# Welcome
# to the DKFZ!

**dkfz.** GERMAN CANCER RESEARCH CENTER IN THE HELMHOLTZ ASSOCIATION

Research for a Life without Cancer

# News from PCF at DKFZ

Martin Schneider

# Topics from DKFZ

- missRanger

- AlphaPept on Linux

- New Software how do you test?

- The field needs to talk about FDR

- Additional thoughts

dkfz.

# missRanger

- ## missForest

  Stekhoven, D.J. and Buehlmann, P. (2012),
  'MissForest - nonparametric missing value imputation for mixed-type data',
  Bioinformatics, 28(1) 2012, 112-118, doi: 10.1093/bioinformatics/btr597

- ## Statistics benchmarks

  https://doi.org/10.29220/CSAM.2023.30.3.331

  ### A comparison of imputation methods using machine learning models

  Communications for Statistical Applications and Methods 2023;30:331-341
  Published online May 31, 2023
  © 2023 Korean Statistical Society.

  We can see that missForest takes much longer to impute missing values than missRanger, while there is not much difference in performances.

dkfz.

# missRanger

- Results constantly differ in between missRanger and missForest

- missRanger based on ground truth performs worse

- missForest now used with only 20 trees performs equivalent and is much faster

**dkfz.**

# AlphaPept on Linux

- DKFZ requires local QC solution

- ShinyQC with MaxQuant ~ 60 min

- ShinyQC 2.0 with AlphaPept

  - Fast ~ 8 min
  - Fully open source
  - Fails to run on Linux
  - Some MS level information missing

dkfz.

# AlphaPept on Linux

- Fails to run on Linux
- Docker for ARM from scratch for AMD64
- Not enough disk space

```
2024-04-02 09:47:27> Size check:
2024-04-02 09:47:27> Size of job (raw files) 0.97 Gb
2024-04-02 09:47:27> Required disk space for / - 0.97 Gb, Available 0.02 Gb.
2024-04-02 09:47:27> Not enough disk space for analysis. Please free disk space.

An exception occurred running AlphaPept version 0.5.2:
```

```python
232  def check_size(settings):
233      sizes = [get_size(_) / 1024 ** 3 for _ in settings['experiment']['file_paths']]
234      base_dirs = [os.path.splitdrive(os.path.abspath(_))[0] for _ in settings['experiment']['file_paths']]
235
236      size_gb = sum(sizes)
237      logging.info(f'Size of job (raw files) {size_gb:.2f} Gb')
238
239      required_size_dict = {}
240
241      for base, size in zip(base_dirs, sizes):
242          if base == "":
243              base = "/"
244          if base in required_size_dict:
245              required_size_dict[base] += size
246          else:
247              required_size_dict[base] = size
248
249          #Require at least file size of raw files as disk space (file conversion, search etc.)
250      for base, size in required_size_dict.items():
251          free = psutil.disk_usage(base).free/1024**3
252          if free < size:
253              logging.info(f'Required disk space for {base} - {size:.2f} Gb, Available {free:.2f} Gb.')
254              logging.info('Not enough disk space for analysis. Please free disk space.')
255              raise
256          else:
257              logging.info(f'Required disk space for {base} - {size:.2f} Gb, Available {free:.2f} Gb OK.')
258
259      logging.info("")
```

dkfz.

# AlphaPept on Linux

- Fix check_size

```
234     base_dirs = [os.path.splitdrive(os.path.abspath(_))[0] for _ in settings['experiment']['file_paths']]
234     base_dirs = [os.path.dirname(os.path.abspath(_)) for _ in settings['experiment']['file_paths']]
```
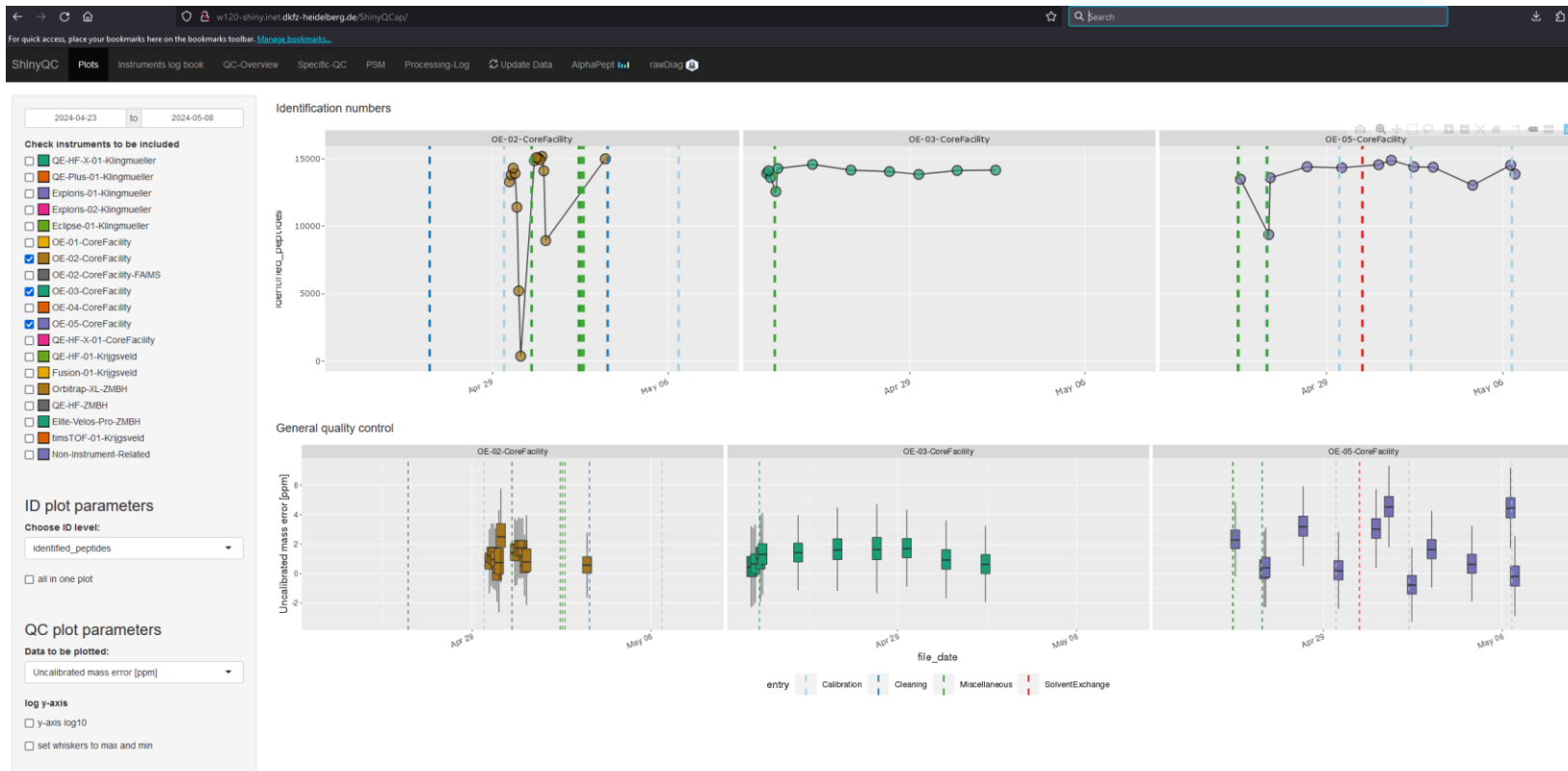
- Additional info on MS level



- Info on LC level
  - Thermo Raw File Reader
  - Own EXE for extracting LC info

dkfz.

# AlphaPept on Linux

- TIC, Injection Time, und AGC-Fill for MSMS
  - rawrr / rawDiag crashes
  - Table for all individual scans too large

- Thermo Raw File Reader
  - Own implementation to extract
  - Visual Studio 2022 to use required NuGet
  - Visual Studio Code did not get it to work
  - Got experience how to run it with local NuGet?

**dkfz.**

# AlphaPept on Linux

# New Software how do you test?

- DIA-NN 1.9.2

- Spectronaut 19.4.241104.62635

- MSAID 1.3.2 (CHIMERYS)

- AlphaDIA 1.8.2

- AlphaPept 0.5.0

- FragPipe 22.0

**dkfz.**

# The field needs to talk about FDR

- DIA-NN 1.9.2 → FDR control
  - …that ensures that each neural network in an ensemble is only used for prediction on samples it has not been trained on.
- Spectronaut default:
  - false discovery rate not in line with reality
- target-decoy; mathematical FDR?; machine learning FDR?

- FDR should match reality
- Which experiment, fading of ratios etc.?

**dkfz.**

## Additional thoughts

- Based on washes tests most abundant peptides are carried over

- Custom internal heavy standard?

  - QC for individual samples (different concentrations of synthetic peptides)

  - Distributed over RT

  - Normalization based on standard

Thank you
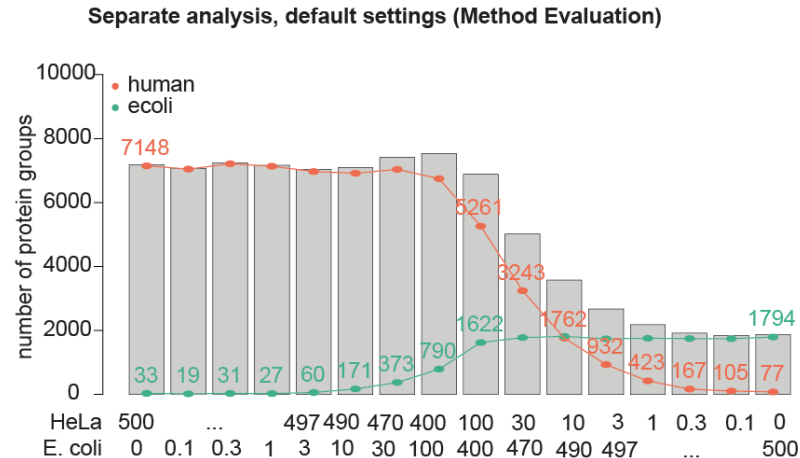for your attention!

Further information on www.dkfz.de

**dkfz.** GERMAN
CANCER RESEARCH CENTER
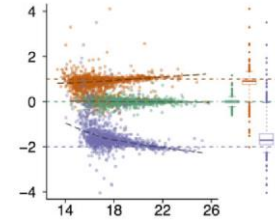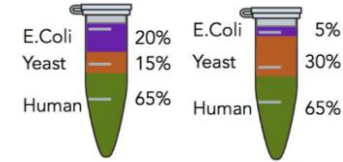IN THE HELMHOLTZ ASSOCIATION

Research for a Life without Cancer

# Optimization HYE, fade-in fade-out, thoughts

- What is carried over?

- How much is in the washes?



Separate analysis, default settings (Method Evaluation)

# Optimization HYE, fade-in fade-out, thoughts
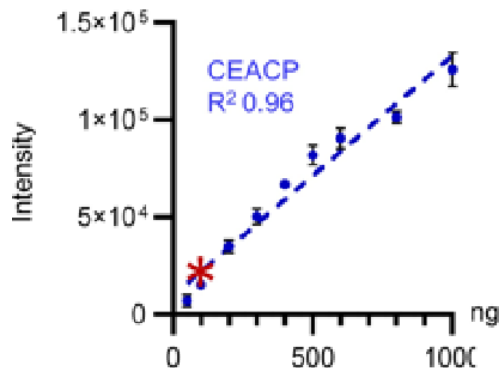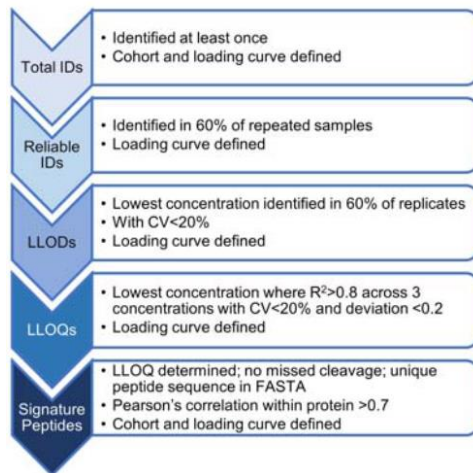
- HYE with dilution series
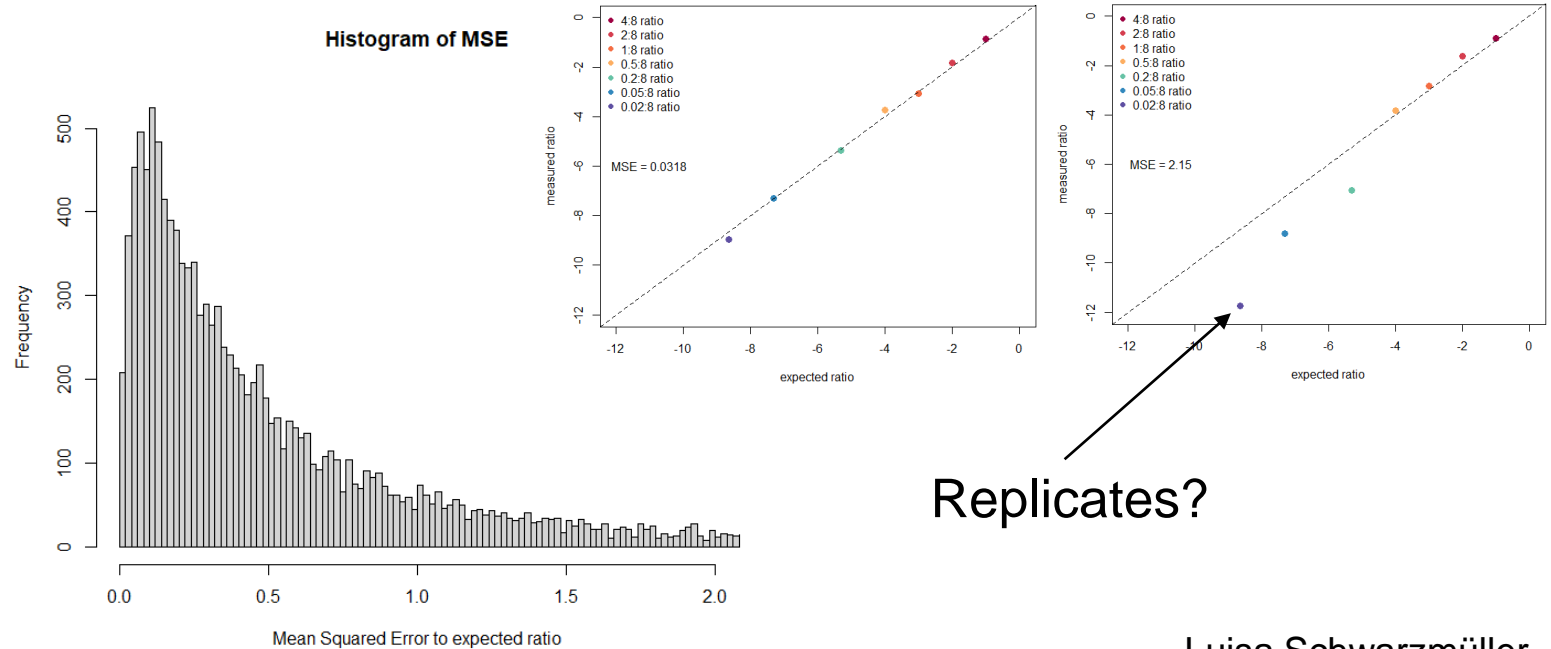


- Human with A. thaliana dilution

# Optimization HYE, fade-in fade-out, thoughts

**Paradigm shift in biomarker translation: a pipeline to generate clinical grade biomarker candidates from DIA-MS discovery**

[1]Qin Fu, [1]Manasa Vegesna, [1]Niveda Sundararaman, [2]Eugen Damoc, [2]Tabiwang N. Arrey, [2]Anna Pashkova, [1]Emebet Mengesha, [1]Philip Debbas, [1]Sandy Joung, [1]Dalin Li, [1]Susan Cheng, [1]Jonathan Braun, [1]Dermot P.B. McGovern, [1*]Christopher Murray, [2*]Yue Xuan, and [1#]Jennifer E. Van Eyk

# Optimization HYE, fade-in fade-out, thoughts



Replicates?

Luisa Schwarzmüller
luisa.schwarzmueller@dkfz.de

# Optimization HYE, fade-in fade-out, thoughts

- Info on all proteins in different abundance

- Pseudo ground truth → they should be on the line

- Feed models HMM, random forest, deep learning

- Which parameter combination good proxy for quality of:

  - Quantification

  - Detection

**dkfz.**