

Biofluid Proteomics: Potential, Pitfalls, and Solutions



1:40 p.m. - 2:05 p.m.

Biofluid Proteomics: Potential, Pitfalls, and Solutions

Niels Heegaard, M.D., D.Sc. (Med)

Director, Department of Clinical Biochemistry and Autoimmunology

Statens Serum Institut (Denmark)

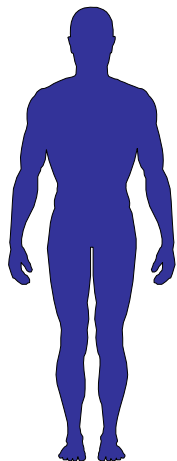
2009 BIOSPECIMEN RESEARCH
NETWORK SYMPOSIUM

March 16-18, 2009

Potential:

- Enormous complexity of the protein and peptide universe present in humans are beginning to be matched by analytical methods
- Promises to give us new biomarkers for diagnosis & individualized medicine

Proteome potential



*~ 20 k genes =>
~ 1,000 k proteins*

*Concentrations vary
(log 7-9)*

*Different over time
and location.*

*Reflects events
in the body*

Biomarker potential

Stratification

*Predisposition to
disease*

*Early indicator of
disease*

Disease type

Drug selection

Dose selection

Toxicity avoidance



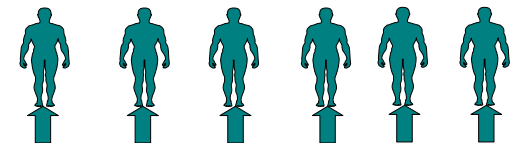
Monitoring

Disease initiation

Progression

Severity

Drug choice/effect



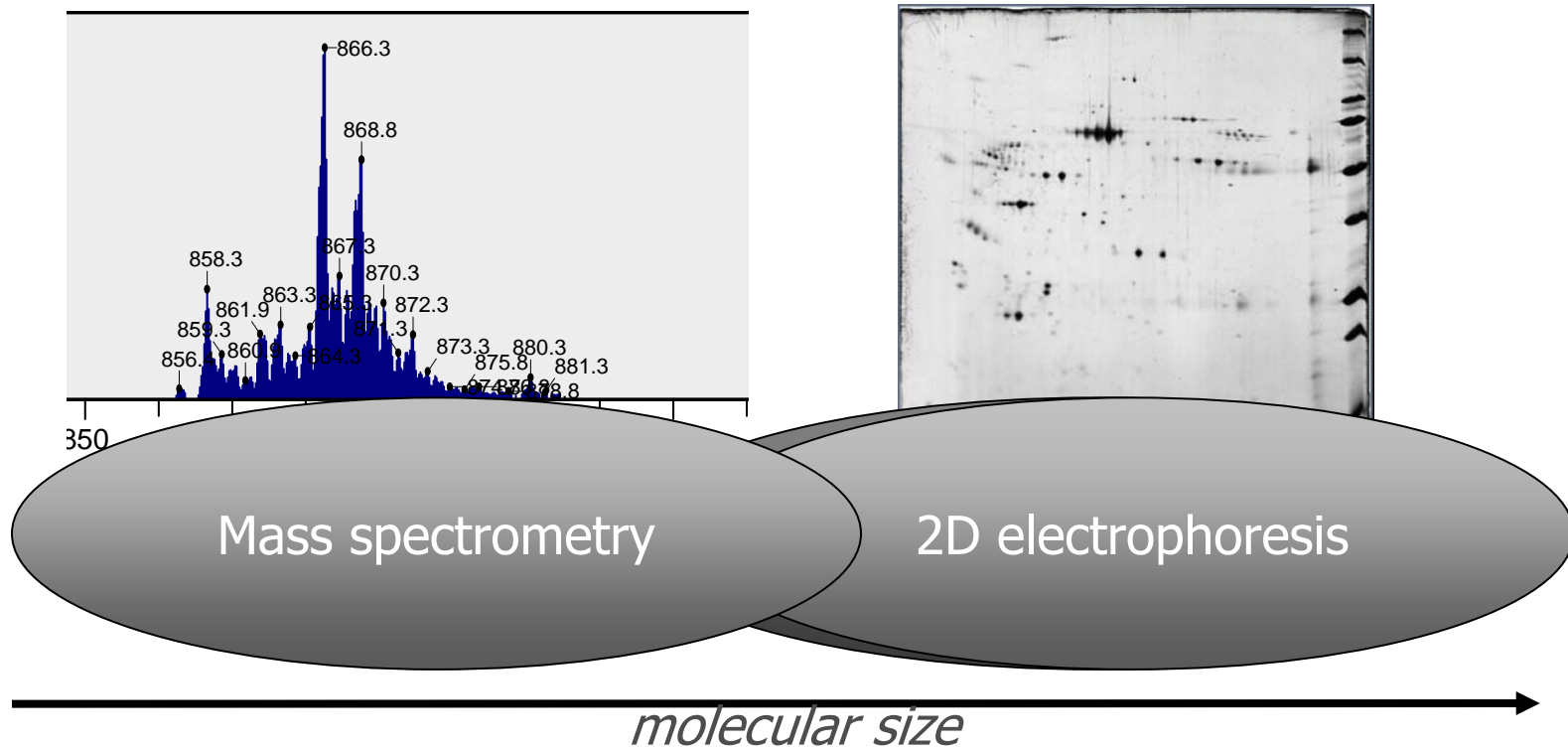


Evolution of protein analytical methods

- Zone gel electrophoresis
 - *5-10 bands*
- Immunoelectrophoresis
 - *20-50 proteins*
- SDS-polyacrylamide gel electrophoresis
 - *50-100 protein bands*
- Two-dimensional gel electrophoresis – *more than 1000 proteins*
- Mass spectrometry - *~30,000 data points/sample; 100 of peptides/proteins*
- Two-dimensional on-line separations (e.g. LC-mass spec)
 - *1,000 of proteins*

Biofluid Proteomics: Potential, Pitfalls, and Solutions

- New biomarkers in the peptide and protein universe (proteomes) of biological fluids?



<i>Output:</i>	30.000 m/z intensity values	up to 2.000 spots
<i>Years in use:</i>	~10	34
<i>New markers in clinical use:</i>	None	1-2



Issues in multiparameter analyses of biofluids for biomarker discovery

- Seeing differences when there are no relevant differences
- Not seeing differences when there are relevant differences
- Seeing differences but not the relevant ones

Variability in biofluid analysis

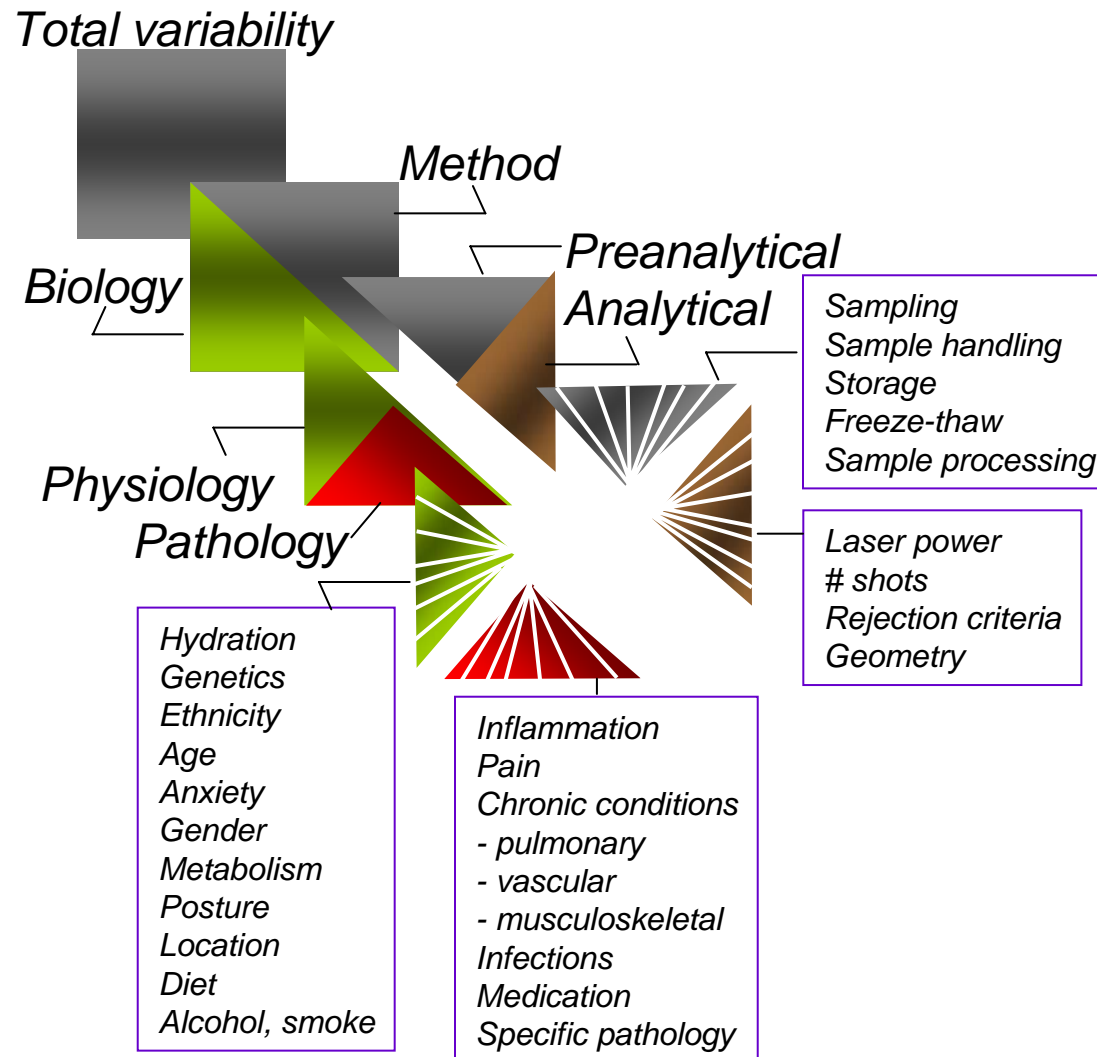
- We are interested in the biological variables associated with a specific pathological condition
- This is obscured by tremendous noise

Biofluid Proteomics: Potential, Pitfalls, and Solutions



STATENS
SERUM
INSTITUT

Variability contributors in high-resolution analysis of biofluids



Seeing differences between groups using analyses yielding multiple variables

- Overfitting is an obvious risk – too many data and too few samples
 - We need at least 10 times more samples than features!
 - Look for univariate effects
 - Bonferroni-correct significance levels (p/N)
 - Use PCA first or other unsupervised method
 - Do permutation analysis
 - Independently train, cross-validate, and validate

Seeing differences

- Minimize noise
 - Study design from patient to analysis: focus on reducing variability
 - Reduce complexity of material
 - Optimize analytical procedures for reproducibility
 - Reduce amount of data (=variable selection)

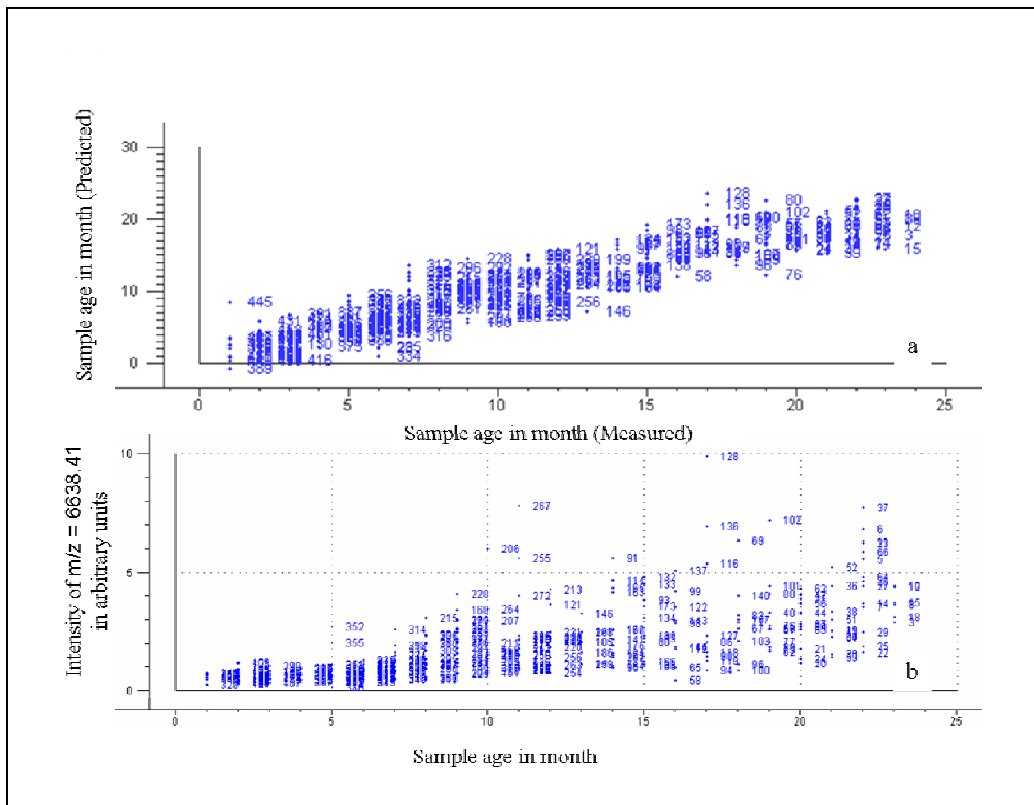
Minimize noise - Study Design

Ovarian Cancer proteomics Pelvic Mass Study

- All samples (disease and disease-control) collected in a specialized clinic following uniform and strict guidelines. Only women with a risk-of-malignancy index $RMI > 150$, and age > 18 years are included.
- >400 sera were collected. 25 % expected malignant (borderline to stage IV). 75 % are expected to be benign conditions
- Samples are thawed on ice and handled on a fractionation robot using WCX beads and IMAC-Cu beads at 22°C and humidity set to 35%



Minimize noise. The introduction of new noise – systematic changes caused by sample handling



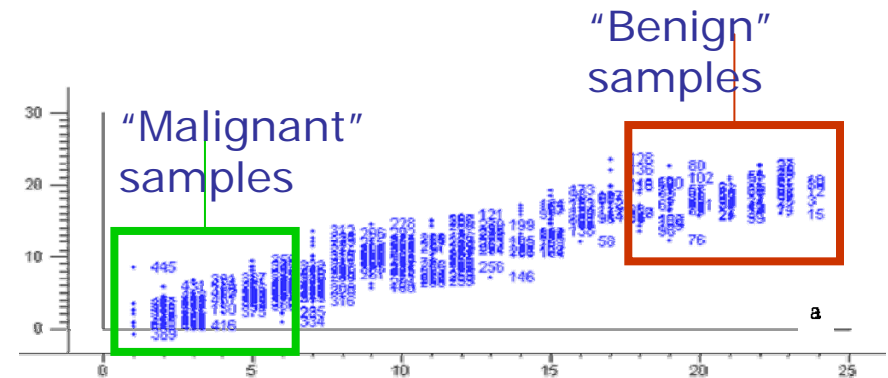
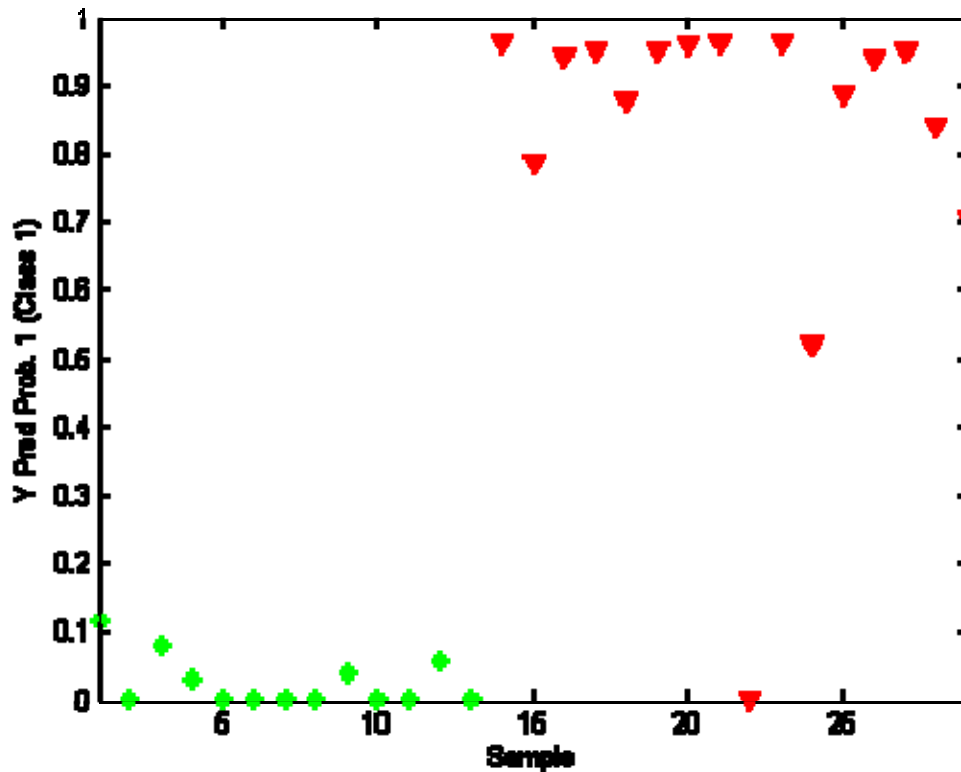
a) Long-term storage effects at -20°C. a) Modelling sample age using PLS. Samples were collected over a period of 23 months.

b) The figure shows a variable with m/z 6638.41 with intensity increasing as a function of storage time. The variance between samples becomes larger with increased storage time.

Biofluid Proteomics: Potential, Pitfalls, and Solutions

Possible consequence of systematic noise

Example of spurious classification model caused by storage artifacts.





Variability arising *in vitro*. Table of typical peaks involved in post-sampling changes in normal sera

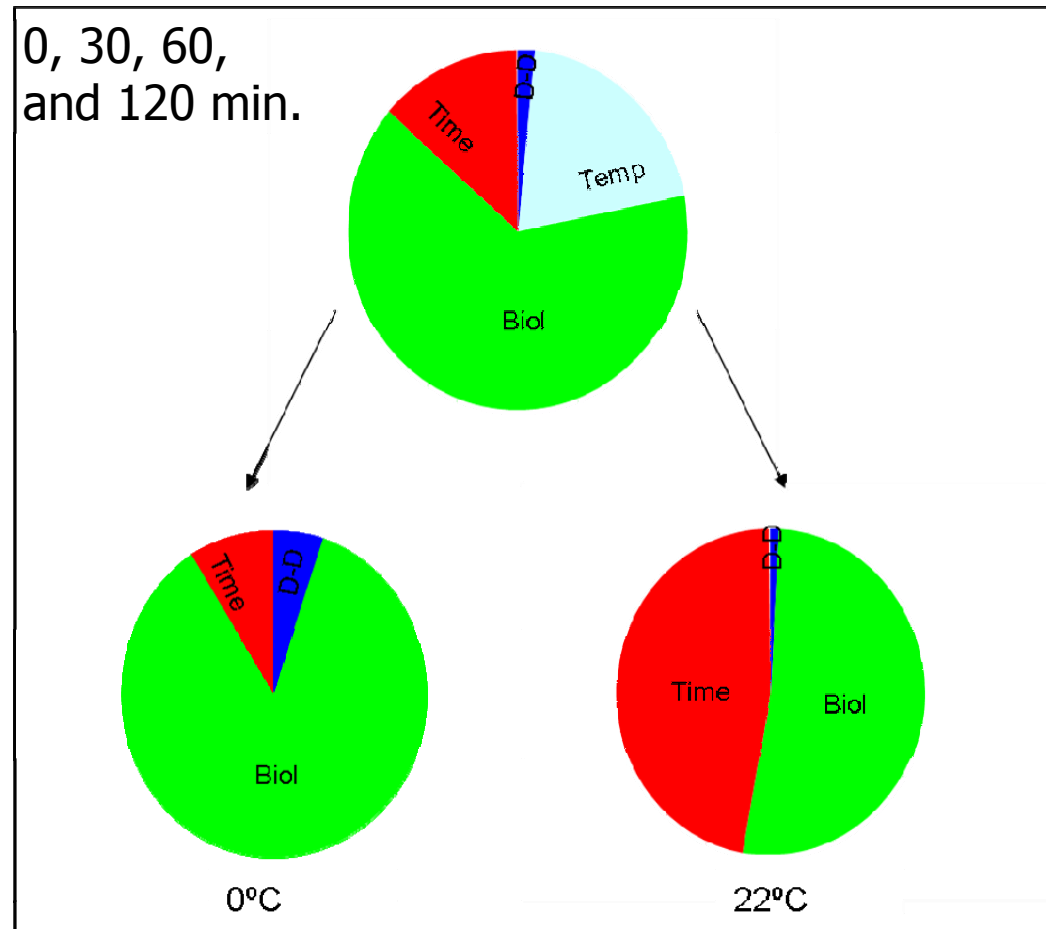
Name/sequence	Fragm. mass
Complement C3f: aa 2-16 (SKITHRIHWESASLL)	1777.9
Complement C3f: aa 1-16 (SSKITHRIHWESASLL)	1864.9
Complement C3f: aa 1-17 (SSKITHRIHWESASLLR)	2021.7
Fibrin alpha C term fragment: aa 81-105 (SSSYSKQFTSSTSYNRGDSTFESKS)	2767.4
Fibrin alpha C term fragment: aa 81-106 (SSSYSKQFTSSTSYNRGDSTFESKSY)	2931.5
Fibrinopeptide A: aa 1-12 (EGDFLAEGGGVR)	1206.5
Fibrinopeptide A: aa 3-16 (SGEGFLA EGGGVR)	1350.7
Fibrinopeptide A: aa 2-16 (DSGEGFLAEGGGVR)	1465.5
Fibrinopeptide A (Modifications: 3 Phosphorylated): aa 1-16 (ADSGEGFLAEGGGVR)	1616.9
Kininogen: aa 439-456 (HNLGHGHKHERDQGHGHQ)	2080.9

West-Nørager *et al.* J Chromatogr B 2007

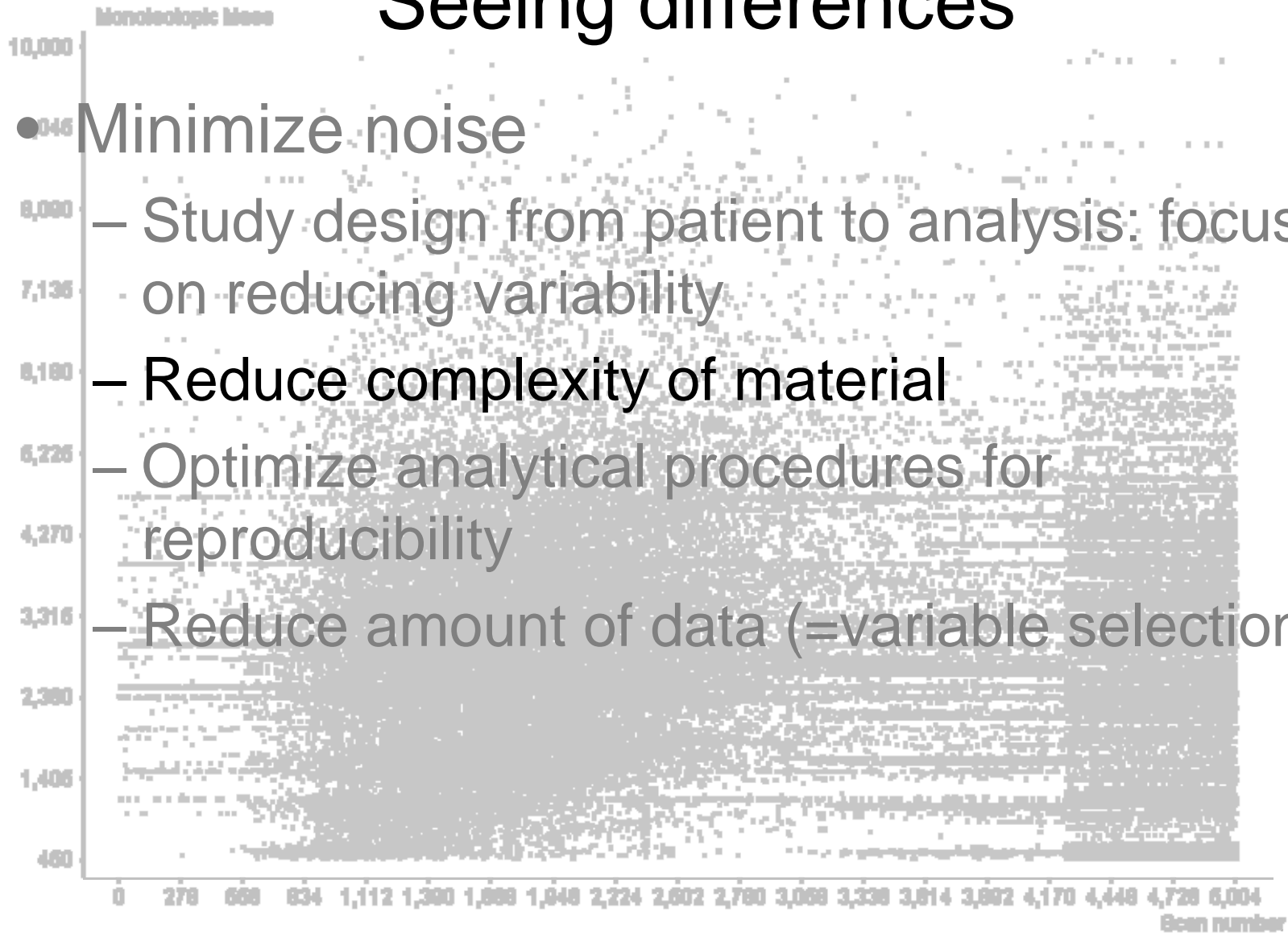
Minimizing noise.

Serum proteomics, WCX-beads

Time, temperature, biology, and day-to-day variation



Seeing differences



● Minimize noise

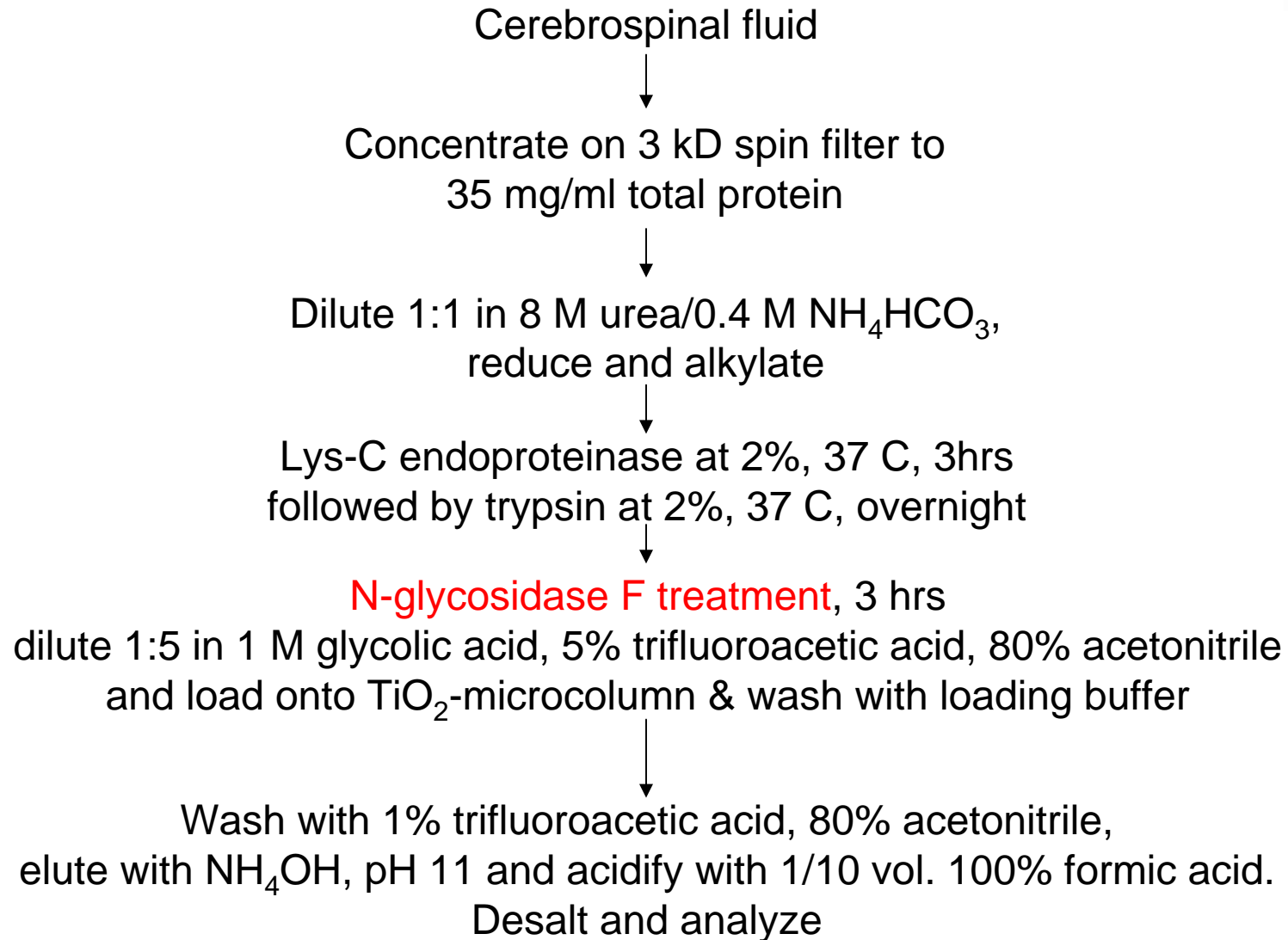
- Study design from patient to analysis: focus on reducing variability
- Reduce complexity of material
- Optimize analytical procedures for reproducibility
- Reduce amount of data (=variable selection)

Simplify input to analytical machines

Sample fractionation

- Chromatographic – “non-biological”
 - Charge, hydrophobicity, size
 - Removal of abundant species
- Biological (*Intelligent Proteomics*)
 - Enrich for post-translational modifications
 - Phosphorylations
 - Glycosylations
 - Others
 - Microparticle isolation
 - Subfractionate on surface markers, size
 - Highly targeted affinity isolation
 - Antibodies
 - Lectins
- COMBINATIONS

Reducing biofluid complexity: CSF-phosphoproteomics.



Biofluid Proteomics: Potential, Pitfalls, and Solutions

Phosphorylation sites identified in this study are underscored. Previously unknown phosphorylation

sites are highlighted.

Bahl JMC et al., Anal Chem 2008

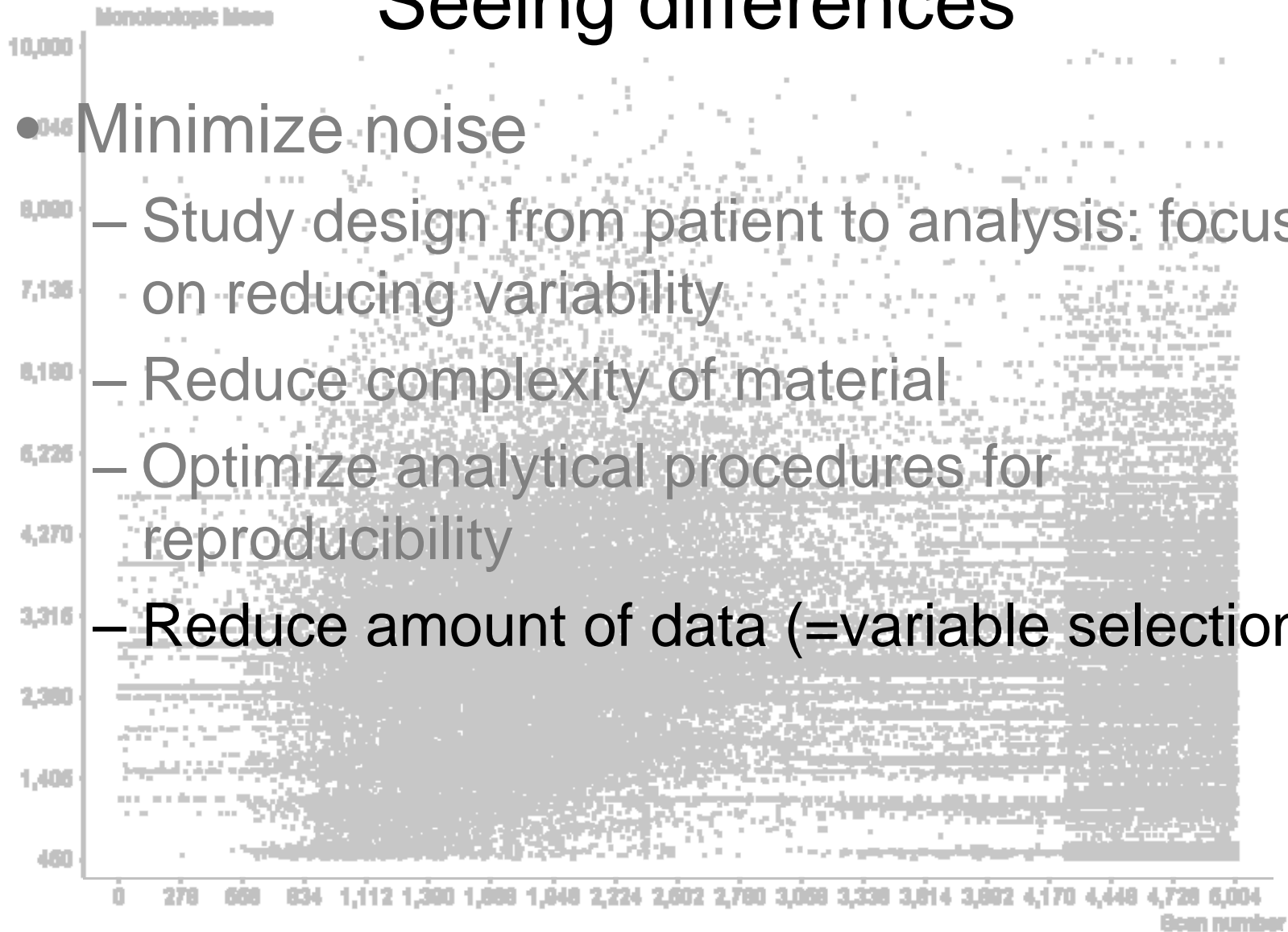


STATENS
SERUM
INSTITUT

Phosphoproteins in control CSF - 56 new phosphosites in 38 proteins

Protein	Peptide	Sequence	Protein	Peptide	Sequence	
Osteopontin	52-70	QNLAPQNAVp <u>SS</u> EETNDFK	Cadherin 2	121-145	LSLKPTLT <u>EE</u> SVK <u>ES</u> AEVEEIVFPR	
	176-203	RPDIQYPDApTDEDITp <u>SH</u> MEp <u>SE</u> ELNGAYK		Neuron Navigator 2	1-9	MPAILV <u>AS</u> K
	204-220	AIPVAQDLNAPpSDWDp <u>SR</u>		Antithrombin III	61-78	KATEDEG <u>SE</u> QKIPEATNR
	223-241	Dp <u>SY</u> ETpSQLDDQpSA <u>ETH</u> SHK		Serum Albumin	76-88	TCVADE <u>SA</u> ENCDK
	250-268	ANDEp <u>S</u> NEHSDVIDp <u>SQ</u> ELp <u>SK</u>		Versican core protein	2108-2124	QEIESETT <u>SEE</u> QIQEEK
	269-290	Vp <u>S</u> REFHpSHEFHpSHEDMLVVDPK		Sushi domain-containing protein 5	290-298	G <u>SE</u> GEQQIMR
300-314	FRIp <u>S</u> HELDpS <u>AS</u> p <u>SE</u> VN	Proenkephalin A		237-252	FAEALPSDEEGESV <u>SK</u>	
SPARC like protein 1	76-89	SKEESHEQ <u>SA</u> EQ GK		Cocaine- and amphetamine-regulated transcript protein	37-51	ALDIYSAVDDA <u>SH</u> EK
	90-98	SS <u>Q</u> ELGLK		Selenoprotein P	262-274	DMPA <u>SE</u> DLQDLQK
	187-220	DQGNQE <p>DPN<u>IS</u>NGEEEEKEPGEV GTHNDNQER</p>		Apolipoprotein L	306-320	VTEP <u>IS</u> AESGEQVER
	256-286	MQDEFDQGNQEEDN <u>NA</u> EMEEENASNVNK		Fructose-bisphosphate aldolase A	29-42	GILAADESTG <u>SI</u> AK
	287-299	HIQETEWQ <u>SE</u> Q GK		Reticulocalbin-1	71-81	TFDQLTPDE <u>SK</u>
	409-426	KAEN <u>S</u> SNEEETSSEG NMR		Proprotein convertase	683-692	HLAQ <u>AS</u> QELQ
Secretogranin 1	123-131	WAEGGGH <u>SR</u>		Matrix remodelling-associated protein 7	125-142	GP <u>SE</u> SGPEEEDGEGFSFK
	134-153	ADEPQWSLYPDSQVp <u>SE</u> EVK		Amyloid precursor protein	439-450	VE <u>SE</u> LEQEAANER
	179-202	GEDp <u>S</u> p <u>SE</u> EKHLEEPGETQNAFLNER		Plasminogen	349-386	IPSCDSSPVST <u>TE</u> QLAPTAPPELTPVWQDCYHGDGQSYR
	235-254	<u>SS</u> Q <u>ES</u> GEAG <u>SE</u> QENHPQESK		Follistatin-related protein 1	163-170	LD <u>SS</u> EFLK
	259-276	p <u>SQ</u> EEp <u>SE</u> EGEEDATSEVDK		Kallikrein-6	81-91	ES <u>SQ</u> EQSSVVR
	305-324	GHPQEEp <u>SE</u> ESNVSMASLGEK		Cofilin-1	2-13	AS <u>G</u> VAVSDGVIK
	372-387	APRPQp <u>SE</u> ESWDEEDKR		Myristoyl alanine-rich C-kinase substrate	12-30	GEAAAERPGEAAVASS <u>PS</u> K
	397-409	MAHGYGEEp <u>SE</u> EER		Golgi phosphoprotein 2	253-267	EE <u>T</u> NEIQV VNEEPQR
617-640	p <u>SA</u> EFPDFYDp <u>SE</u> EPVSTHQEAENEK	Latent-transforming growth factor beta binding protein isoform 1L		1280-1297	GFVPA <u>GES</u> SEAGGENYK	
Secretogranin 2	257-269	IESQTQEEVRD <u>SK</u>		Testican -3	33-42	<u>SD</u> GGNFLDDK
	549-561	EHLNQGS <u>SE</u> QETDK		Ubiquitin	55-63	TL <u>SD</u> YNIQK
Receptor-type tyrosine-protein phosphatase N2	355-368	AALGE <u>SE</u> GEQADGPK		Thy-1 membrane protein	88-97	VLYLSAFT <u>SK</u>
Dickkopf related protein 3	428-448	SEHPSSLS <u>SE</u> EETAGVENVK				
Apolipoprotein E	81-103	AS <u>SE</u> VNLANLPPSYHNETNTDTK				
Fibulin 1	138-152	GEVQAMLGQ <u>SE</u> EELR				
Kininogen 1	147-163	<u>SE</u> QETGDLDV GGLQETDK				
Cystatin C	325-343	ET <u>CS</u> KE <u>SE</u> NEELTESCETK				
Secretogranin 3	35-51	LVGGPMDA <u>SE</u> VEEIGVRR				
Trans Golgi network integral membrane protein 2	35-53	EL <u>SA</u> ERPLNEQIAEAEEDK				
Extracellular matrix protein 2	65-83	DSPSK <u>SA</u> EAQTPEDTPNK				
	209-227	EALQ <u>SE</u> EDEEVKEEDTEQK				

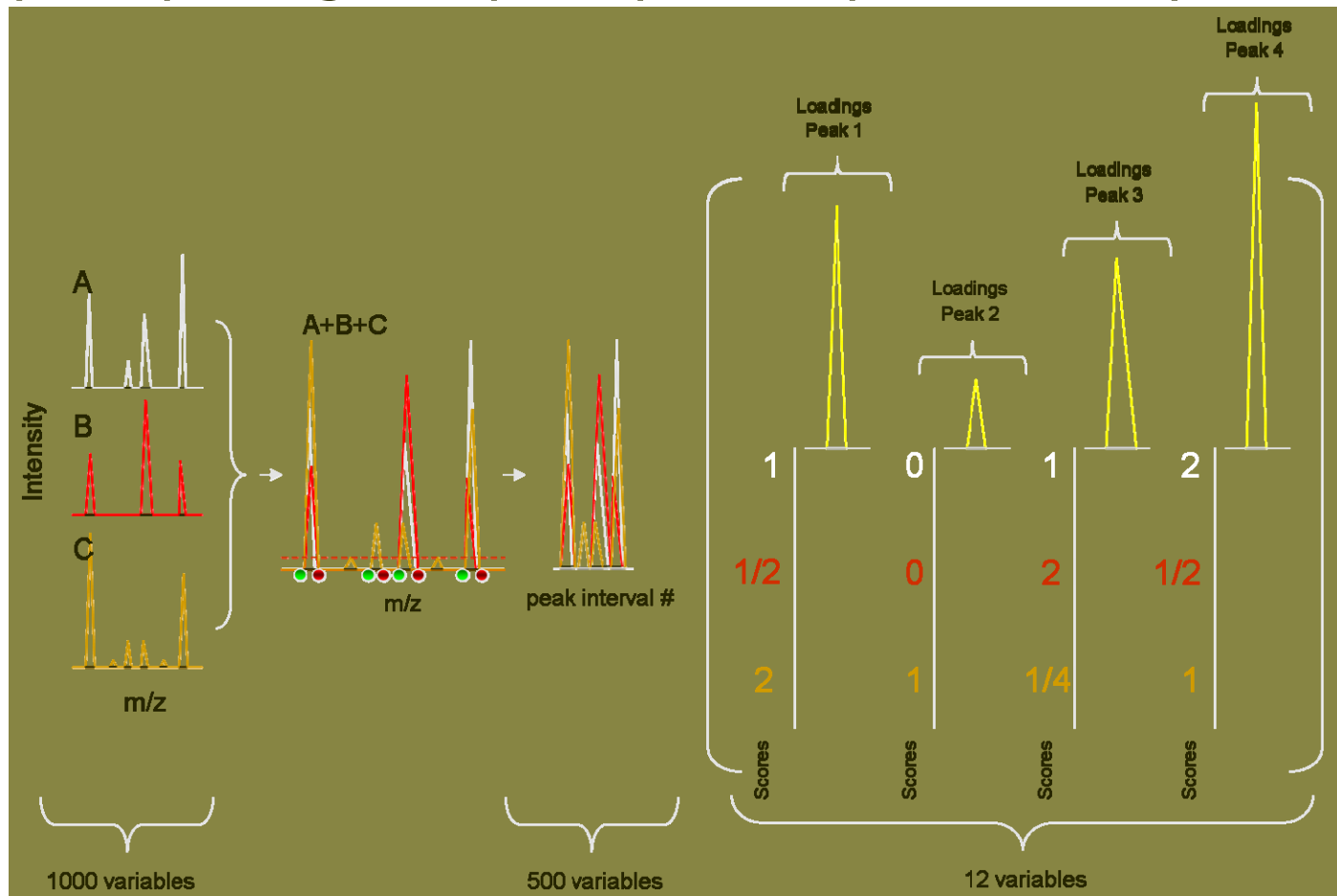
Seeing differences



● Minimize noise

- Study design from patient to analysis: focus on reducing variability
- Reduce complexity of material
- Optimize analytical procedures for reproducibility
- Reduce amount of data (=variable selection)

A method for data reduction that involves both peak picking and principal component analysis

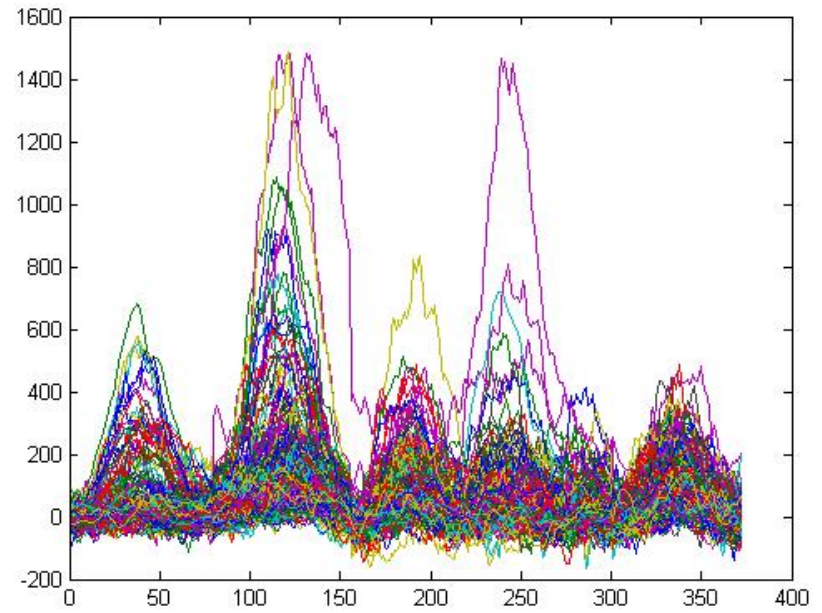
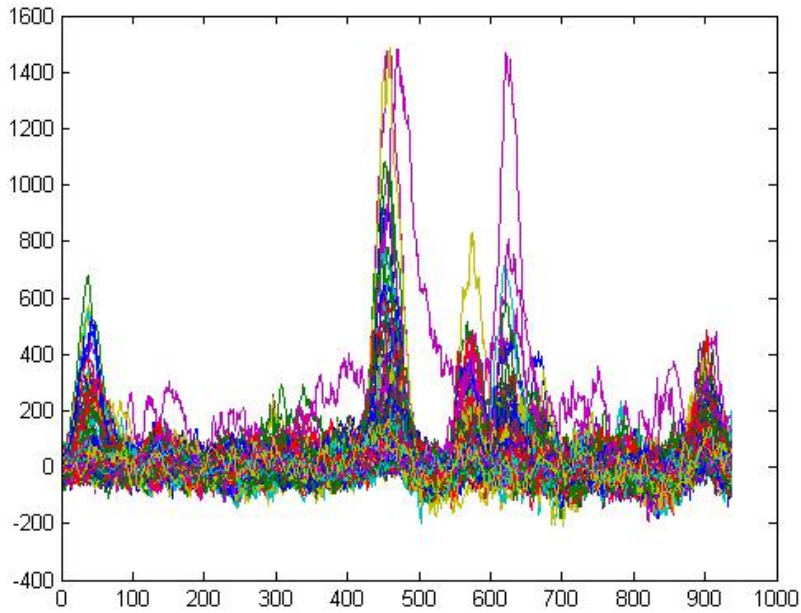


300 samples → 33.000 points = matrix with 9.900.000 points

Through data reduction, matrix reduced to 165.000 points

98% in reduction with less than 1% loss of data

Data Reduction I. Peak picking (aka feature selection).

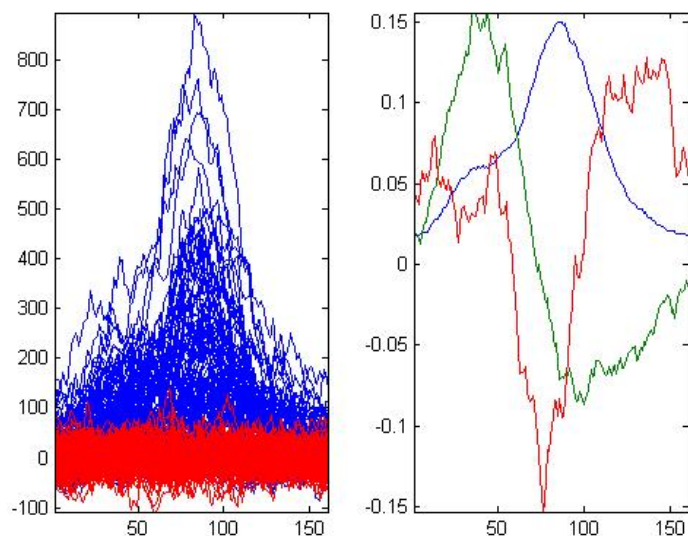


Removal of data between peaks. Yields 50% variable reduction

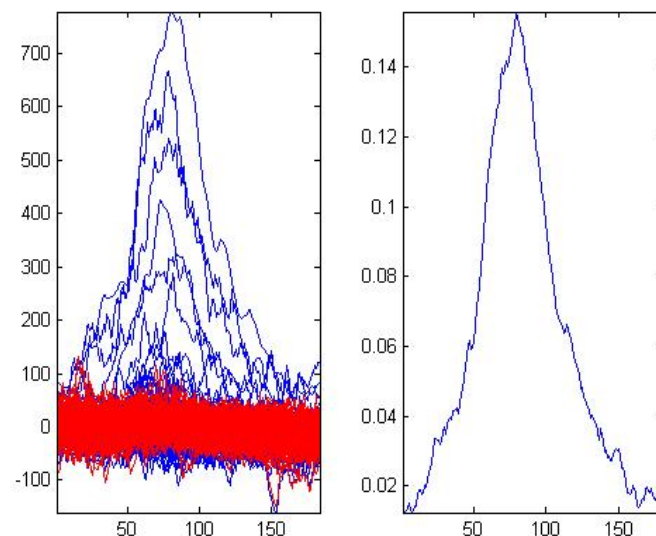


Data Reduction II.

PCA (principal component analysis) performed in each peak interval (MatLab-script available)



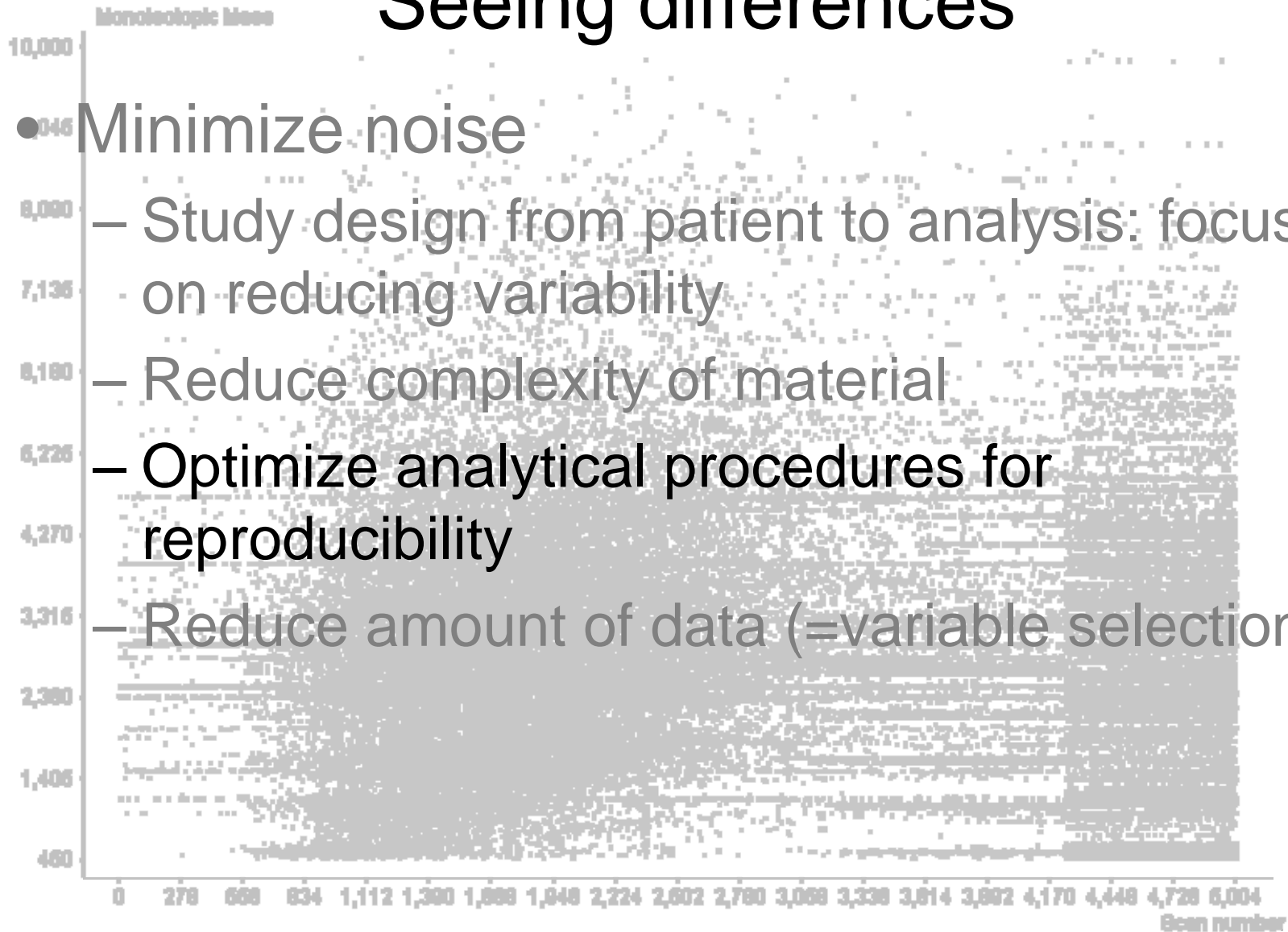
Peak #107



Peak #114

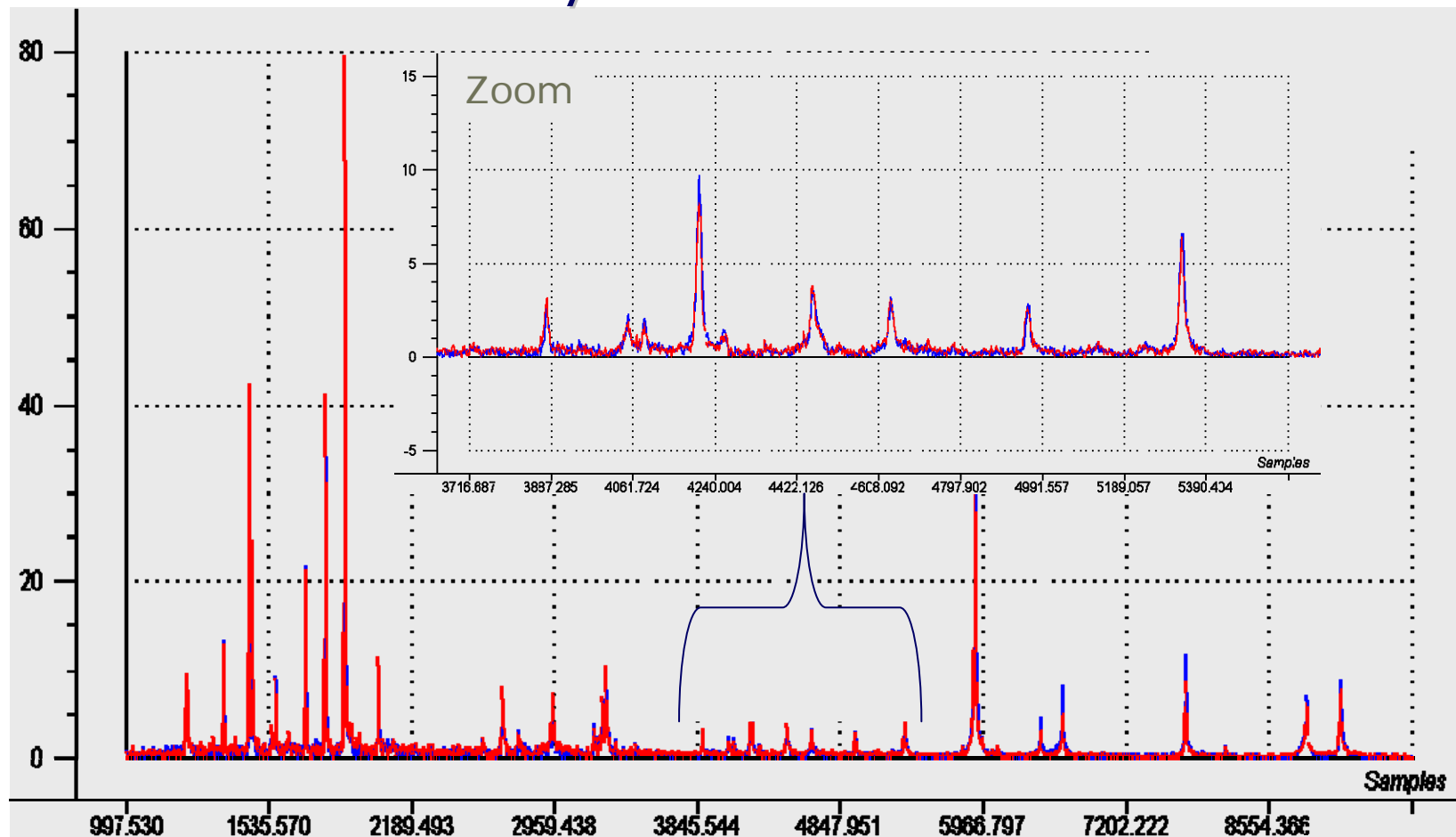
Procedure reduces the cumulated data matrix from 8,745,000 to 137,000 data sets. Thus, every sample is represented by 520 variables, - a reduction of 98% without losing information

Seeing differences

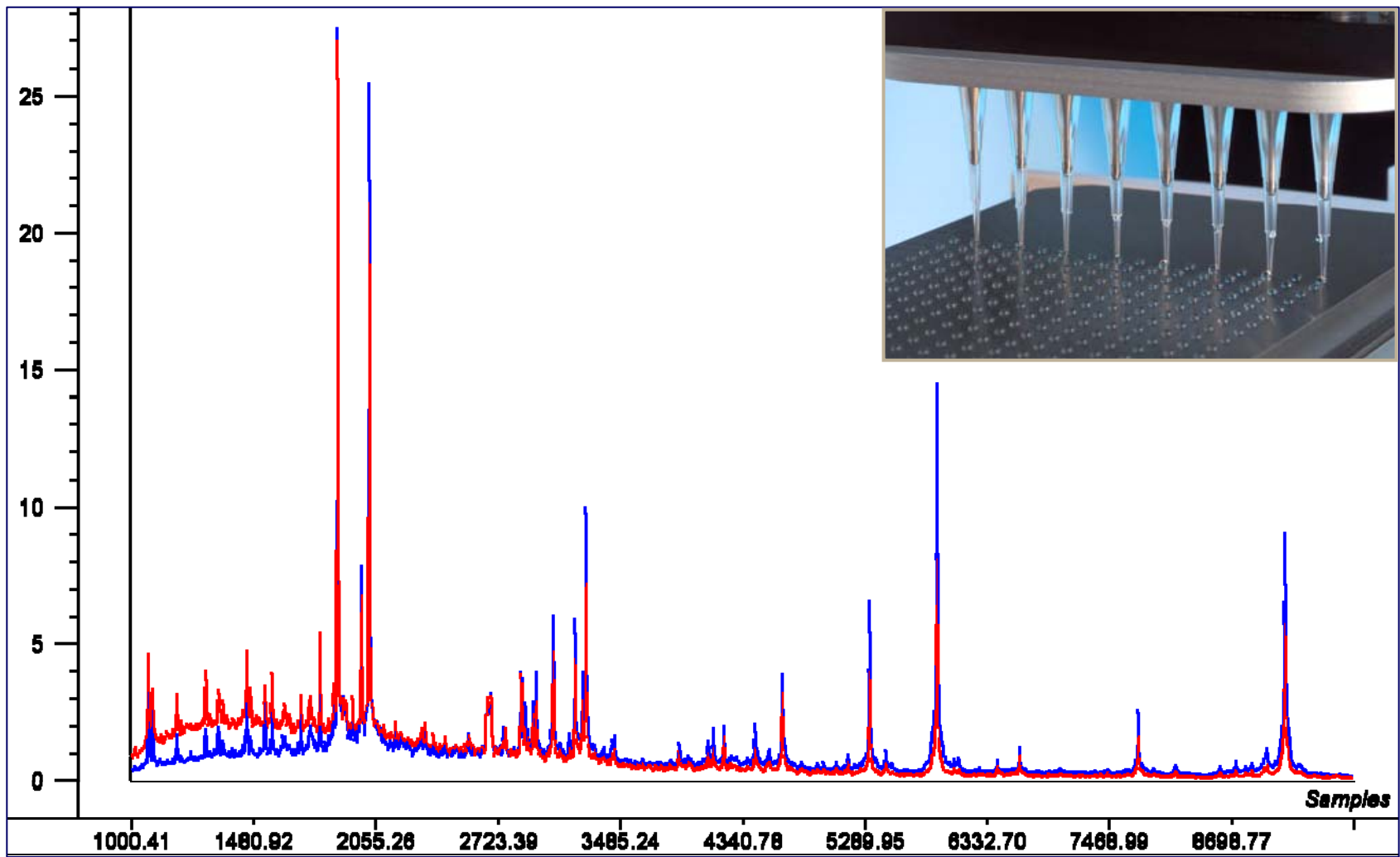


MALDI-TOF MS reproducibility

- Spot to spot variation in same run
 - R^2 usually ≥ 0.99

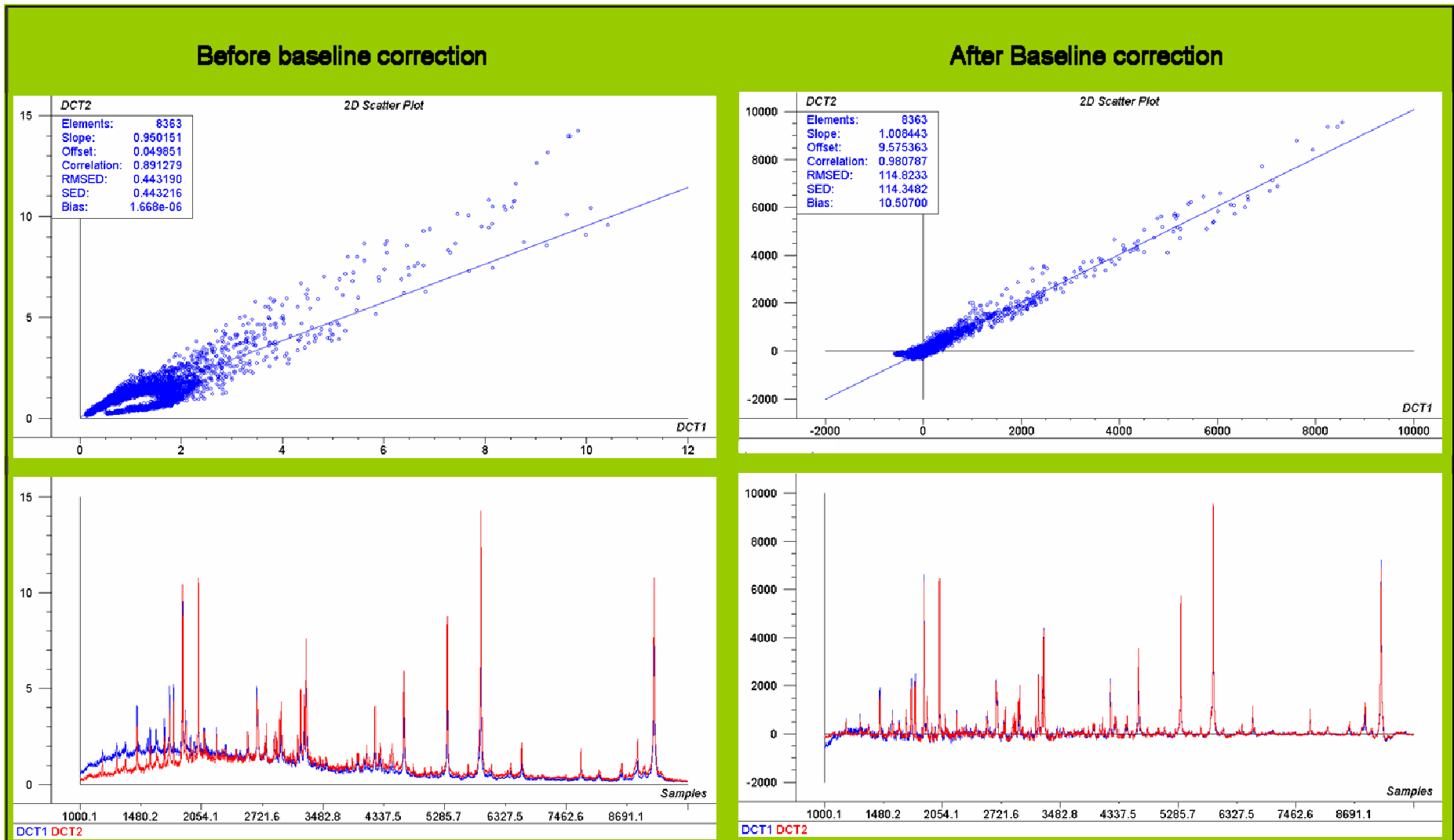


Same sample, different targets

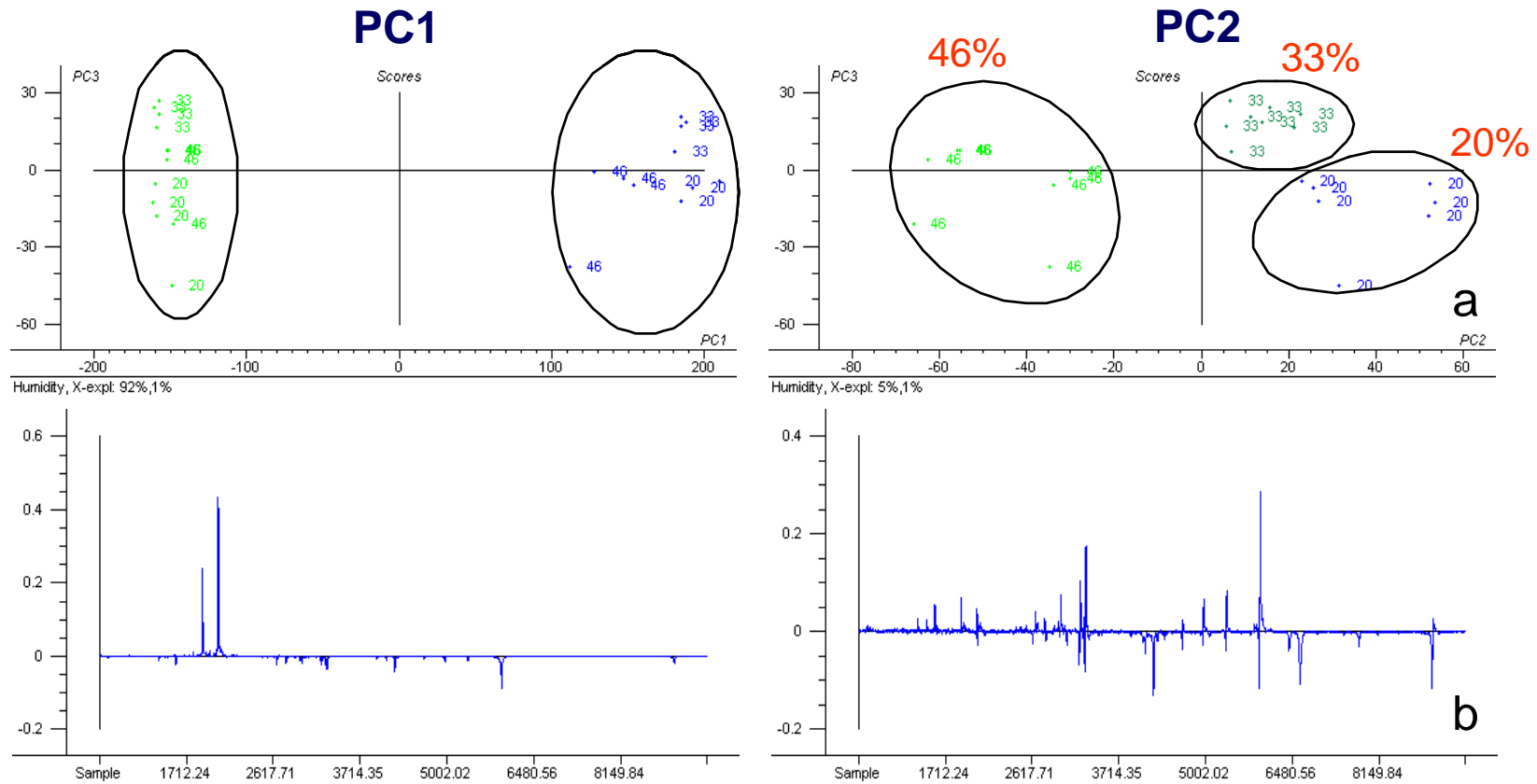


Same sample, different targets

Baseline correction removes most of the target-induced variation



Ambient humidity during sample crystallization affects results



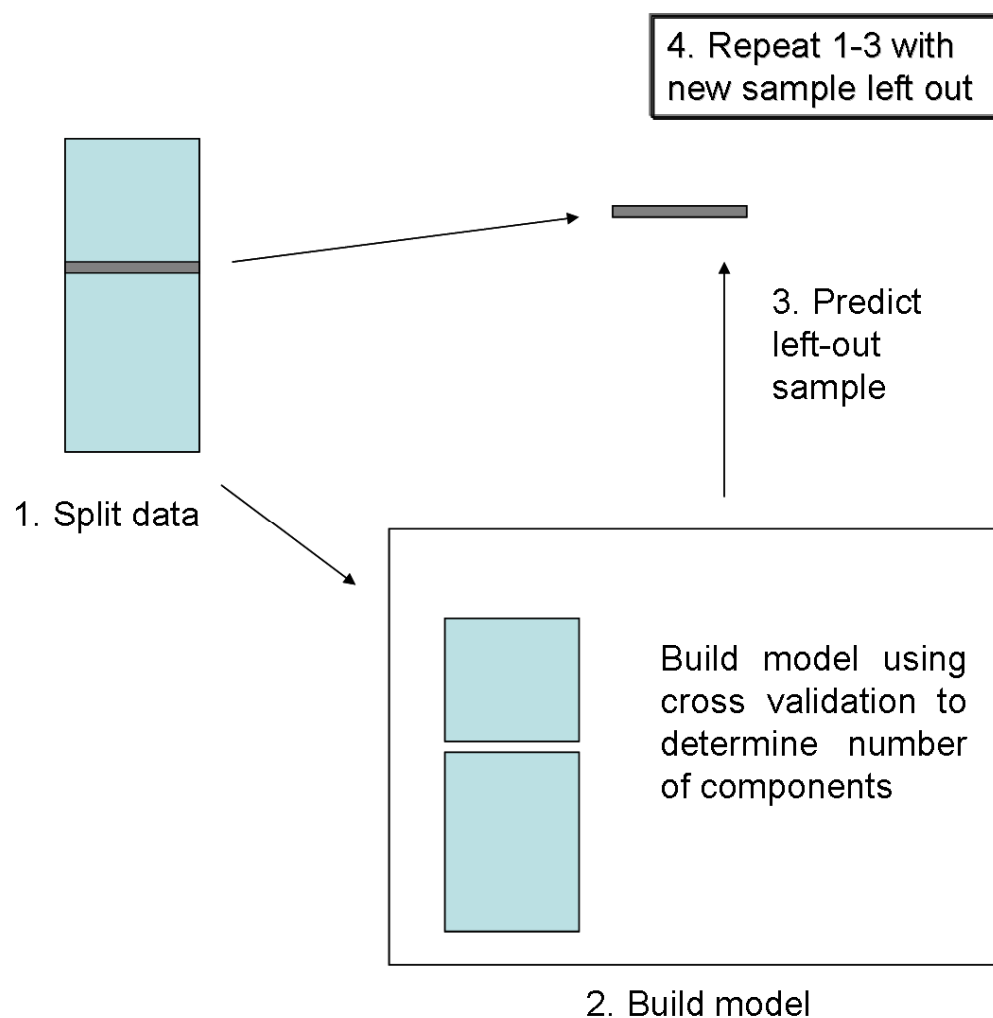
PCA analysis of the crystallization process at three different humidity levels (20%, 33%, and 46%) using replicates of two different samples. Score plots show both the biological difference (PC1) and the different crystallization dependent on the humidity level (PC2). The related loading plots outline the variables accountable for the object pattern in the score plot.

Seeing differences between groups using analyses yielding multiple variables

- Overfitting is an obvious risk – too many data and too few samples
 - We need at least 10 times more samples than features!
 - Look for univariate effects
 - Bonferroni-correct significance levels (p/N)
 - Use PCA first or other unsupervised method
 - Do permutation analysis
 - Independently train, cross-validate, and validate

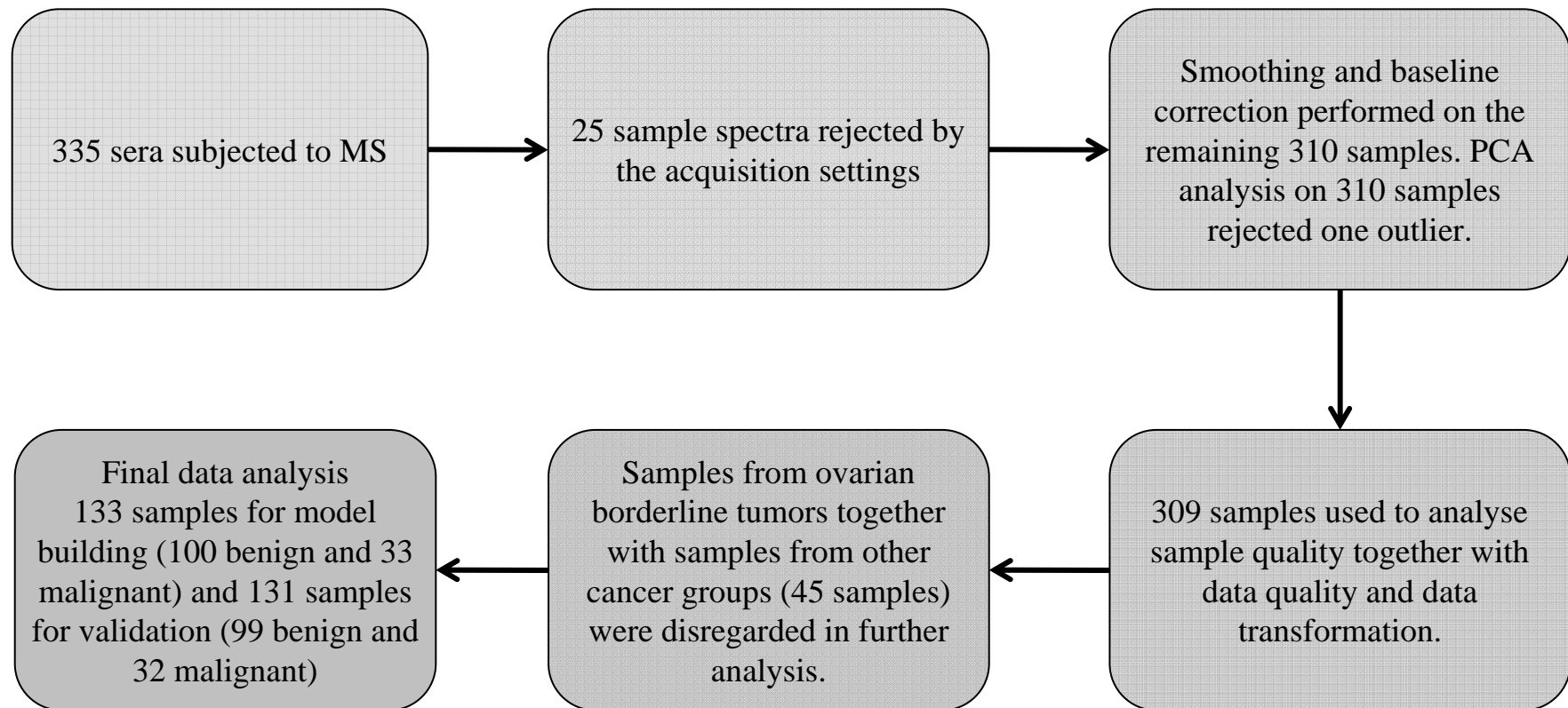
Biofluid Proteomics: Potential, Pitfalls, and Solutions

Cross-Model-Validation is more stringent than usual *leave-one-out* cross validation for testing model





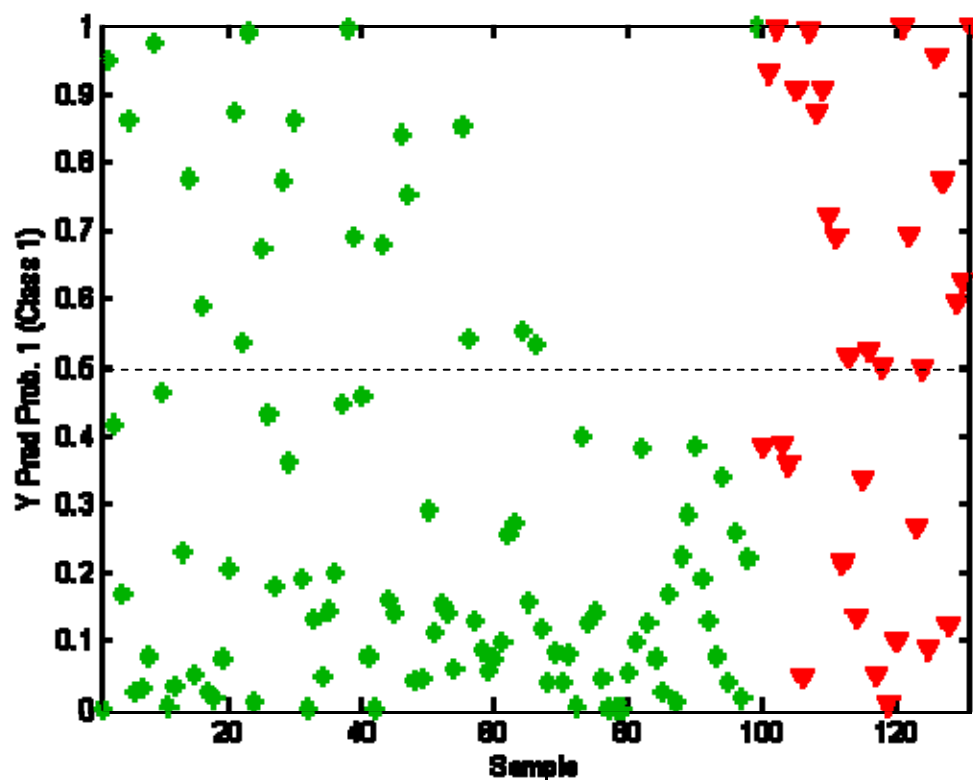
Pelvic Mass Study Serological biomarkers for ovarian cancer



Pelvic Mass Study Serological biomarkers for ovarian cancer

Overview of the performance of the final model

Independent validation sample set of 99 benign (green asterisks) and 32 malignant (red triangles) from the PLS-DA model – 29 variables (26 peaks)



The final model discriminates malignant from benign conditions in 68% of the cases with 56% sensitivity and 79% specificity.

Model / Correctly predicted	Benign	Malignant
Cross validation (CV)	86 %	74 %
Model prediction (Final validation)	79 %	56 %

Clinical proteomics

Goal: Measurements that discriminate and characterize disease vs. health

What the Analyst wants: Comprehensible, quantitative, and reproducible measurements of the contents of complex biosamples

What the Statistician wants: The best descriptors for maximizing the distinction between two or more groups of samples

What the Physician wants: Laboratory tests that increase life expectancy of patients

Acknowledgements

Mikkel West-Nørager

Estrid Høgdall

Claus Høgdall

Rasmus Bro

Martin R.Larsen

Søren Skov Jensen

Justyna Czarna Bahl

Keld Poulsen

Mark Lim, NIH

