**RUHR-UNIVERSITÄT** BOCHUM

# USING PROTEIN-GRAPHS TO GENERATE FASTA-DATABASES WITH VARIATIONALLY PEPTIDES

IOW: HOW TO ENABLE SEARCH-ENGINES TO SEARCH FOR VARIATIONAL PEPTIDES?
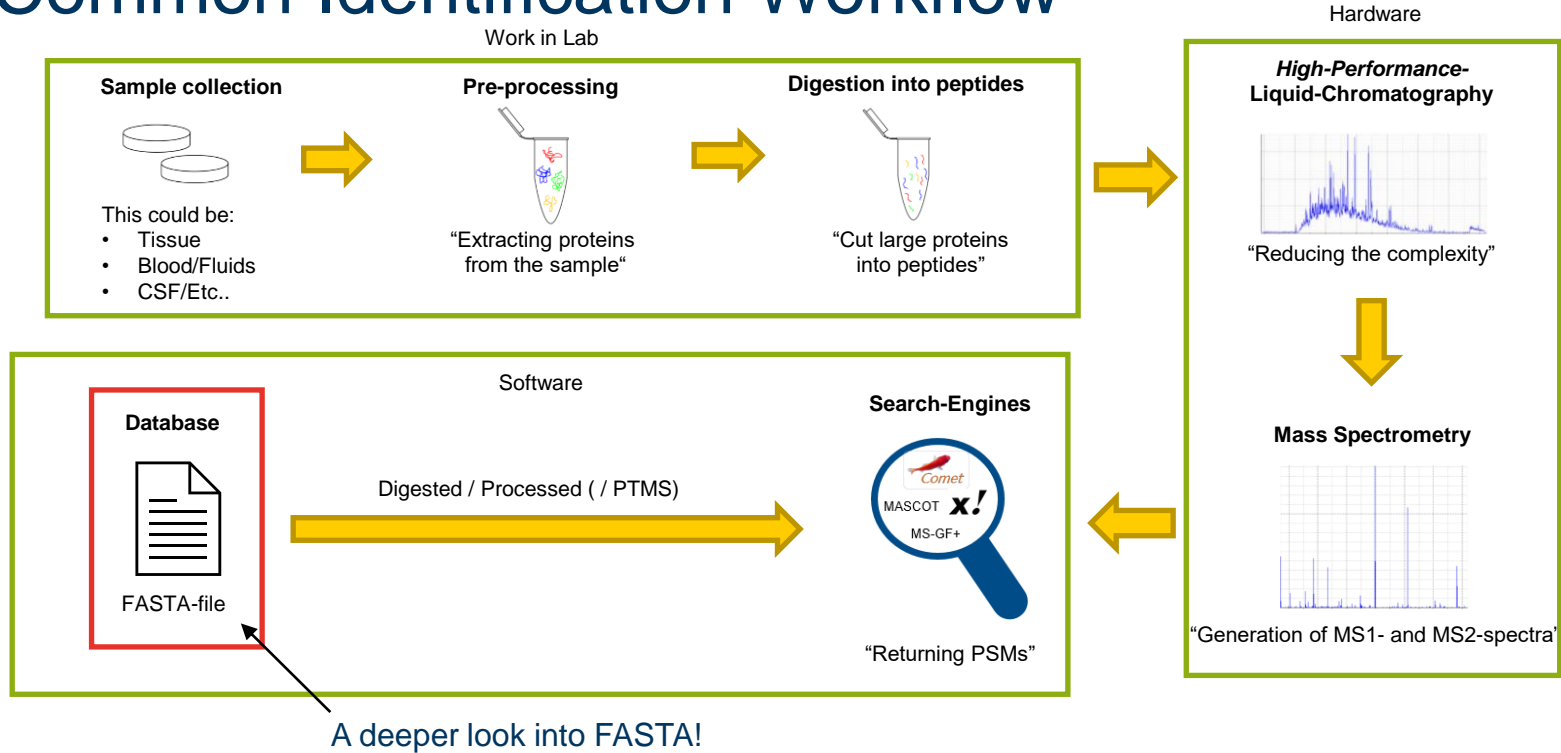
Dominik Lux

This is a graph!

# Common Identification Workflow

Work in Lab

Hardware

**Sample collection**

This could be:
- Tissue
- Blood/Fluids
- CSF/Etc..

**Pre-processing**

"Extracting proteins from the sample"

**Digestion into peptides**

"Cut large proteins into peptides"

***High-Performance-Liquid-Chromatography***

"Reducing the complexity"

**Mass Spectrometry**

"Generation of MS1- and MS2-spectra"

Software

**Database**

FASTA-file

Digested / Processed ( / PTMS)

**Search-Engines**

Comet
MASCOT **X!**
MS-GF+

"Returning PSMs"

A deeper look into FASTA!

**mpc**
MEDIZINISCHESPROTEOMCENTER

**RUHR UNIVERSITÄT BOCHUM**

**RU**B

# FASTA-files and their potential!

**Database**

FASTA-file

- Usually contains only canonical sequences (sometimes isoforms, cRAP)
- "Used unprocessed" by search engines

```
...
>sp|ACCESSION|GENE_NAME  Protein XYZ Some protein functions...
MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGP
DEAPRMPEAAPPVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAK
SVTCTYSPALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHE
RCSDSDGLAPPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCTTIHYNYMCNS
SCMGGMNRRPILTIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEEENLRKKGEPHHELP
PGSTKRALPNNTSSSPQPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPG
GSRAHSSHLKSKKGQSTSRHKKLMFKTEGPDSD
...
...
...
```
Example-Entry in FASTA

→ Large potential:

  → Include additional information (e.g., variants)

  → Organize by precise header-information

  → Enable search engines to search, which was not possible before

**UniProt**

Unique/Shared-entries (by proteins)

- Search with "infinite many" miscleavages
- 2 (3 or more) digestion enzymes at once
- ....

→ Sophisticated FASTA-generator would be interesting!

**mpc** MEDIZINISCHES**PROTEOM**CENTER

**RUHR UNIVERSITÄT BOCHUM**

**RUB**

# Parsing and Encoding feature Information

→ Large potential:
→ Include additional information (e.g., variants)
→ Organize by precise header-information
→ Enable search engines to search, which was not possible before



**ProtGraph**

Available on:
- BioConda*
- PyPI*
- GitHub*

Python (CLI)

SP-EMBL (UniProt, species)

Contains:
- Canonical- (Isoform-) sequence
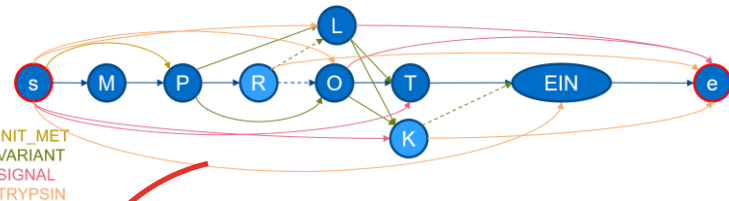- Variants  (/Mutagens/Conflicts)
- Peptides (Pro- Signalpeptides)
- …

Produces many protein-graphs:
- Encodes feature information
- Allows (multiple) digestion(s)
- Allows to add PTMs
- …

INIT_MET
VARIANT
SIGNAL
TRYPSIN

1 protein-graph per entry
(this example contains 46 peptides)

Protein-Graph
- Feature-Information represented in graph-format
- Only contains valid paths from s to e
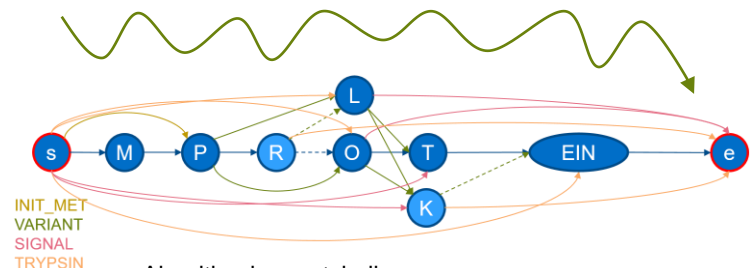- Very compact representation

**Database**

FASTA-file

## Next Step:
## → Convert Protein-Graphs into FASTA-entries

* conda install -c bioconda protgraph
* pip install protgraph
* https://github.com/mpc-bioinformatics/ProtGraph

Using Protein-graphs to generate FASTA-Databases with variationally peptides
EuBIC September 2022 | Dominik Lux

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Naïve Approach: Depth-First-Search



→ Large potential:
  → Include additional information (e.g., variants)
  → Organize by precise header-information
  → Enable search engines to search, which was not possible before

INIT_MET
VARIANT
SIGNAL
TRYPSIN

Algorithm in a nutshell:
- Begin at s
- Report every path to e
- Constrain by specific thresholds
  - Peptide-/Path-Mass
  - #Variants
  - #Miscleavages
  - …

**(peptide-) Database**

...
>pg|ACCESSION|Protein ...
RPILTIITLEDSSGNLLGR
>pg| ACCESSION |Protein ...
TCPVQLWVDSTPPPGTR
>pg| ACCESSION |Protein ...
CSDSDGLAPPQHLIR
>pg| ACCESSION |Protein ...
KPLDGEYFTLQIR
>pg| ACCESSION |Protein ...
GEPHHELPPGSTK
...

FASTA-file

A shared peptide between Protein A and B!

...
>pg|ACCESSION|**ProteinA(100:119, mssclvg:1), ProteinB(90:109, mssclvg:1)**
RPILTIITLEDSSGNLLGR

>pg|ACCESSION|**ProteinA(55:72, mssclvg:0, VARIANT[56:56, L->C]**
TCPVQLWVDSTPPPGTR
...

A unique peptide of Protein A with a variant

→ Organize headers by traversed path (concatenate same FASTA-entries)

→ Search-Engines would need to search peptide-FASTAs as is (peptidomics)
  → IOW: "Digestion turned off"

Using Protein-graphs to generate FASTA-Databases with variationally peptides
EuBIC September 2022 | Dominik Lux

**RUHR
UNIVERSITÄT
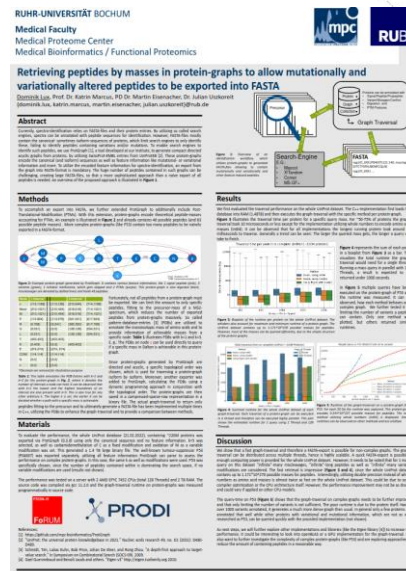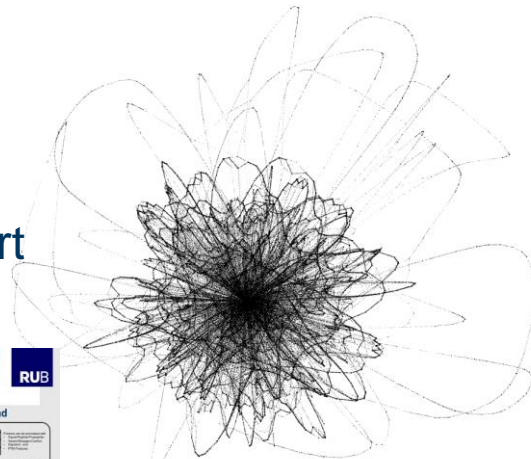BOCHUM**

**RU**B

# Drawbacks of the Naïve Approach

→ "Some" protein-graphs are too complex, for a naïve export
    → Too complex to save on disk (even in large scale!)

Solution:

- Search Engines do not need all entries
    - Export only peptides, fitting to the MS2-precursor

No details*,
just an overview how it was "solved"

\* This is not a secret! If you are interested how it was "solved", please ask!

Example: P04637 (P53_Human)
- 1363 Variants
- 36 Mutagens
- 9 Isoforms
- 13988 Nodes, 93924 Edges

- Encodes 1.7E+224 peptides (computed via ProtGraph in a few seconds)

Poster presented at
Proteomic Forum 2022

# Sophisticated Approach: Target-Value-Search

Problem:        Return all Paths, where the mass of the peptide equals to the MS2-precursor
Target-Value:   MS2-Precursor (+- Xppm)



INIT_MET
VARIANT
SIGNAL
TRYPSIN

Algorithm (C++) in a nutshell:
- Begin at s
- Use so-called PDBs* to traverse only through possible solution paths**
- Reach target e and report a peptide fitting the MS2-Precursor
- Repeat until all solutions are reported
- Constrain by:
  - #Variants

→ Executable on all proteins! (whole UniProtKB, ~230M Proteins)

→ Only a few Proteins, have long running time during traversal***

* Pattern databases, Schmidt et al, A depth first search approach to target value search
** Early expansion-prevention in branches with no solutions
*** Guess: ~200 proteins in the whole UniProtKB

Using Protein-graphs to generate FASTA-Databases with variationally peptides
EuBIC September 2022 | Dominik Lux

# Generation of peptide-FASTA-Databases

ISA*-RAW-files containing 28123 distinct MS2-Precursors (+/-5ppm, Oxidation M (variable), Carbamidomethylation C (fixed))

| | E.Coli | Mus Musculus | ~1000 Species (excluding homo sapiens) | Any other species |
|---|---|---|---|---|
| Organism: | | | | |
| # Proteins: | 4448 | 55319 | 11 973 189 | |
| Restrictions: | None | None | None | (probably feasible) |
| **(peptide-)FASTA** Generation** | ⬇ | ⬇ | ⬇ | (if no very complex protein-graphs are present) |
| **(peptide-) Database** # Entries | 66 680 808 (peptides) (21 GB) | 85 702 533 (peptides) (46 GB) | 126 462 579 (peptides) (27 GB) | |
| Generation time: | 1 h 30m | 3 h 54m | 1d 14h 4m | |

→ Generation of a MS2-specific-peptide FASTA with features***
  is feasible

\* Internal Standard (for benchmarking)
\*\* Generated on a server with 64 threads
\*\*\* All features ProtGraph can parse

Using Protein-graphs to generate FASTA-Databases with variationally peptides
EuBIC September 2022 | Dominik Lux

RUHR
UNIVERSITÄT
BOCHUM

RUB

# What about Homo Sapiens?

→ Most well researched species (→ most annotated proteins)

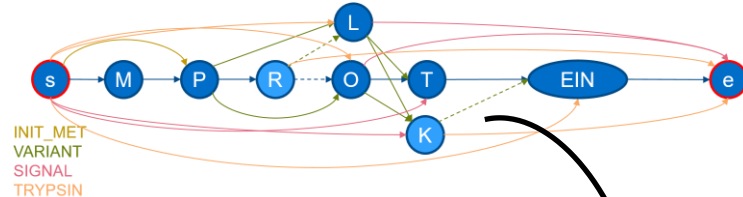| | |
|---|---|
| Complex protein: | P04637 (P53_Human) |
| Benchmark/Target-Value: | Every 50 Dalton |
| Configuration: | All Features, Digested, PTMs: Carbamidomethylation of C (fixed) Oxidation of M (variable) |

| | |
|---|---|
| Complex protein: | P68871 (HBB_Human) |
| Benchmark/Target-Value: | Every 50 Dalton |
| Configuration: | All Features, Digested, PTMs: Carbamidomethylation of C (fixed) Oxidation of M (variable) |

Using Protein-graphs to generate FASTA-Databases with variationally peptides
EuBIC September 2022 | Dominik Lux

**RUHR UNIVERSITÄT BOCHUM**

**RUB**

# Summary



- Generation of protein-graphs via ProtGraph
- Protein-graphs can contain up to infinite
  - Miscleavages
  - Features (Variants/Peptide/Signal-/Propeptide/Mutagens …)
  - Digestion enzymes

  → Can be exported into a peptide-FASTA-database
  → Search-Engines benefit from it!
    → see above
  → Easier interpretation with the new header-format

- Unique/shared PSMs
- (trivial) Inference
- More precise information (how the peptide was generated)

Using Protein-graphs to generate FASTA-Databases with variationally peptides
EuBIC September 2022 | Dominik Lux

**RUHR UNIVERSITÄT BOCHUM**  **RU**B

# Acknowledgement

**And Thank you for your attention!**

@ProteomCenter
@lululuxii

mpc-bioinformatics
Luxxii



(MPC-Bioinformatics-Team)
(We need to update this!)

s with variationally peptides

EuBIC September 2022 | Dominik Lux