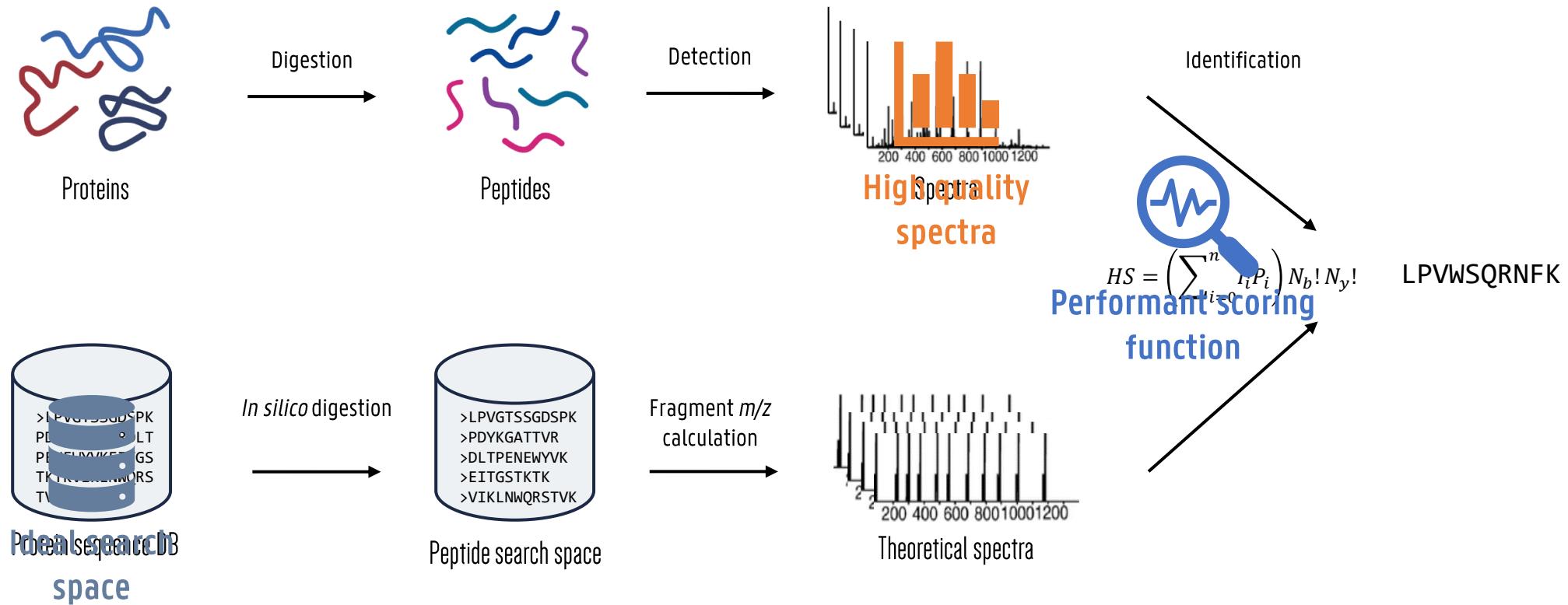


$$y = a \cdot x_1 + b$$

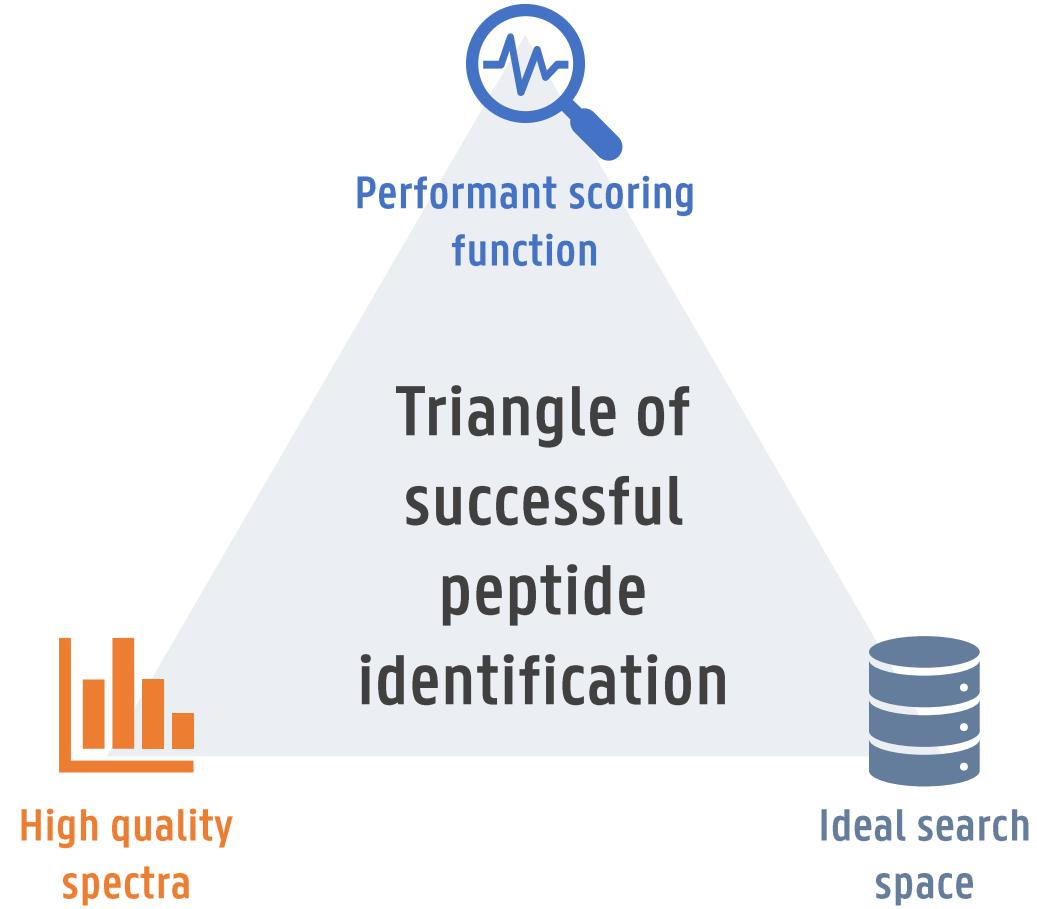


**proteomics
+ mass spectrometry
+ machine learning
= ?**

Key elements to successful peptide spectrum identification?



Key elements to successful peptide spectrum identification?



Key elements to successful peptide spectrum identification?



Triangle of
successful
peptide
identification



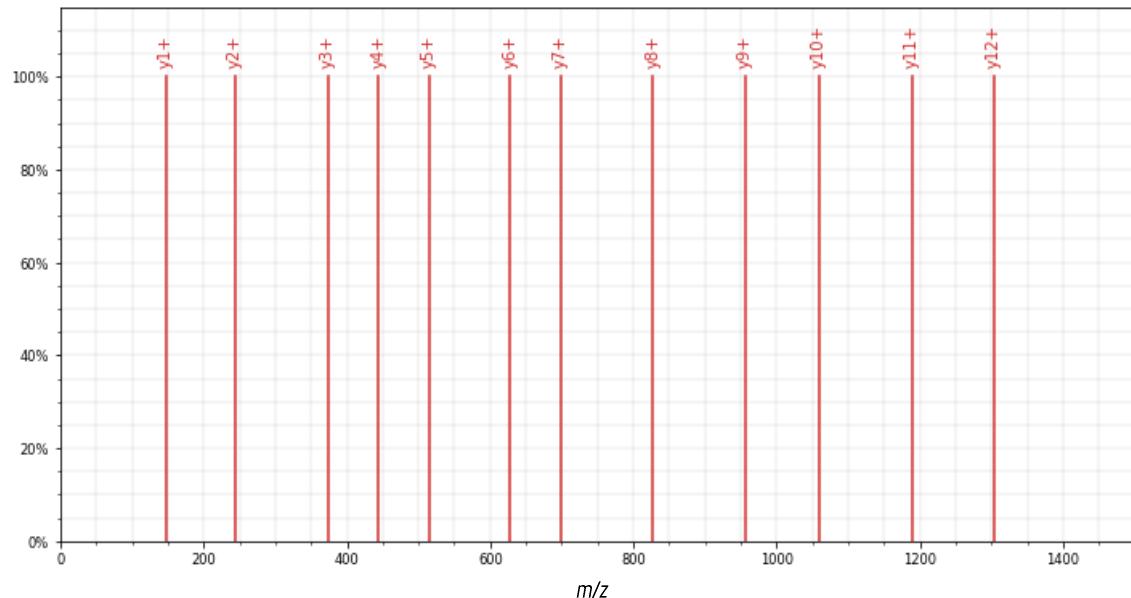
High quality
spectra



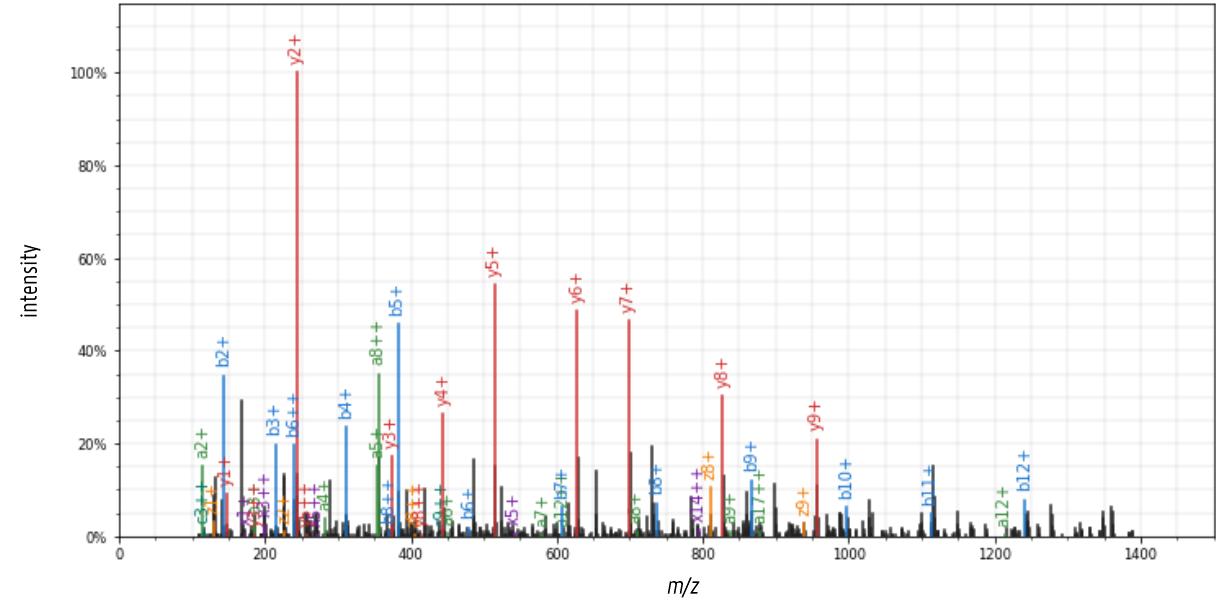
Ideal search
space



Machine learning to the rescue!

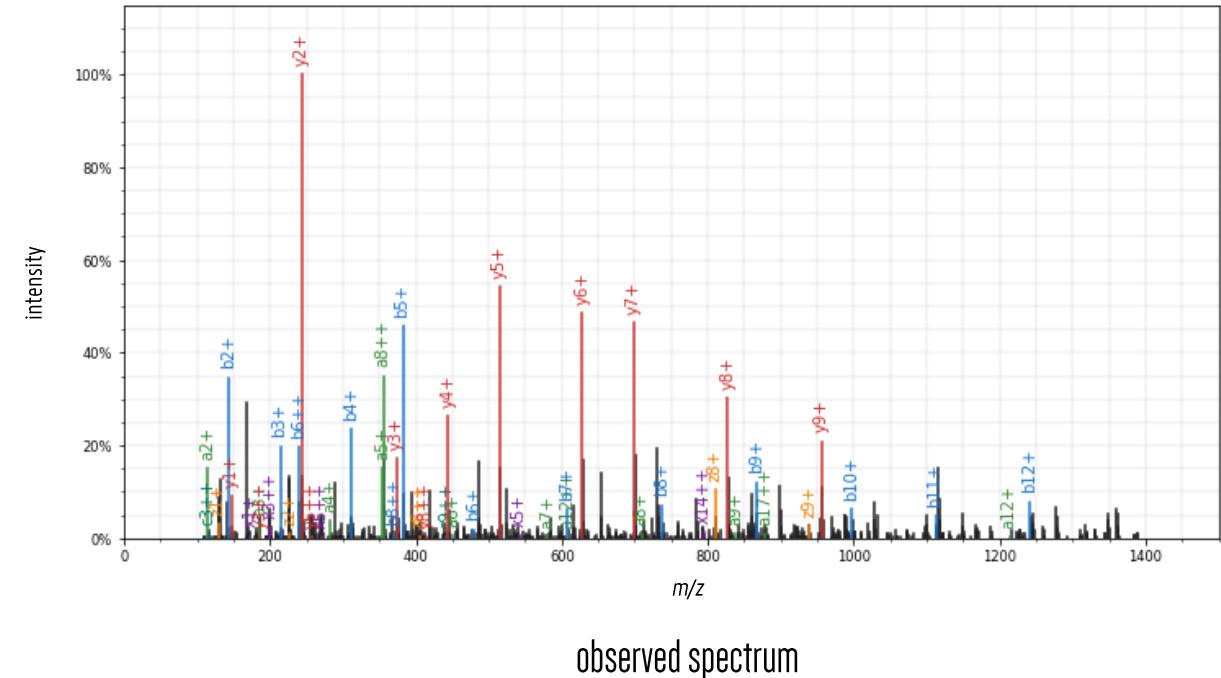
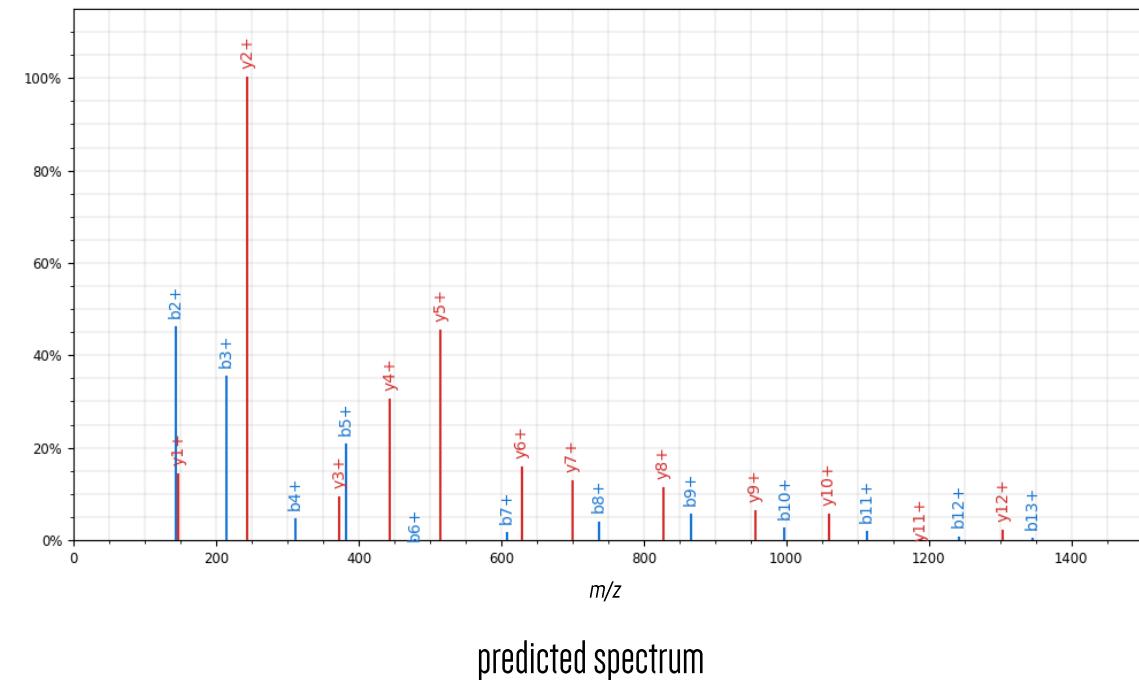


theoretical spectrum

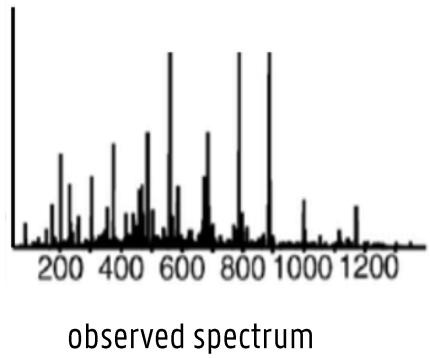


observed spectrum

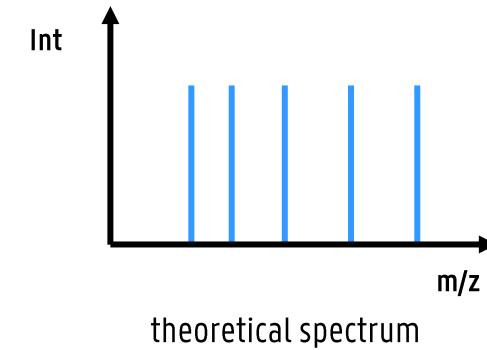
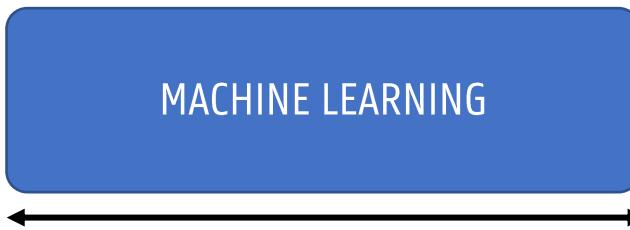
Machine learning to the rescue!



Machine learning to the rescue!



observed spectrum



theoretical spectrum

Introduction (Dutch)
proteins and proteomics
mass spectrometry
machine learning

Results (English)

MS²PIP: Peptide spectrum prediction for multiple fragmentation methods,
instruments, and labeling techniques

Removing the hidden data dependency of DIA with predicted spectral libraries

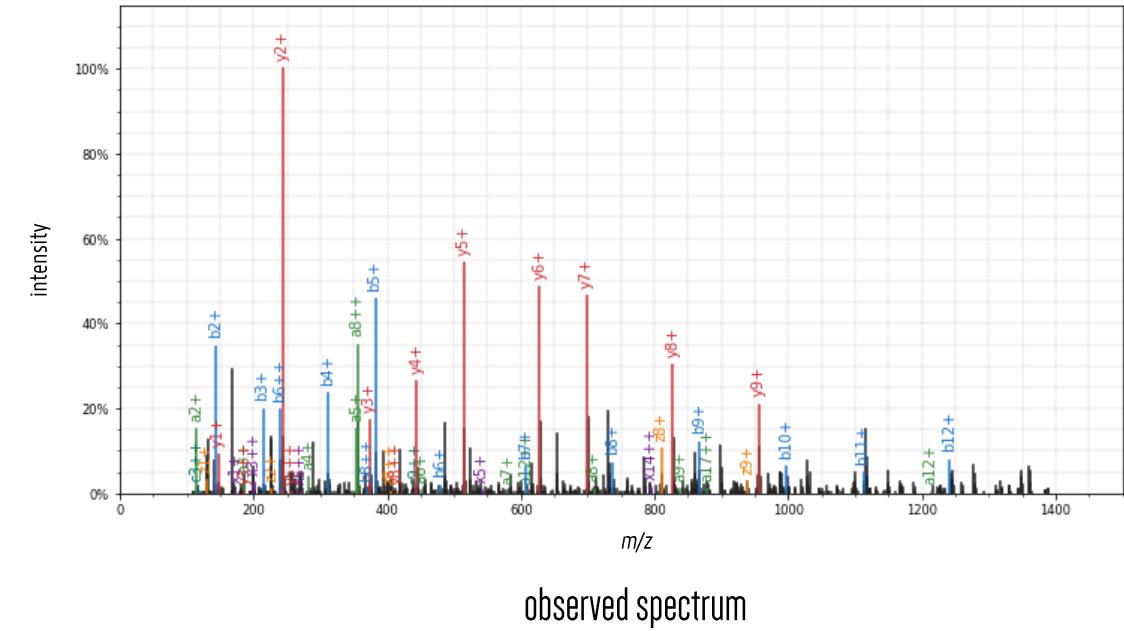
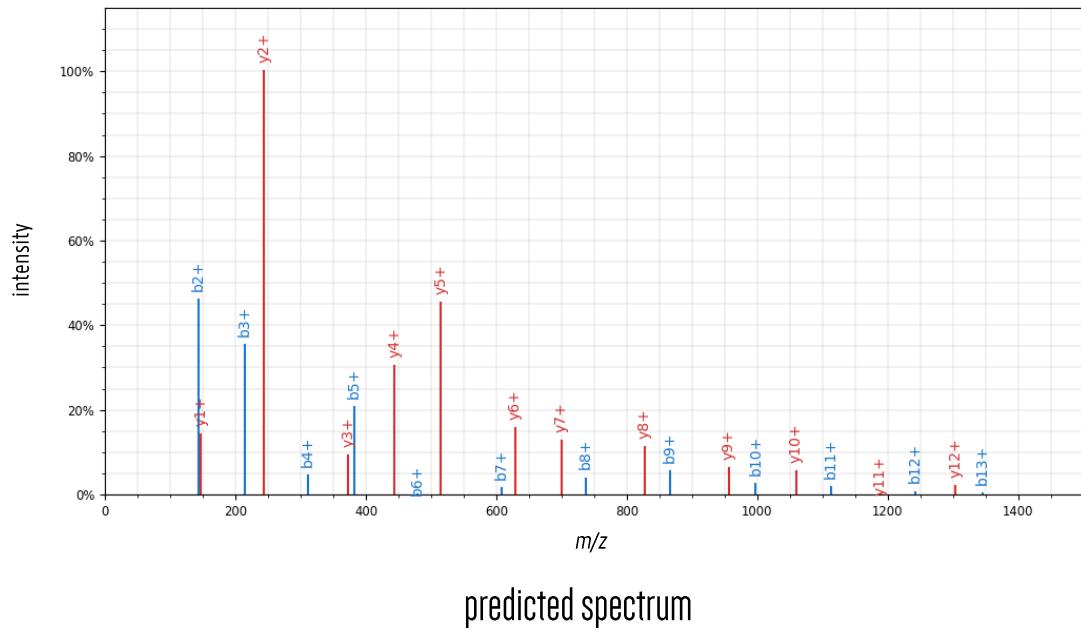
MS²Rescore: Leveraging spectrum predictions to enable novel proteomics workflows

MS²DIP: spectrum prediction for modified peptides

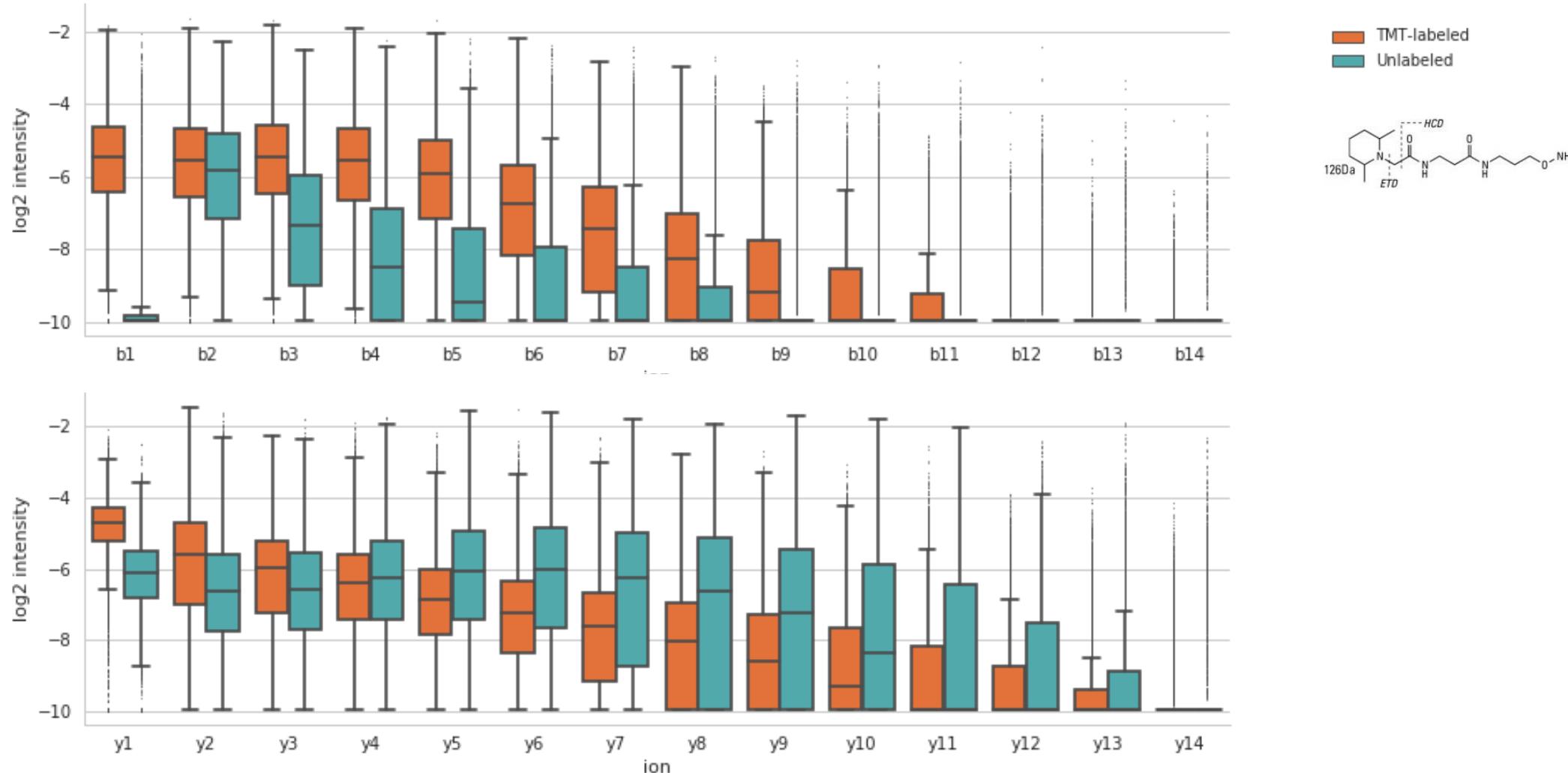
Discussion and future perspectives (English)

MS²PIP: Peptide spectrum prediction for multiple fragmentation methods, instruments, and labeling techniques

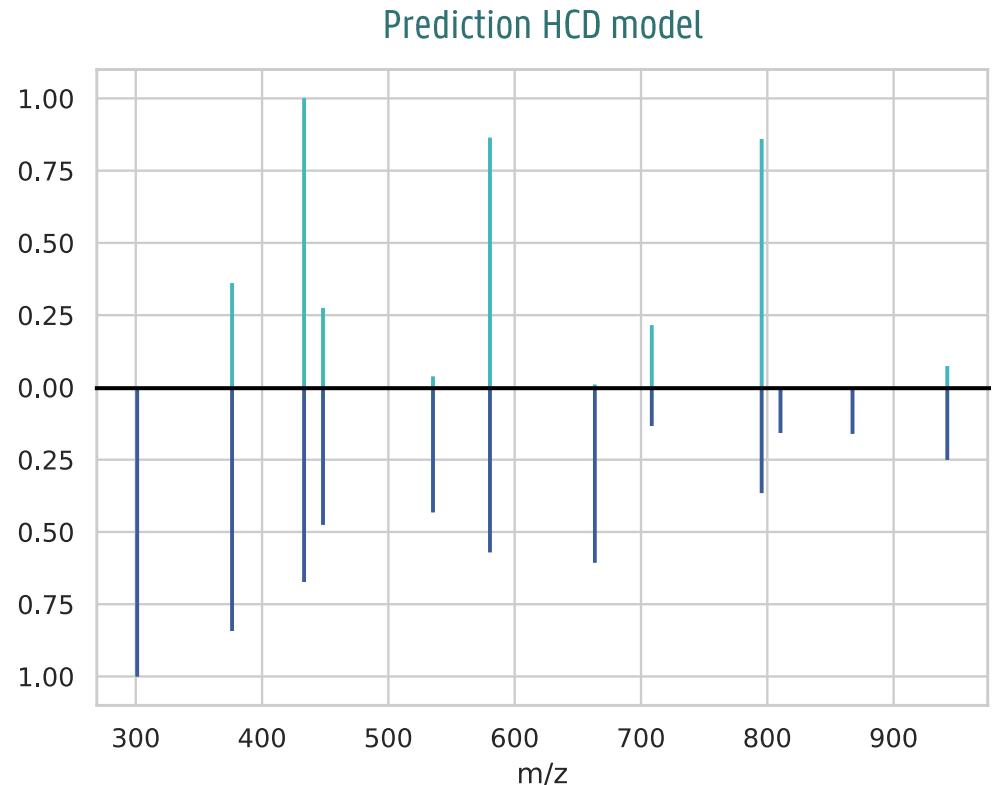
MS²PIP can accurately predict peptide fragmentation spectra



However, peptide fragmentation can drastically differ between experimental setups

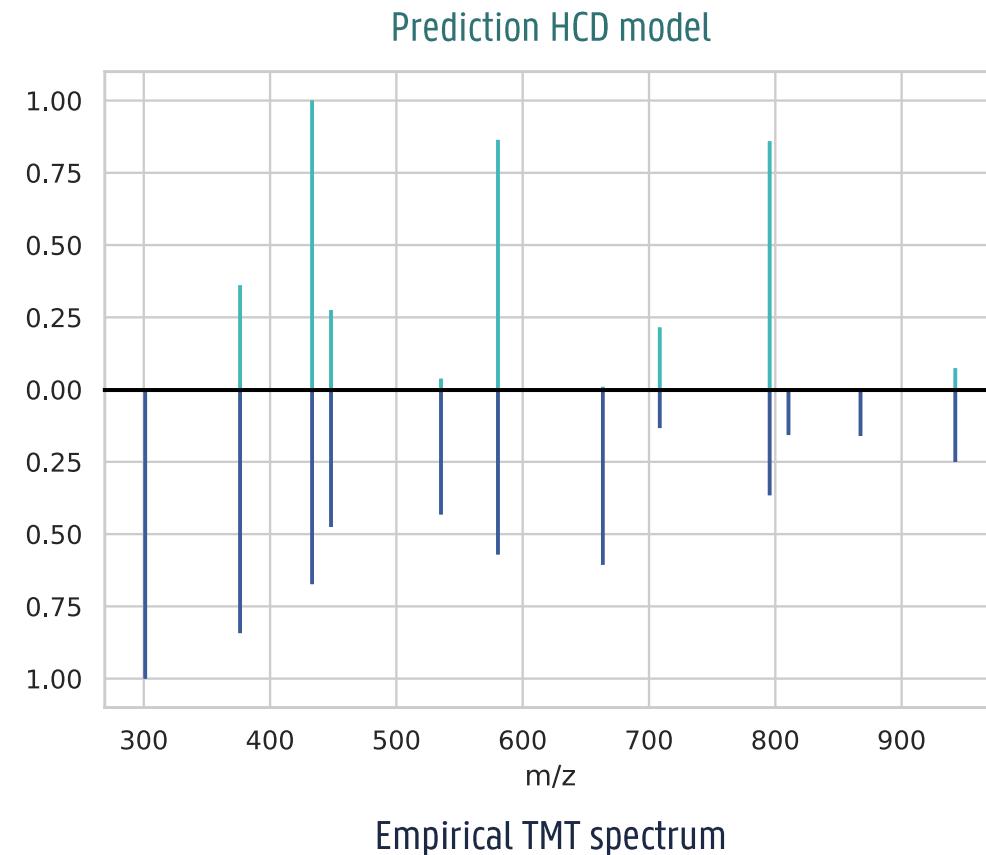
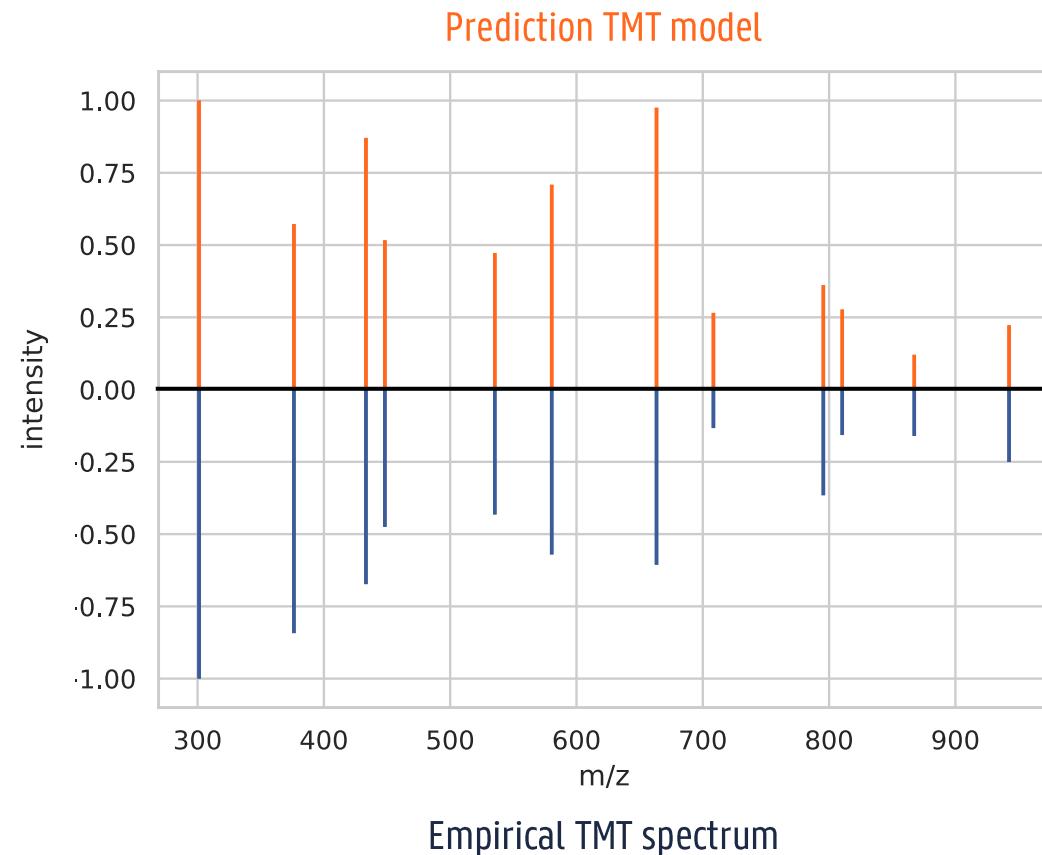


Specialized setups require specialized MS²PIP models

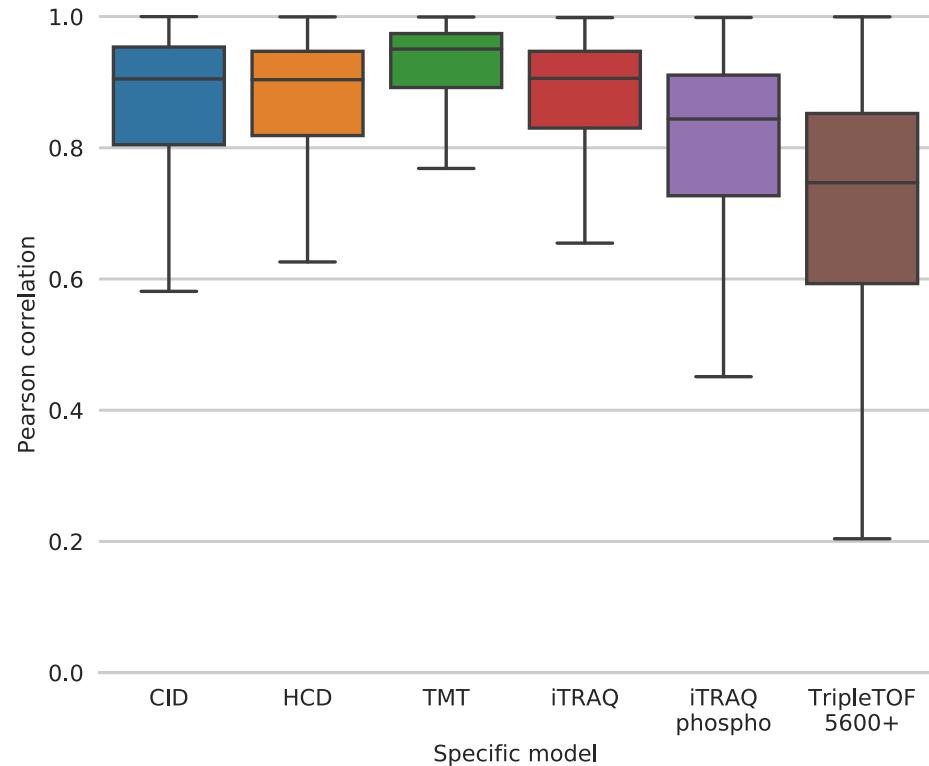


→ Pearson correlation = 0.20

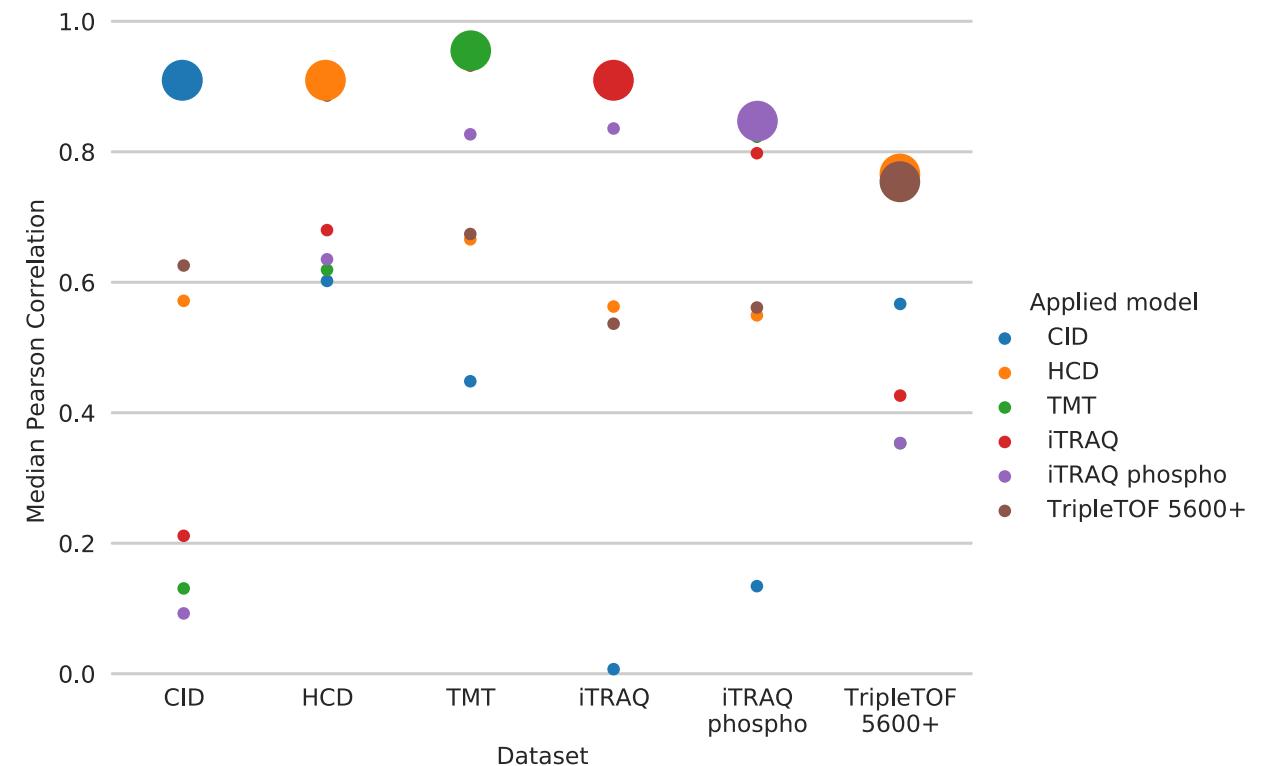
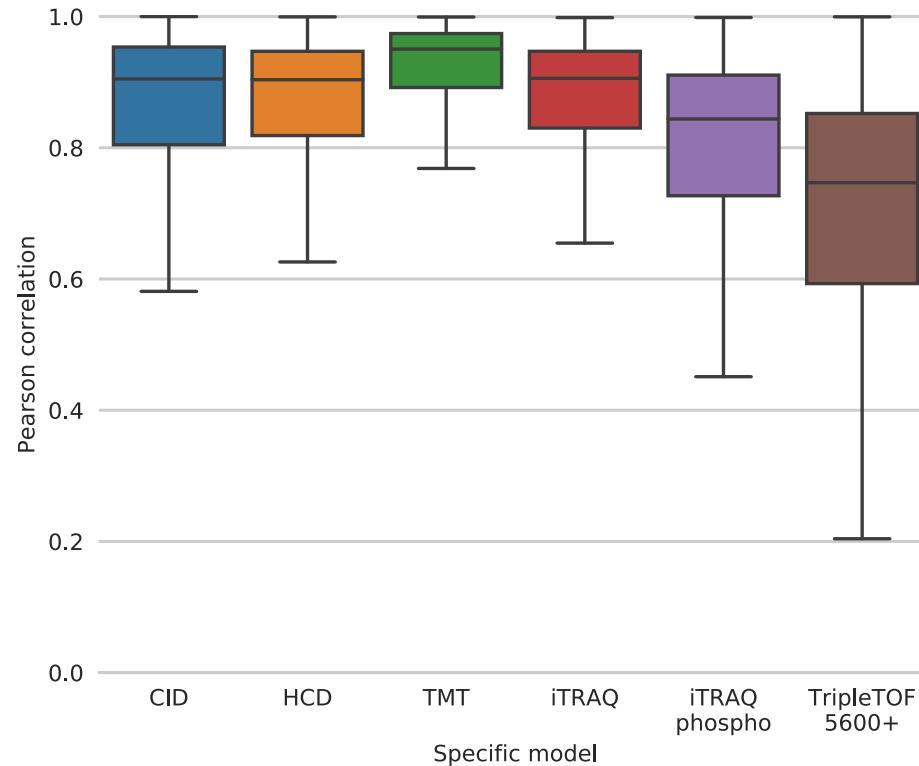
Specialized setups require specialized MS²PIP models



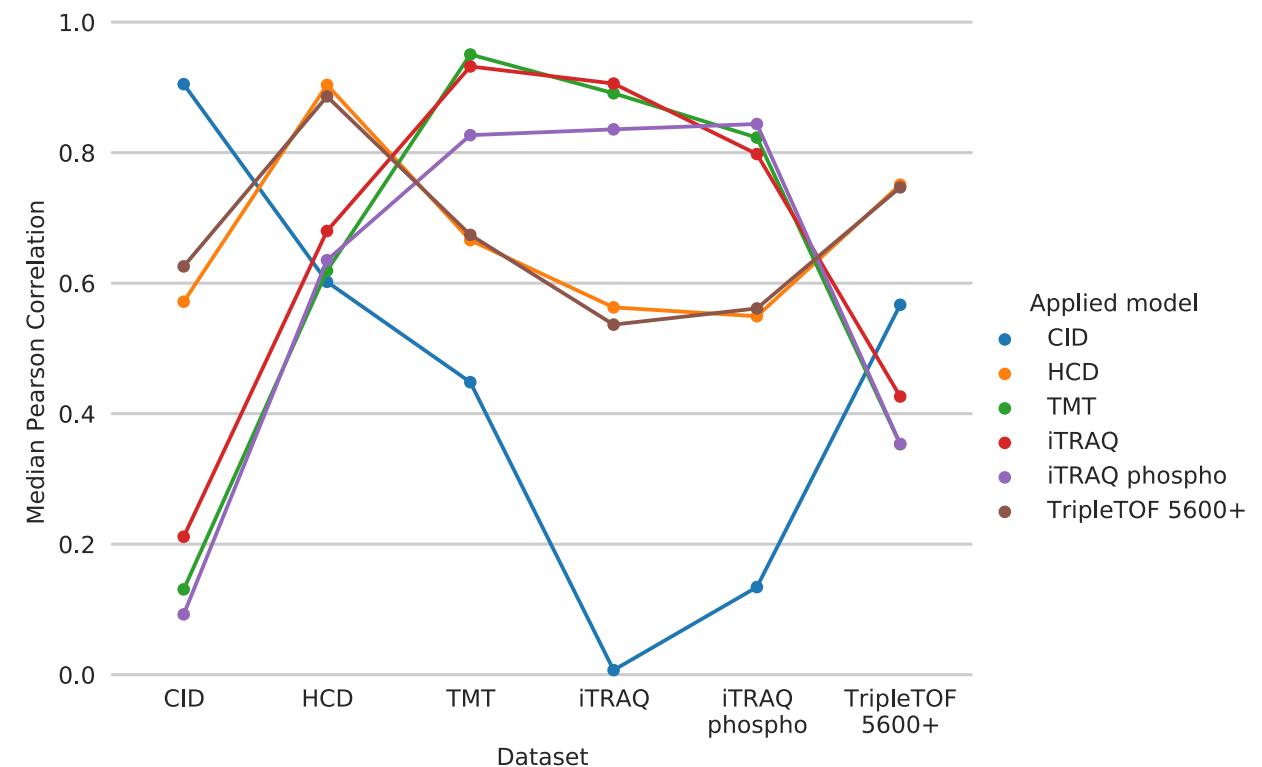
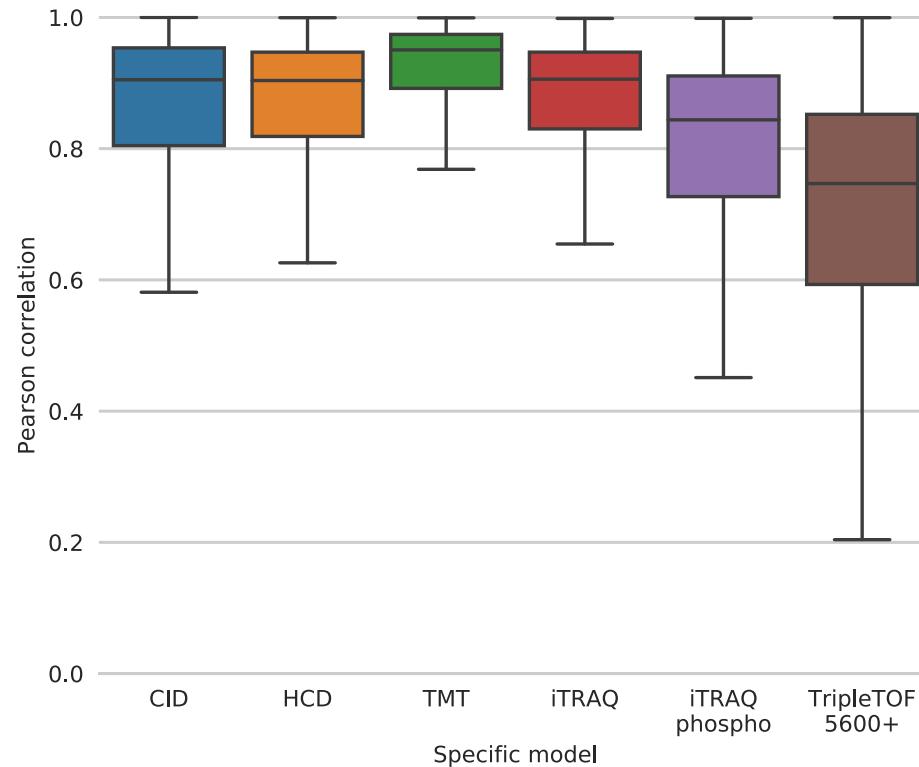
I trained new models for multiple fragmentation methods, instruments and labelling techniques



As expected, training data-specific models substantially improves the predictions

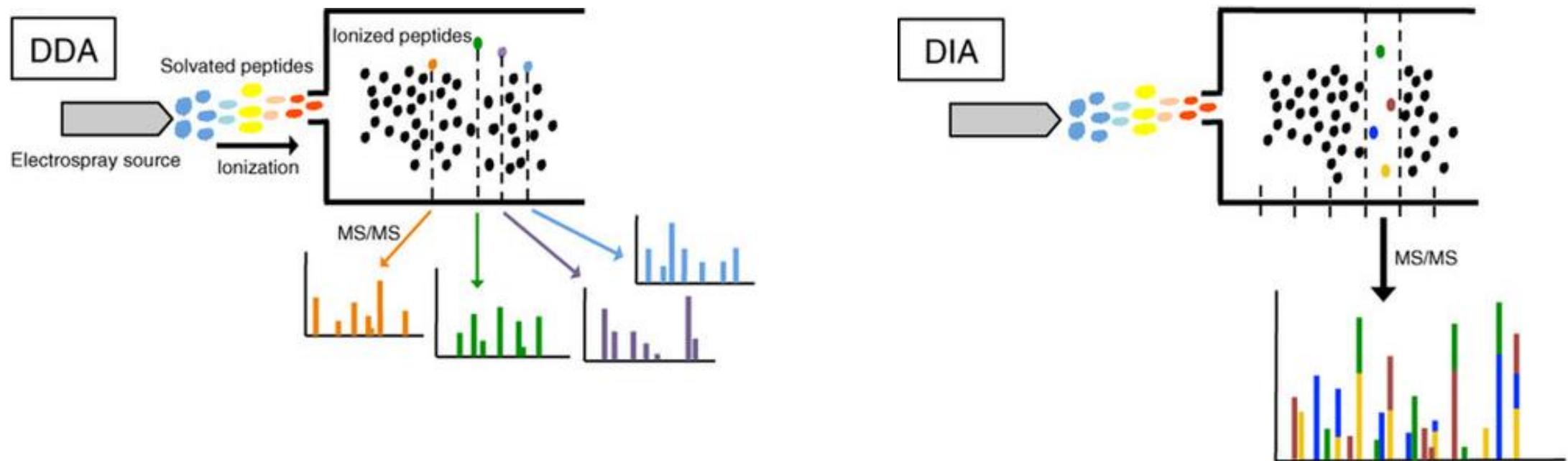


Interestingly, similar fragmentation effects across setups are highlighted by the model performances

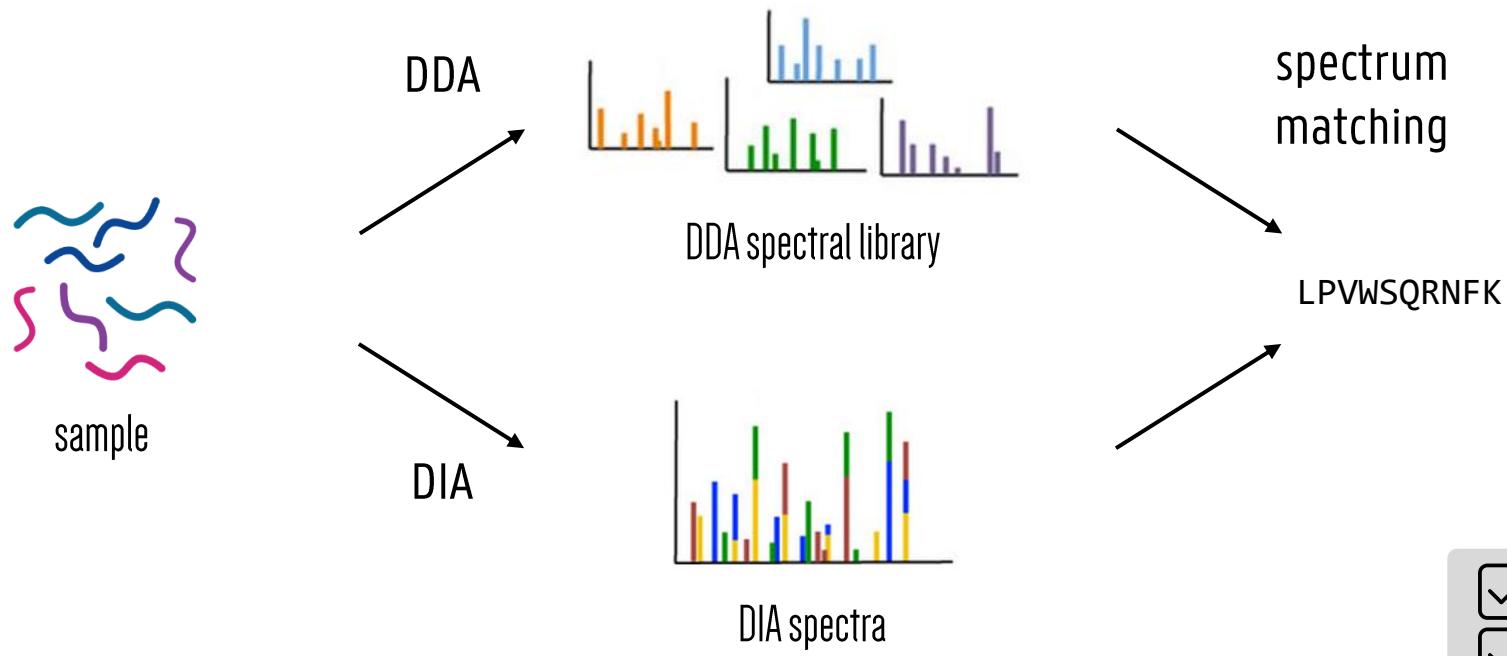


Removing the hidden data dependency of DIA with predicted spectral libraries

DIA uses wide isolation windows, resulting in complex chimeric spectra

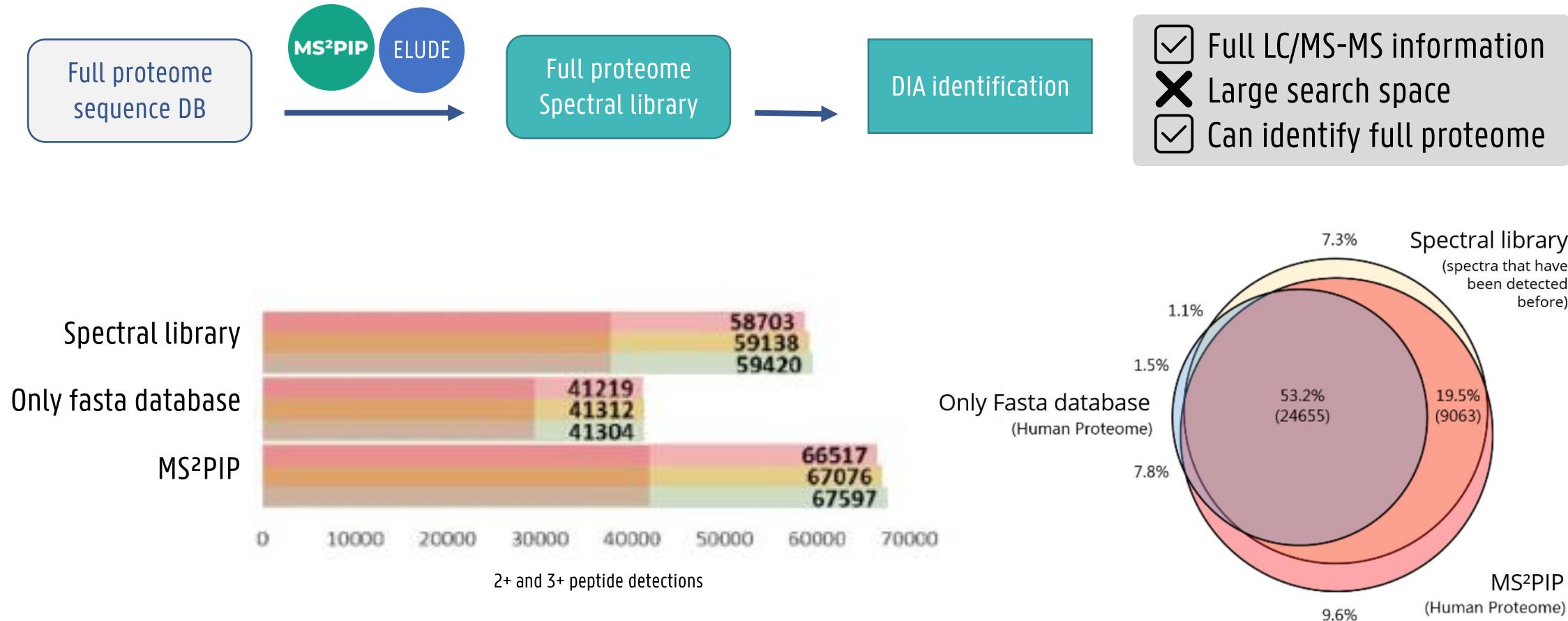


Therefore, DDA spectral libraries are often used for spectrum identification



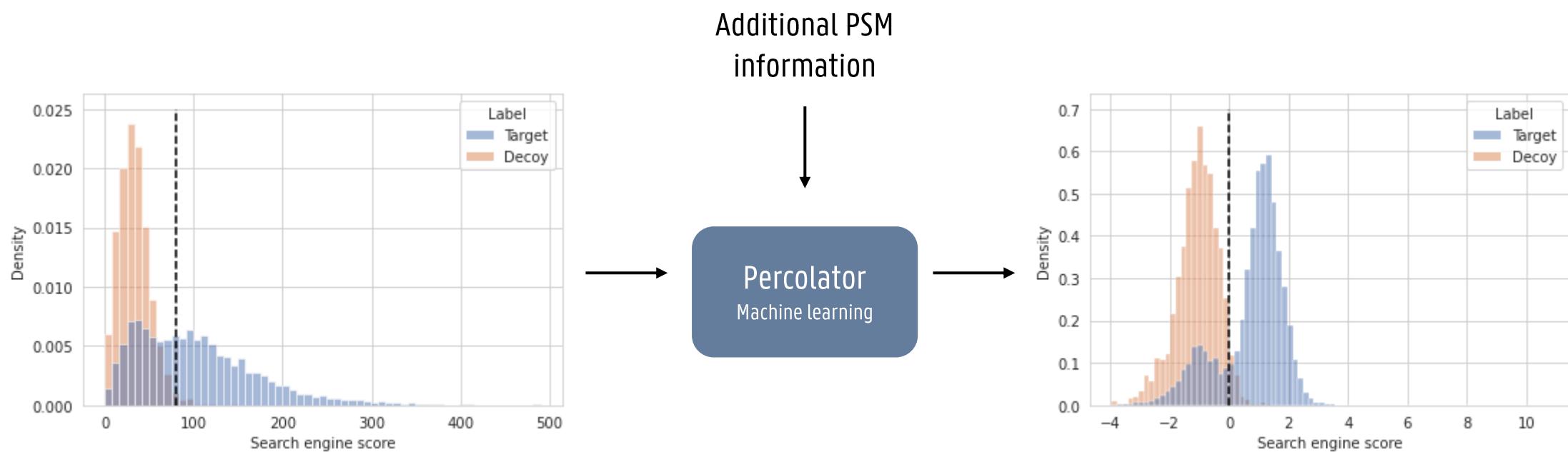
- Full LC-MS/MS information
- Manageable search space
- Limited to what DDA can identify

In silico predicted spectral libraries can free DIA from this DDA data dependency

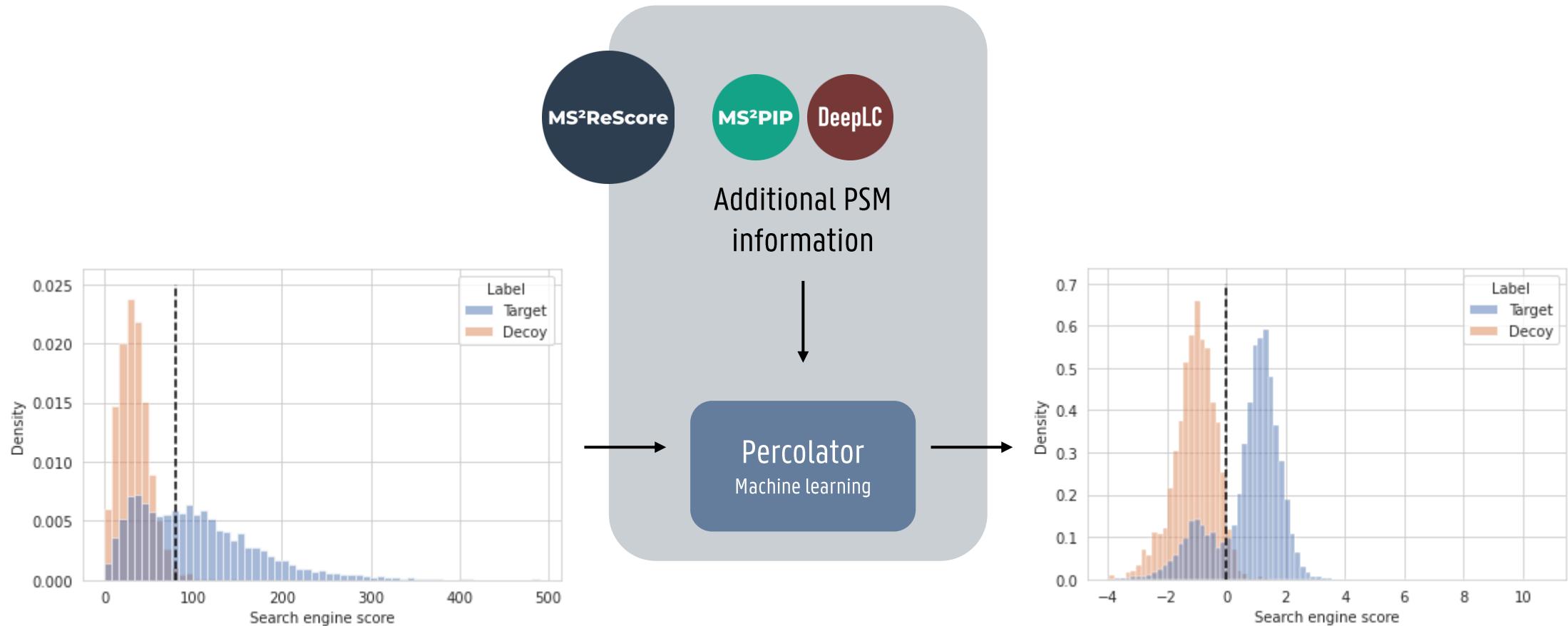


MS²Rescore: Leveraging spectrum predictions to enable novel proteomics workflows

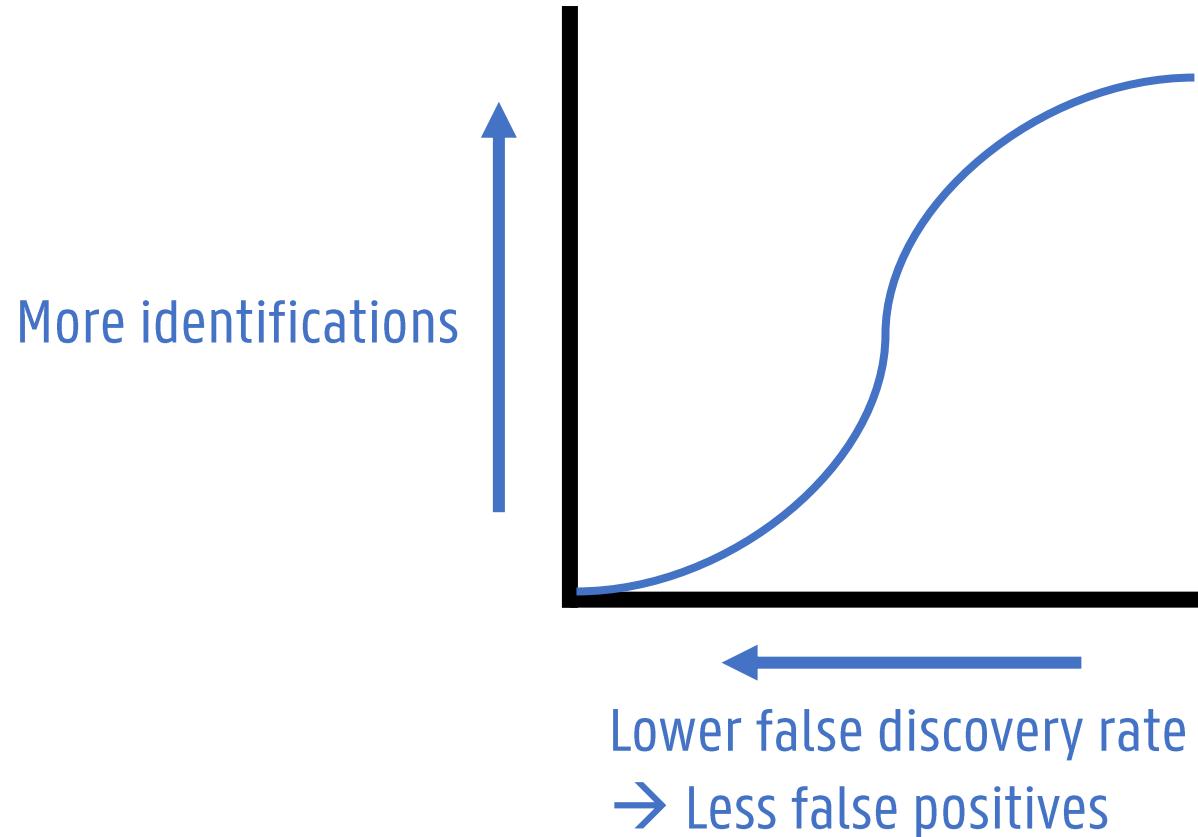
Percolator combines different peptide match information sources and learns how to separate true from false identifications



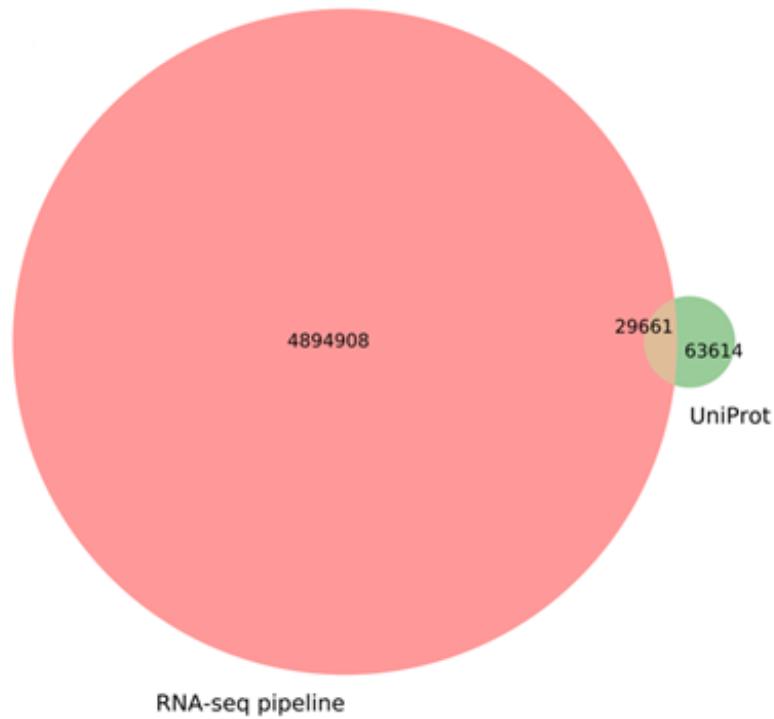
Machine learning-based information can be passed to Percolator for improved rescoring



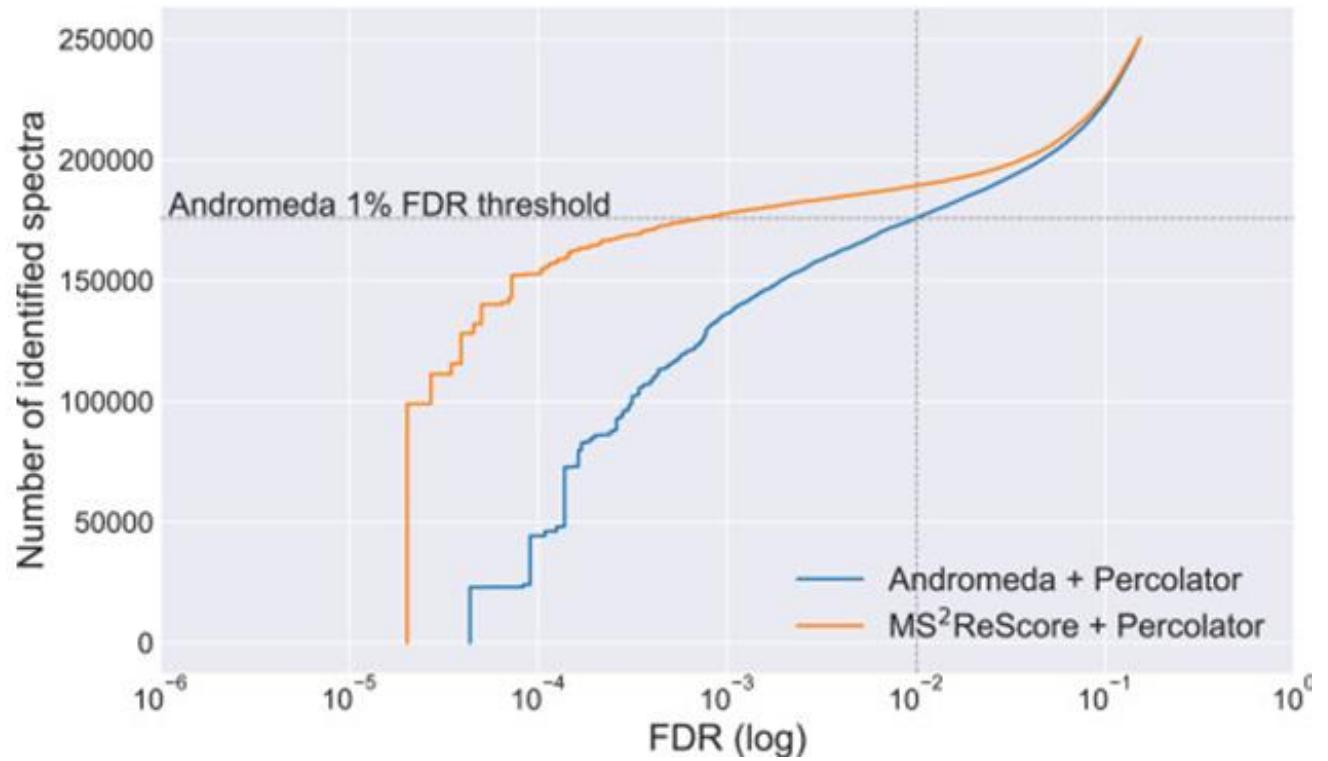
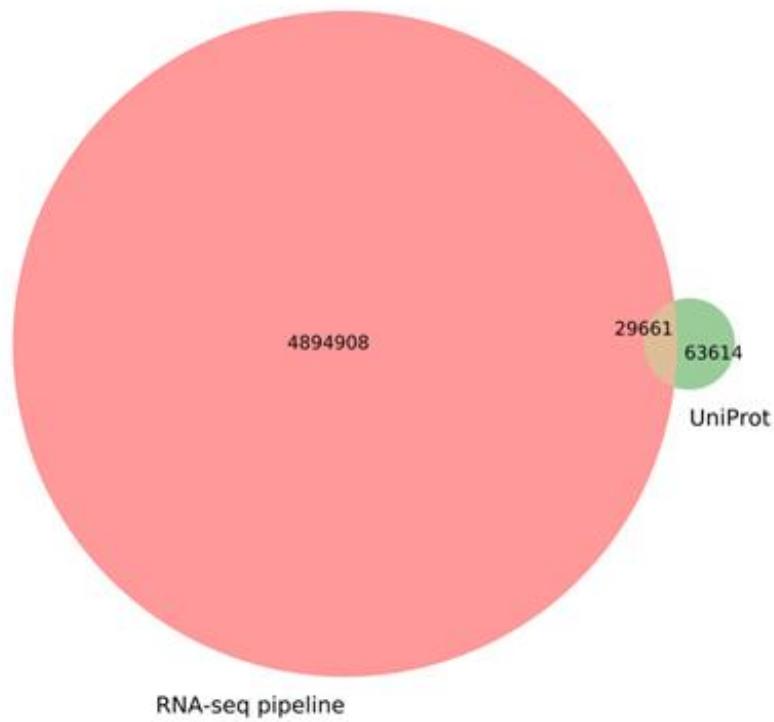
Explaining FDR/identification rate plots



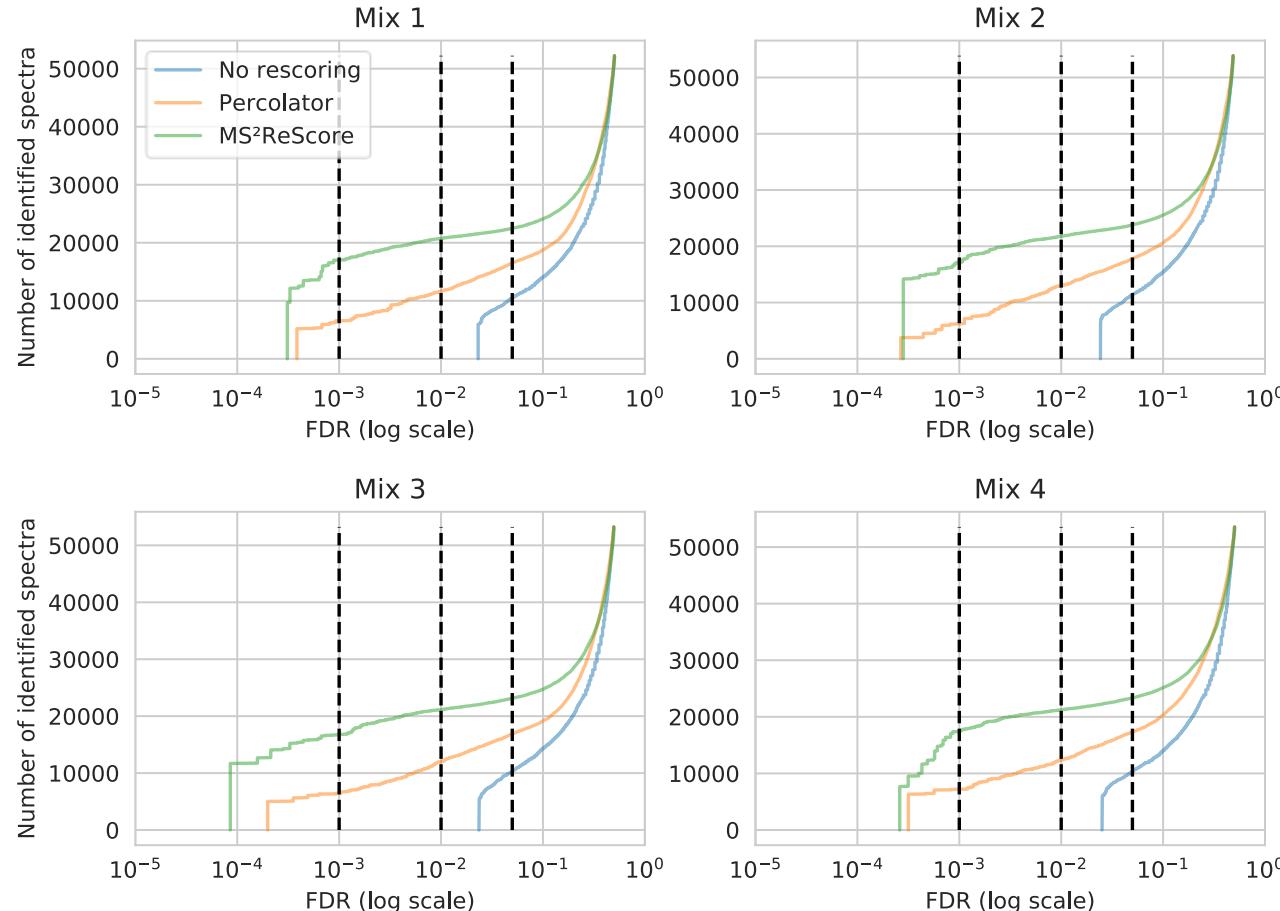
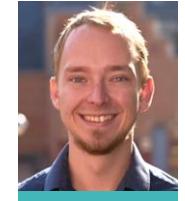
MS²Rescore in proteogenomics: Higher ID rate, at a 10-fold lower FDR threshold



MS²Rescore in proteogenomics: Higher ID rate, at a 10-fold lower FDR threshold



MS²Rescore in metaproteomics: From 0 to 20 000 identified spectra

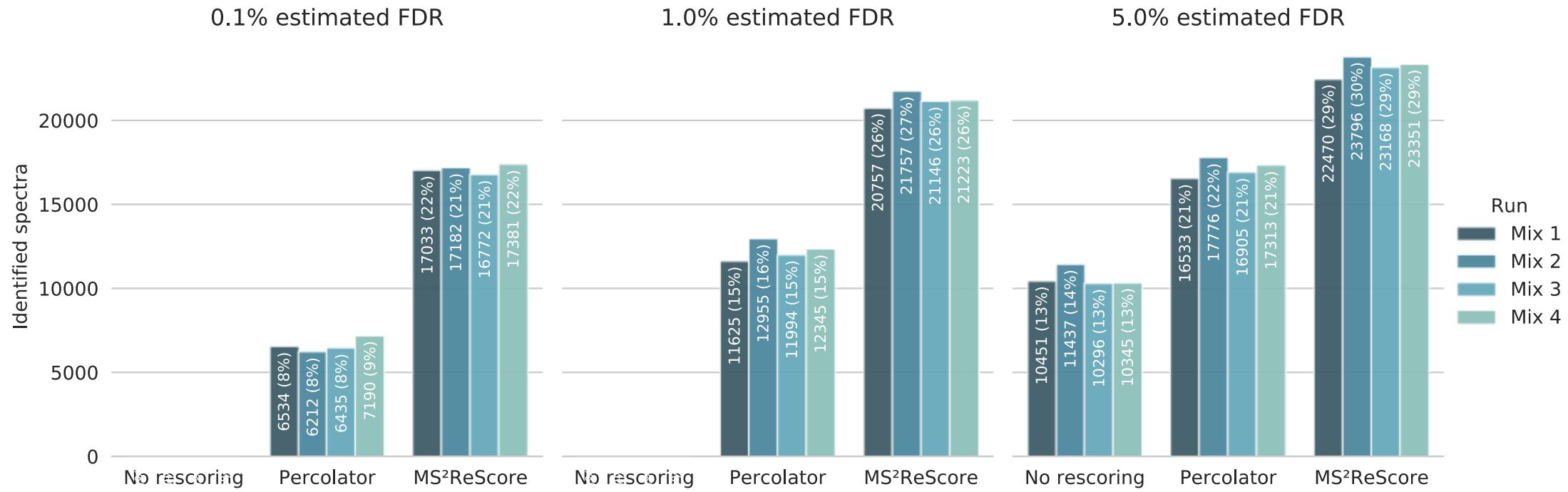


iPRG 2020 study:

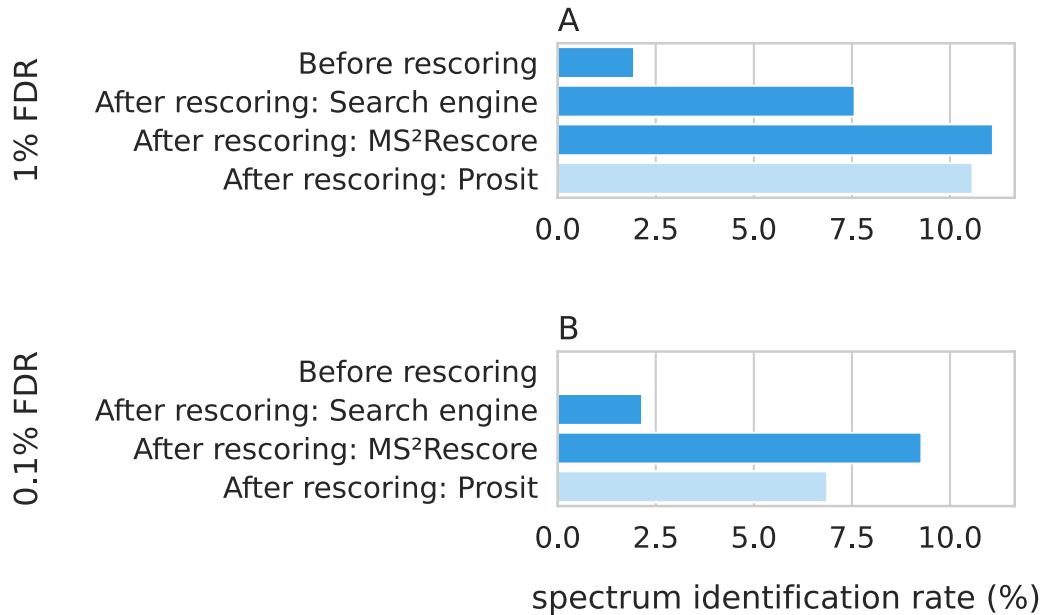
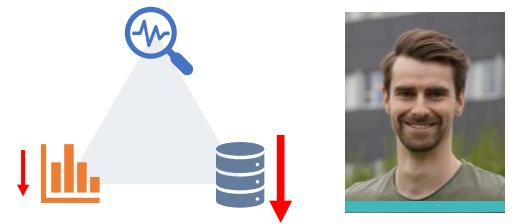
- 4 microbial mixes, with unknown sequences
- Extremely large search space
- X!Tandem
- CID ion trap acquisition



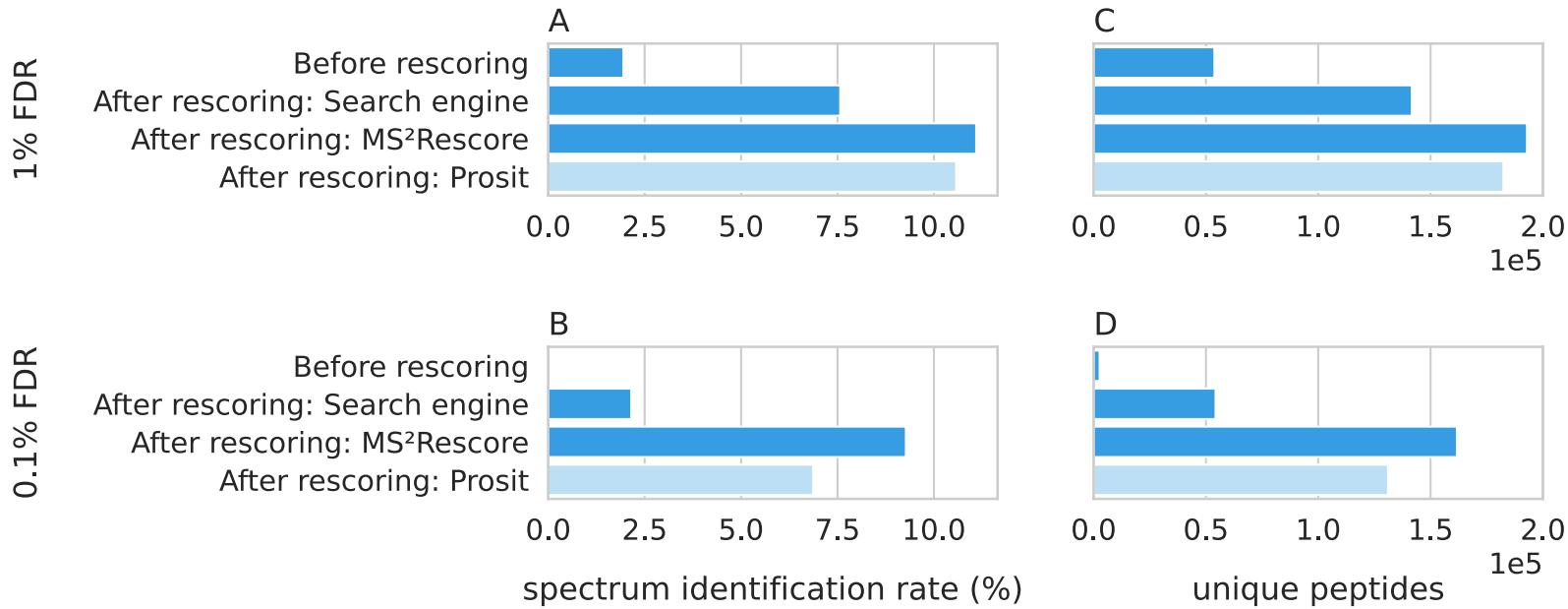
MS²Rescore in metaproteomics: From 0 to 20 000 identified spectra



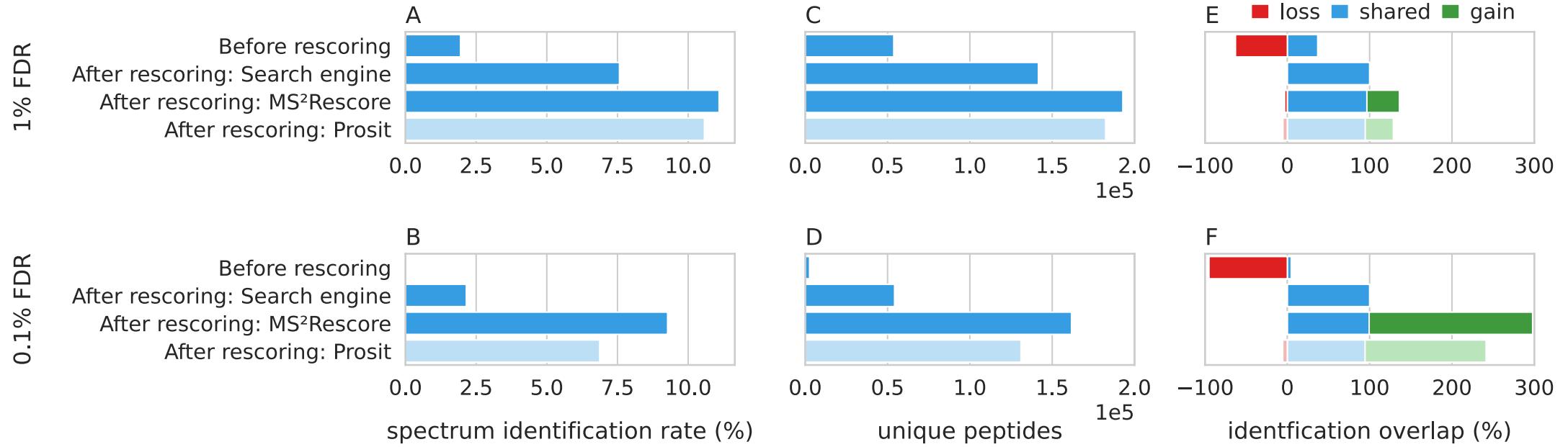
MS²Rescore in immunopeptidomics: +46% spectrum identification rate



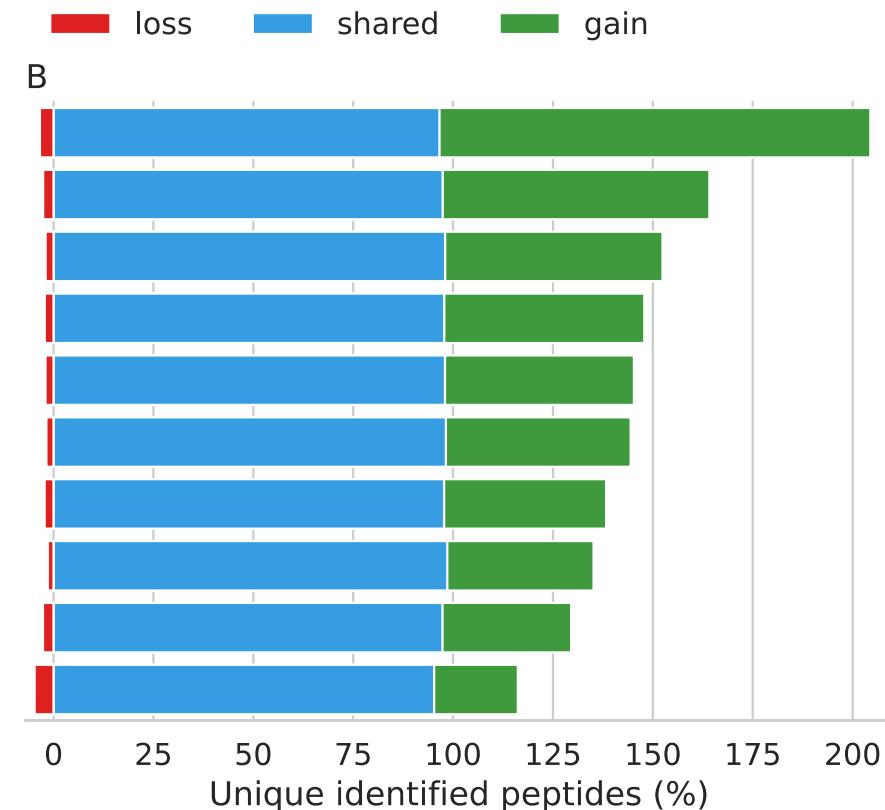
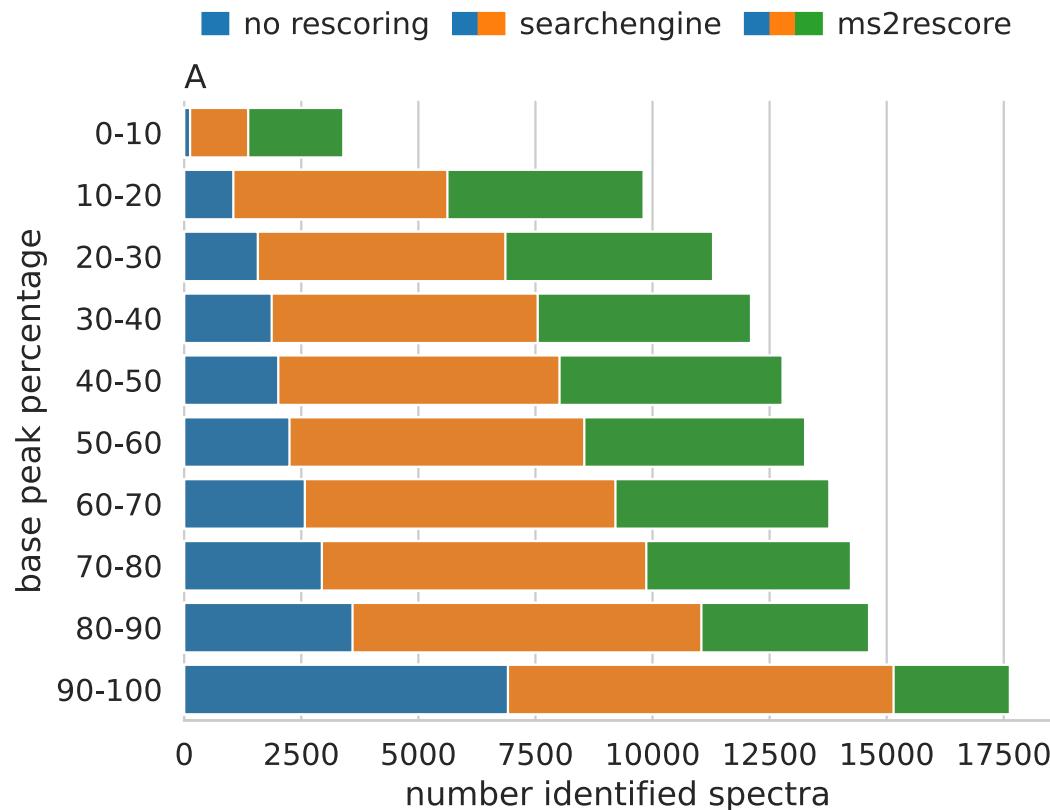
MS²Rescore in immunopeptidomics: +36% unique identified peptides



MS²Rescore in immunopeptidomics: +36% unique identified peptides

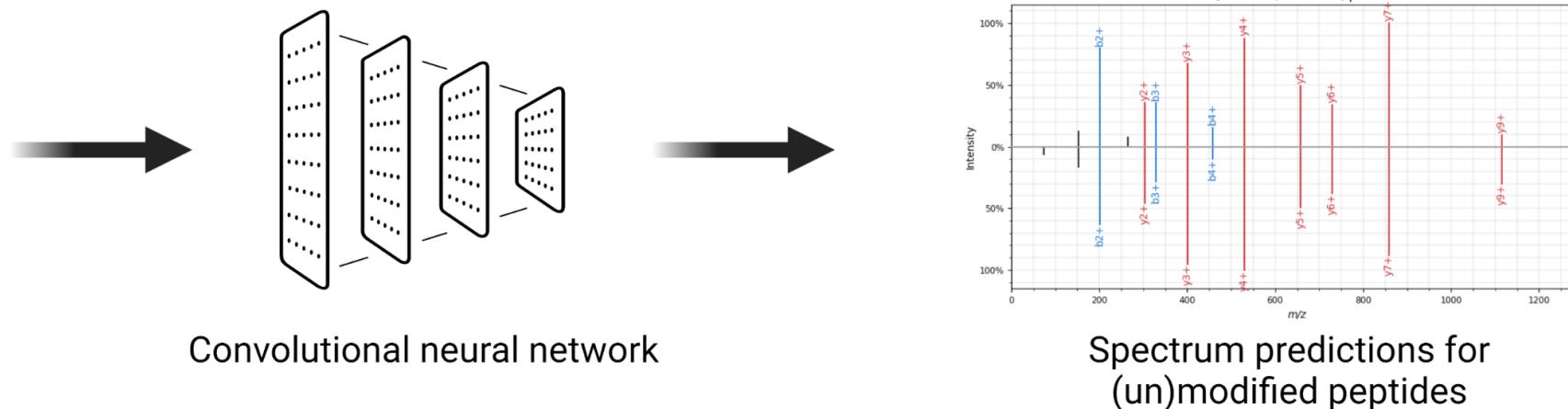
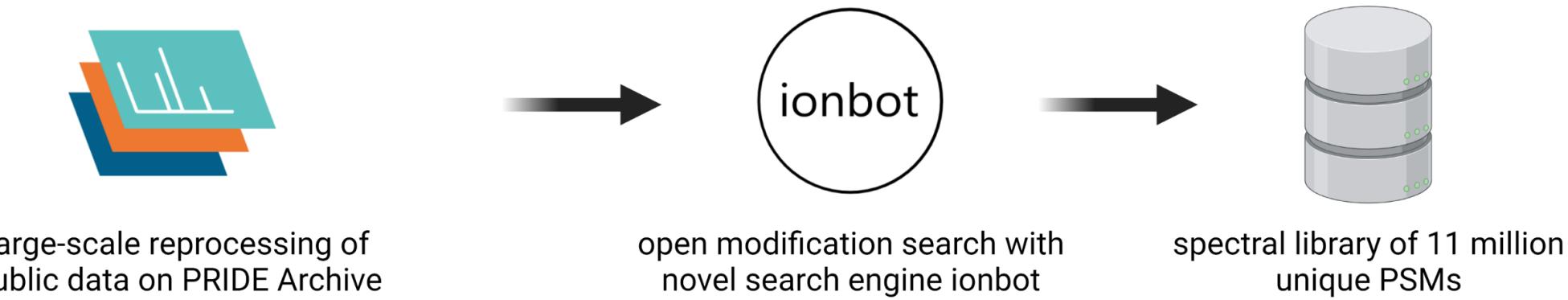


MS²Rescore in immunopeptidomics: +46% identification rate

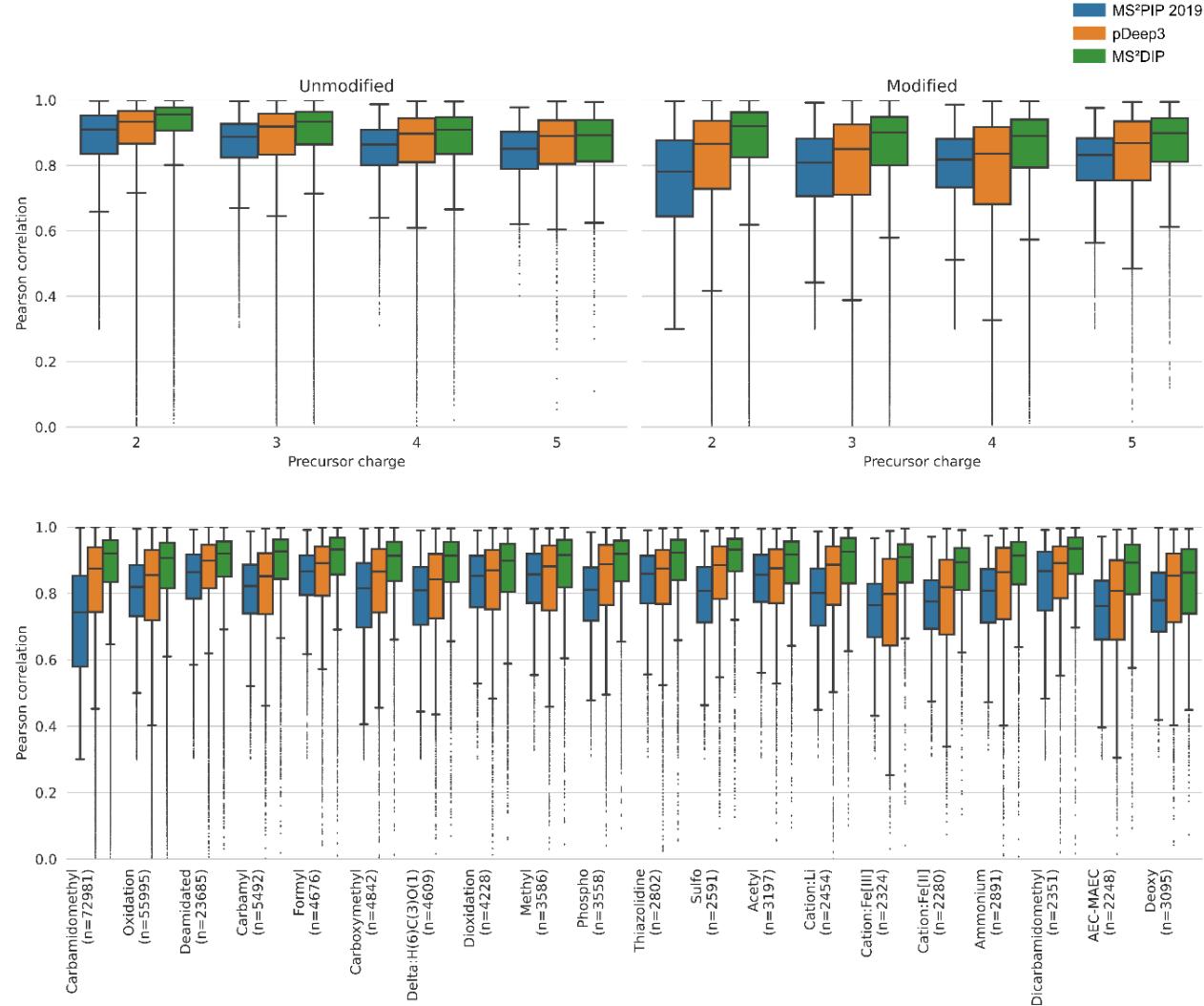


MS²DIP: **spectrum prediction for modified peptides**

Deep learning allows a generalization across unmodified and modified peptides



Preliminary results already outperform MS²PIP and out-of-the-box pDeep3



Introduction (Dutch)
proteins and proteomics
mass spectrometry
machine learning

Results (English)

MS²PIP: Peptide spectrum prediction for multiple fragmentation methods,
instruments, and labeling techniques

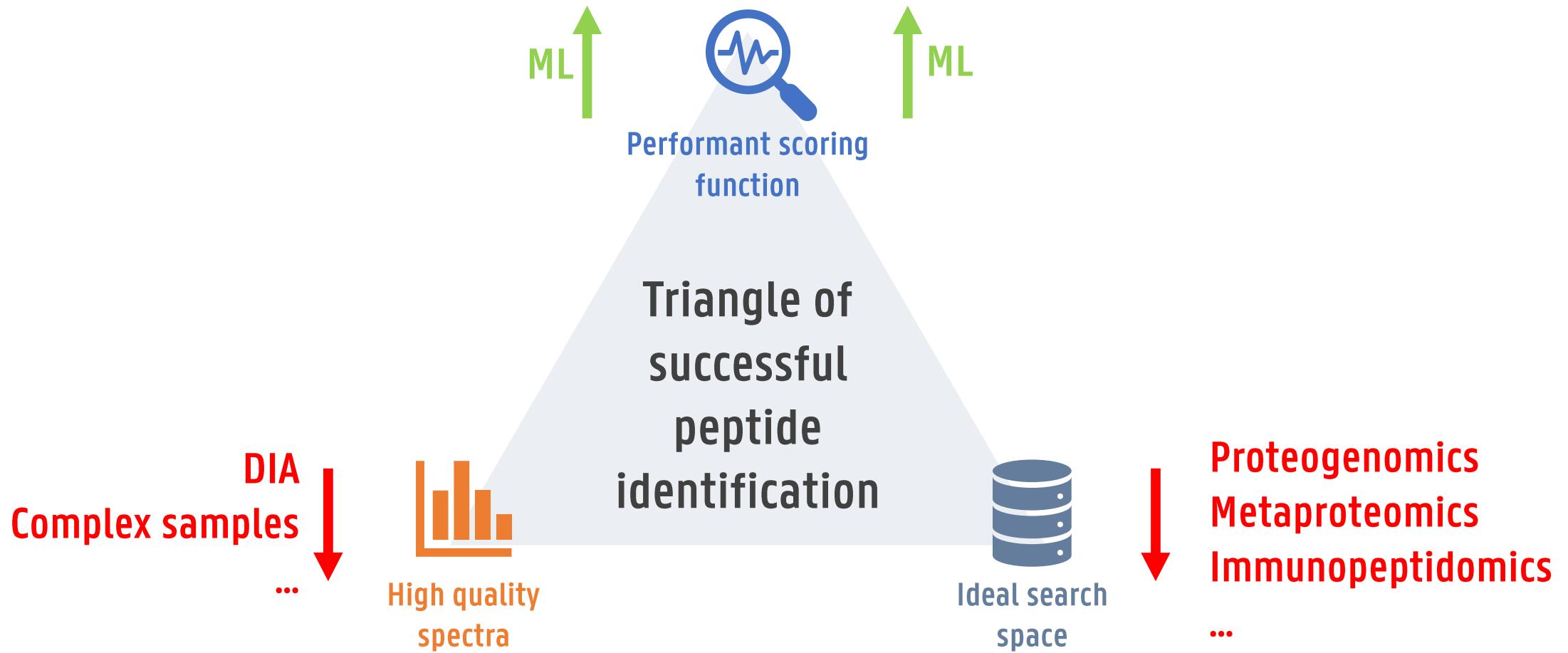
Removing the hidden data dependency of DIA with predicted spectral libraries

MS²Rescore: Leveraging spectrum predictions to enable novel proteomics workflows

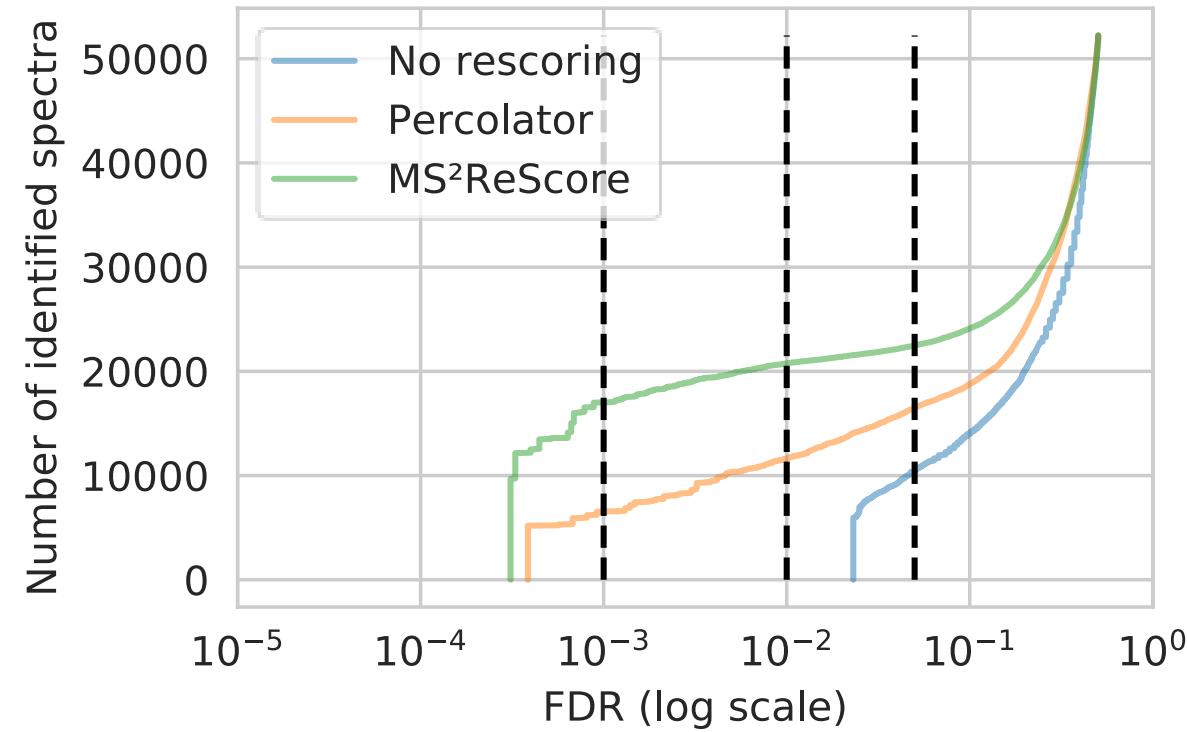
MS²DIP: spectrum prediction for modified peptides

Discussion and future perspectives (English)

In challenging experimental workflows, machine learning can rescue peptide identification rates



More sensitive scoring methods also allow for more specificity: 0.1% FDR to replace 1% FDR?



My “ML in proteomics wish list” for the (near) future

Better metadata in public data

Use of standardized file formats

More modularity and interoperability in bioinformatics software

Modularity and sharing of deep learning networks

High quality, community-curated training data

**The future of MS-based proteomics lies in the discovery of the unexpected
requiring more sensitive methods
where machine learning will take a central role**

MACHINE LEARNING TO THE RESCUE: ENABLING NOVEL PROTEOMICS WORKFLOWS WITH DATA-DRIVEN BIOINFORMATICS METHODS

Ralf Gabriels

Dissertation submitted to obtain the degree Doctor in Health Sciences
Academic year 2021-2022

