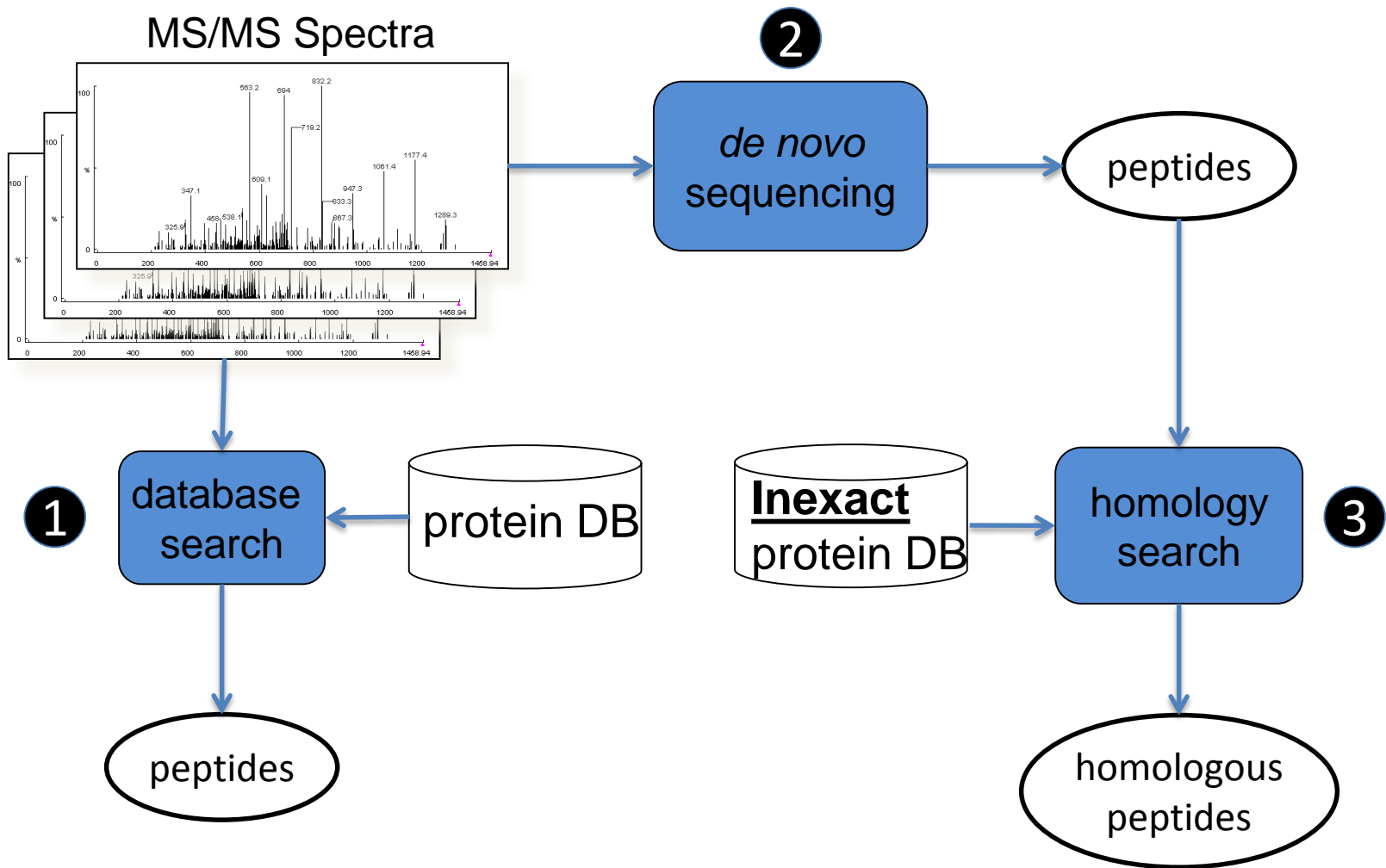# Integrating database search and de novo sequencing to improve the peptide identification

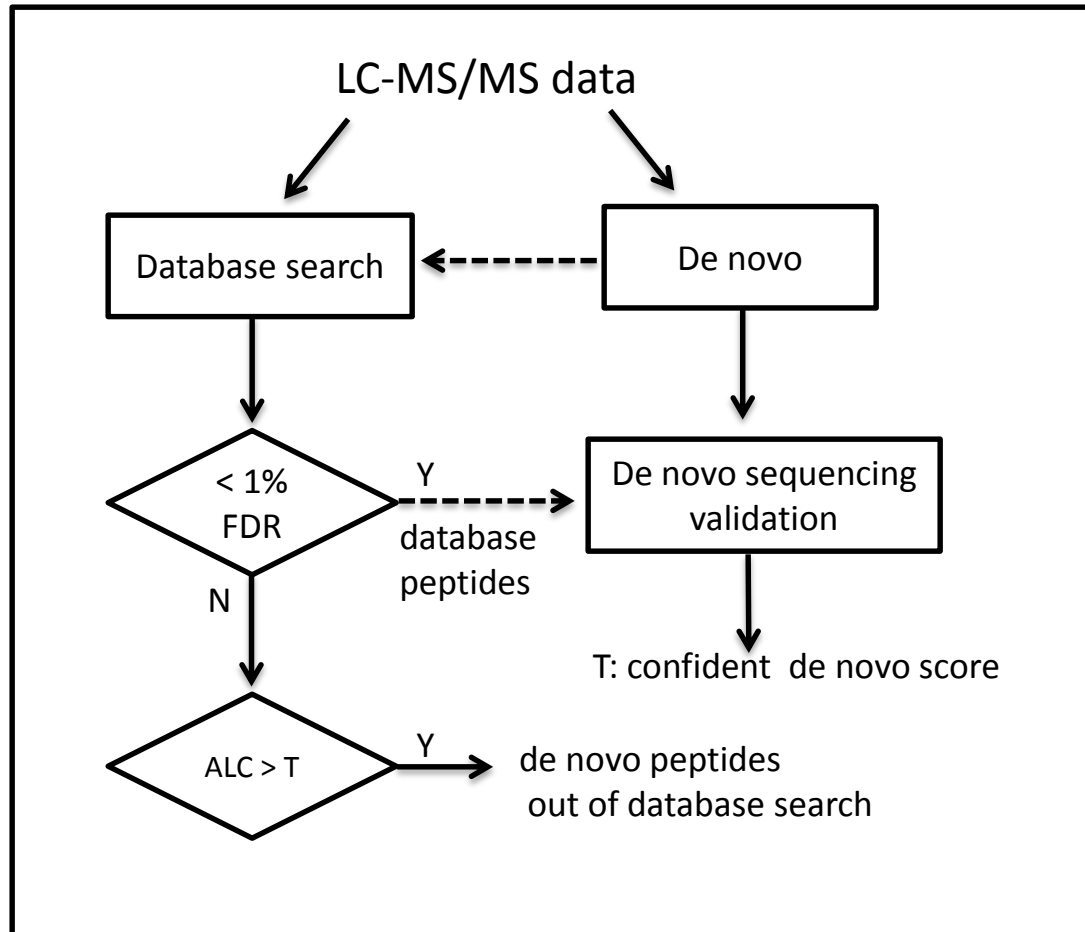Baozhen Shan

Bioinformatics Solutions Inc.

# 1. Approaches for LC/MSMS data analysis
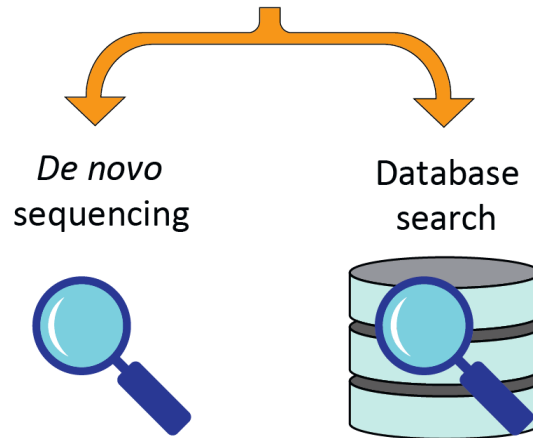


Ma and Johnson. *MCP* (2012) 11: O111.014902

# 2. Integrating de novo and database search

# 3. De novo sequencing improves DB search

- Problem
  - Coverage
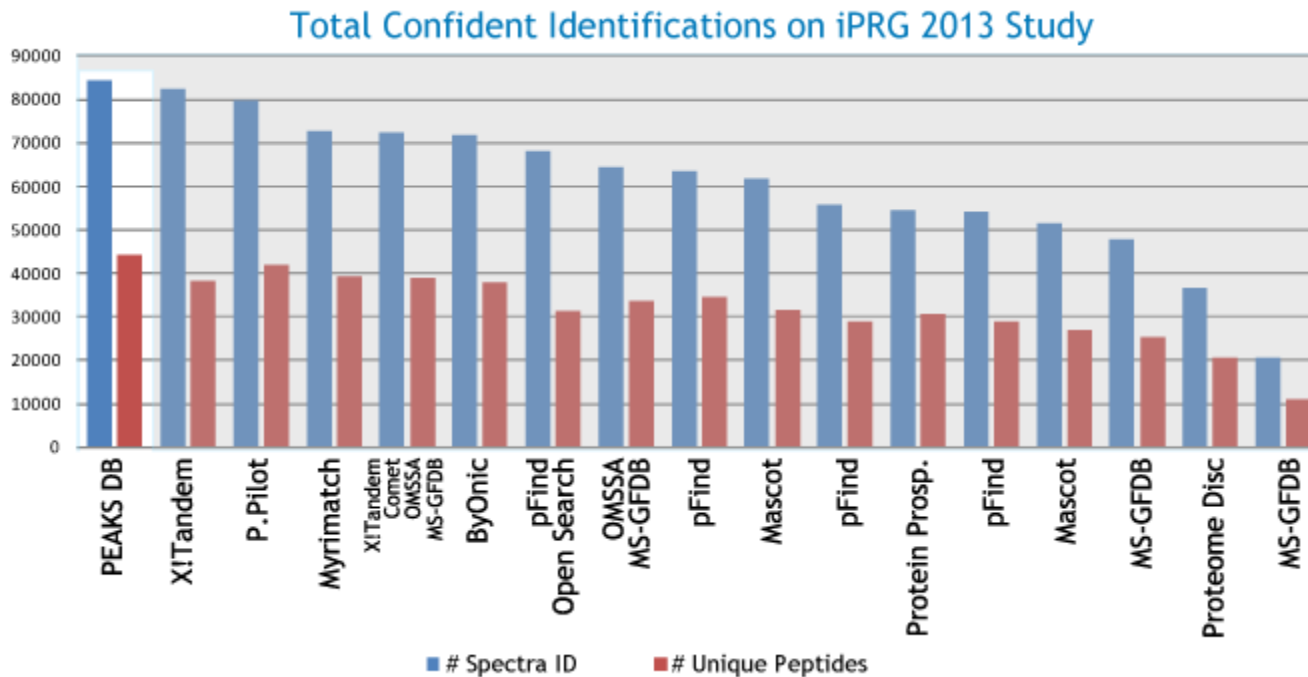  - Modifications
  - Incomplete database

- Solution



De novo
sequencing

Database
search

# 3.1. Good scoring function

- Uses many more factors than other algorithms
  - particularly the similarity between *de novo* and DB sequence
  - many other scoring features considered
- Better separation of true and false means better accuracy and sensitivity.



Zhang et al, *MCP* (2012) 11, M111.010587.2011.

# High sensitivity and accuracy

- ABRF iPRG 2013 study

# 3.2. De novo assisted PTM "Blind Search"
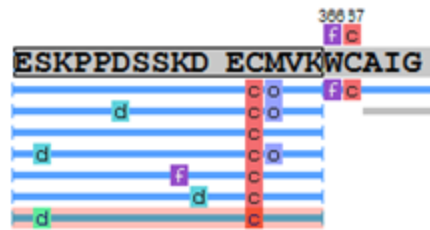
- Search for PTM when there is a tag match.

PEAKS DB



| | ΔM | PTM |
|---|---|---|
| c | +58.01 | Carboxymethyl |
| d | +0.98 | Deamidation (NQ) |
| o | +15.99 | Oxidation (M) |

PEAKS PTM

| | ΔM | PTM |
|---|---|---|
| c | +58.01 | Carboxymethyl |
| d | -18.01 | Dehydration |
| d | +0.98 | Deamidation (NQ) |
| o | +15.99 | Oxidation (M) |
| a | -17.03 | Ammonia-loss (N) |
| p | -18.01 | Pyro-glu from E |
| f | +27.99 | Formylation (TS) |
| f | +27.99 | Formylation |
| a | +42.01 | Acetylation (N-term) |
| s | +21.98 | Sodium adduct |
| p | +68.06 | Piperidination |
| r | +53.92 | Replacement of 2 proton.. |

X. Han et al., *JPR* (2011), 10, 2930-2936.

# ABRF iPRG 2012 study

One spiked peptide

PRDX1_HUMAN



DISLSDYK

3,7-phospho  *vs*  5,7-phospho

# 3.3. SPIDER homology search

*Problem*：de novo errors、database mutations

```
(denovo)   X:     LSCFAK
                       |
(homolog)  Z:     SLAAFK
```
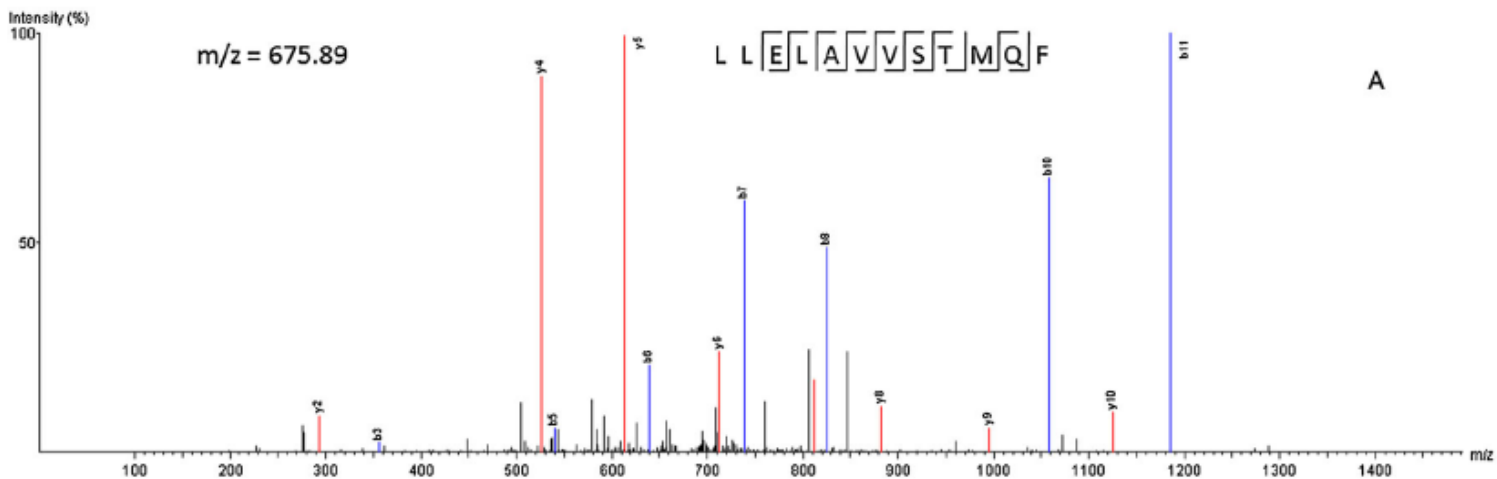
de novo error

```
(denovo)   X:     [LS]C[FA]K
(real)     Y:     [SL]C[AF]K
                   ||    || |
(homolog)  Z:     [SL]A[AF]K
```

mutation

Solution：minimize de novo errors and mutations

Y. Han *et al*, *JBCB* (2005) 3, 697-716.

# Peptides in camel milk

# 4. DB search validates *De novo* sequencing
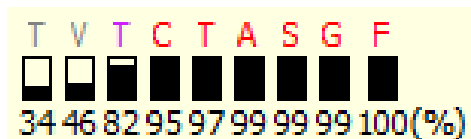
- ## Problem

   De novo sequencing

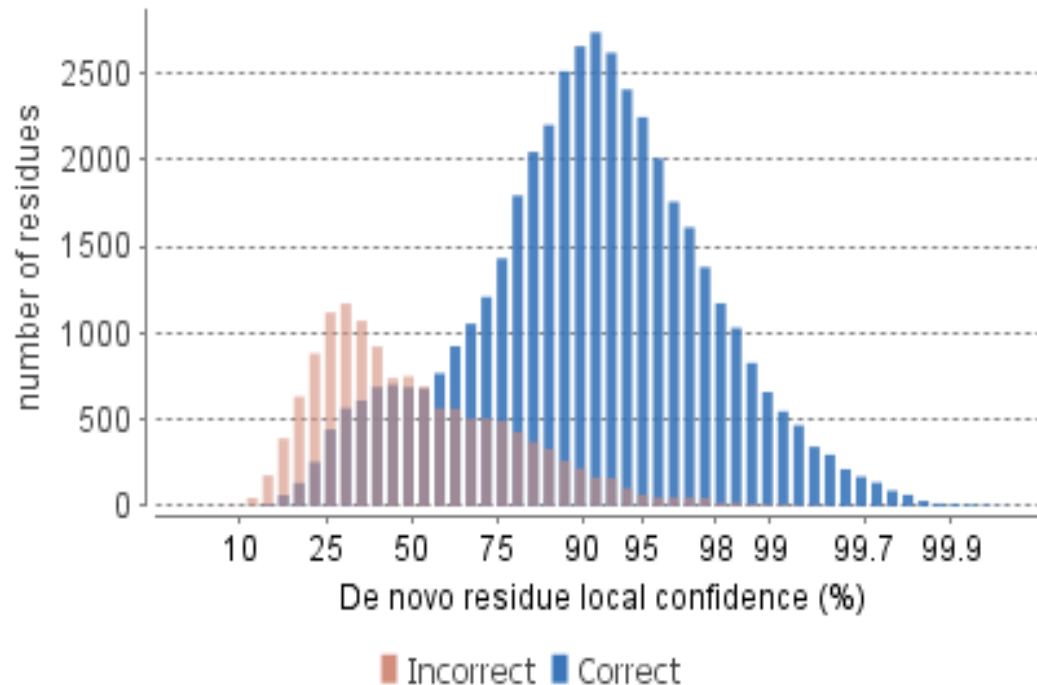   - Ambiguity of de novo sequence
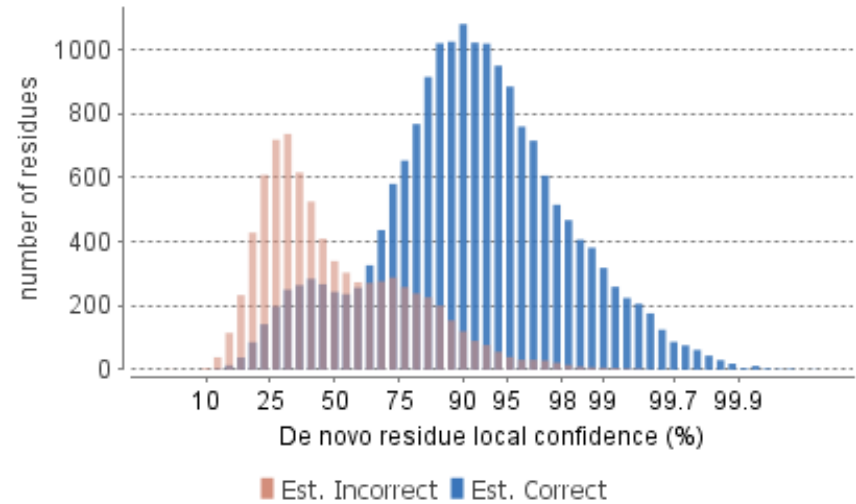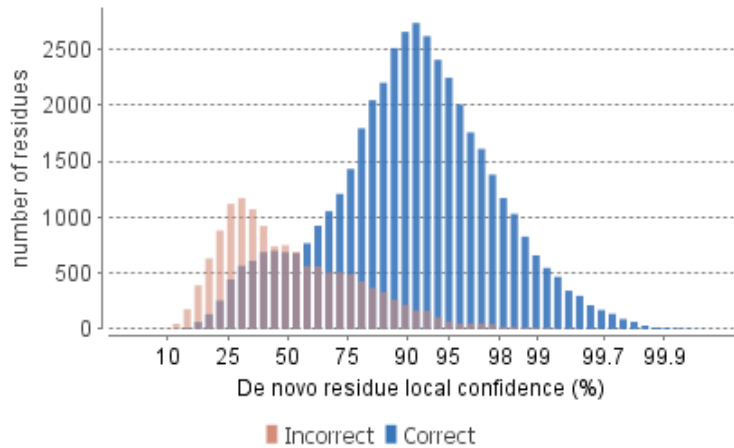   - Partially correct tags



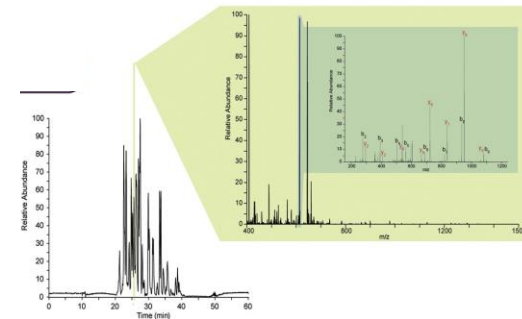- ## Solution

   Local confidence score

# Validation with DB peptides
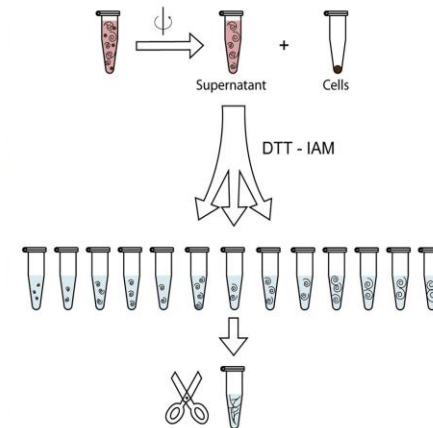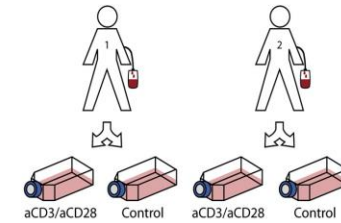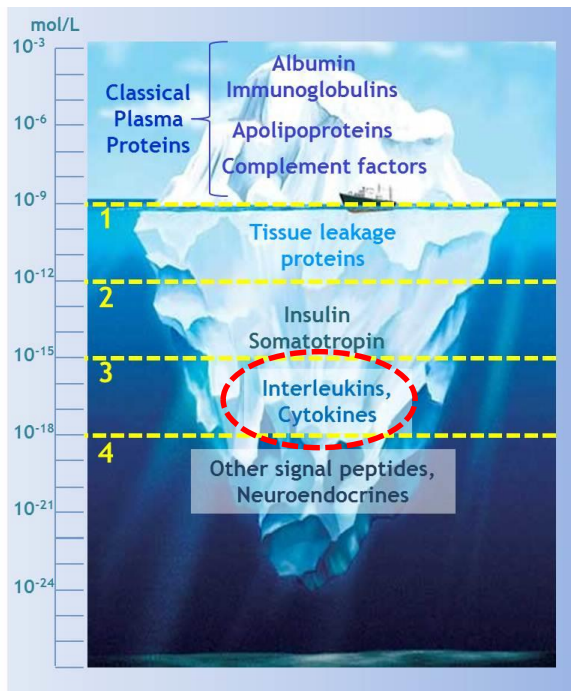
# De novo - only peptides
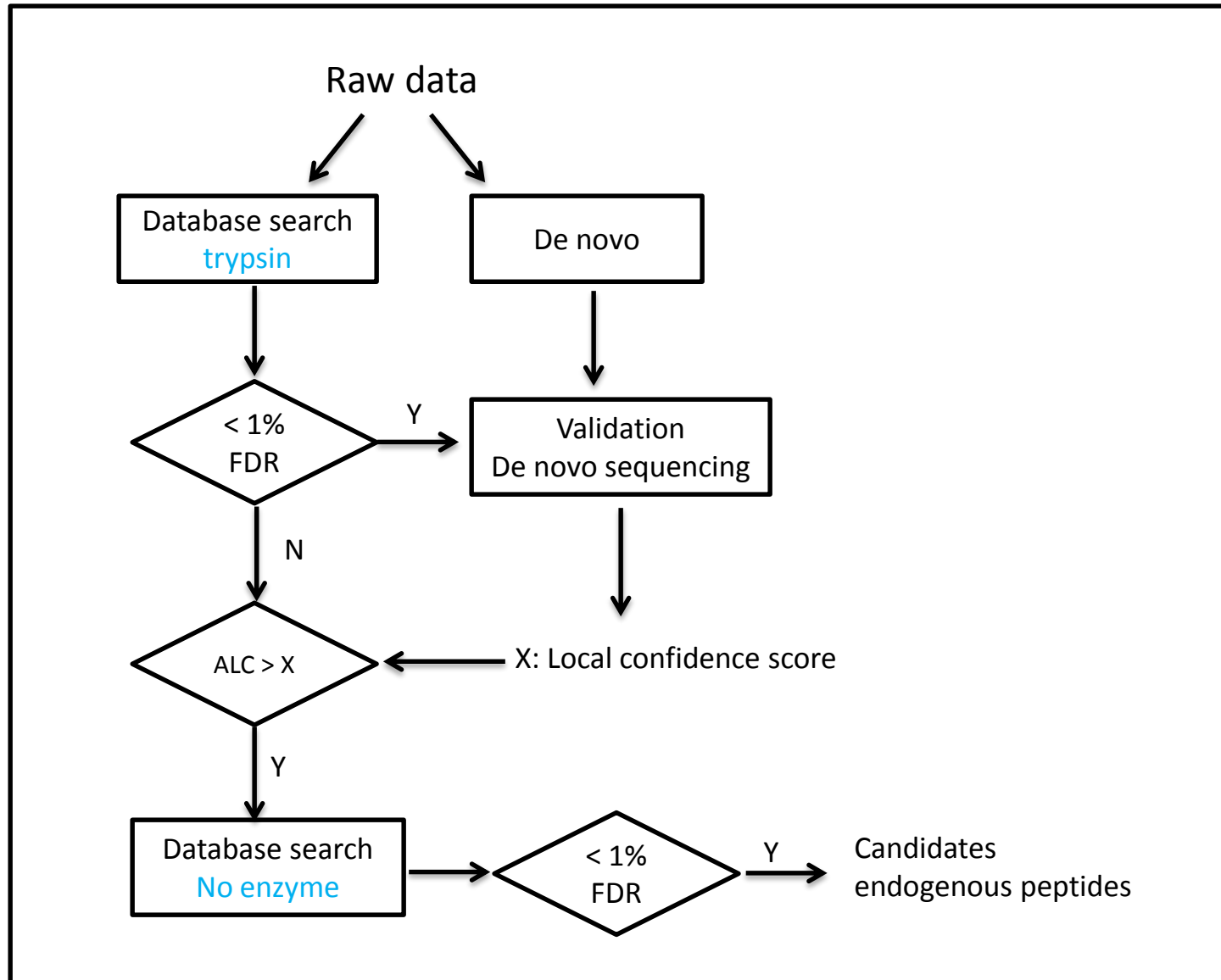


*de novo* peptides validated by DB



score distribution
of *de novo* "only" peptides
with estimated correctness

# 5. Finding endogenous peptides

- **Extracellular proteome**

  low abundance, esp. signaling peptides

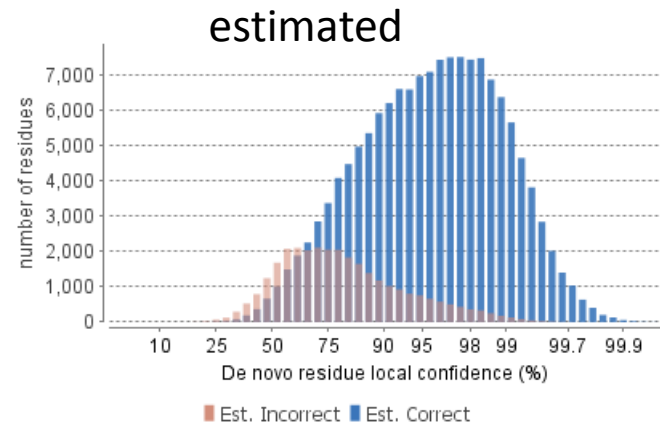  contaminated by intracellular proteins



Finoulst et al. J. Proteomics (2012)

# Workflow for endogenous peptides identification

# Identification of peptides

| | # MS/MS | # peptides |
|---|---|---|
| LC-MS/MS | 1954303 | |
| Database search | 584614 | 18625 |
| De novo sequencing | 15597 | 987 |

validated

estimated



## 70 Human non-tryptic peptides

# Example of an endogenous peptide



Q08554|DSC1_HUMAN

a member of the desmocollin subfamily

extracellular region

# Implement in PEAKS

Raw MS data

*de novo* seq.

PEAKS DB → DB peptides

PEAKS PTM → PTMs & mutations

SPIDER → More mutations

**de novo peptides**

# Acknowledgement

- PEAKS R&D team at BSI
- Collaborations
  - Prof. Bin Ma at UW
  - Prof. Gilles A. Lajoie at UWO
  - Prof. Kaizhong Zhang at UWO
  - Prof. Peter Verhaert at Delft University of Technology
  - Cuijie Zhang at Samuel Lunenfeld Research Institute