# Standards for Proteomics Data Generated by LC-MS-MS

*Ruedi Aebersold, Ph.D.*
*ETH Zurich, Switzerland and*
*Institute for Systems Biology*
*Seattle, Washington*

# Theses:

- Different requirements for data processing, dissemination and storage apply for mass spectrometry applied to the analysis of proteins and proteomes.

- Proteomics is a genomic science and needs to develop "genomics" data analysis/dissemination strategies

# LC-MS/MS as a protein analysis tool

- Relatively low number of proteins analyzed per experiment

- Extensive (biological, manual) validation of data

- Projects centered in single group and focused on specific question

- Data stored in notebook or local computer

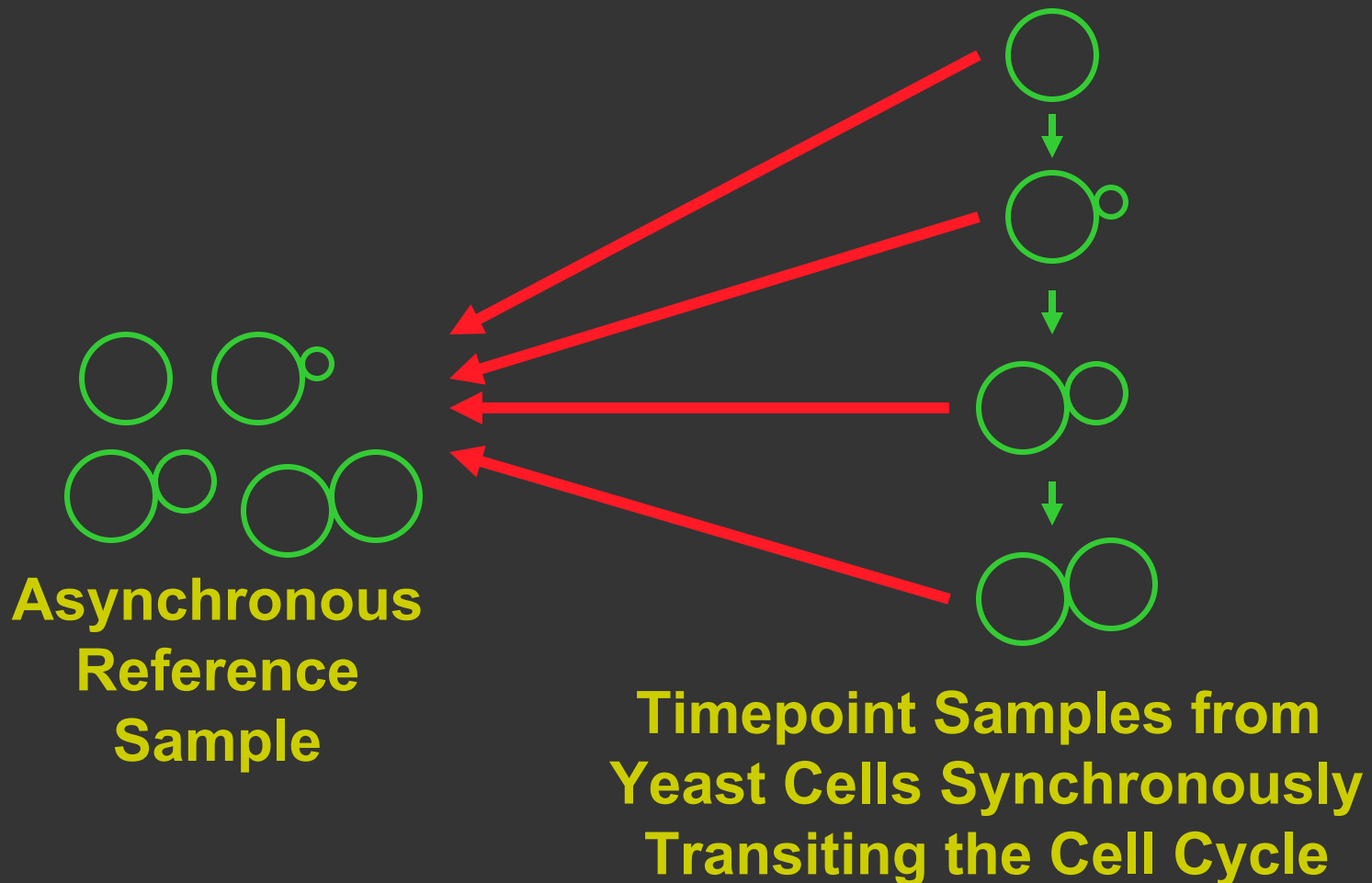- Reports focused on the biological meaning of the data

# LC-MS/MS as a genomic technology

- Many – ideally, all – proteins in a proteome analyzed repeatedly

- Extensive and consistent biological or manual validation of all data impossible

- Value of information increases if data from multiple experiments/groups can be integrated and collectively mined

- Proteomics is a community effort

- Data are collected and organized in relational databases

- Whole data sets should be made accessible/published

# Discussion Points

Many – ideally, all – proteins in a proteome analyzed repeatedly, generating large volumes of data

# Data Summary

|       | T0   | T30  | T60  | T90  | T120 |
|-------|------|------|------|------|------|
| T0    | 1648 | 1095 | 1184 | 1112 | 892  |
| T30   |      | 1523 | 1055 | 1140 | 921  |
| T60   |      |      | 1448 | 1051 | 871  |
| T90   |      |      |      | 1713 | 960  |
| T120  |      |      |      |      | 1229 |

- 2735/6562 proteins quantified across all timepoints (42%)
- 696 proteins quantified in every experiment
- 1513 proteins quantified in at least one timepoint
- 34,400 peptides quantified on average per timepoint
- >1 million mass spectra collected

# Discussion Points

Many – ideally, all – proteins in a proteome analyzed repeatedly, generating large volumes of data

Current status:

- Large volumes of data are being generated to identify relatively small numbers of proteins

- Information from prior experiments is not used, making the process relatively inefficient

Recommendations:

- Improved strategies for more efficient data collection and analysis are required

- To develop those, access to data is essential

# Discussion Points

Extensive biological and/or manual validation of all data impossible

# Protein Identification by MS/MS



protein sample

protein identifications

A B C D

A B C

peptide mixture

peptide identifications

MS/MS spectra

# Output from search algorithm



sort by search score

# Threshold Model



sort by search score

threshold

SEQUEST:
*Xcorr* > 2.0
$\Delta C_n$ > 0.1

MASCOT:
Score > 47

# Difficulty Interpreting Protein Identifications based on MS/MS

- Different search score thresholds used to filter data

- Unknown and variable false positive error rates

- No reliable measures of confidence

# Protein Identification by MS/MS



**MS/MS spectra**

# Amplification of False Positive Error Rate from Peptide to Protein Level



Peptide Level: 50% False Positives

Protein Level: 71% False Positives

# Protein ID False Positive Rate: Control Dataset Examples



Data Filters:
— Publ. threshold model #1
— Publ. threshold model #2
— Statistical model (p ≥ 0.5)
— Statistical model predicted

Control Datasets:
1   18 purified proteins vs. 18+Human (22 runs)
2   Halobacterium vs. Halo+Human (4 runs)
3   Halobacterium vs. Halo+Human (45 runs)

# False Positive Error Rates among Single-hit Proteins

| Data Filter | Control Dataset | | |
|---|---|---|---|
| | 1 | 2 | 3 |
| Publ. Threshold model #1 | 11% | 37% | 67% |
| Publ. Threshold model#2 | 14% | 32% | 82% |

Control Datasets:
1    18 purified proteins vs. 18+Human (22 runs)
2    Halobacterium vs. Halo+Human (4 runs)
3    Halobacterium vs. Halo+Human (45 runs)

# Serum Protein Identifications from Large-scale (~375 run) Experiment

| Data Filter | # ids | # non-single hits | # single-hits |
|---|---|---|---|
| Publ. Threshold model#1 | 2257 | 359 | 1898 |
| Publ. Threshold model #2 | 2742 | 441 | 2301 |
| Statistical model, p$\geq$ 0.5 (*predicted error rate: 7%*) | 713 | 511 | 202 |

# Consistency of Manual Validation of SEQUEST Search Results

Manual Authenticators →

Search Results



Correct Validation   Incorrect Validation   Validation Withheld

# Discussion Points

Extensive (biological, manual) validation of all data impossible

Current status:

- Peptide and protein identifications are largely made based on threshold model
- Manual validation is often used as "gold standard"

Recommendations:

- Develop, validate and use statistical models that calculate accurate false positive and false negative error rates for peptide AND protein identifications
- Discourage manual validation of spectra as "gold standard".
- Tools should be transparent and generally available

# Discussion Points

- Value of information increases if data from multiple experiments/groups can be integrated and collectively mined

- Proteomics is a community effort

- Data are collected and organized in relational databases

IFN-treated

Mock-treated

C12

ICAT label

C13

C12 C13

HPLC-MS/MS

Wei Yan et al. MCP 2004

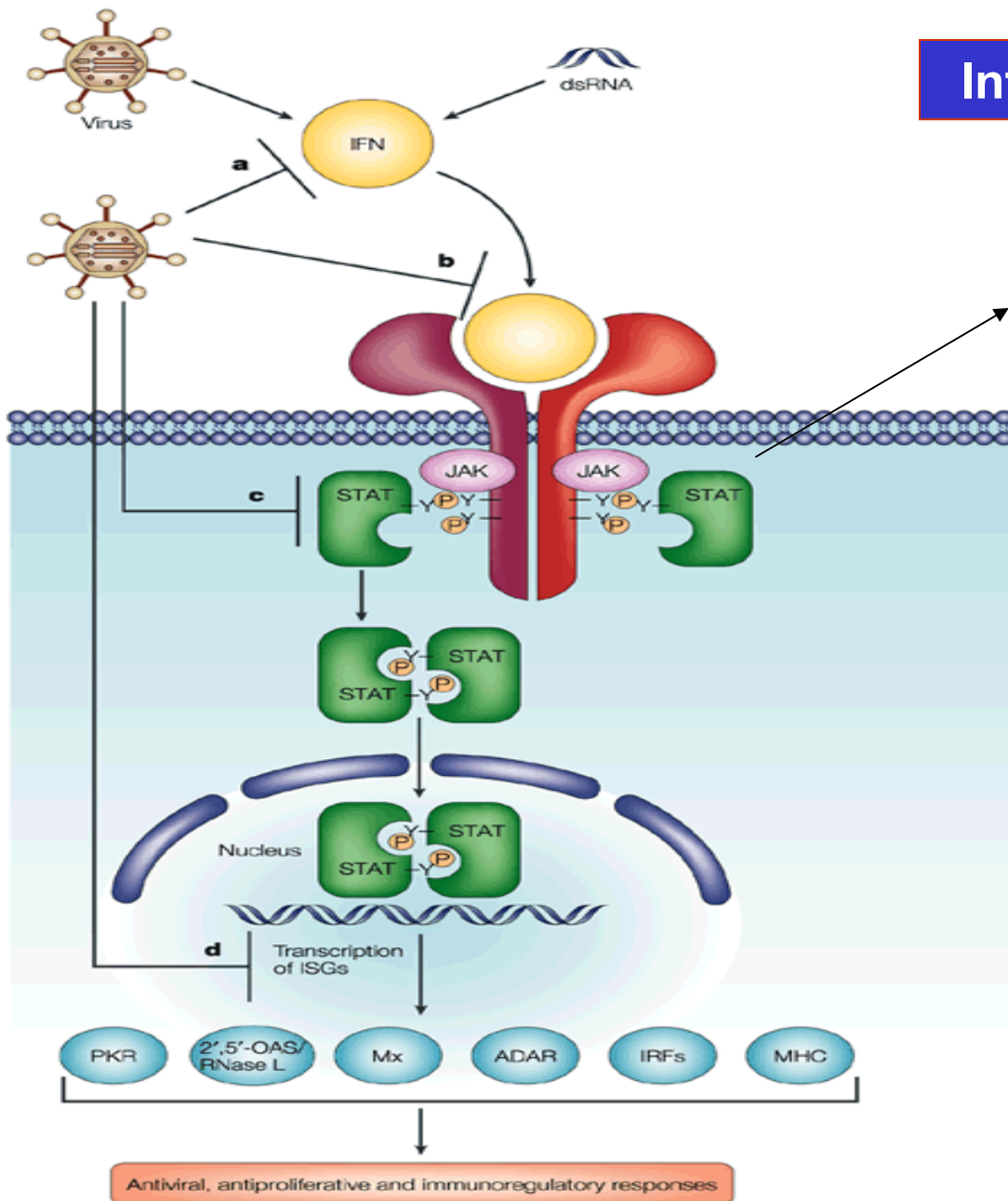| Name | Cellular pathway | Probability | ASAPRatio Mean | ASAPRatio Std. |
|---|---|---|---|---|
| DNAH11: dynein, axonemal, heavy polypeptide 11 | moto protein complex | 0.94 | 9999 | -1 |
| UBE2L6: ubiquitin-conjugating enzyme E2L 6 | ubiquitination and protein degradation | 0.57 | 9999 | -1 |
| **IFIT1: interferon-induced protein with tetratricopeptide repeats 1** | unknown and ESIs | 0.48 | 9999 | -1 |
| GPR111: G protein-coupled receptor 111 | G-protein coupled receptor and G-protein signaling | 0.63 | 21.270 | 4.741 |
| PASK PAS domain co | | | | 1.024 |
| **ADRM1: adhesion re** | | | | 1.043 |
| CSA_PPIasePEPTIDY | | | | 1.070 |
| AHCY: S-adenosylhom | | | | 0.936 |
| **IFIT4: interferon-indu** | | | | 0.794 |
| FLJ32915: hypothetica | | | | 4.883 |
| GNB1: guanine nucleo | | | | 0.133 |
| **G1P2: interferon, alp** | | | | 0.661 |
| MTP: microsomal trigly | | | | 0.751 |
| PLCD1: phospholipase | | | | 0.116 |
| CD7: CD7 antigen (p4 | | | | 2.204 |
| **PRKR protein kinase** **dependent** | | | | 0.659 |
| KIAA1276: KIAA1276 | | | | 0.058 |
| NUDT2: nudix (nucleo | | | | 0.224 |
| CABC1: chaperone, A | | | | 1.659 |
| ACACA: acetyl-Coenz | | | | 0.259 |
| KNS2: kinesin 2 60/70 | | | | 0.335 |
| LOC151636: rhysin 2 | cytoskeleton and intracellular transport? | 1 | 2.975 | 0.231 |
| M96: likely ortholog of mouse metal response element binding transcription factor 2 | transcription | 0.98 | 2.923 | 0.390 |
| ETFA: electron-transfer-flavoprotein, alpha polypeptide (glutaric aciduria II) | electron transfer | 0.45 | 2.890 | 0.484 |
| **NMI: N-myc (and STAT) interactor** | signaling pathway; transcription | 0.57 | 2.875 | 0.138 |
| GSA7: ubiquitin activating enzyme E1-like protein | ubiquitination and protei | 0.98 | 2.844 | 0.663 |
| MGC3207: hypothetical protein MGC3207 | | 0.61 | 0.499 | 0.071 |
| SPK: symplekin | | 1 | 0.496 | 0.029 |
| KRT10: keratin 10 (epidermolytic hyperkeratosis; kerat plantaris) | | 0.97 | 0.495 | 0.055 |
| SARDH: sarcosine dehydrogenase | | 0.98 | 0.484 | 0.008 |
| TRA1: tumor rejection antigen (gp96) 1 | | 1 | 0.452 | 0.165 |
| GPS1: G protein pathway suppressor 1 | | 0.98 | 0.455 | 0.138 |
| SRRM2: serine/arginine repetitive matrix 2 | | 0.82 | 0.434 | 0.224 |
| KIAA0007: KIAA0007 protein | | 1 | 0.426 | 0.014 |
| FACL4: fatty-acid-Coenzyme A ligase, long-chain 4 | | 0.98 | 0.416 | 0.081 |
| FXR2: fragile X mental retardation, autosomal homolog | | 0.95 | 0.391 | 0.074 |
| TUBA6: tubulin alpha 6 | | 1 | 0.383 | 0.165 |
| CPSF4: cleavage and polyadenylation specific factor 4 | | 0.96 | 0.378 | 0.154 |
| MAPRE1: microtubule-associated protein, RP/EB fami | | 0.98 | 0.339 | 0.016 |
| OAT: ornithine aminotransferase (gyrate atrophy) | | 0.98 | 0.331 | 0.018 |
| PPGB: protective protein for beta-galactosidase (galact | | 1 | 0.323 | 0.084 |
| WNT9A: wingless-type MMTV integration site family, mem | | 0.99 | 0.316 | 0.091 |
| FASN: fatty acid synthase | lipid and fatty acid metabolism | 0.99 | 0.304 | 0.100 |
| Ig lambda chain C regions | immune response | 0.98 | 0.265 | 0.110 |
| G2AN: alpha glucosidase II alpha subunit | carbohydrate metabolism | 1 | 0.198 | 0.033 |
| Hypothetical protein FLJ21140 | unknown | 0.71 | 0.043 | 0.064 |
| KRT6: keratin 6 | cytoskeleton and intracellular transport | 1 | 0.003 | 0.008 |
| MIG-6: Gene 33/Mig-6 | signaling pathway | 0.99 | 0.000 | -1.250 |
| HIC1: hypermethylated in cancer 1 | transcription suppression | 0.94 | 0.000 | -1.250 |

| | S100 | P100 | P3 | Sum | Unique ID |
|---|---|---|---|---|---|
| $P \geq 0.9$ | 523 | 270 | 671 | 1464 | 1113 |
| $P \geq 0.4$ | 590 | 330 | 748 | 1668 | 1272 |

54 IFN-induced proteins (2-fold)

15 previously reported

39 novel

23 IFN-repressed proteins (0.5-fold)

# Lots of data -what does it mean?

# Interferon (IFN) Pathway

2.215 ± 0.079

| | IFN / Mock |
|---|---|
| PKR | 3.963 ± 0.659 |
| 2',5'-OAS | 2.460 ± 0.076 |
| Mx | 2.359 ± 0.149 |
| ADAR | 1.398 ± 0.118 |
| IRFs | Not identified |
| MHC | |
| β-2-microglobulin (MHC I) | 2.768 ± 0.583 |
| IFI-30 (MHC II) | 2.219 ± 0.183 |

Katze et al (2002) 2: 675

Nature Reviews | Immunology
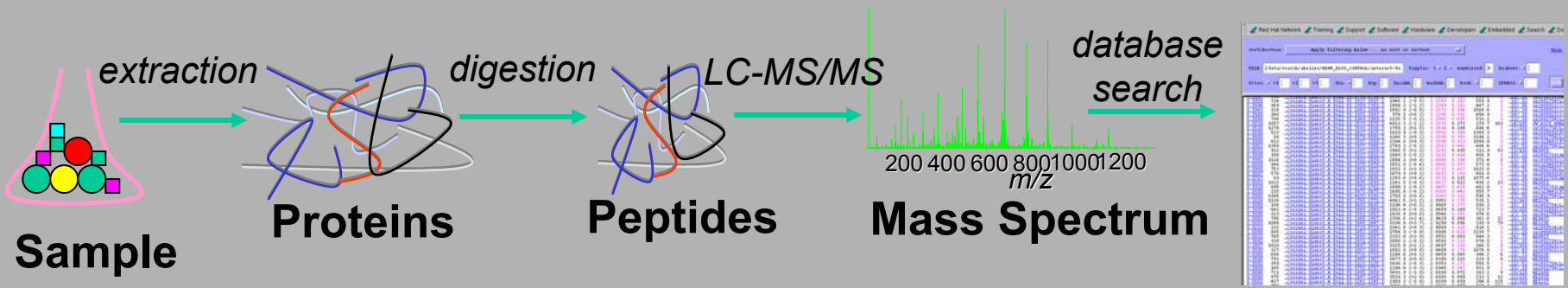
GO Analysis of Interferon regulated proteins

- Single proteomics datasets tend to rediscover the known

- New insights can be made from the comparison of many datasets

# From Peptides to Genome Annotation

# Discussion Points

- Value of information increases if data from multiple experiments/groups can be integrated and collectively mined

- Proteomics is a community effort

- Data are collected and organized in relational databases

Current status:

- Very little proteomics data publicly accessible

- Publications usually only show conclusions but not data

Recommendations:

- Develop and support infrastructure for data sharing and mining

- Make data access condition for publication

# Summary

    If proteomics is to truly operate as a discipline of the genomic sciences, data processing, management and dissemination strategies proven in other fields of genomics must be applied. These include:

- Statistical validation of large data sets

- Providing community access to all data (not just selected data points)

- Providing transparent tools for data processing to community