# Data Independent Acquisition Strategies for Quantitative Proteomics: The Challenges of Scaling Up to Meet Demand
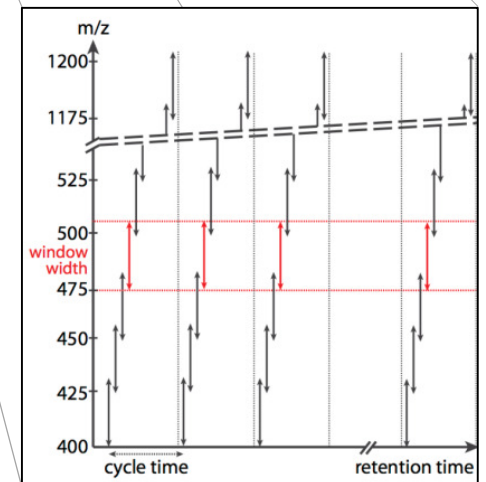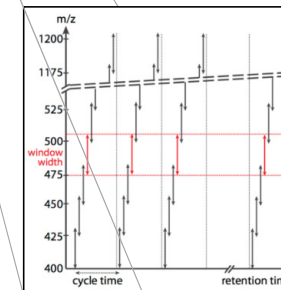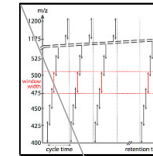
Data Independent Acquisition Interest Group
Presiding:  Ben Collins (ETH Zurich), Isabell Bludau (ETH Zurich)

65th ASMS Conference on Mass Spectrometry and Allied Topics
June 7th, 2017
Room: 235 – 238 -- Indiana Convention Center, Indianapolis, IN

Real time poll

[bit.ly/2sEPrpu](bit.ly/2sEPrpu)

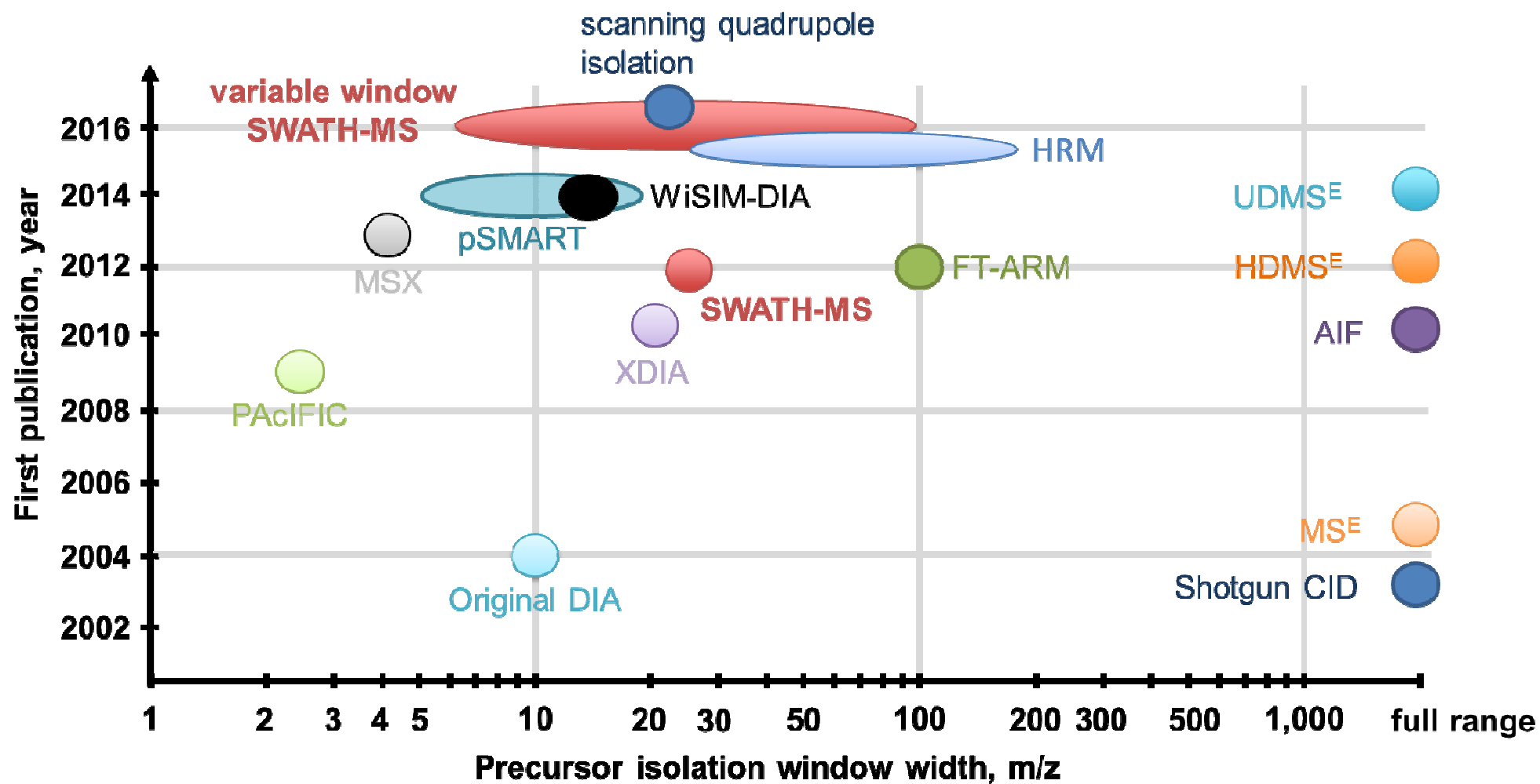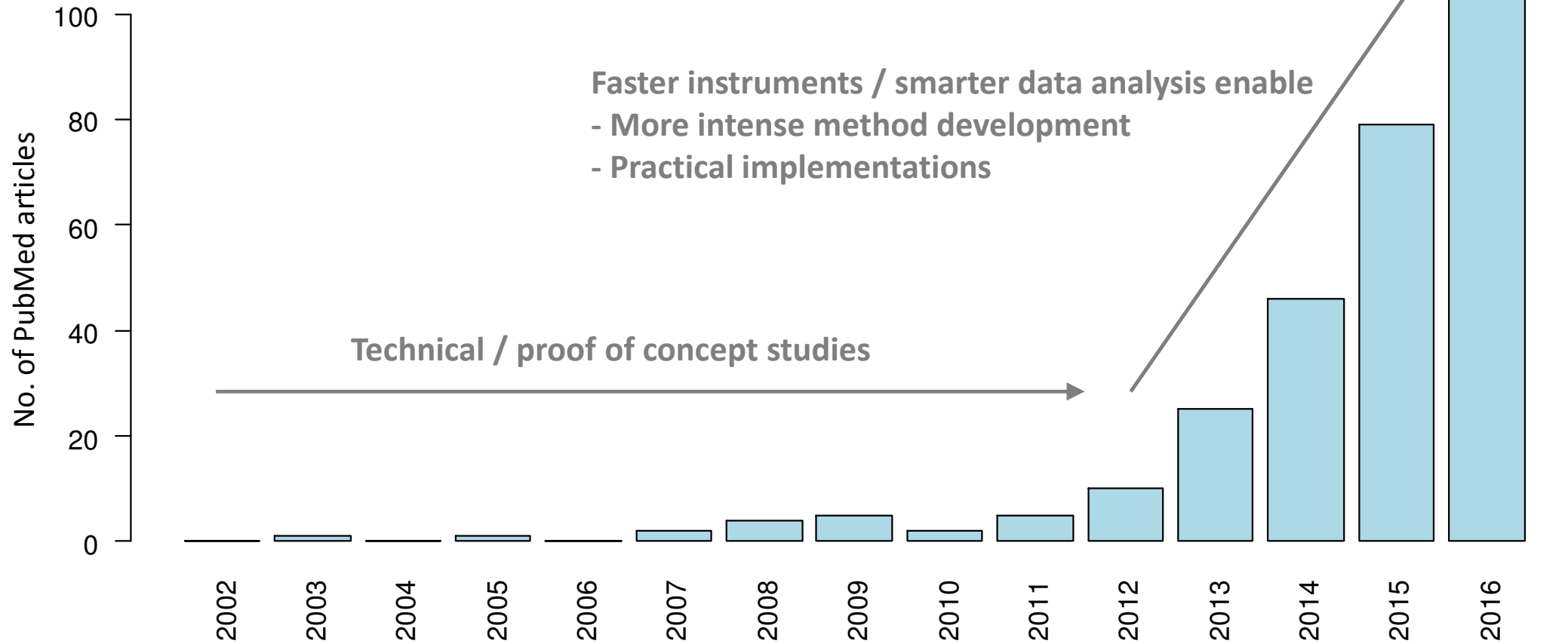# The breadth of DIA methodology is increasing
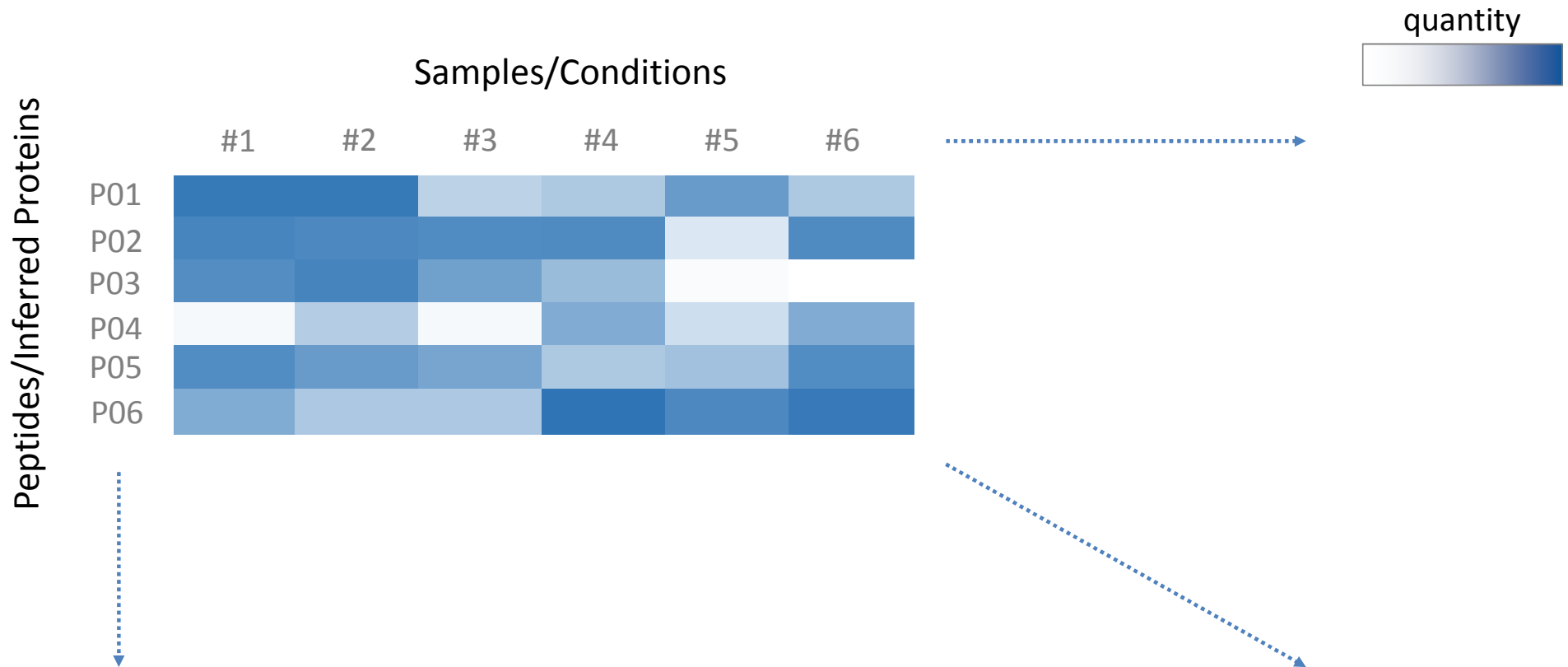


Figure from Tina Ludwig – SWATH-MS tutorial (in prep)

# The development and application of DIA to life science research is increasing

PubMed query:
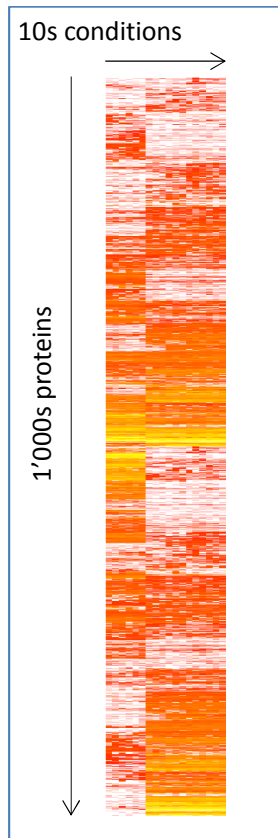((proteom* OR protein) OR peptide) AND ("data independent acquisition" OR SWATH)

Faster instruments / smarter data analysis enable
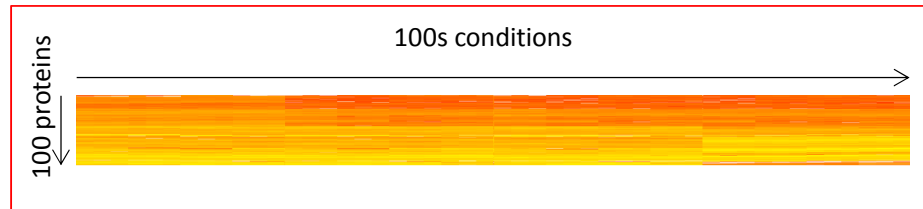- More intense method development
- Practical implementations

Technical / proof of concept studies

No. of PubMed articles

100

80

60

40

20

0

2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016

# Scaling up – in which dimensions?

# The primary goal of DIA is data completeness

**Historical problem of DDA based approaches**

10s conditions

1'000s proteins

**Targeted proteomics (SRM/PRM)**

100s conditions

100 proteins

**DIA**

100s conditions

1'000s proteins

Broad/deep coverage

Flexibility in targets

Figures from Gillet, L. C., et al. *Annual Review of Analytical Chemistry* **9,** (2016).

# Scaling up DIA – no. of samples/conditions

**Population-based analysis (personalized med., etc)**

Published online: February 4, 2015

*Article*

TRANSPARENT PROCESS    OPEN ACCESS    molecular systems biology

Quantitative variability of 342 plasma proteins in a human twin population

Yansheng Liu[1,*,†], Alfonso Buil[2,†], Ben C Collins[1], Olga Vitek[3], Jeppe Mouritsen, Genevieve Lach Ruedi Aebersold[1,5,**]

**232 human** 

**Large-scale knock-out screens**

bioRχiv beta

Precise label-free quantitative proteomes in high-throughput by microLC and data-independent SWATH acquisition

hael Mülleder[1,2]

**Genetic ass**

RESEARCH ARTICLE

PROTEOMICS

Systems proteomics
mitochondria function

Evan G. Williams,[1*] Yibo Wu,[2*] Pooja Jha,[1] Sébastien Dubuis,[2] Peter Blattmann,[2] Carmen A. Argmann,[3] Sander M. Houten,[3] Tiffany Amariuta,[1] Witold Wolski,[2] Nicola Zamboni,[2] Ruedi Aebersold,[2,†] Johan Auwerx[1]†

**386 mouse liver samples**

**ple numbers**

---

**Issues - scaling up samples?**    **Workshop**
- Error control (FDR)    ✔
- Data completeness    ✔
- Quantitative precision/accuracy    ✔ ✘
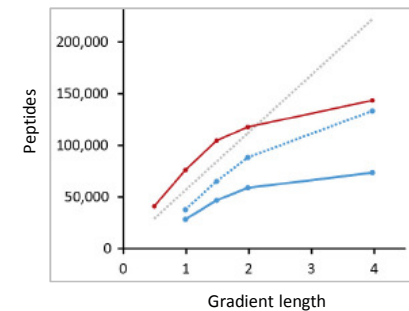- Quality control    ✘
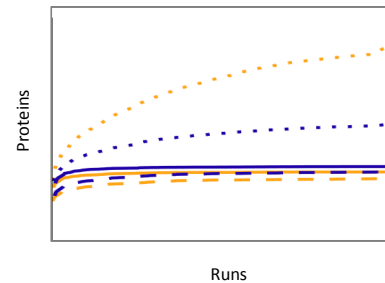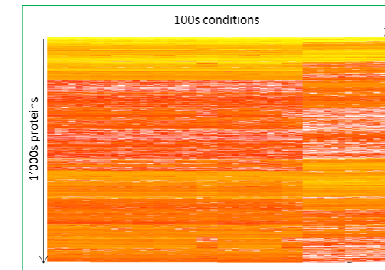
# Scaling up DIA – no. of peptides/inferred proteins
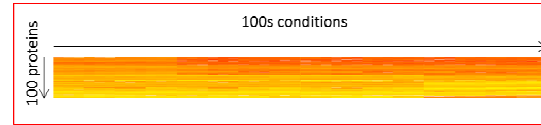
# Topics for discussion

1. Mike MacCoss (Univ. of Washington)
   - DIA as targeted proteomics

2. Alexey Nesvizhskii (Univ. of Michigan)
   - DIA as discovery proteomics

3. Isabel Bludau (ETH Zurich)
   - Error rate control at various levels

4. Lukas Reiter (Biognosys)
   - Depth of proteome coverage

5. Ben Collins (ETH Zurich)
   - Repeatability and data completeness
   (optional - depending on time)

# The goal is a community discussion!

- If you have a question or comment
    1. Raise your hand
    2. Shout
    3. Throw something
    4. Use the ASMS app
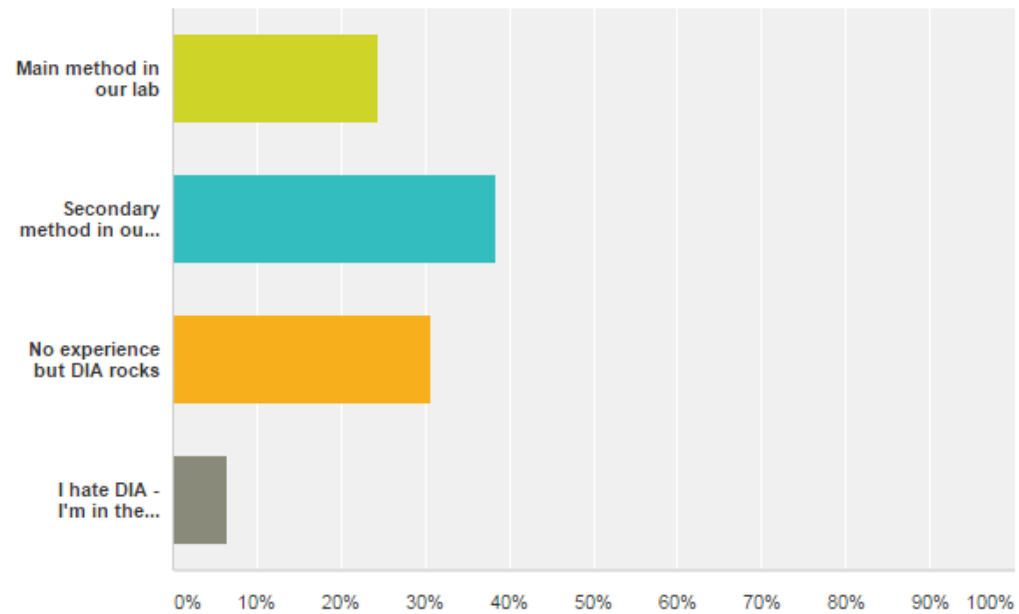
- Let's answer Question 1 in the poll

# Survey question 1

## 1. What's your experience with DIA?

○ Main method in our lab

○ Secondary method in our lab

○ No experience but DIA rocks

○ I hate DIA - I'm in the wrong room!

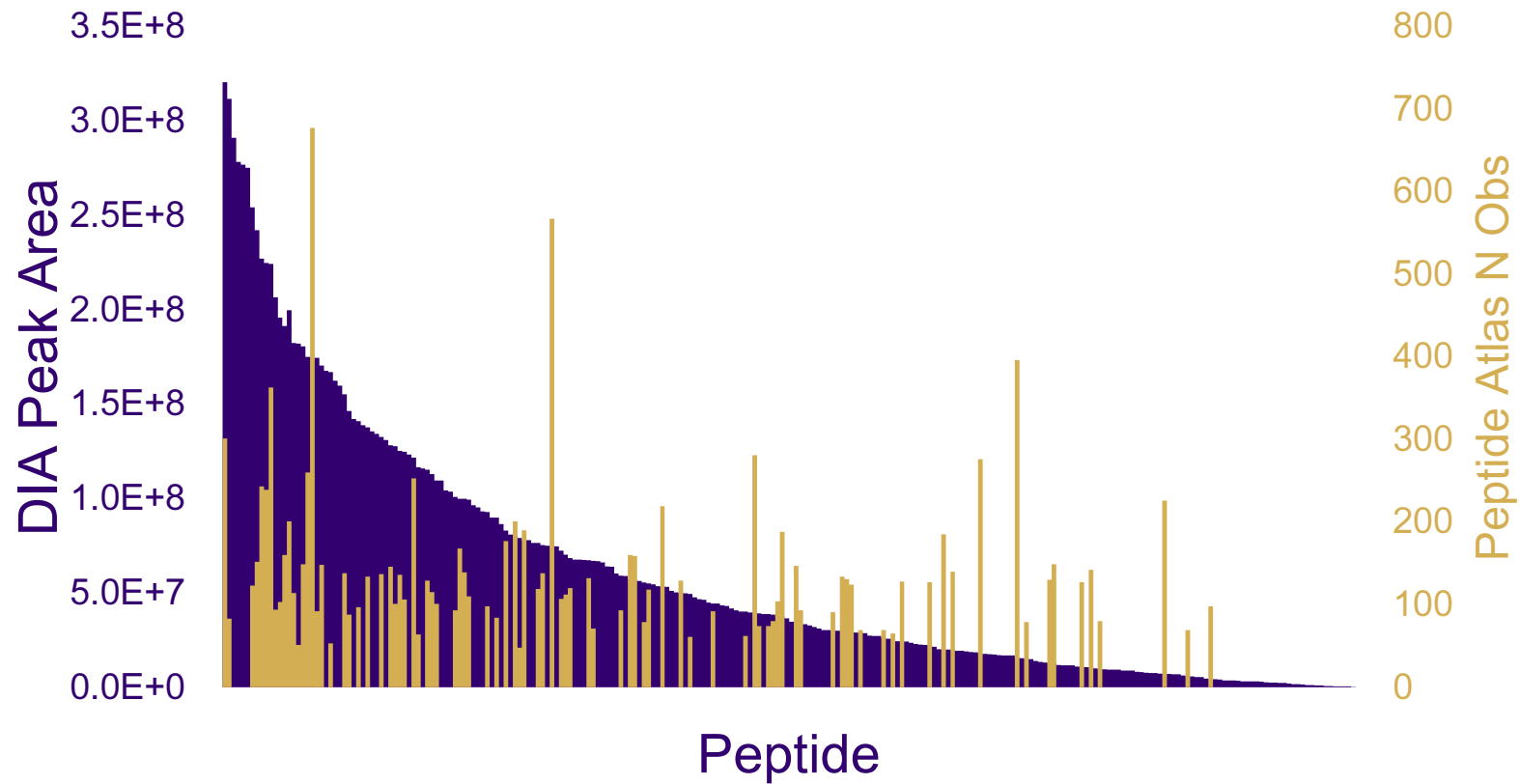# What's your experience with DIA?

Answered: 78    Skipped: 0



| Answer Choices | Responses | |
|---|---|---|
| Main method in our lab | 24.36% | 19 |
| Secondary method in our lab | 38.46% | 30 |
| No experience but DIA rocks | 30.77% | 24 |
| I hate DIA - I'm in the wrong room! | 6.41% | 5 |
| Total | | 78 |

# Mike

# Improving DIA Assay Design using Lessons Learned from Targeted Assay Development
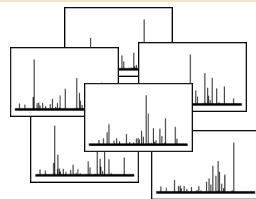
Apolipoprotein B100
DIA vs. PeptideAtlas
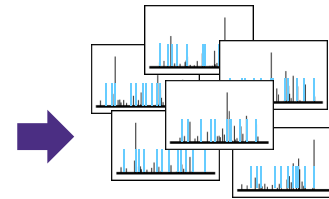
# Pecan: Detecting Peptides Directly from DIA Data

**Spectrum-centric analysis**

**What peptides best explain the data?**



**MS/MS spectra**

example tools:
SEQUEST, MASCOT ... etc.

**Protein sequence DB**

**Peptide spectrum matches (PSMs)**

p-value
q-value

**Peptide-centric analysis**

**Which peptides are detected in our data?**

**Peptides of interest**

APTNVTCILK

example tools:
OpenSWATH, Skyline

**Extracted MS/MS data**
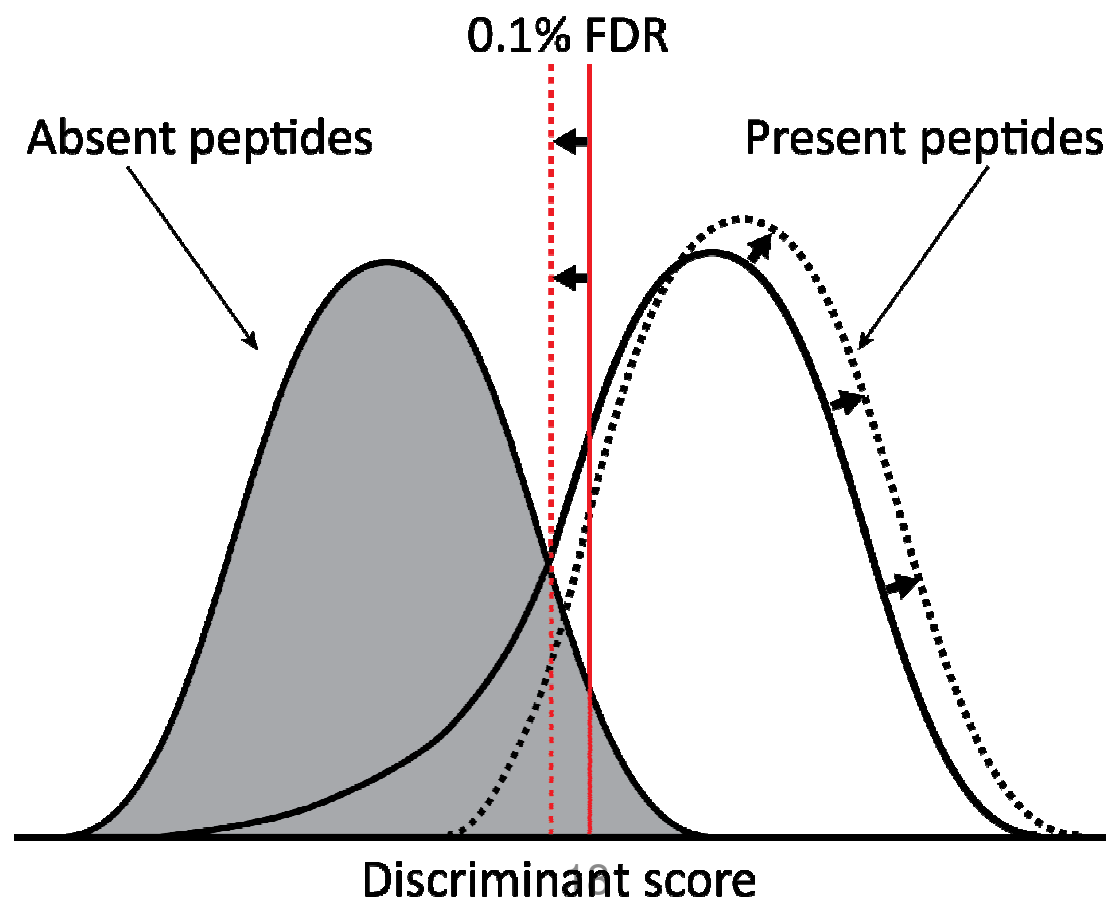
**Evidence for detection**

p-value
q-value

UNIVERSITY *of* WASHINGTON

**Problem #1: Detection of Peptides in DIA Data is Inversely Proportional to Isolation Width**

Narrow Isolation Windows Improves the Score Discrimination Between Absent and Present Peptides

0.1% FDR

Absent peptides

Present peptides

Discriminant score

# Picking Peptides Directly from DIA Data
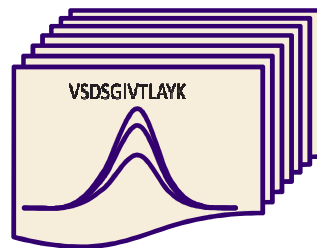## 12 LC-MS/MS Runs: ~1 µL plasma



UNIVERSITY of WASHINGTON

# ~~Spectral~~ Chromatogram Library Workflow



**Build a Chromatogram Library**

12 DIA Runs
**2 *m/z* Isolation**
Pooled Sample

Pecan Peptide
Detection

VSDSGIVTLAYK
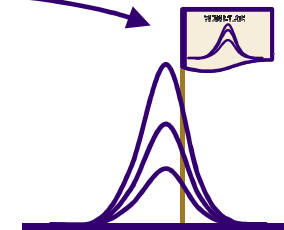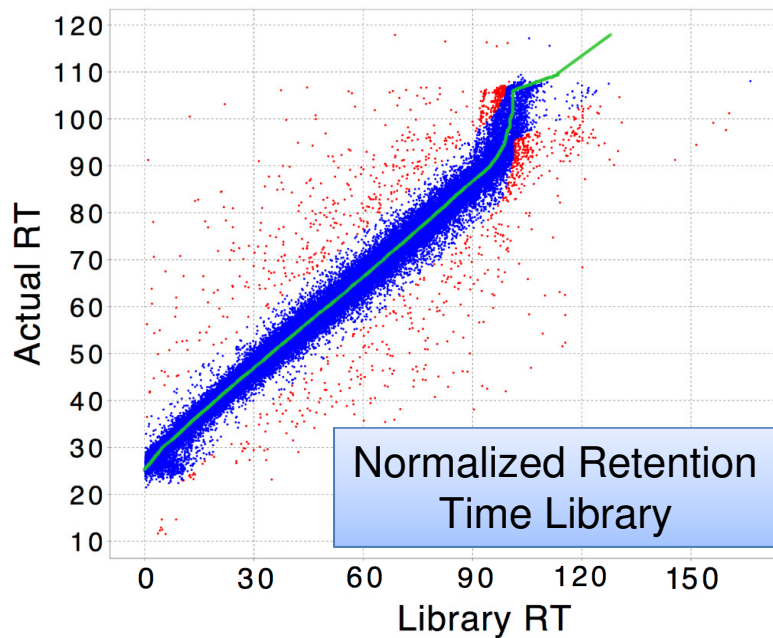
**Acquire DIA Data**

Precursor *m/z*

Time

1 DIA Run
**20 *m/z* Isolation**
/per sample

Chromatogram Extraction
Quantification

UNIVERSITY *of* WASHINGTON
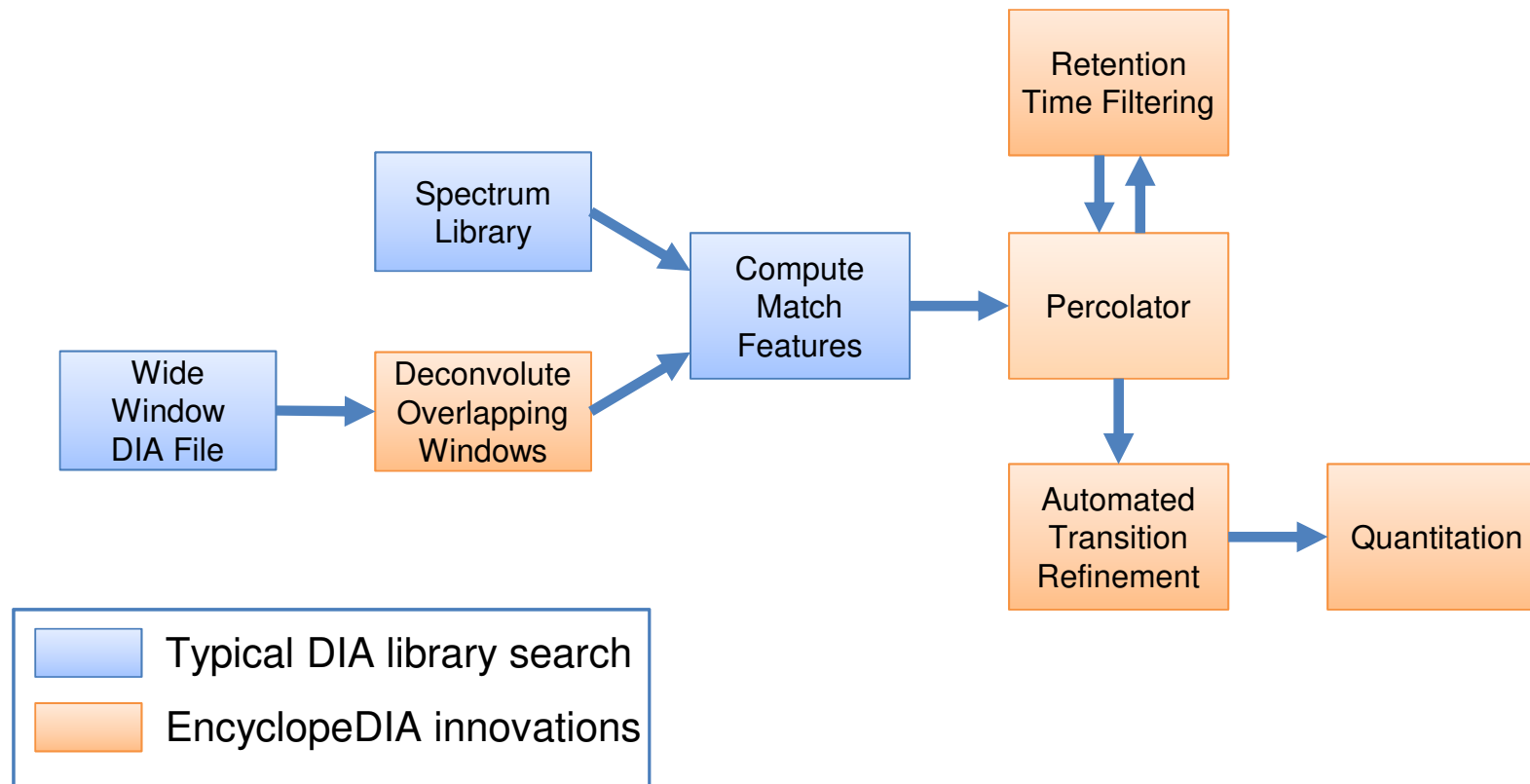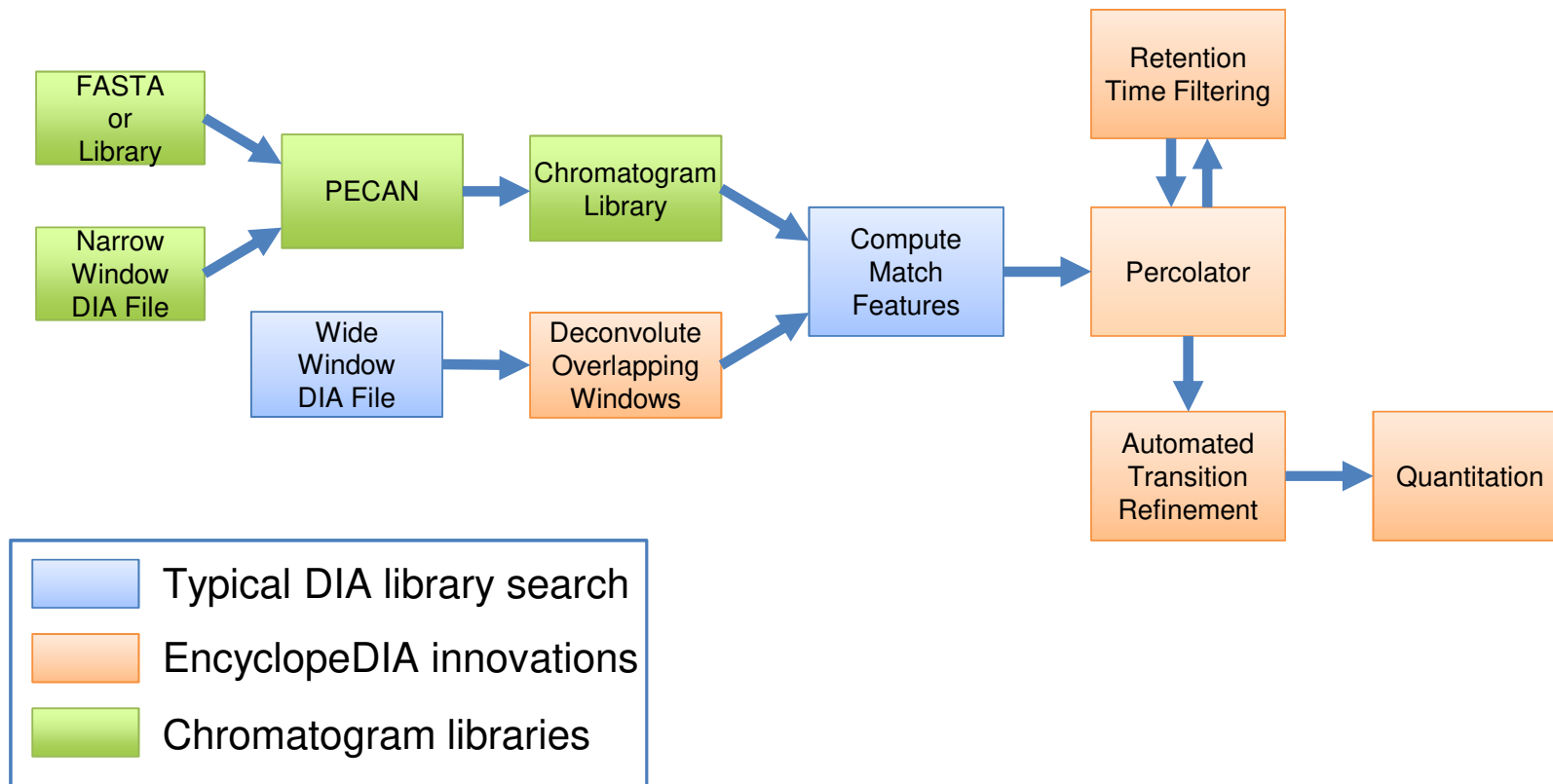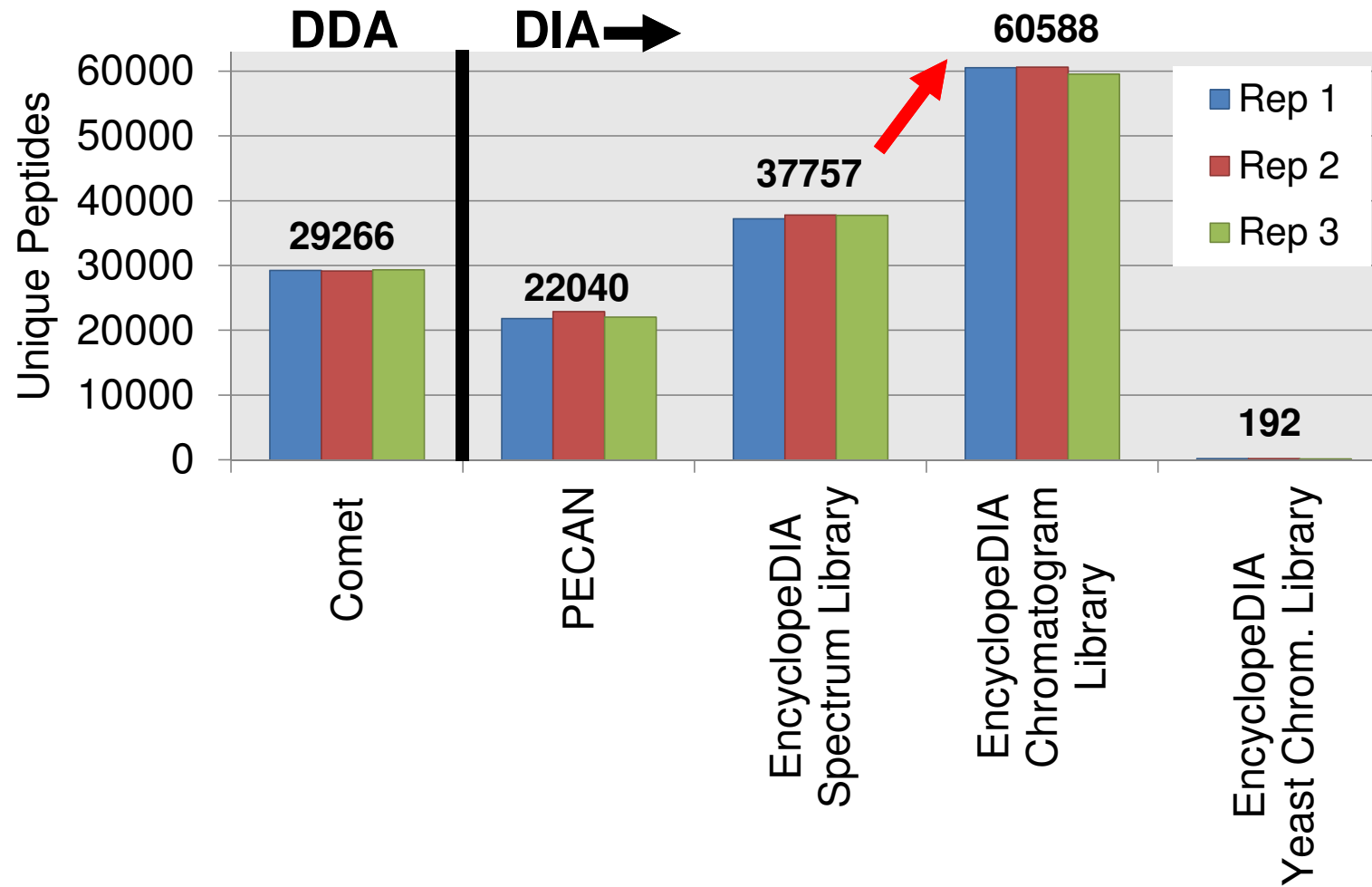
# EncyclopeDIA workflow



Percolator from Käll L et al, Nat Methods. 2007 Nov;4(11):923-5.

# EncyclopeDIA workflow

# Chromatogram libraries are significantly more powerful than spectrum libraries

# Human Plasma Chromatogram Library

**4,244** unique peptides  |  **495** proteins

**3,880** peptides mapping to one protein



Detections are made directly in the sample matrix

0.1% FDR

Incorrect PSMs

Correct PSMs

Discriminant score

# Lesson 3: Not all Peptides are Stable



🕐 3 days  > 4.5 x

**Apolipoprotein B100** | TTLTAFGFASADLIEIGLEGK @ 4°C

**Jim Bollinger** | ASMS 2014

# Assessing Peptide Stability with DIA



UNIVERSITY of WASHINGTON

Peptide Stability
Apolipoprotein B100
*top 50 peptides

Peak Area

3.5E+8
3.0E+8
2.5E+8
2.0E+8
1.5E+8
1.0E+8
5.0E+7
0.0E+0

Peptide

0 hr
72 hr

UNIVERSITY of WASHINGTON

# Lesson 4: Know your Digestion
## Alpha-1-acid glycoprotein

# Digestion Time Course by DIA

Digestion Time →



15 min      1 hr      2 hr      4 hr      10 hr      18 hr

1 DIA Run
**20 *m/z* Isolation**
/per sample

VSDSGIVTLAYK

# Lesson 5: Evaluate Linearity and LoQ



6 LC-MS/MS Runs
(not including replicates)

Grant RP, Hoofnagle AN, Clinical
Chemistry 2014

UNIVERSITY *of* WASHINGTON

# Are the measurements quantitative or just differential?



LOD: Limit of Detection
LLOQ: Lower Limit of Quantitation
ULOQ: Upper Limit of Quantitation

# Method to measure both LOQ and Linearity

**Diluent Sample Matrix**  **Pooled Reference Sample Matrix**

- Possible Samples to Use as a Diluent Matrix
  - Stable isotope labeled version of the matrix.
    - 15N or SILAC labeled cells
  - A diverged species
    - For human plasma we use chicken plasma.

# Reference Yeast BY4742 Diluted in 15N Yeast (S288c)



**Protein SSA1** (2.69E+05 copies/cell)

**Protein LSM2** (2.21E+03 copies/cell)

# Assess Reproducibility (5x5)



Replicate

Day

Grant RP, Hoofnagle AN, Clinical Chemistry 2014

UNIVERSITY *of* WASHINGTON

# DIA Assay Workflow

Peptide Detection and Stability
24 x 2 hr LC-MS/MS

Digestion Time Course
8 x 2 hr LC-MS/MS

Linearity / LoQ Assessment
10 x 2 hr LC-MS/MS

Reproducibility
25 x 2 hr LC-MS/MS

## 67 LC-MS/MS Runs

# Alexey

YEAH, WELL, THAT'S JUST, LIKE, YOU KNOW, YOUR OPINION, MAN.

# DIA

Is it a shotgun proteomics method?

*Yes, DIA as not "less shotgun" than DDA*

Is it a discovery proteomics method?

*Yes, DIA is a untargeted data acquisition method. It is even "less*

*targeted" than DDA*

# Lessons from History

High throughput proteomics methods that rely on the previously generated (proteomics) data have not been very successful

- Reference databases of 1D and 2D SDS Page gels (SWISS-2DPAGE database)
- AMT (accurate mass and time) approach
- Spectral library searching as replacement for database search

1) Match spot
2) Make map active
3) Click on spot
4) Putative ID pops up

# Hybrid (Direct+ Targeted) Strategy



C.C. Tsou *et al.* DIA-Umpire: comprehensive computational framework for data independent acquisition proteomics *Nature Methods*, 2015

# Violation of the Target-Decoy Assumption in Closed Searches



- ▶ Selected spectra corresponding to common modifications identified in open search and examined their identifications in closed search

- ▶ Under target-decoy assumptions, these spectra should all be incorrect and match equally to target and decoy sequences

- ▶ Target-decoy assumption is violated: 6X difference for carbamylation, 9X for oxidation

# Clothes And Body Parts

Match the images on the left to their corresponding images on the right.

# Clothes And Body Parts

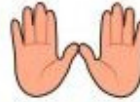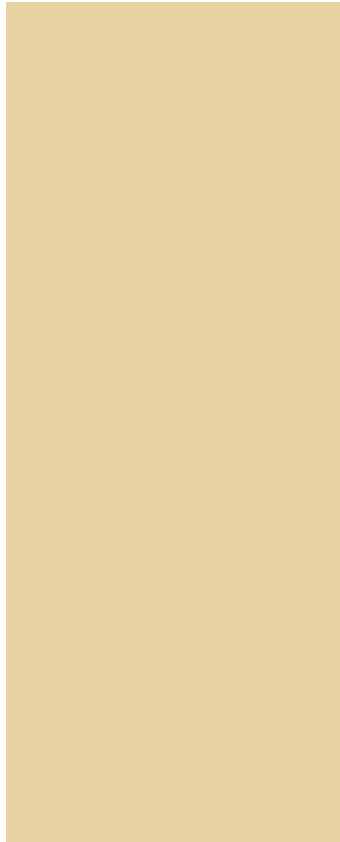Match the images on the left to their corresponding images on the right.

# Clothes And Body Parts

Match the images on the left to their corresponding images on the right.

# Survey question 2

## 2. In which mode should DIA data be analyzed?

- ◯ Targeted
- ◯ Discovery
- ◯ Both (depends on experiment)

# In which mode should DIA data be analyzed?

Answered: 77    Skipped: 1



| Answer Choices | | Responses | |
|---|---|---|---|
| ▼ Targeted | | 7.79% | 6 |
| ▼ Discovery | | 14.29% | 11 |
| ▼ Both (depends on experiment) | | 77.92% | 60 |
| Total | | | 77 |

# Isabell

# Scaling up DIA – Error rate control


focused ← DIA 'targeting' scale → global

**Sample-specific spectral library**

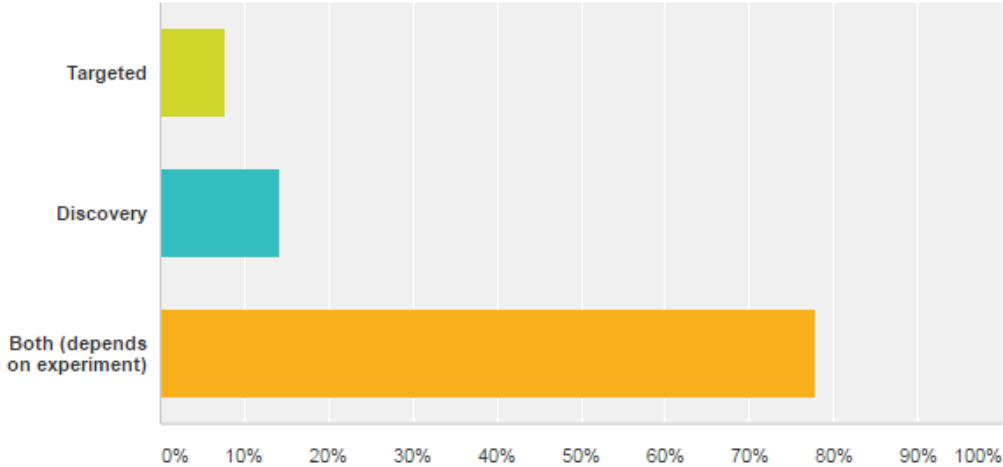- Built from same samples as DIA data
- Majority of peptides/proteins detectable

**Comprehensive deep spectral library**

- Built from multiple different samples
- Majority of peptides/proteins NOT detectable

# Scaling up DIA – Error rate control

focused ← DIA 'targeting' scale → global

**Sample-specific spectral library**

- Built from same samples as DIA data
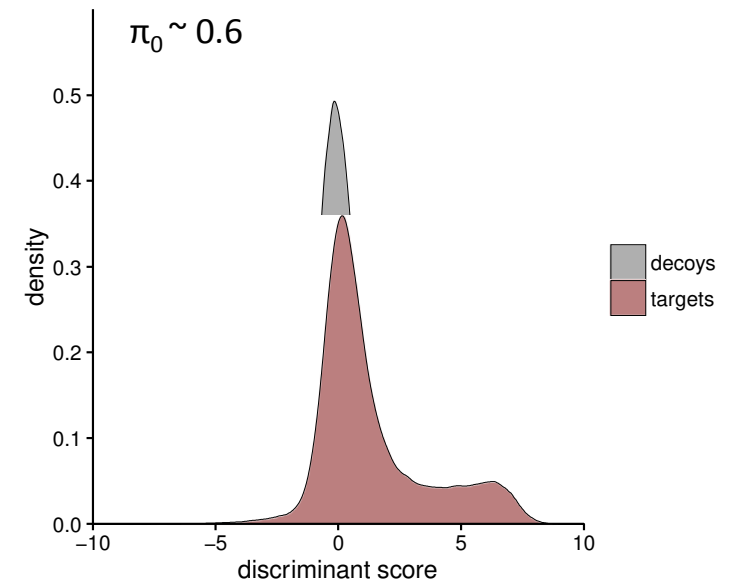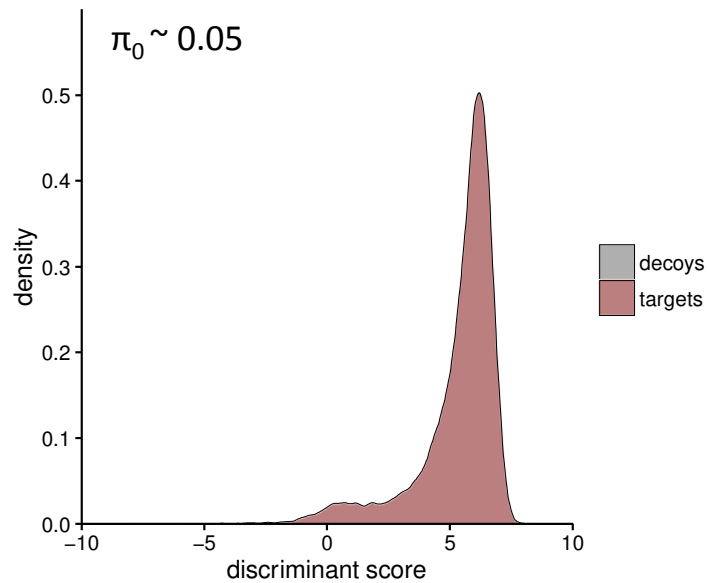- Majority of peptides/proteins detectable

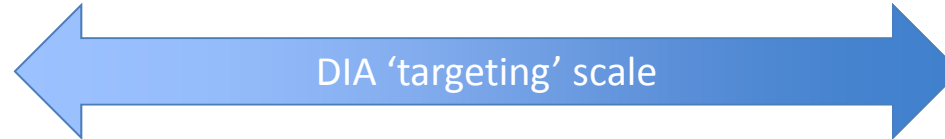**Comprehensive deep spectral library**

- Built from multiple different samples
- Majority of peptides/proteins NOT detectable

Protein-level error can b handled on library generation level

Protein-level error accumulates if not carefully controlled

# Scaling up DIA – Error rate control



focused ← DIA 'targeting' scale → global

**Comprehensive deep spectral library**
- Built from multiple different samples
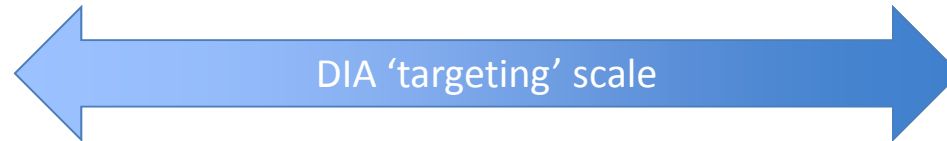- Majority of peptides/proteins NOT detectable

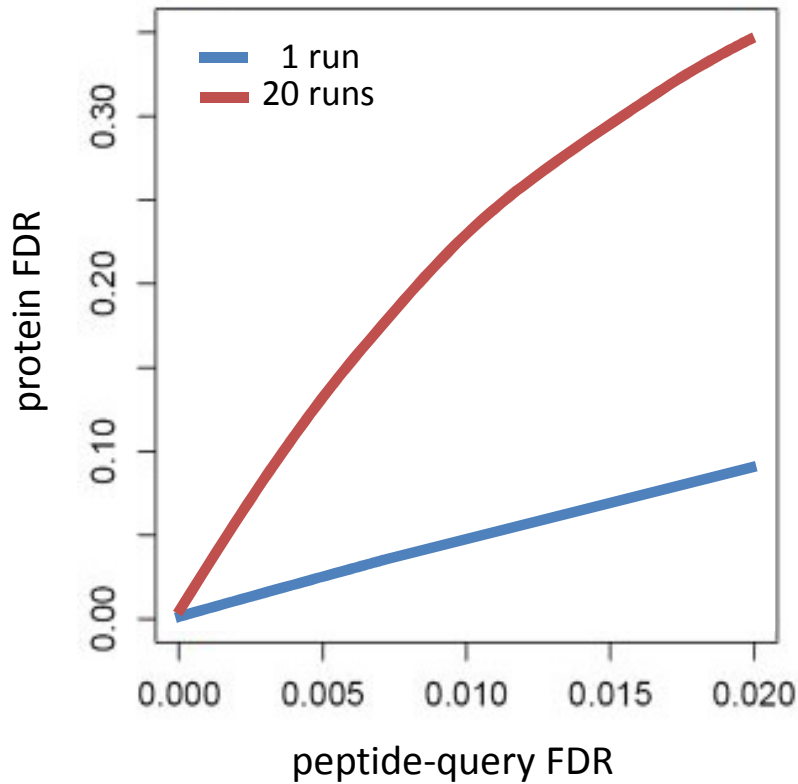Protein-level error accumulates if not carefully controlled

# Scaling up DIA – Error rate control

1. **Control error rate on protein level**:
   Take best peptide peak group per protein for FDR / q-value estimation on protein level

2. **Control error rate globally across all samples within a study: Protein master list**
   Take best peptide peak group per protein across all samples in a study to generate a protein master list at 1% FDR

Extended version of PyProphet: https://github.com/PyProphet
Rosenberger & Bludau *et al.* (submitted)

# Scaling up DIA – Error rate control

**DIA data:**

- Inter-laboratory study:
  229 DIA measurements of same HEK-293 cell lysate
  Collins et al. (2017)

**Spectral library:**

- Combined assay library (CAL):
  331 DDA injections of different human tissues and cell types including HEK293
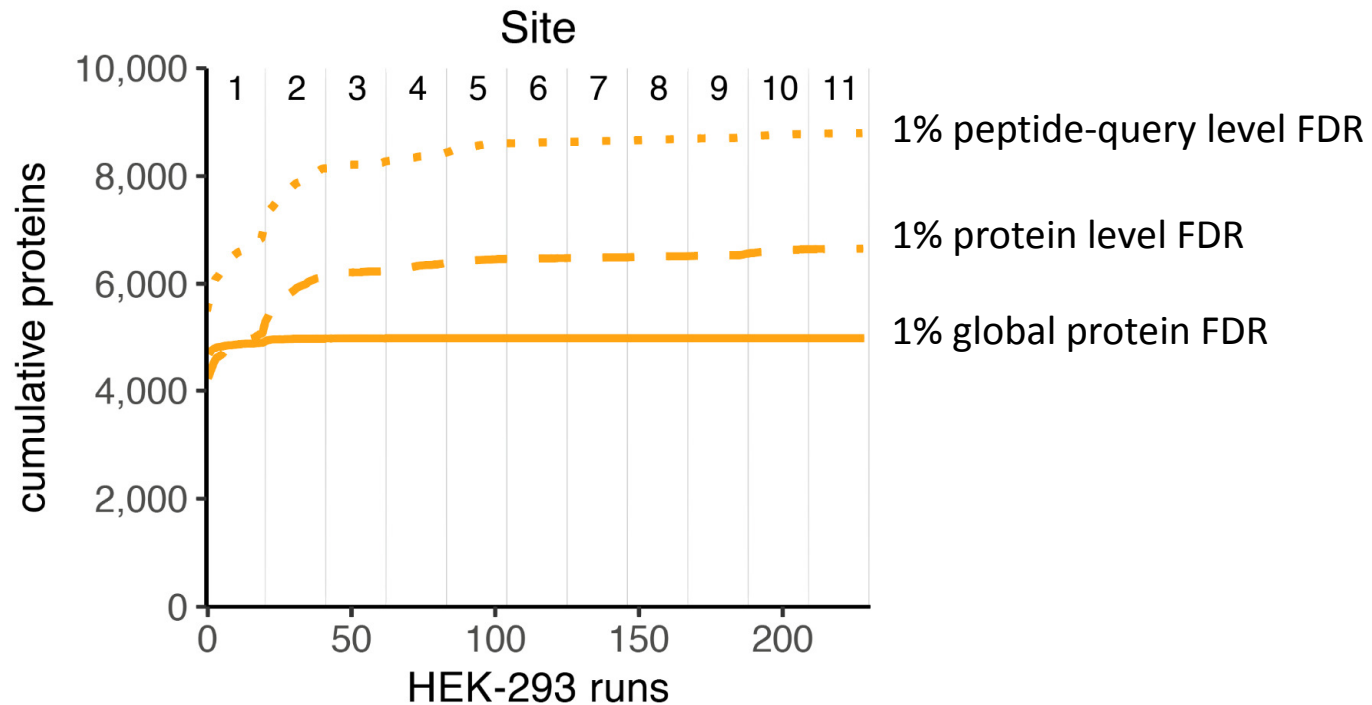  Rosenberger, G. *et al.* (2014)

# Scaling up DIA – Error rate control

**DIA data:**

- Inter-laboratory study:
  229 DIA measurements of same HEK-293 cell lysate
  Collins et al. (2017)

**Spectral library:**

- Combined assay library (CAL):
  331 DDA injections of different human tissues and cell types including HEK293
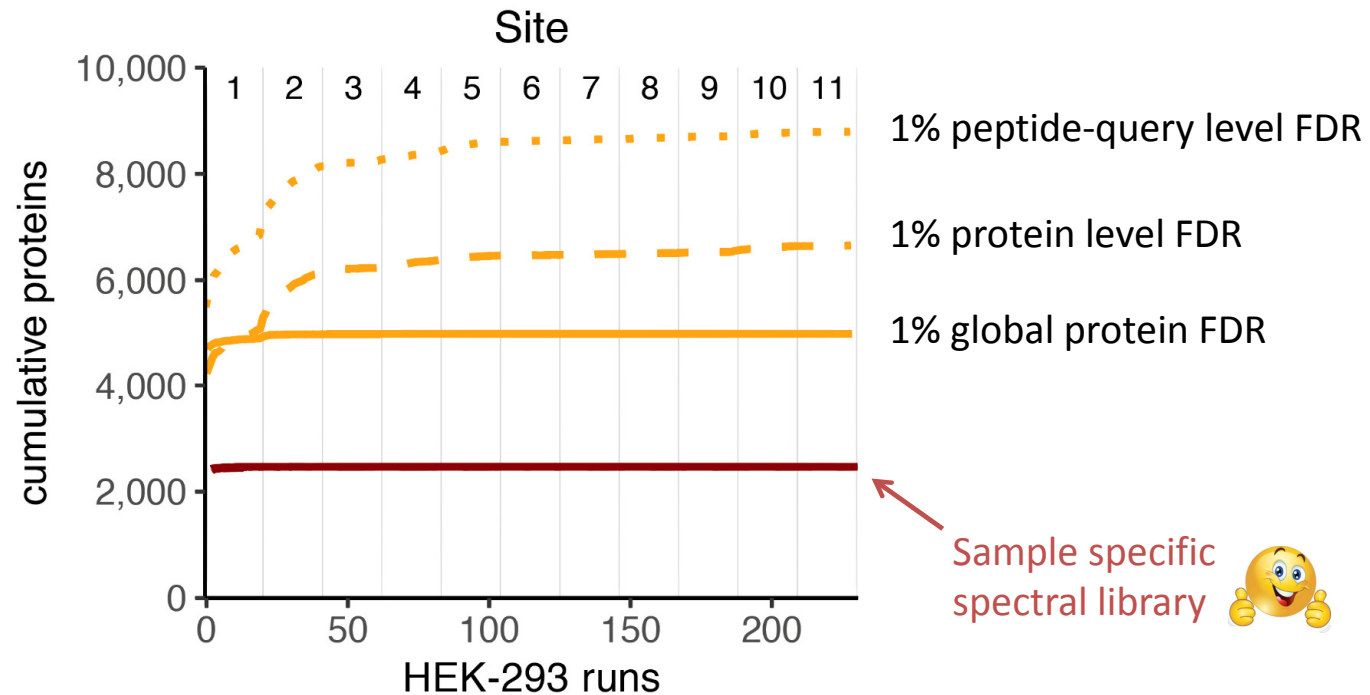  Rosenberger, G. *et al.* (2014)

# Scaling up DIA – Error rate control

**DIA data:**

- Blood plasma dataset with 246 samples
  Liu et al. (2015)

**Spectral library:**

- Combined assay library (CAL):
  331 DDA injections of different human
  tissues and cell types including HEK293
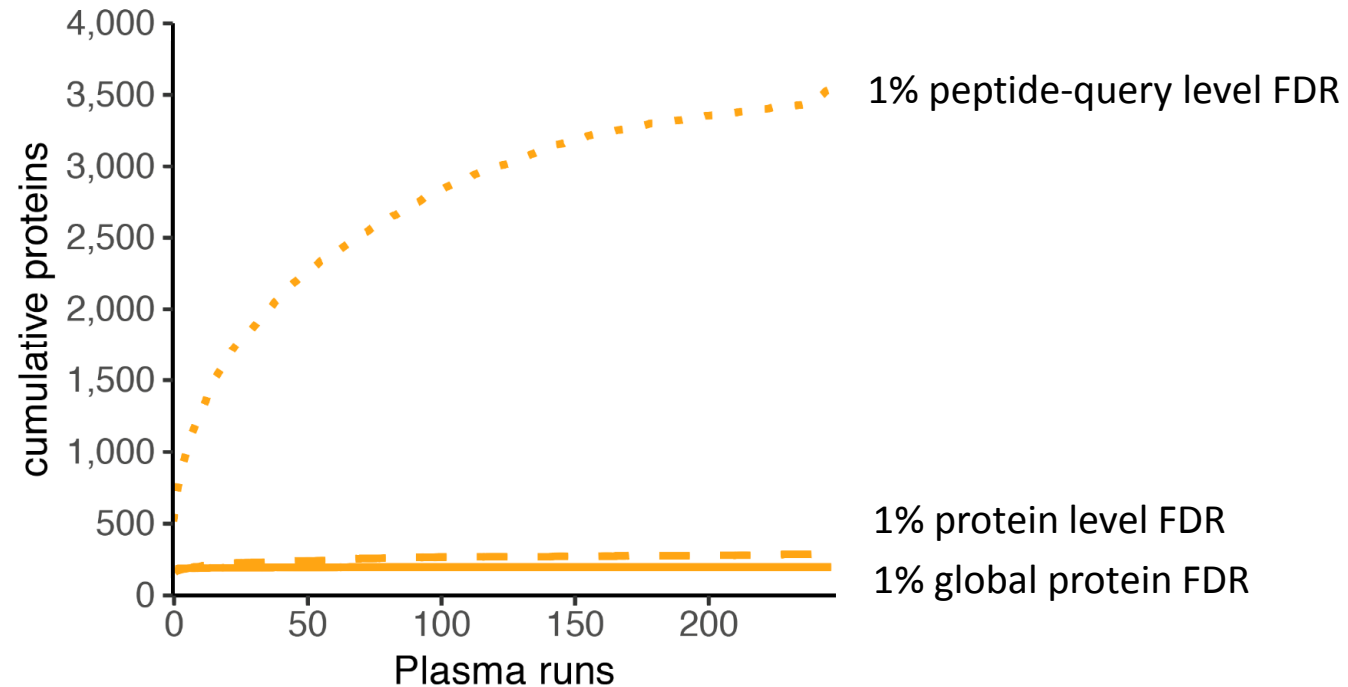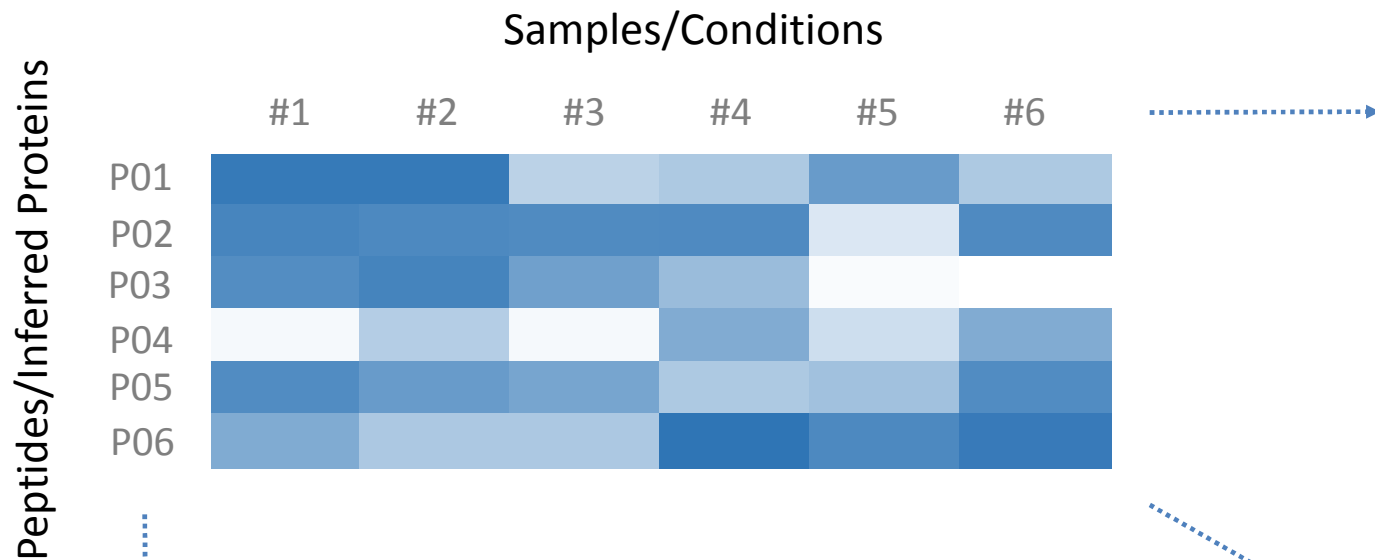  Rosenberger, G. *et al.* (2014)

# Scaling up DIA – Error rate control

Samples/Conditions



For large comprehensive
spectral libraries:
Use protein-level FDR

&

Context matters!
Use global protein master list
constraint to avoid error
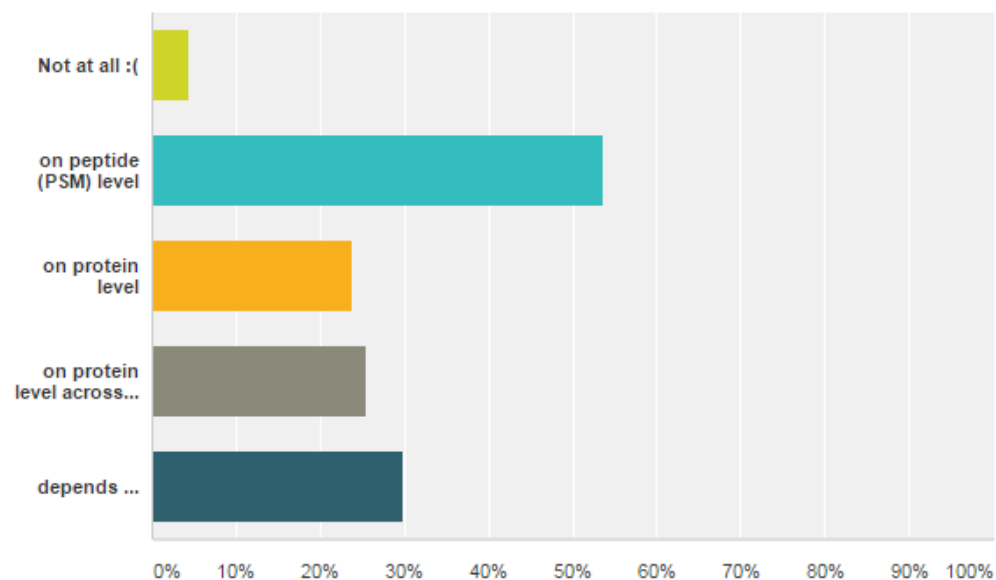accumulation across samples!

# Survey question 3

## 3. At which level do you control your FDR?

- [ ] Not at all :(

- [ ] on peptide (PSM) level

- [ ] on protein level

- [ ] on protein level across the entire dataset

- [ ] depends ...

# At which level do you control your FDR?

Answered: 67    Skipped: 11



| Answer Choices | Responses | |
|---|---|---|
| Not at all :( | 4.48% | 3 |
| on peptide (PSM) level | 53.73% | 36 |
| on protein level | 23.88% | 16 |
| on protein level across the entire dataset | 25.37% | 17 |
| depends ... | 29.85% | 20 |
| Total Respondents: 67 | | |

# Lukas

ASMS 2017

# DIA Workshop –
# "Depth of Proteome Coverage"

Lukas Reiter,  Biognosys

# Depth of Proteome Coverage

Why having a large proteome coverage?

- Discovery
    - E.g. drug target deconvolution with Limited proteolysis (LiP) *
- Low abundant wish list proteins combined with discovery
- Multi OMICS -> increase overlap with other data sets

How can the proteome coverage be increased?

- Sample (prep)
- **Chromatography**
- Instrumentation
- **DIA Method**
- **Spectral library**
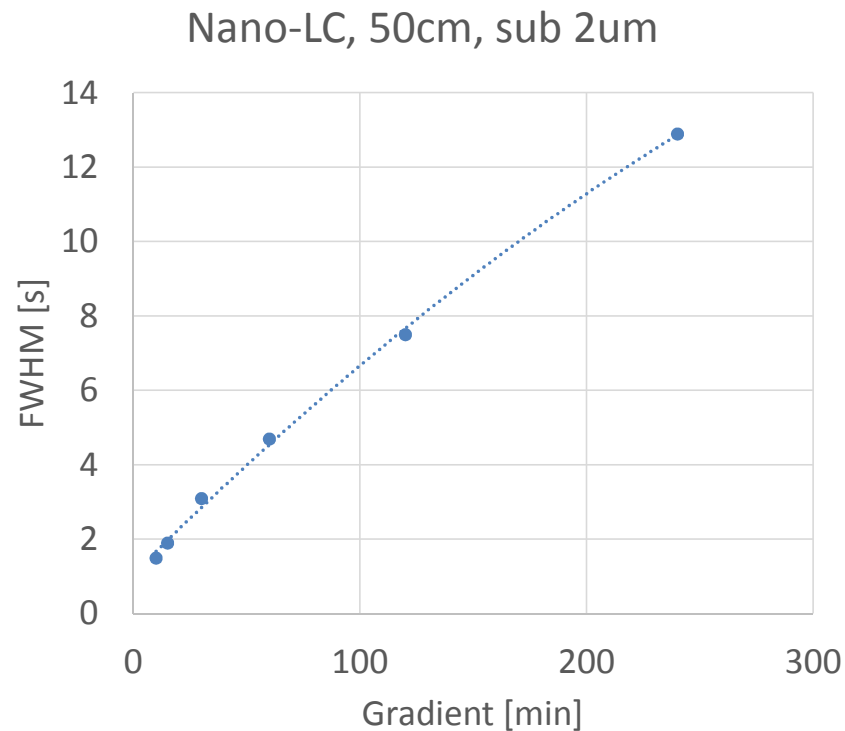- **Precision iRT**

*Leuenberger et al. Cell-wide analysis... Science (2017)*

66

# Chromatography

How to get a high peak capacity?

- UHPLC

- Nano-LC

- 75um ID long columns

- Sub 2um beads

- Long gradients

- Low dead volumes

**1m column, 4h gradient**

- Peak capacity > 700

FWHM:                    measured as median for all peptides identified in a HeLa digest

Peak capacity:           $1 + g / (FWHM*1.7)$



Nano-LC, 50cm, sub 2um

# DIA Method

**Step by step optimization**

## 1) Data points per peak

- Generate a scouting method which over samples the peaks
- Scale the number of MS2 segments from the scouting method to result in 5, 8, 11 and 14 data points per peak
- Pick the best method

## 2) MS1 resolution

- Vary the MS1 resolution from 30'000 to 240'000 (balance the MS2 segments to keep data points per peak constant)
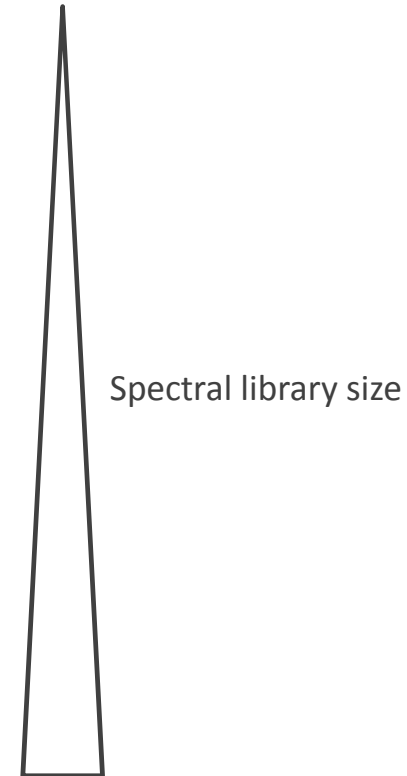- Pick the best method

## 3) MS2 resolution

- Vary the MS2 resolution from 15'000 to 120'000 (balance the MS2 segments to keep data points per peak constant)
- Pick the best method



Data points per peak
14
11
8
5

MS1 resolution
240,000
120,000
60,000
30,000

MS2 resolution
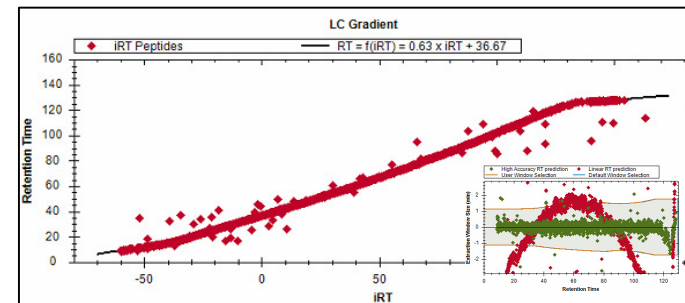120,000
60,000
30,000
15,000

# Spectral Library

**Evolution over time for Biognosys**

- Replicate injections and DDA on project samples

- Mild fractionation

  - 6 high pH reversed phased fractions

- Deep fractionation

  - Two condition pools

  - Pooled micro fractions from UHPLC

  - 10 fractions each

…

- Resource spectral libraries or from large sets of synthetic peptides

Spectral library size

# Precision iRT

- Extends the indexed retention time (iRT) concept
- Allows very precise targeting in retention time dimension
- Dynamically adapts
- XIC windows of 1-3% of gradient length can be achieved
  - Especially when using spectral libraries acquired on exactly the same setup



*Bruderer et al. High-precision iRT prediction… Proteomics (2016)*

# Depth of Proteome Coverage

*Some example data*

- Setup

  - Deep project-specific spectral library (MaxQuant & Spectronaut)

  - Analyzed with Spectronaut

  - Peptide and protein FDR 1%

  - 4h gradients, 1m column

  - HEK-293 sample

- Single run results

  - 7'060 protein groups, 154'643 precursors

  - Median XIC width 8.5 min

  - Peak capacity 710, median FWHM 13 s

- Technical triplicates

  - 6'534 proteins with CVs < 20%

  - 123'700 precursors with CVs < 20%

  - Data completeness for precursors: 91%

# Survey Question 4

## 4. What do you consider a deep proteome coverage?

○ 100 + proteins

○ 1,000 + proteins

○ 5,000 + proteins

○ 10,000+ proteins

# What do you consider a deep proteome coverage?

Answered: 65    Skipped: 13



| Answer Choices | Responses | |
|---|---|---|
| 100 + proteins | 4.62% | 3 |
| 1,000 + proteins | 9.23% | 6 |
| 5,000 + proteins | 61.54% | 40 |
| 10,000+ proteins | 24.62% | 16 |
| Total | | 65 |

# Ben

# Study design -- Inter-lab SWATH-MS

# Peptide/protein detection rates (HEK293 lysate)



Total -- 4,960 proteins / 39,928 peptide PGs
Site median – 4,691 / 34,286 PGs

(1% protein FDR / 1% peptide PG FDR)

- - accumulated proteins detected

# Data completeness



No alignment

No ID propagation between runs

4,960 inferred proteins

80 % complete
4,064 inferred proteins

log2 protein abundance
10 15 20 25

1 2 3 4 5 6 7 8 9 10 11
Site
(229 SWATH files)

0 20 50 80 100
% missing values

# Repeatability of identification



Inter-lab SWATH-MS

Tabb et al. – DDA (2010)

3 reps

21 reps

229 reps

No alignment!

Tabb, D. L. *et al.* Repeatability and Reproducibility in Proteomic Identifications by Liquid Chromatography–Tandem Mass Spectrometry. *J. Proteome Res.* **9,** 761–776 (2010).

# Alignment can only improve completeness…



**TRIC**
TRansfer of
Identification
Confidence

Röst, H. L. *et al. Nat Meth* **13,** 777–783 (2016)

# Comments? Questions?

# Linearity, dynamic range, and response differences



**All sites**

average

Normalize (based on HEK293 medians)

**1 site**

**All sites**

Legend:
- AEN.LVE
- HNRPC.GFA
- SFPQ.FAQ
- SMD3.VAQ
- ZMAT5.APP
- ANRA2.SVQ
- DDX5.FVI
- ESRP2.EAS
- RUXF.C[CAM]NN
- SFSWA.EAQ
- SMD1.EPV
- SNRPA.EVS
- APR.APA
- E2F7.LDF
- NCBP2.SDS
- ROA2.GGN
- SMD2.NNT
- SRSF7.AFS
- FL2D.YTD
- PHLA3.SGG
- SF3B2.VGE
- SFR19.TPE
- SRSF3.AFG
- WDR83.DGQ
- AK17A.HDW
- HNRPU.SSG
- RBM5.GLP
- RBMX.LFI
- SRRM4.LGQ
- TRA2A.TGP

Reproducibility (30 x SIL peptides)

# Global similarity of quantitative protein abundance profiles (HEK293 lysate)



Median Pearson correlation (overall) = 0.940
Median Pearson correlation (within sites) = 0.971
Minimum Pearson correlation = 0.868

Collins BC*, Hunter C*, Liu Y* et al, Nature Communications (in press)

# If you want learn more about DIA/SWATH



dia-swath-course.ethz.ch

Registration for this year is closed but lecture videos will be posted **late July 2017**

Thanks for participating!!

Ideas for discussion topics for next year to:

collins@imsb.biol.ethz.ch

bludau@imsb.biol.ethz.ch