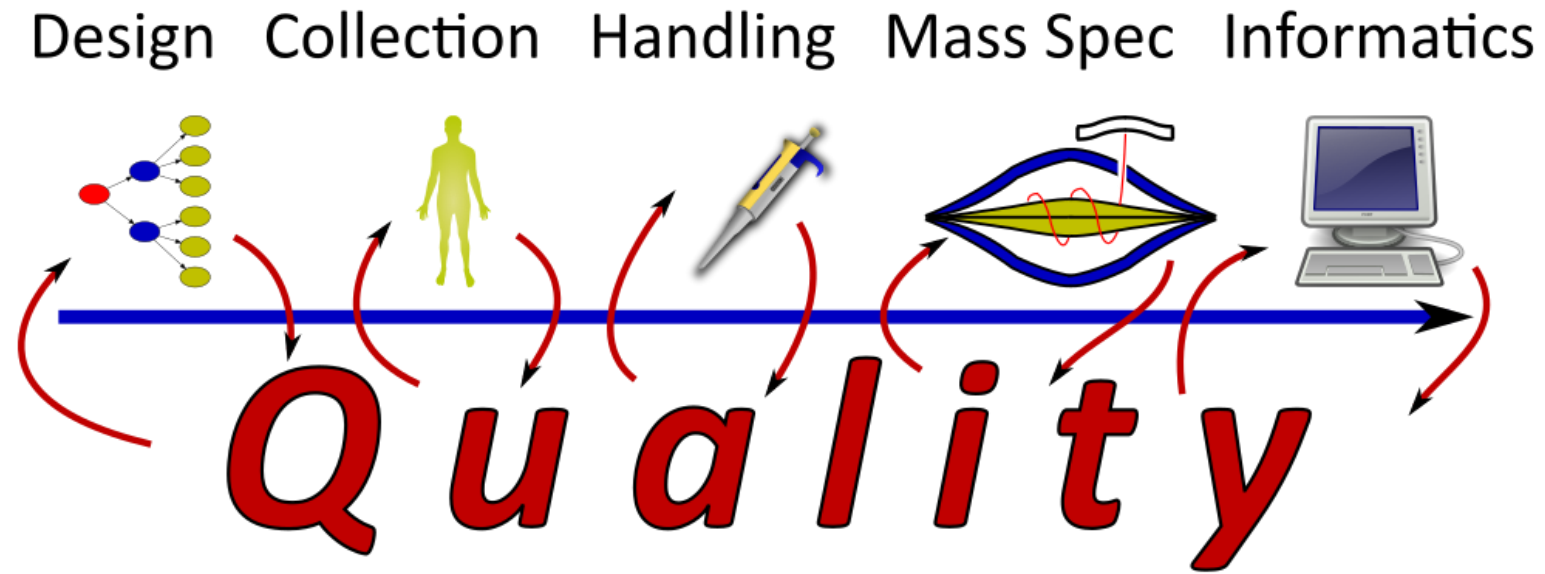# PCA-powered quality control for large-scale proteomics
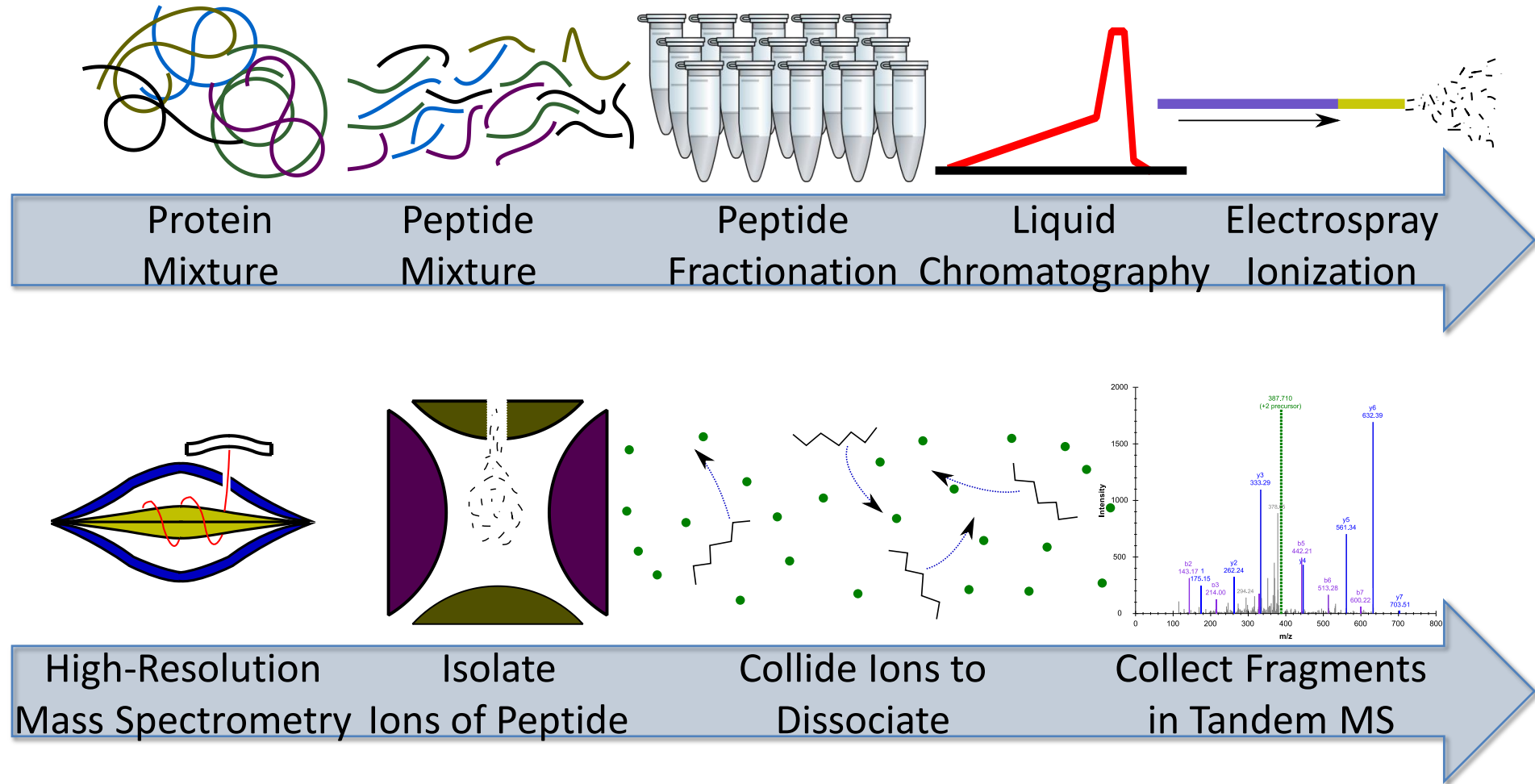
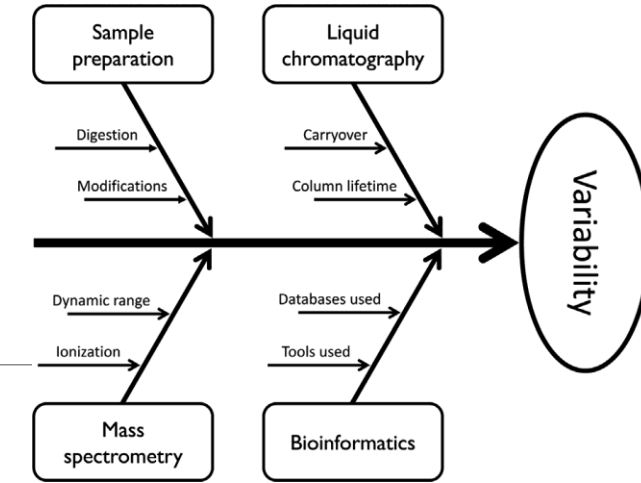DAVID L. TABB, PH.D.

OCTOBER 20, 2020

# Overview

- Proteomics: technical complexity accumulates variability

- Quality metrics quantify variation, but how to interpret?

- Principal Components Analysis: a brief introduction

- Outliers, batch effects, and analysis of variance, oh my!

# Discovery Proteomics



Protein Mixture → Peptide Mixture → Peptide Fractionation → Liquid Chromatography → Electrospray Ionization

High-Resolution Mass Spectrometry → Isolate Ions of Peptide → Collide Ions to Dissociate → Collect Fragments in Tandem MS
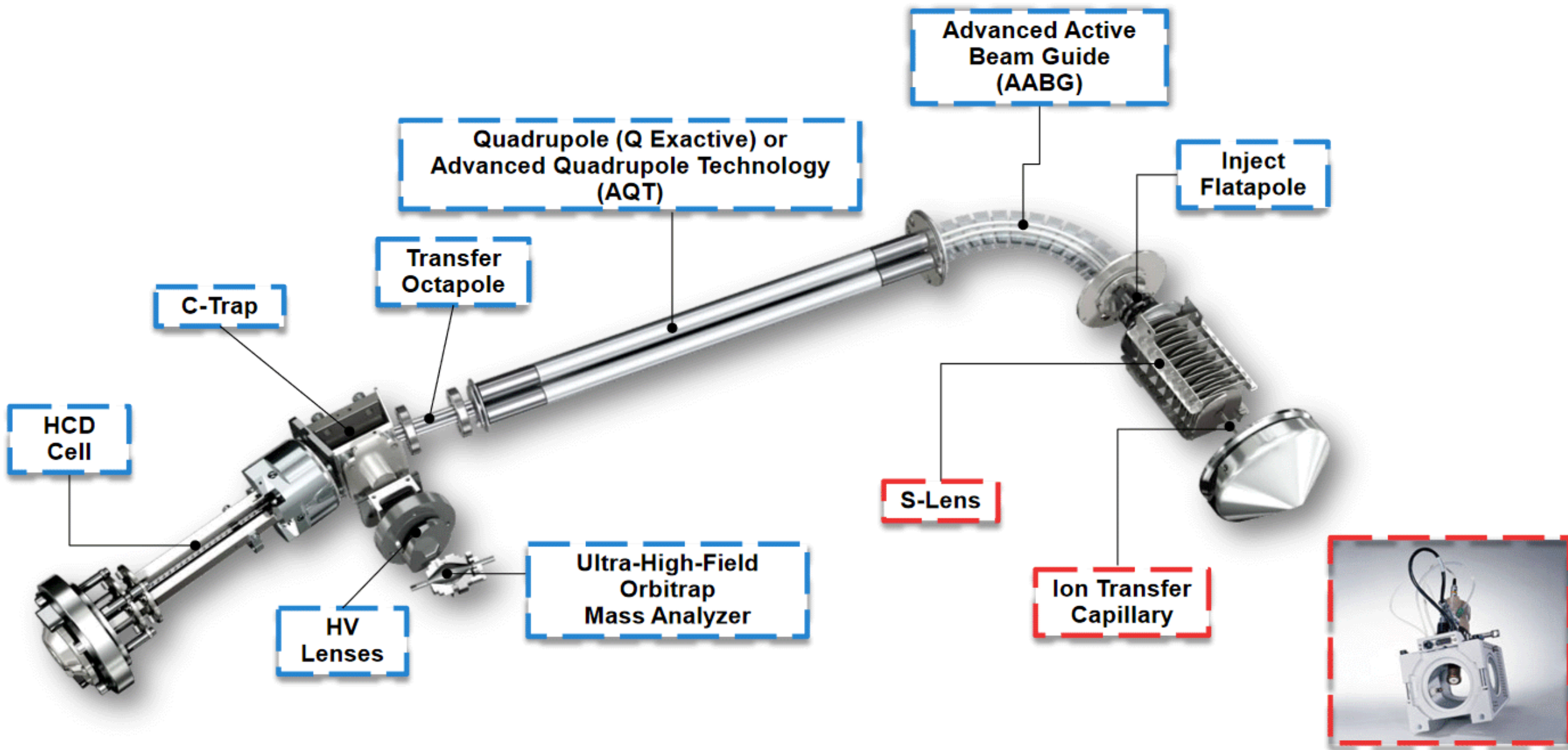
# Factors contributing variability

- Enrichment (Co-IP or PTM)
- Depletion (from biofluids)
- Denaturation and digestion
- Fractionation
- Liquid chromatography
- Electrospray

- Tune response
- Mass calibration
- Source cleanliness
- Ion optics
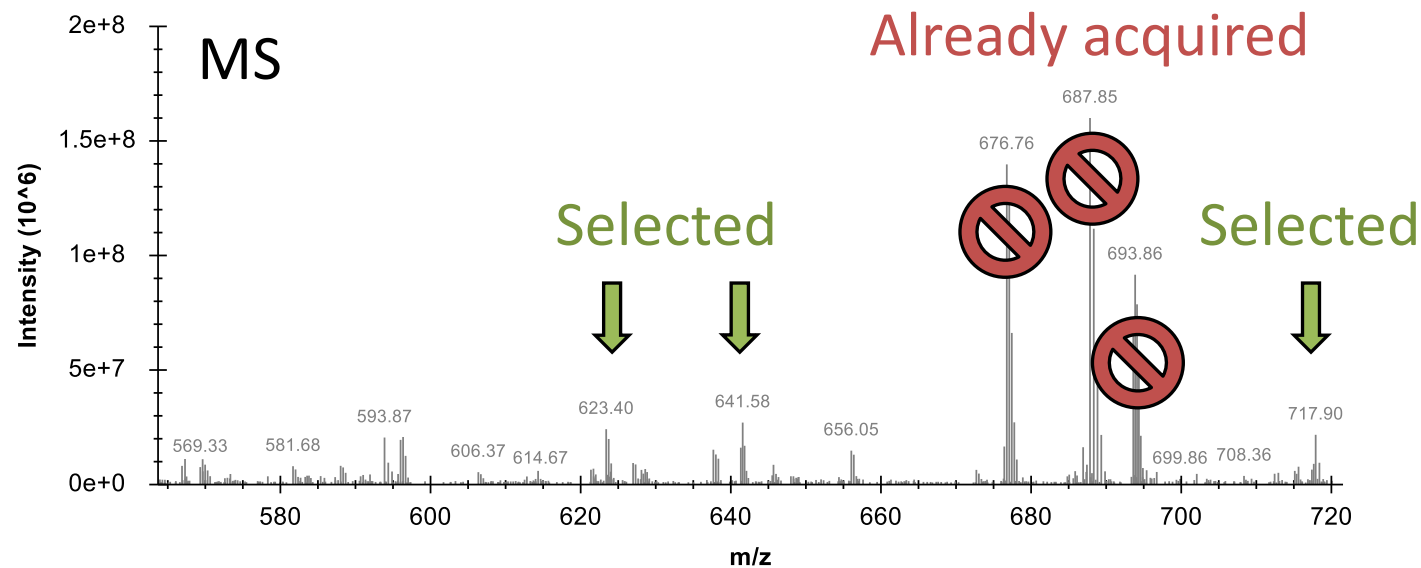- Vacuum pump efficiency
- Detector age

Advanced Active Beam Guide (AABG)

Quadrupole (Q Exactive) or Advanced Quadrupole Technology (AQT)

Inject Flatapole

Transfer Octapole

C-Trap

HCD Cell

S-Lens

Ultra-High-Field Orbitrap Mass Analyzer

HV Lenses

Ion Transfer Capillary

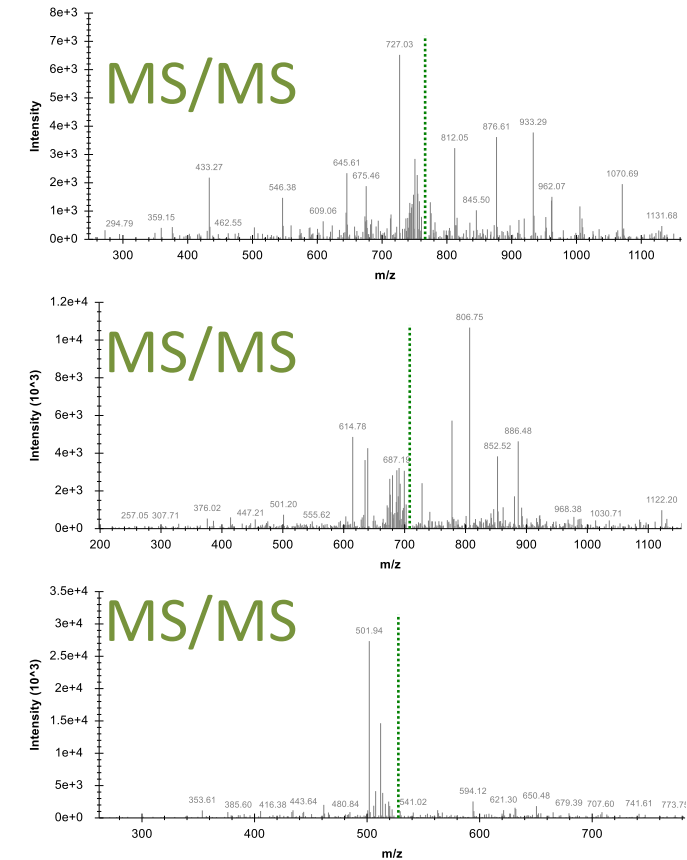ThermoFisher: WS-MS-Q-Exactive-Calibration-Maintenance-iQuan2016-EN.pdf

# MS/MS acquisition prioritizes intense precursor ions in MS1.



In "data-dependent acquisition," the computer chooses high-intensity m/z values from MS that have not recently been subjected to MS/MS.
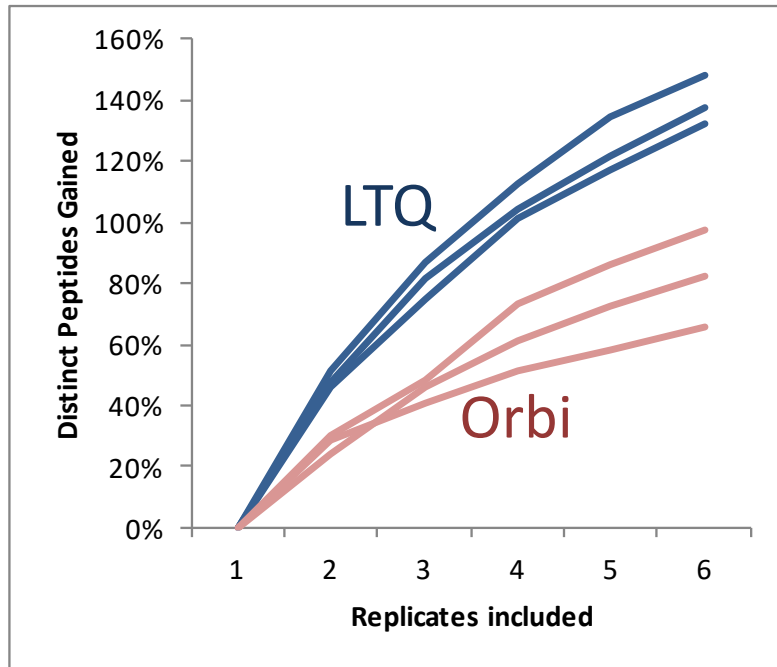
DDA method leads to *stochastic* differences among experiments with the same sample.

# CPTAC: Clinical Proteomic Technology Assessment for Cancer

### Typically stable spectral counts by instrument

| Instrument | Accession | MVH[1] | Rep1 | Rep2 | Rep3 | Rep4 | Rep5 | Rep6 |
|---|---|---|---|---|---|---|---|---|
| LTQ@73 | YDR012W | 1.05 | 28 | 27 | 29 | 23 | 25 | 18 |
| LTQ2@95 | YGL062W | 0.90 | 2 | 2 | 2 | 2 | 5 | 3 |
| LTQc@65 | YLR441C | 1.06 | 12 | 14 | 10 | 8 | 10 | 9 |
| Orbi@86 | YDL124W | 0.64 | 5 | 3 | 3 | 4 | 2 | 3 |
| OrbiP@65 | YLR150W | 0.41 | 4 | 4 | 4 | 4 | 3 | 2 |
| OrbiW@56 | YBR109C | 0.54 | 3 | 4 | 3 | 4 | 2 | 5 |

[1] Natural log of multivariate hypergeometric probability ratio: observed distribution versus expected distribution

DL Tabb et al. *J. Proteome Res*. (2009) 9: 761-776.

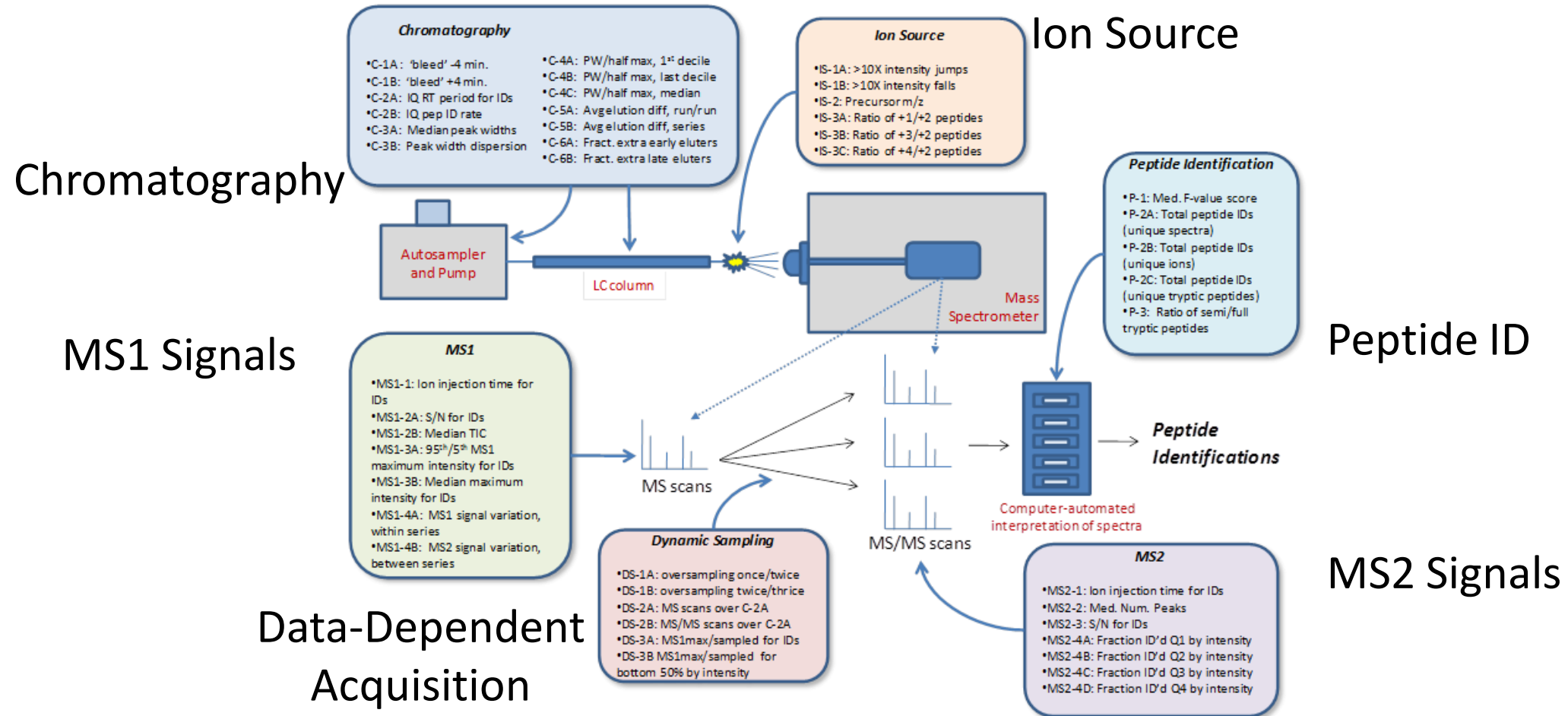# Are the differences we claim reproducible?



- Blue differential genes are found in common by another experiment in this instrument, while orange ones are unique to a single replicate.

- For iTRAQ, the confirmation was required to come from a different multiplex.

DL Tabb et al. *J. Proteome Res*. (2015) 15: 691-706.

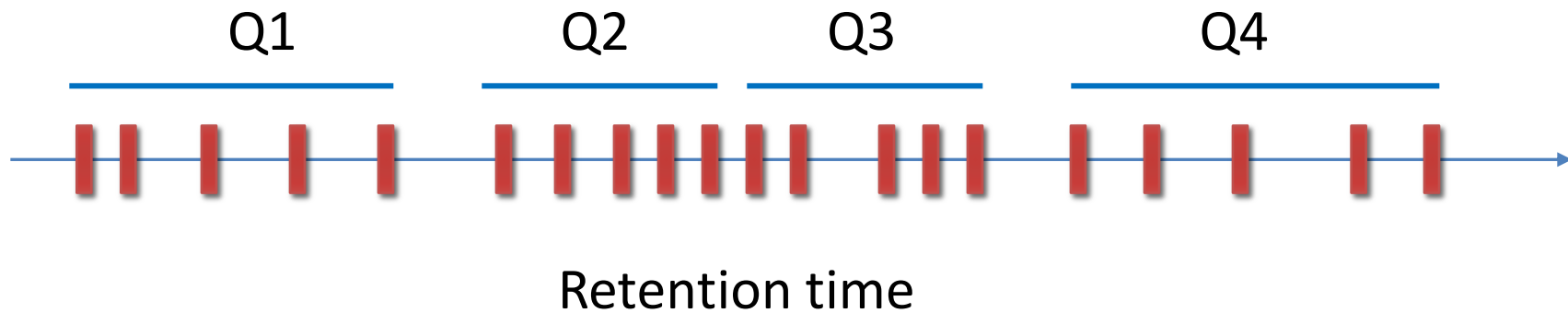# National Institute of Standards and Technology MSQC metrics

# QuaMeter "IDFree" metrics

- 40 metrics categorized into extracted ion chromatograms, retention times, mass spectrometry, and tandem mass spectrometry

- Metrics emphasize quartiles and log ratios to give robust behavior in extreme variation

- Example fields for rapid characterization:
  - XIC-FWHM-Q2: median of peak widths in time
  - RT-TIC-Q$x$: relative duration for TIC integration
  - MS1-Density-Q2: median of peak counts for MS
  - MS2-Freq-Max: Fastest rate of MS/MS acquisition

```
quameter.exe *.raw -MetricsType idfree -OutputFilepath metrics.tsv
```

| RT-Duration | What is the highest scan time observed minus the lowest scan time observed? |
|---|---|
| RT-TIC-Q1 | The interval when the first 25% of TIC accumulates divided by RT-Duration |
| RT-TIC-Q2 | The interval when the second 25% of TIC accumulates divided by RT-Duration |
| RT-TIC-Q3 | The interval when the third 25% of TIC accumulates divided by RT-Duration |
| RT-TIC-Q4 | The interval when the fourth 25% of TIC accumulates divided by RT-Duration |
| RT-MS-Q1 | The interval for the first 25% of all MS events divided by RT-Duration |
| RT-MS-Q2 | The interval for the second 25% of all MS events divided by RT-Duration |
| RT-MS-Q3 | The interval for the third 25% of all MS events divided by RT-Duration |
| RT-MS-Q4 | The interval for the fourth 25% of all MS events divided by RT-Duration |
| RT-MSMS-Q1 | The interval for the first 25% of all MS/MS events divided by RT-Duration |
| RT-MSMS-Q2 | The interval for the second 25% of all MS/MS events divided by RT-Duration |
| RT-MSMS-Q3 | The interval for the third 25% of all MS/MS events divided by RT-Duration |
| RT-MSMS-Q4 | The interval for the fourth 25% of all MS/MS events divided by RT-Duration |



Q1    Q2    Q3    Q4
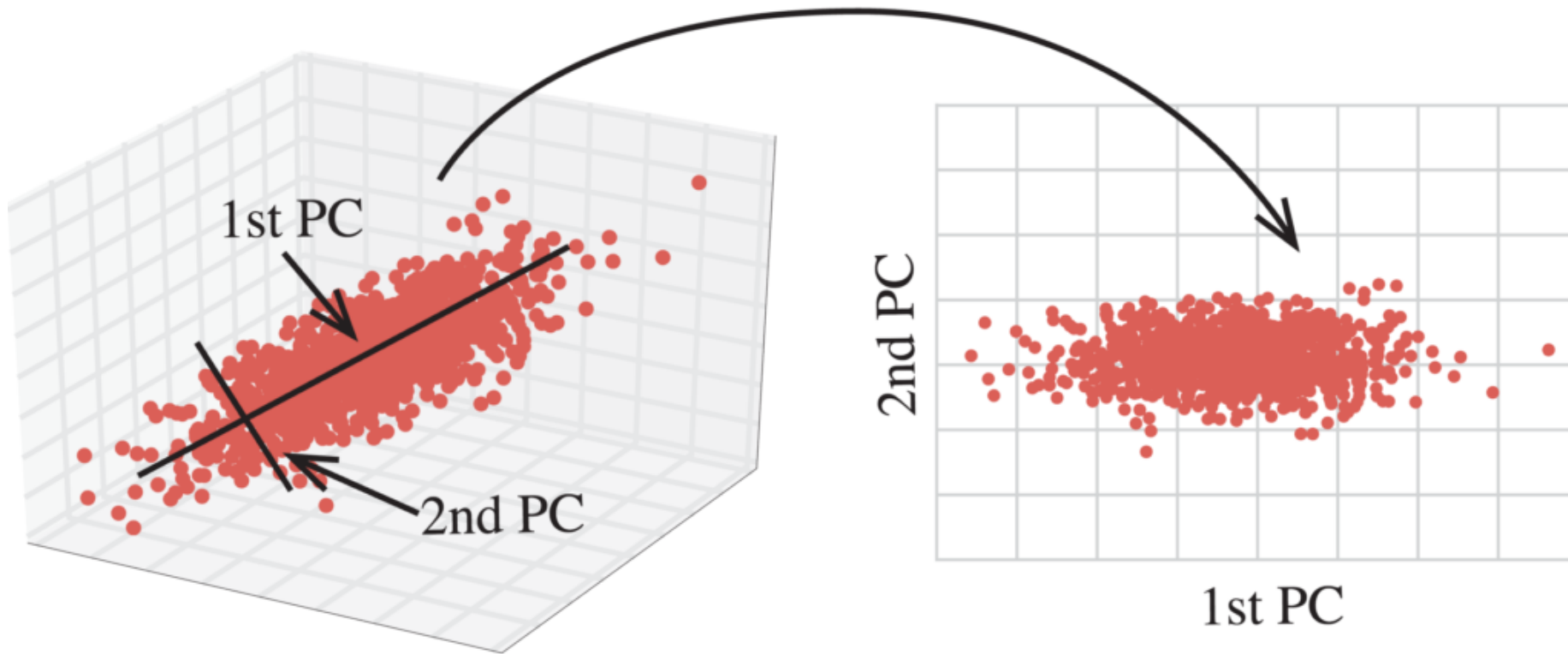
Retention time

# A role for dimensionality reduction

- Measurements may be redundant. They may contain *mutual information* or be *correlated* with each other.

- Principal Components Analysis combines these metrics into components that account for observed variance.
  - "Component:" PCA accepts a table of $n$ metrics for each sample; each component is a linear combination.
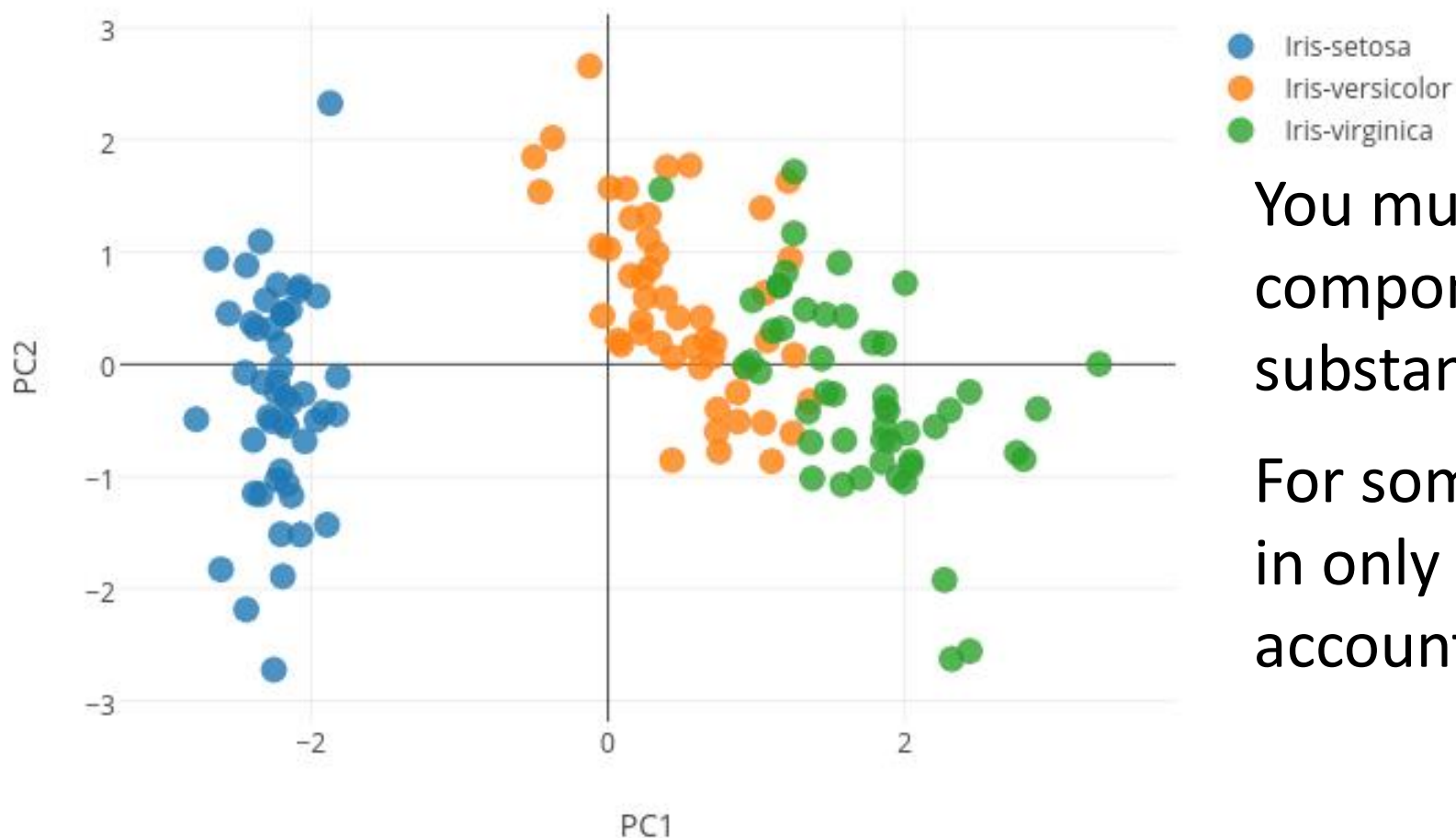  - "Principal:" PCA prioritizes the components by the amount of observed variability that each explains.

http://www.theanalysisfactor.com/tips-principal-component-analysis/

# *Projection pursuit*:
# PCA rotates data in *n*-D space

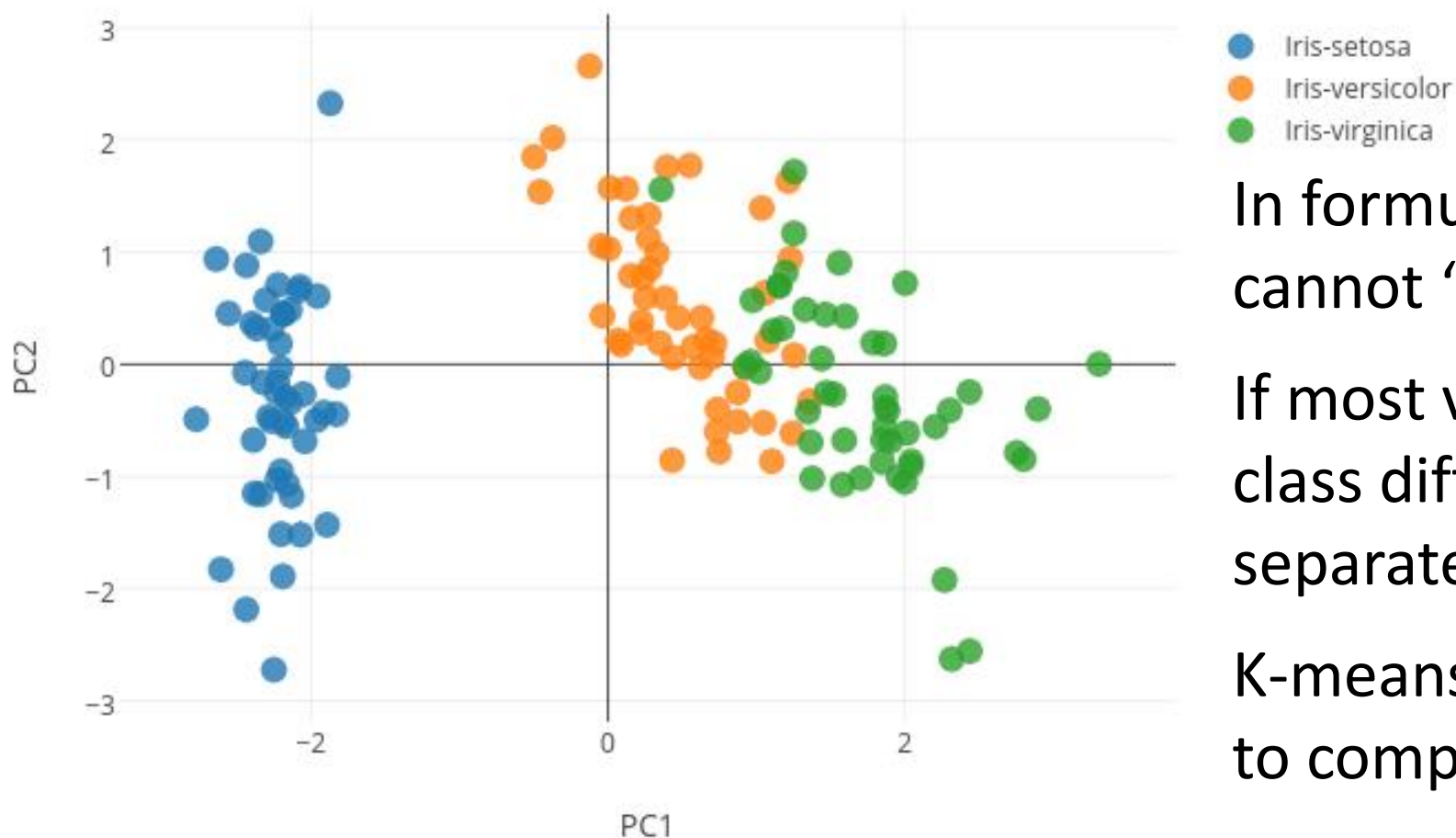# PCA plot typically shows first two components



You must keep enough components to account for substantial variability.

For some experiments, visualizing in only two dimensions may not account for enough variability.

https://plot.ly/ipython-notebooks/principal-component-analysis/

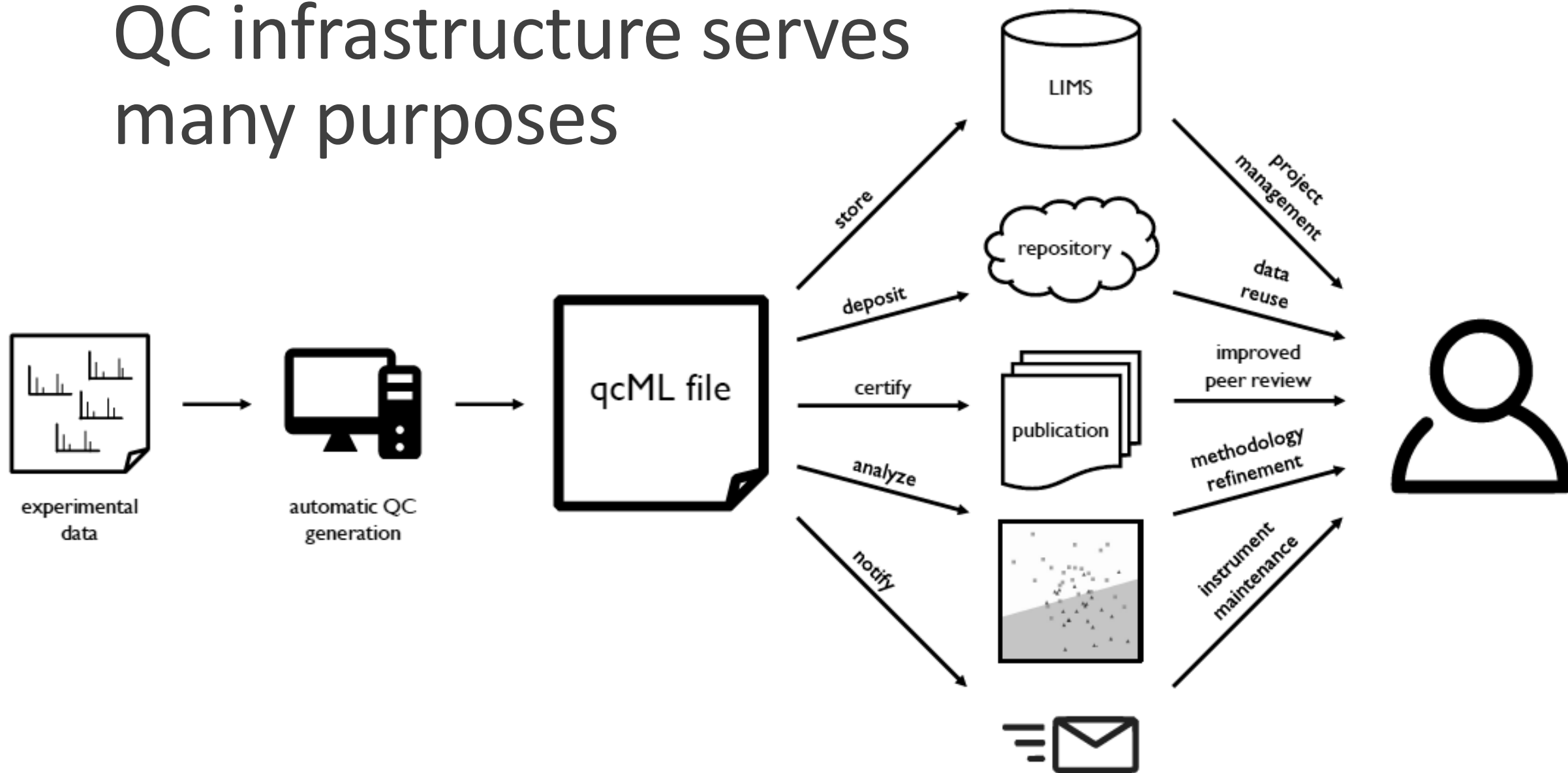# *Unsupervised learning*: PCA may reveal classes



In formulating components, PCA cannot "see" class labels.
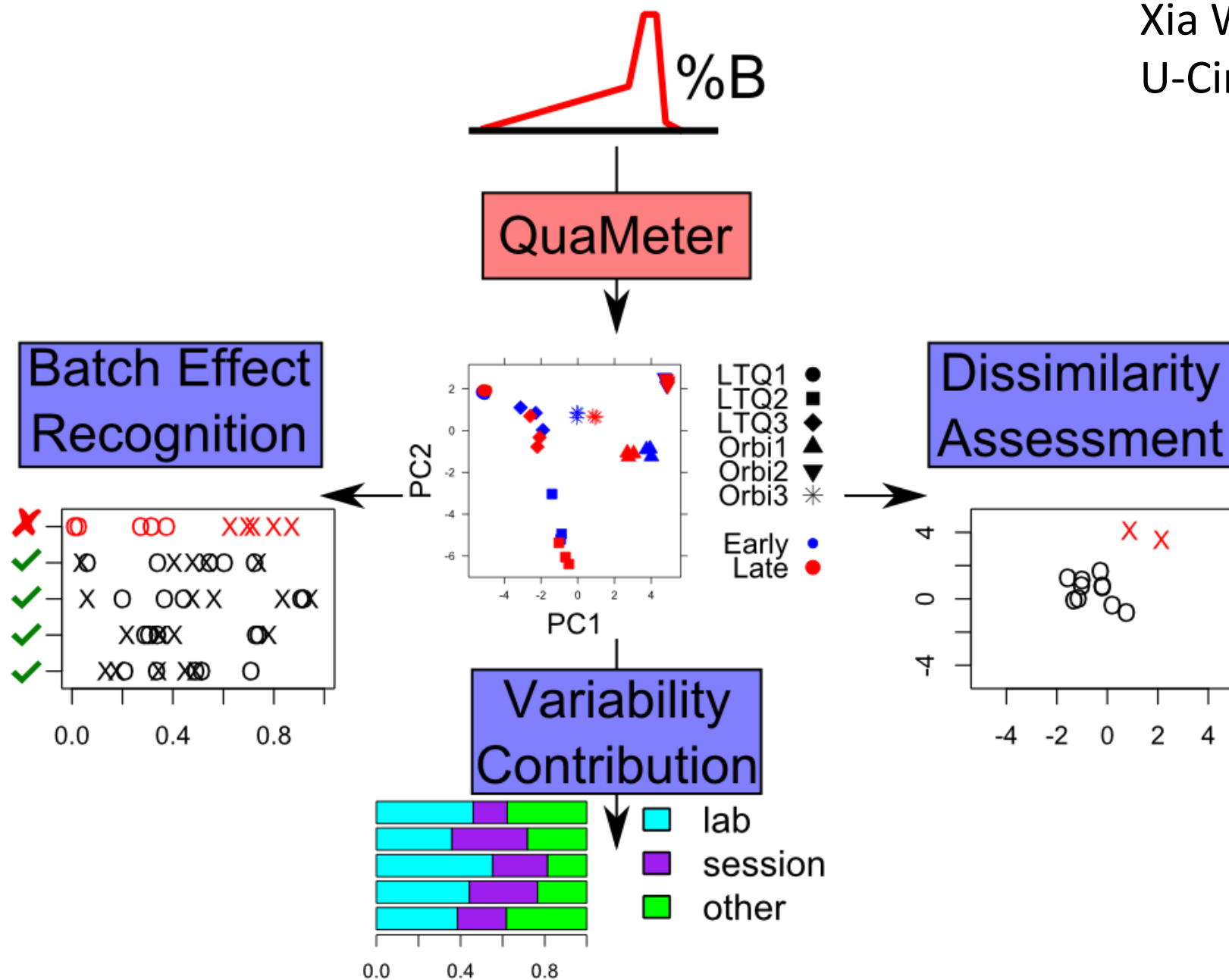
If most variance comes from inter-class differences, the classes will separate in PCA plot.

K-means clustering can be applied to components.

C. Ding and X. He. *ICML '04: Proc. 21$^{st}$ Intl Conf. Machine Learn*
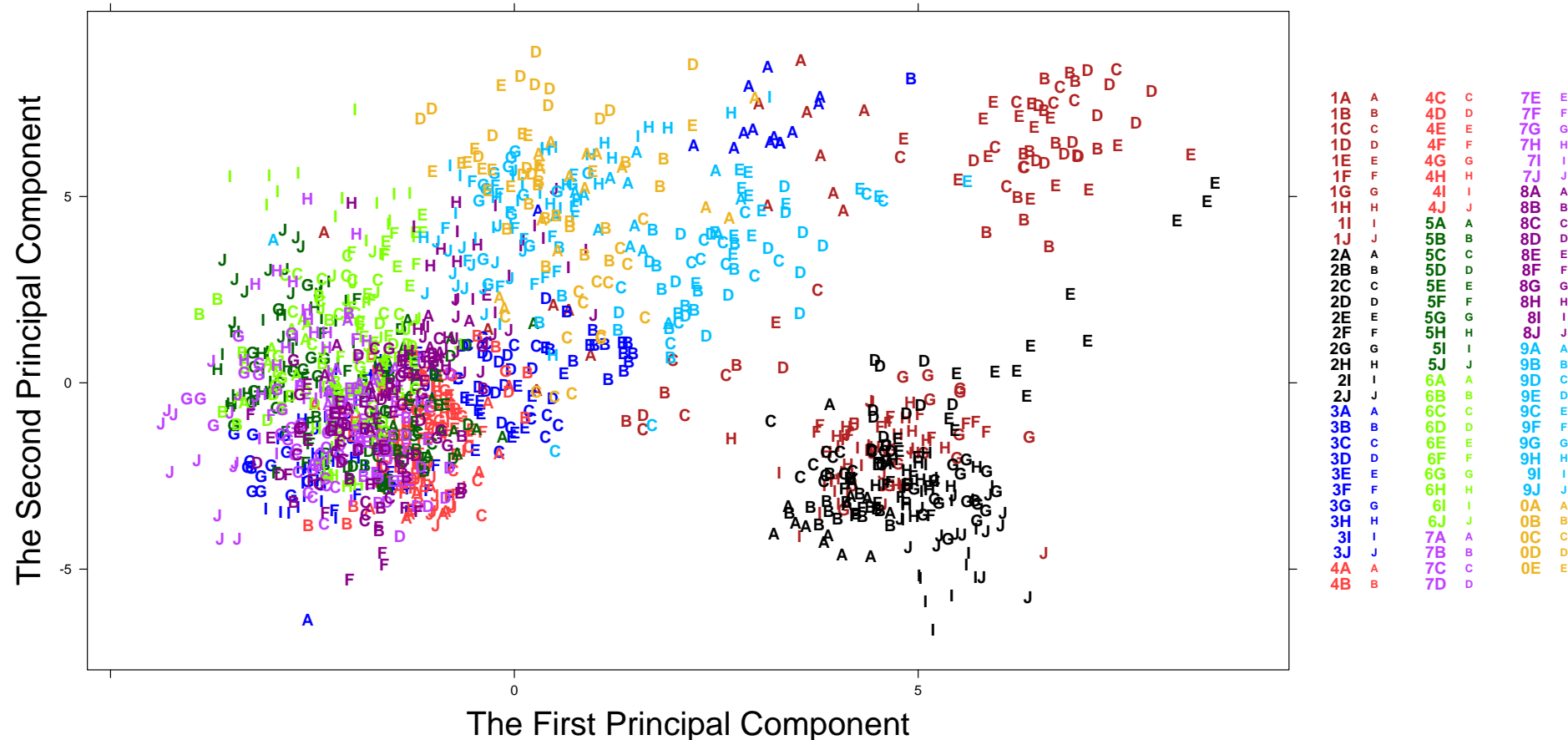
# QC infrastructure serves many purposes



W. Bittremieux et al. *Anal Chem*. (2017) 89: 4474-4479

Xia Wang, U-Cincinnati

Wang et al. (2014) *Analytical Chemistry* 86: 2497-2509.

# CPTAC Colon: 95 bRPLCs of 15 fxns over eight months of operation

Slebbos et al. *Scientific Data* (2015) 2: 150022

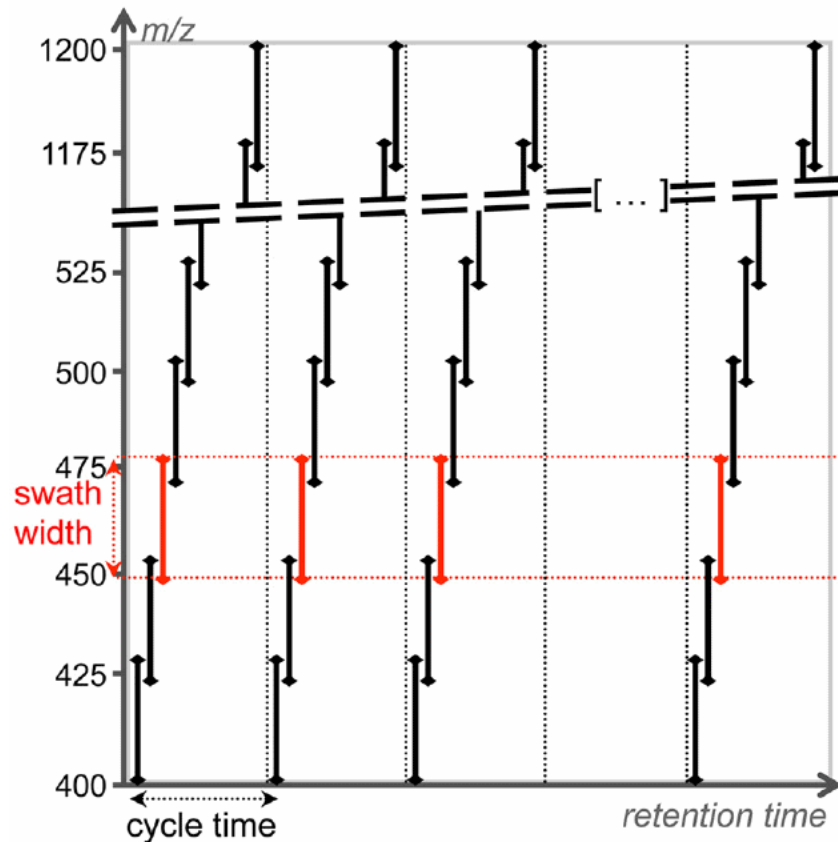Marina Kriek, SUN

# Outliers may dominate variance



- Three *M. tuberculosis* strains were subjected to GeLC-MS.

- Two LC-MS/MS experiments were far from the "cloud."

- We rediscovered that data production was temporarily halted after this pair.

Marina Kriek, SUN

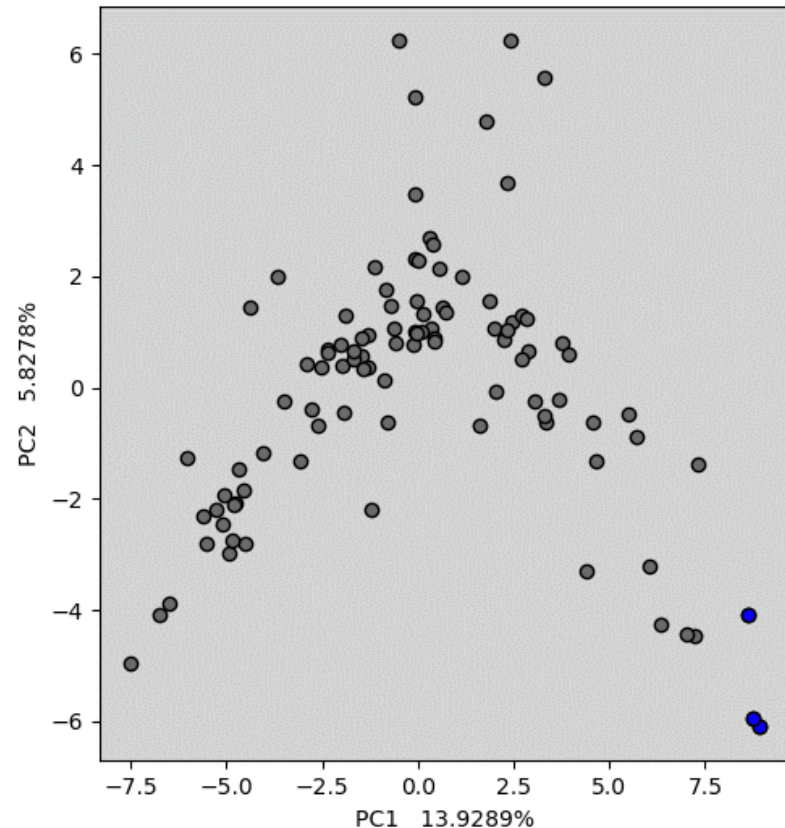# Specializing QC metrics to DIA



The SwaMe project describes SWATH experiments with QC metrics:

- Comprehensive metrics apply to complete experiment.

- Isolation window metrics detail each SWATH individually.

- RT metrics subdivide LC retention to different time intervals.

Marina Kriek, SUN

# Struggling with R? Assurance simplifies QC

PCA of the comprehensive QC metrics



- Assurance enables you to produce QuaMeter or SwaMe metrics and visualize them through a GUI.
  - PCA of metrics highlights outliers.
  - Images interrogate experiments for a particular metric.
  - Machine learning examines longitudinal series data.

https://github.com/marinaPauw/Assurance

# mzQC: designing a HUPO-PSI standard for QC information

▪Tab-separated text is fine for simple tables, but what if a single metric is a matrix of values?

▪Just as mzML represents a text format for mass spec data, an mzQC represents a text format for quality data, easing passage from one tool to another.

▪The use of JSON streamlines the size of mzQC documents.

https://github.com/HUPO-PSI/mzQC

# Takeaways

- Proteomics QC software is growing in maturity.  Bench researchers are no longer required to "eyeball it."

- Statistical models based on QC metrics can recognize outliers, batch effects, and factors contributing variability.

- PCA is ubiquitous and valuable; we can all improve our understanding of how it works!