

# STATISTICAL METHODS FOR CLASS PREDICTION

Olga Vitek

College of Science

College of Computer and Information Science

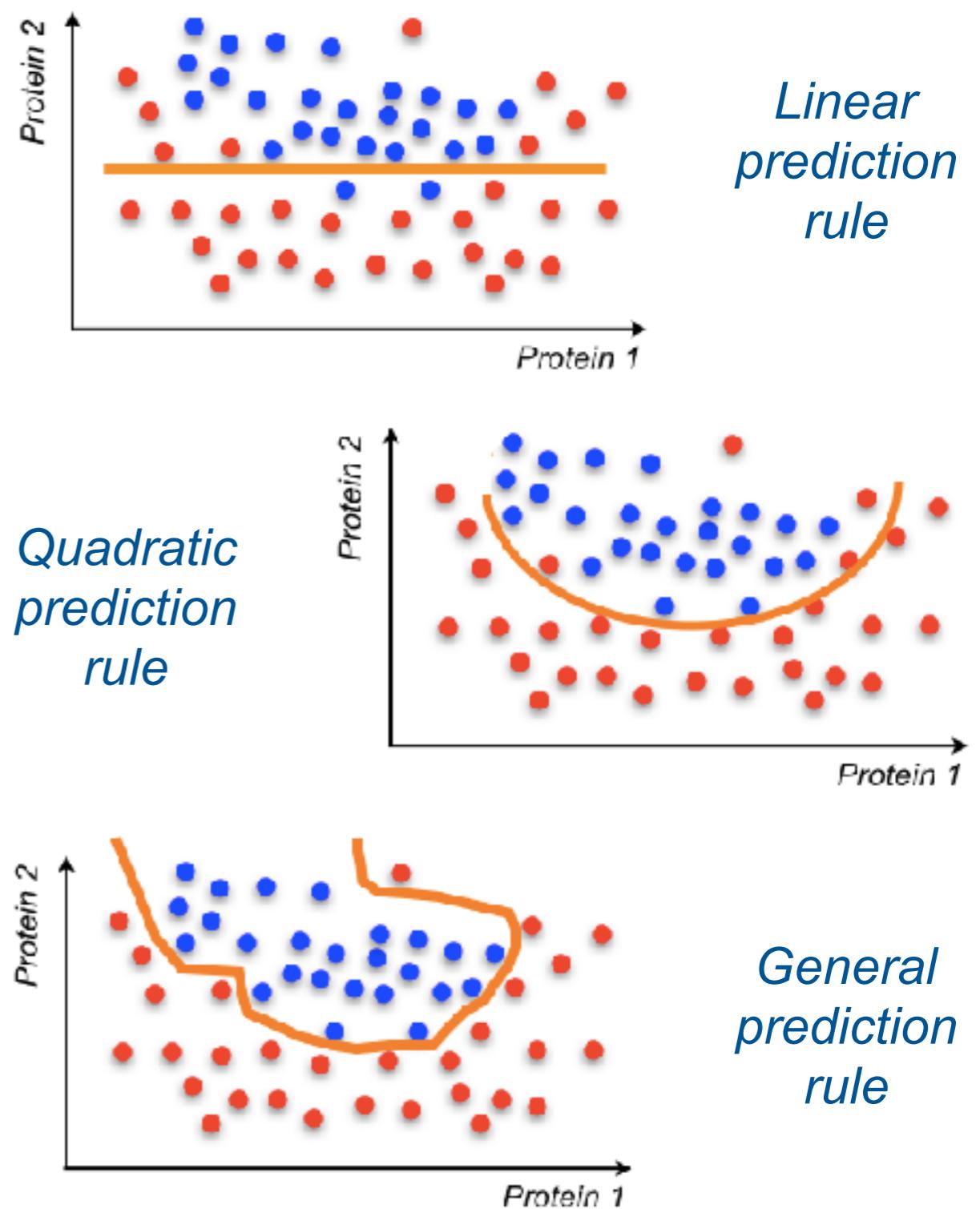


Northeastern University

# STATISTICAL GOAL 3: CLASS PREDICTION

Classify each subject into a known group

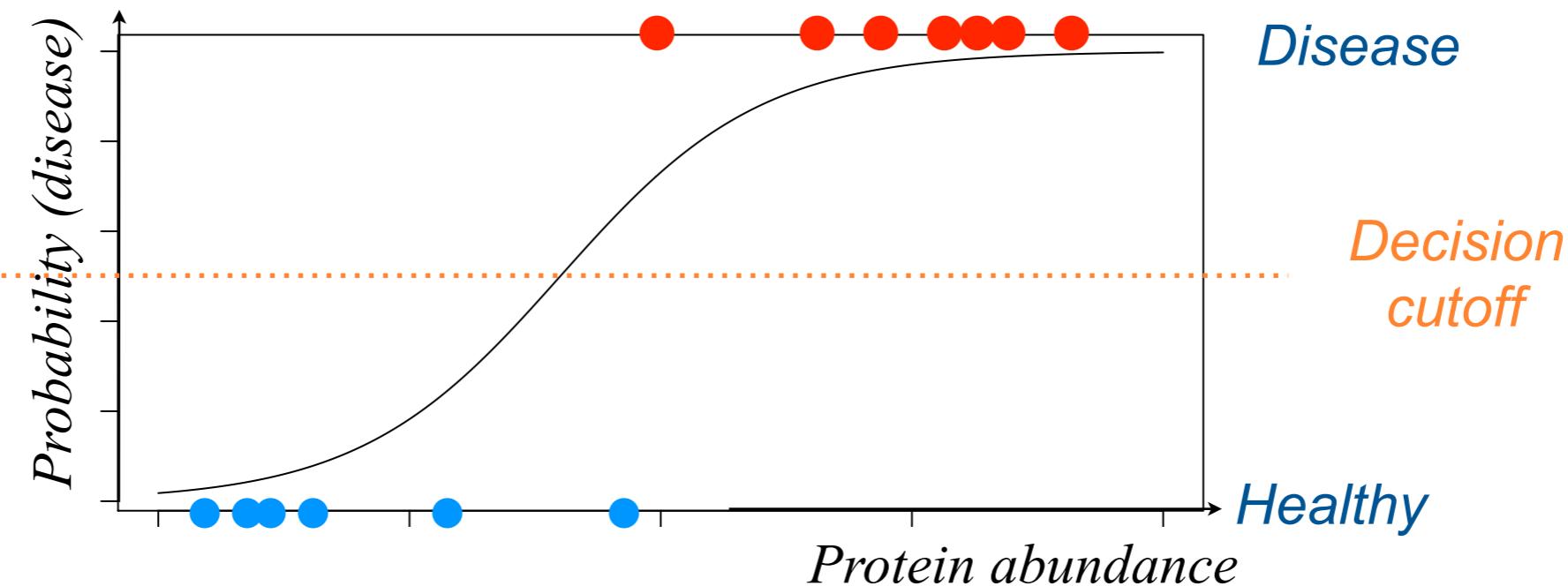
- Known class labels
  - Predict individual subjects
  - Report misclassification error (sensitivity, specificity, predictive value etc)
- Useful when focus on an individual
  - Tier I or Tier II biomarker discovery studies



# OUTLINE

- Supervised classification
  - Classifiers, evaluation, issues
- Case studies
  - Alzheimer's diagnosis with plasma signaling proteins
  - Colorectal cancer diagnosis from plasma proteins
  - MACQ II: common practices in assay validation
- Experimental design for prediction
  - Impacts of dimensionality and noise

# EXAMPLE: LOGISTIC REGRESSION



*The model has a relatively simple form*

$$P(\text{disease}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 \text{Protein})}}$$

*Add terms to reflect multiple proteins and confounding factors*

*The parameters have an interpretation that is familiar to clinicians*

$$\beta_1 = \log \frac{\text{Odds}(\text{Protein} + 1)}{\text{Odds}(\text{Protein})}$$

*Selects a subset of proteins as predictors*

# EVALUATION OF PREDICTIVE ABILITY

Two classes, for a particular probability cutoff

		Decision		Totals
		Negative	Positive	
Truth	Negative	TN	FP	N
	Positive	FN	TP	
				P

## Summaries from the decision table

$$\text{fp rate} = \frac{FP}{N}$$

$$\text{tp rate} = \frac{TP}{P}$$

$$\text{precision} = \frac{TP}{TP+FP} \quad \text{recall} = \frac{TP}{P}$$

$$\text{accuracy} = \frac{TP+TN}{P+N}$$

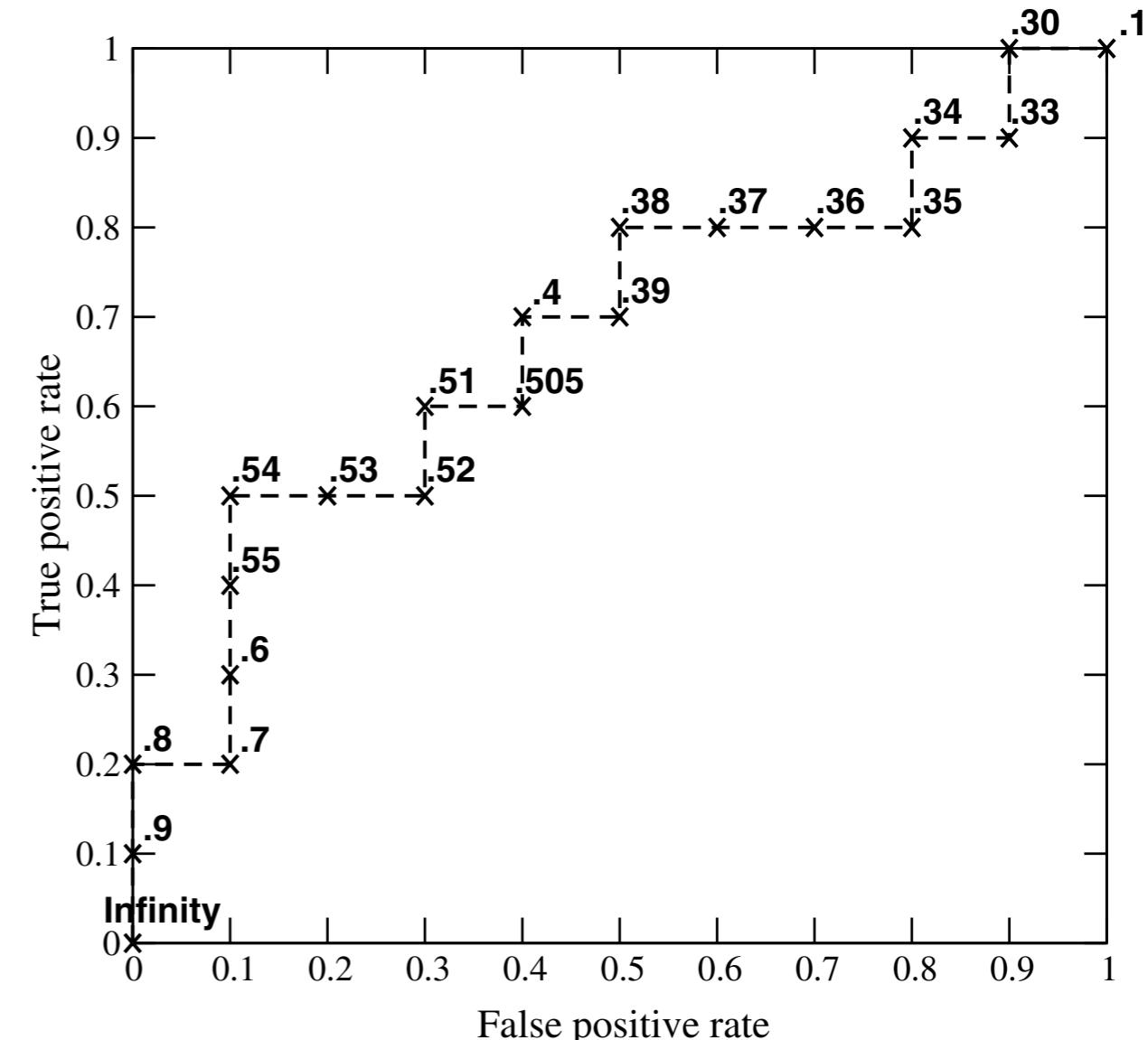
$$\text{F-measure} = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}}$$

One separate table for each probability cutoff

# SUMMARY OF PREDICTIVE ABILITY

## Varying predicted probability cutoff

Inst#	Class	Score	Inst#	Class	Score
1	p	.9	11	p	.4
2	p	.8	12	n	.39
3	n	.7	13	p	.38
4	p	.6	14	n	.37
5	p	.55	15	n	.36
6	p	.54	16	n	.35
7	n	.53	17	p	.34
8	n	.52	18	n	.33
9	p	.51	19	p	.30
10	n	.505	20	n	.1



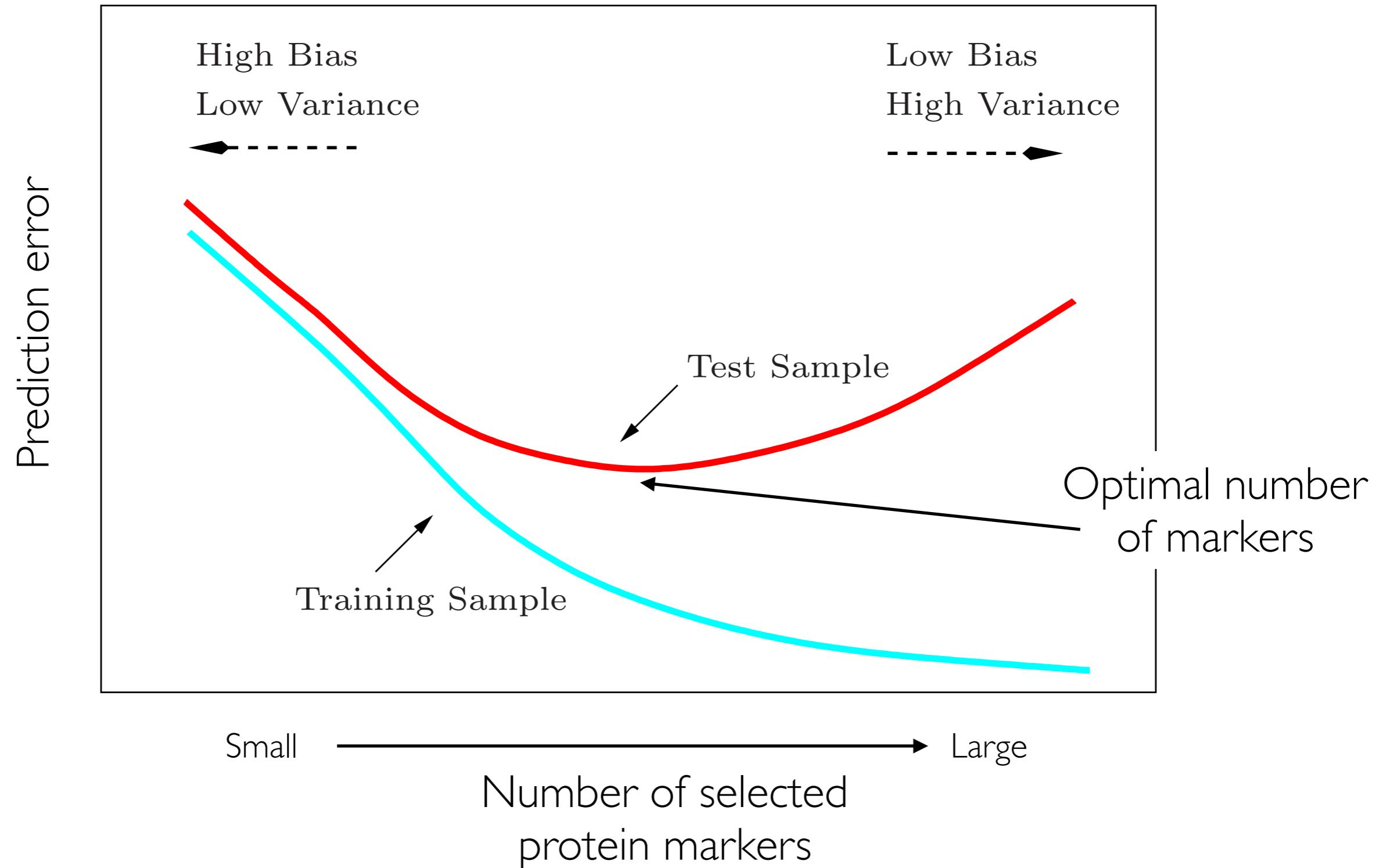
*Subjects in the study and their predicted probability of the disease*

**ROC curve**, summarizing the predictive ability for various predicted probability cutoffs

# SELECTION OF PREDICTIVE PROTEINS

7

## Bias-variance trade-off



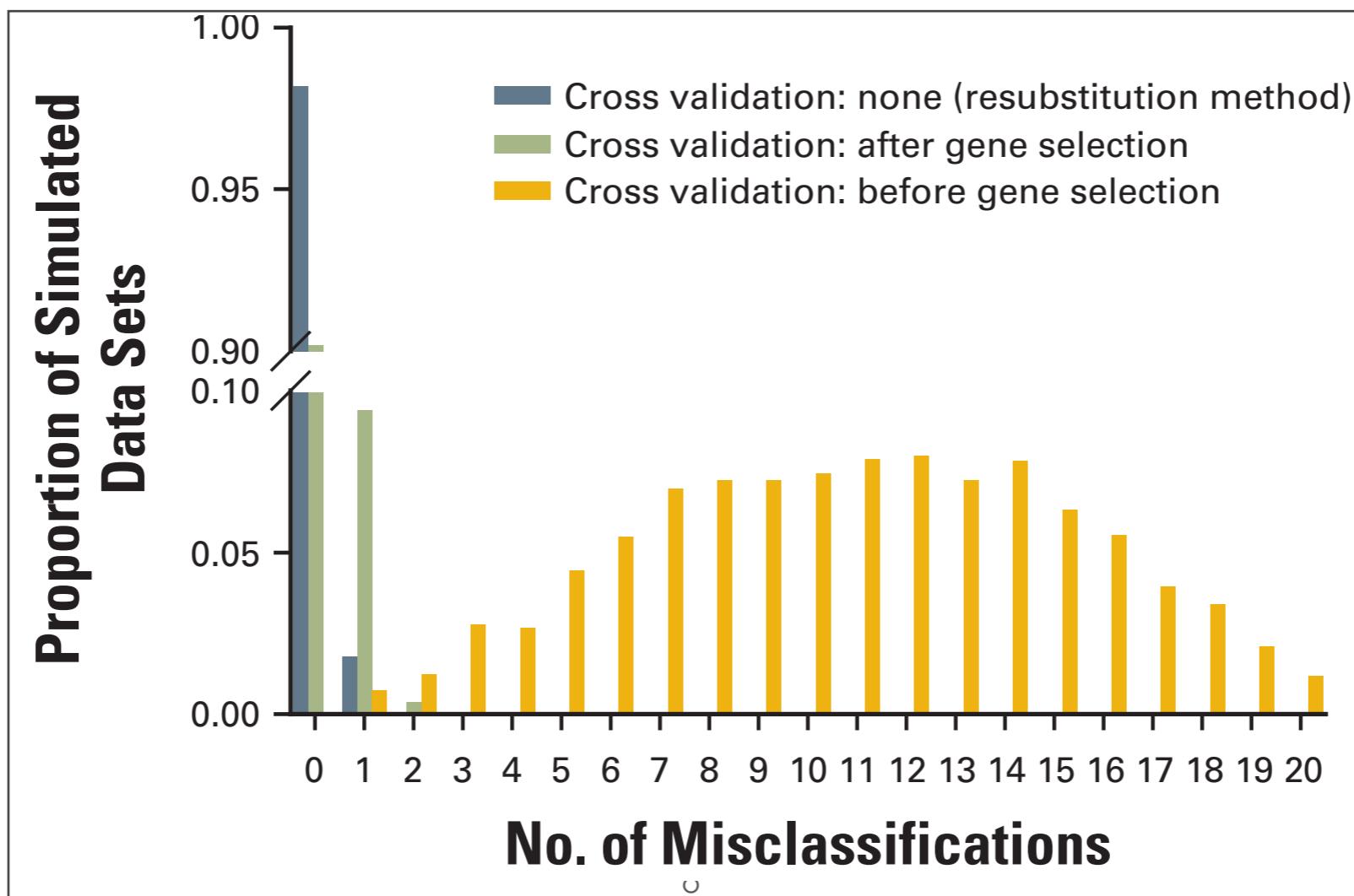
# SELECTION OF PREDICTIVE PROTEINS

Should be done within (cross-) validation

- Simulated data with no structure
  - 20 observations with random labels
  - 6,000 possible but unrelated predictors
  - Repeated 200 times

- Estimated predictive accuracy using
  - no cross-validation
  - selecting features on full dataset, then using cross-validation
  - selecting features at each step of cross-validation

Simon et al., JNCI 2003



# OUTLINE

- Supervised classification
  - Classifiers, evaluation, issues
- Case studies
  - Alzheimer's diagnosis with plasma signaling proteins
  - Colorectal cancer diagnosis from plasma proteins
  - MACQ II: common practices in assay validation
- Experimental design for prediction
  - Impacts of dimensionality and noise

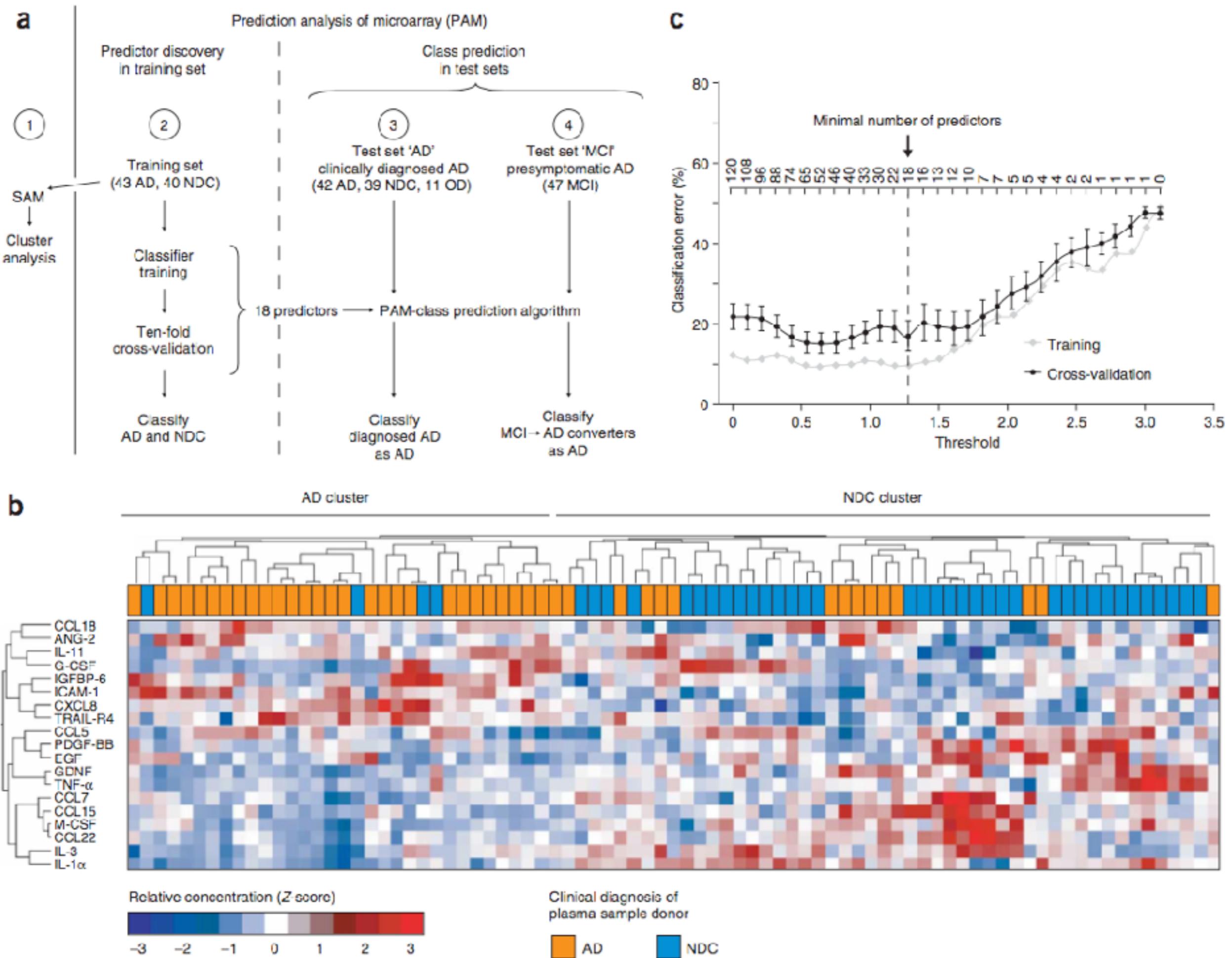
# Case study 2: cytokine arrays

nature  
medicine

## Classification and prediction of clinical Alzheimer's diagnosis based on plasma signaling proteins

Sandip Ray<sup>1,16</sup>, Markus Britschgi<sup>2,16</sup>, Charles Herbert<sup>1</sup>, Yoshiko Takeda-Uchimura<sup>2</sup>, Adam Boxer<sup>3</sup>, Kaj Blennow<sup>4</sup>, Leah F Friedman<sup>5</sup>, Douglas R Galasko<sup>6</sup>, Marek Jutel<sup>7</sup>, Anna Karydas<sup>3</sup>, Jeffrey A Kaye<sup>8</sup>, Jerzy Leszek<sup>9</sup>, Bruce L Miller<sup>3</sup>, Lennart Minthon<sup>10</sup>, Joseph F Quinn<sup>8</sup>, Gil D Rabinovici<sup>3</sup>, William H Robinson<sup>11</sup>, Marwan N Sabbagh<sup>12</sup>, Yuen T So<sup>2</sup>, D Larry Sparks<sup>12</sup>, Massimo Tabaton<sup>13</sup>, Jared Tinklenberg<sup>5</sup>, Jerome A Yesavage<sup>5</sup>, Robert Tibshirani<sup>14</sup> & Tony Wyss-Coray<sup>2,15</sup>

Ray et al., *Nature Medicine*, 2007



# OUTLINE

- Supervised classification
  - Classifiers, evaluation, issues
- Case studies
  - Alzheimer's diagnosis with plasma signaling proteins
  - Colorectal cancer diagnosis from plasma proteins
  - MACQ II: common practices in assay validation
- Experimental design for prediction
  - Impacts of dimensionality and noise

# CASE STUDY

Research Article



EMBO  
Molecular Medicine

## Prediction of colorectal cancer diagnosis based on circulating plasma proteins

Silvia Surinova<sup>1,†</sup>, Meena Choi<sup>2</sup>, Sha Tao<sup>3,‡</sup>, Peter J Schüffler<sup>4</sup>, Ching-Yun Chang<sup>2</sup>, Timothy Clough<sup>2</sup>, Kamil Vysloužil<sup>5</sup>, Marta Khoylou<sup>6</sup>, Josef Srovnal<sup>6</sup>, Yansheng Liu<sup>1</sup>, Mariette Matondo<sup>1</sup>, Ruth Hüttenhain<sup>1</sup>, Hendrik Weisser<sup>1</sup>, Joachim M Buhmann<sup>4</sup>, Marián Hajdúch<sup>6</sup>, Hermann Brenner<sup>3,7</sup>, Olga Vitek<sup>2,8,9,\*</sup> & Ruedi Aebersold<sup>1,10,\*\*</sup>

Research Article

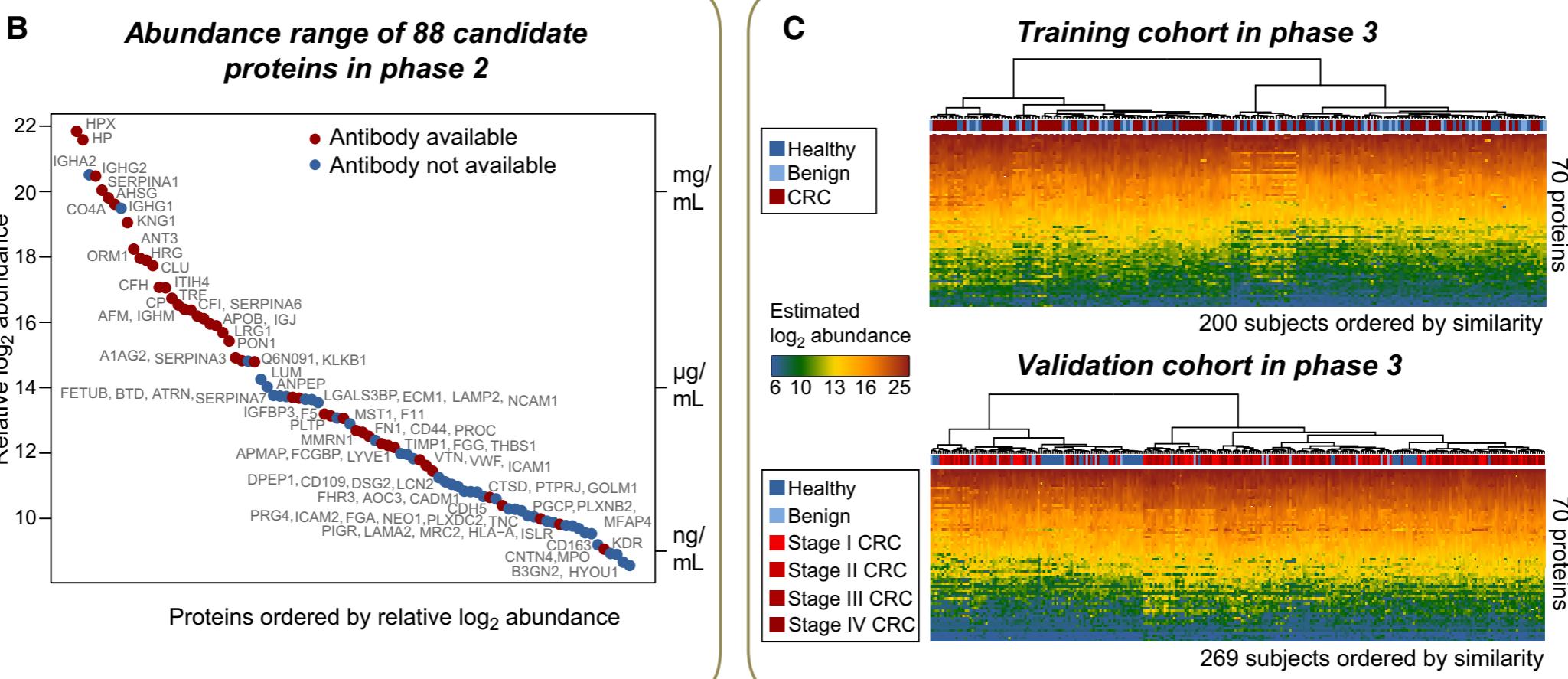
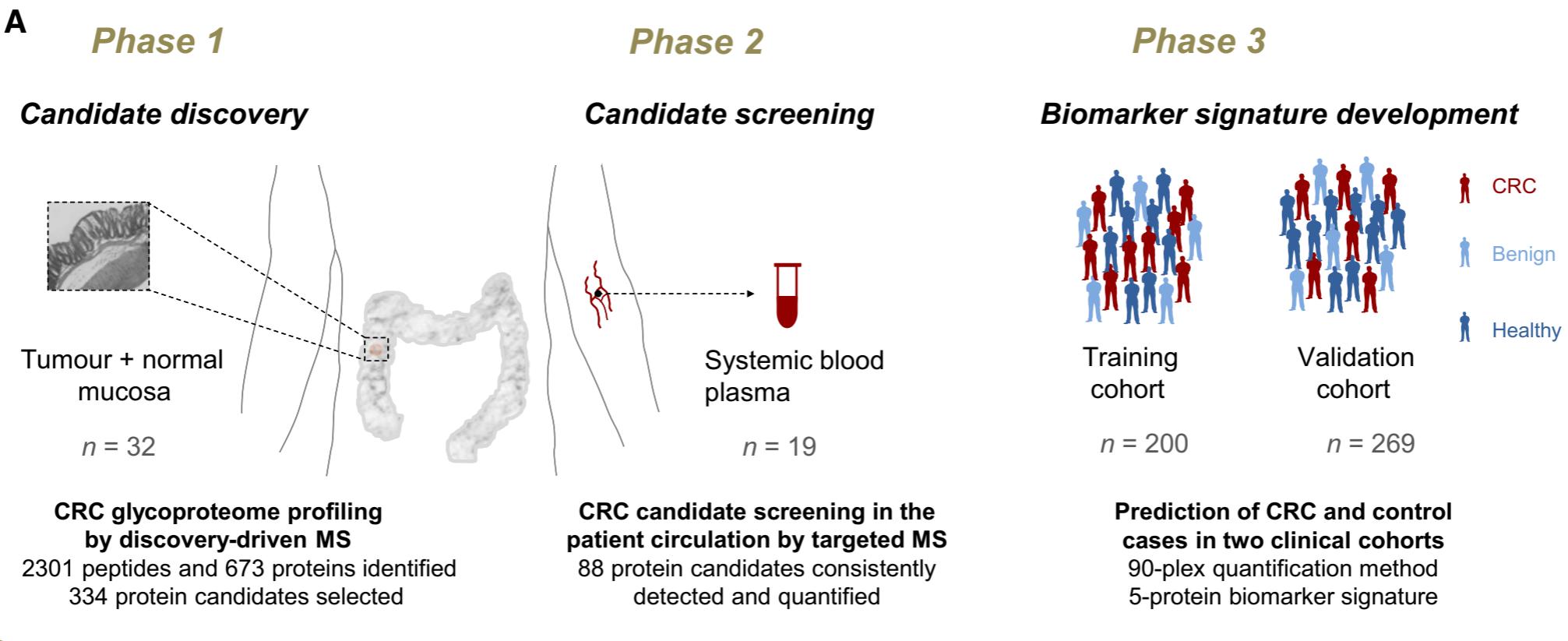


EMBO  
Molecular Medicine

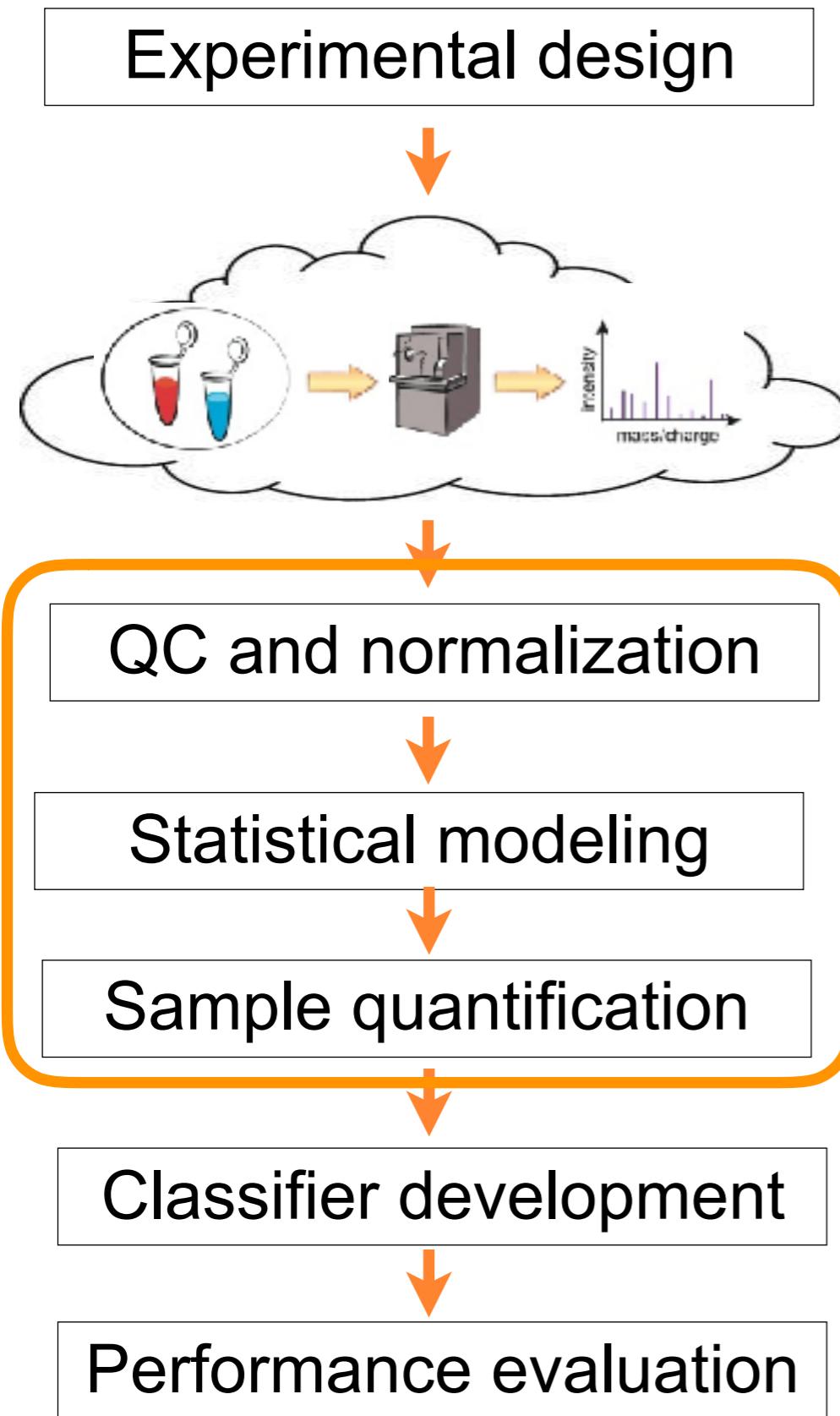
## Non-invasive prognostic protein biomarker signatures associated with colorectal cancer

Silvia Surinova<sup>1,†</sup>, Lenka Radová<sup>2,‡</sup>, Meena Choi<sup>3</sup>, Josef Srovnal<sup>2</sup>, Hermann Brenner<sup>4,5</sup>, Olga Vitek<sup>3,6,7,\*</sup>, Marián Hajdúch<sup>2</sup> & Ruedi Aebersold<sup>1,8,\*\*</sup>

# OVERVIEW



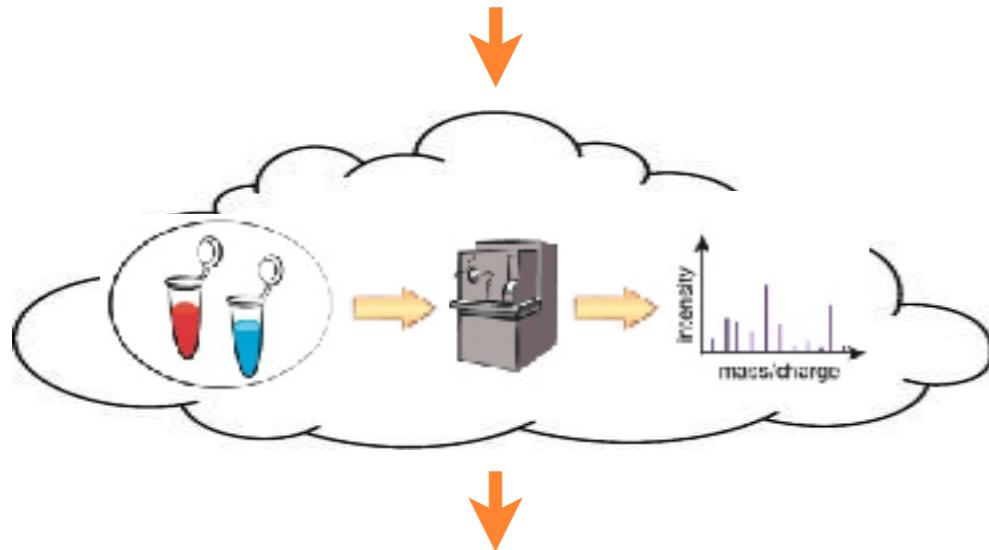
# ANALYSIS WORKFLOW WITH MSSTATS



**Bioconductor**

# ANALYSIS WORKFLOW WITH MSSTATS

## Experimental design



## QC and normalization

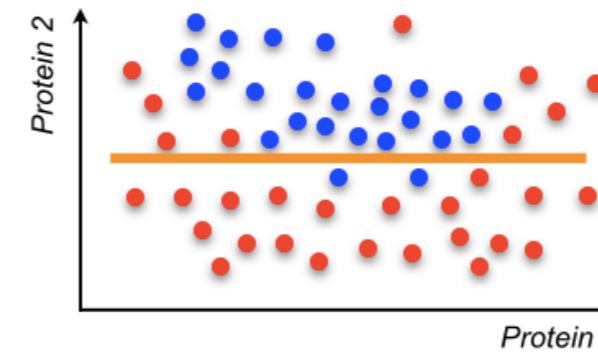
## Statistical modeling

## Sample quantification

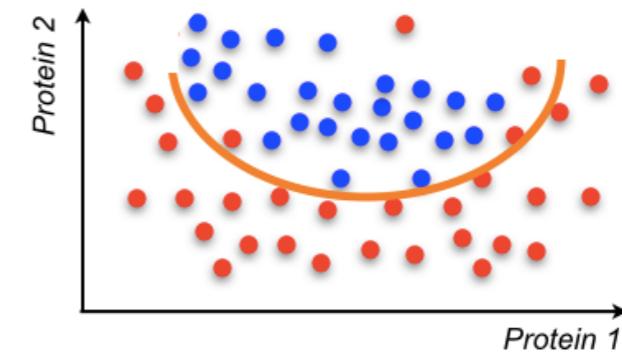
## Classifier development

## Performance evaluation

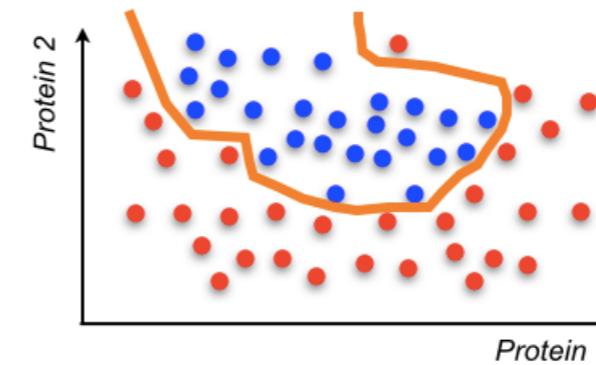
- Objective: Train a model that can classify each subject into correct class.



*Linear  
prediction  
rule*



*Quadratic  
prediction  
rule*



*General  
prediction rule*

# SELECTION OF BIOMARKER PANEL

Phase 3, training cohort: 100 CRC and 100 control subjects

## Step 1 Discovery of a multivariate signature within 10-fold cross-validation

- Systematically rotate samples between 10 folds
  - Select differentially abundant proteins on 9/10 of samples ► about 20 proteins per fold
  - Fit logistic regression models with stepwise protein selection ► 4-5 proteins per model
  - Predict the disease status of 1/10 ‘left-out’ samples ► AUROC (mean=0.62, SE=0.04)
- Generate a consensus logistic regression model with proteins selected in  $\geq 5$  folds:  
 $\ln\{\text{Prob}(\text{CRC})/(1-\text{Prob}(\text{CRC}))\} = -15.15 + 0.73 \times \text{CP} - 1.06 \times \text{PON1} + 0.69 \times \text{SERPINA3} + 0.37 \times \text{LRG1} + 0.41 \times \text{TIMP1}$

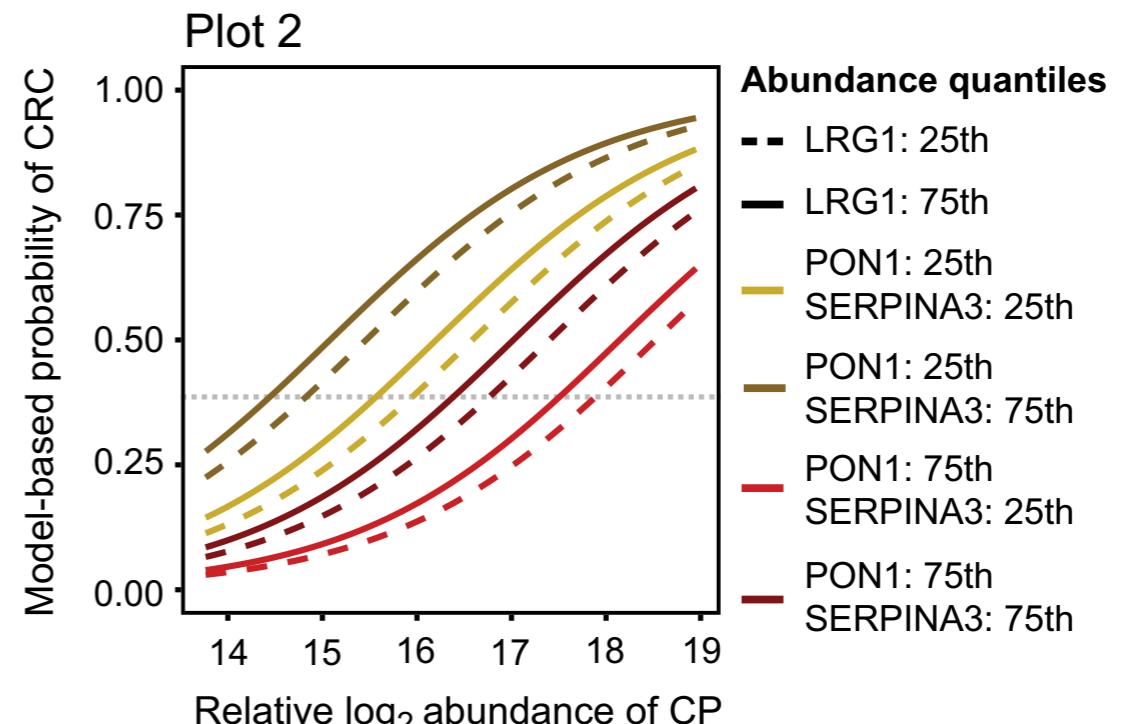
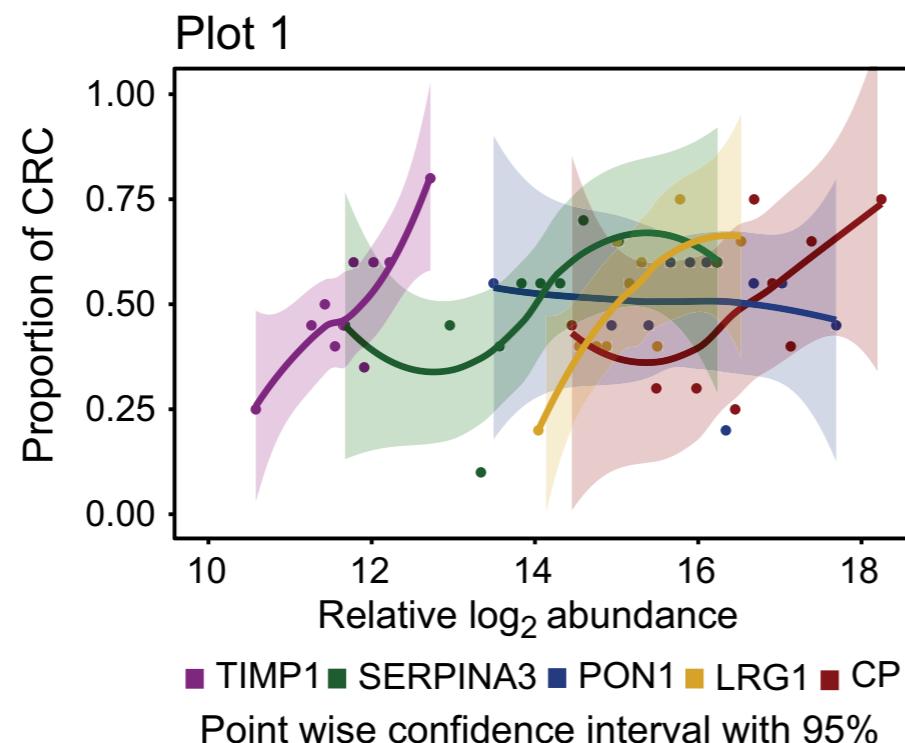
10x

## Step 2 Characterization of selected candidate biomarkers on the full training cohort

- Quantify protein fold changes and standard errors

	CP	PON1	SERPINA3	LRG1	TIMP1
Fold change (SE)	1.3 (0.05)	0.9 (0.04)	1.3 (0.06)	1.2 (0.03)	1.2 (0.03)

- Illustrate the range of protein intensities and their univariate contribution to disease (plot 1)
- Illustrate model-based estimated probability of disease (plot 2)



# EVALUATION OF BIOMARKER PANEL

Phase 3, validation cohort: 202 CRC and 67 control subjects

**Step 3 Assessment of the predictive ability of the biomarker signature**



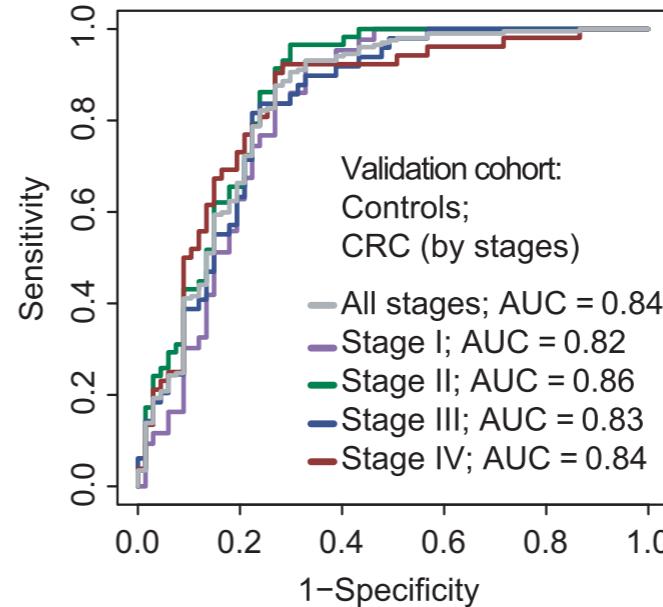
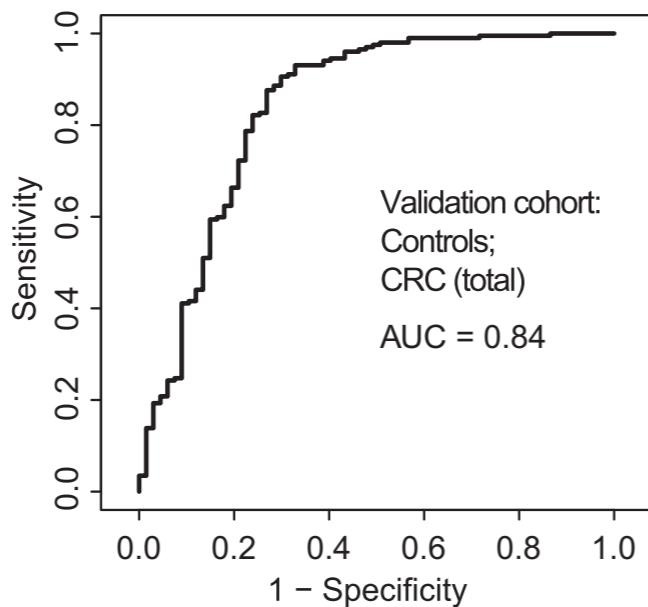
- Predict CRC and controls on the full validation dataset ►AUROC=0.84

**Characterization of selected candidate biomarkers on the validation cohort**

**Step 4**

- Evaluate the reproducibility of fold changes and standard errors

	CP	PON1	SERPINA3	LRG1	TIMP1
Fold change (SE)	1.6 (0.03)	0.8 (0.03)	1.2 (0.04)	1.5 (0.04)	1.4 (0.04)



**Summary statistics**

		99% CI
AUC	0.84	0.75 – 0.92
Specificity	0.79	0.64 – 0.91
Sensitivity	0.70	0.38 – 0.95
Accuracy	0.72	0.65 – 0.79

**AUC comparisons**

Stage	p-value				99% CI
	II	III	IV	All	
I	0.49	0.87	0.75	0.74	0.72–0.92
II		0.59	0.71	0.68	0.76–0.94
III			0.88	0.87	0.71–0.92
IV				1.00	0.73–0.93

# OUTLINE

- Supervised classification
  - Classifiers, evaluation, issues
- Case studies
  - Alzheimer's diagnosis with plasma signaling proteins
  - Colorectal cancer diagnosis from plasma proteins
  - MACQ II: common practices in assay validation
- Experimental design for prediction
  - Impacts of dimensionality and noise

# Case study 1: MACQCII

Shi et al., 2010

nature  
biotechnology

ARTICLES

The MicroArray Quality Control (MAQC)-II study of common practices for the development and validation of microarray-based predictive models

MAQC Consortium\*

Gene expression data from microarrays are being applied to predict preclinical and clinical endpoints, but the reliability of these predictions has not been established. In the MAQC-II project, 36 independent teams analyzed six microarray data sets to generate predictive models for classifying a sample with respect to one of 13 endpoints indicative of lung or liver toxicity in rodents, or of breast cancer, multiple myeloma or neuroblastoma in humans. In total, >30,000 models were built using many combinations of analytical methods. The teams generated predictive models without knowing the biological meaning of some of the endpoints and, to mimic clinical reality, tested the models on data that had not been used for training. We found that model performance depended largely on the endpoint and team proficiency and that different approaches generated models of similar performance. The conclusions and recommendations from MAQC-II should be useful for regulatory agencies, study committees and independent investigators that evaluate methods for global gene expression analysis.

**Negative controls:**  
randomly assigned endpoints I & M

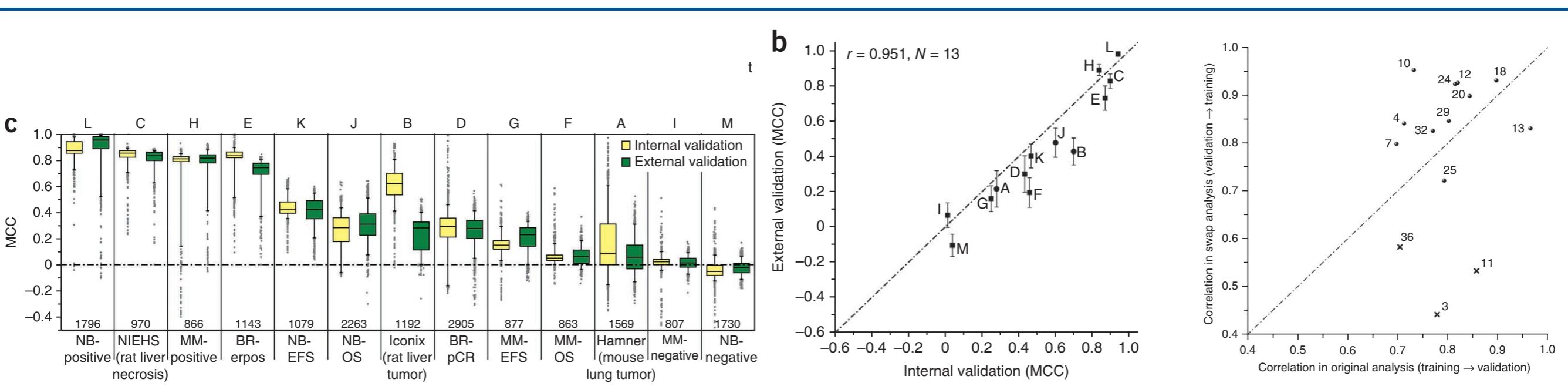
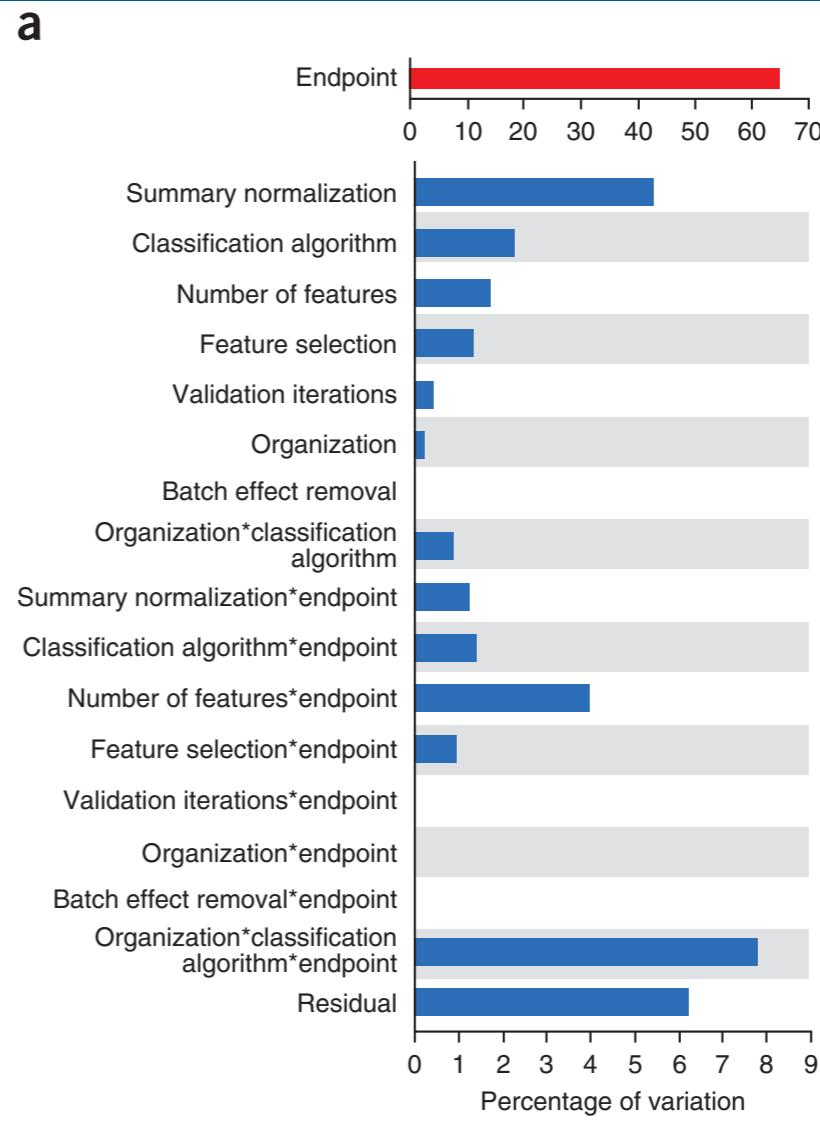
**Positive controls:**  
gender H & L

# Summary of the results

**Table 2 Modeling factor options frequently adopted by MAQC-II data analysis teams**

		Original analysis (training => validation)		
Modeling factor	Option	Number of teams	Number of endpoints	Number of models
Summary and normalization	Loess	12	3	2,563
	RMA	3	7	46
	MAS5	11	7	4,947
Batch-effect removal	None	10	11	2,281
	Mean shift	3	11	7,279
Feature selection	SAM	4	11	3,771
	FC+P	8	11	4,711
	T-Test	5	11	400
	RFE	2	11	647
	0~9	10	11	393
Number of features	10~99	13	11	4,445
	≥1,000	3	11	474
	100~999	10	11	4,298
Classification algorithm	DA	4	11	103
	Tree	5	11	358
	NB	4	11	924
	KNN	8	11	6,904
	SVM	9	11	986

Analytic options used by two or more of the 14 teams that submitted models for all endpoints in both the original and swap experiments. RMA, robust multichip analysis; SAM, significance analysis of microarrays; FC, fold change; RFE, recursive feature elimination; DA, discriminant analysis; Tree, decision tree; NB, naive Bayes; KNN, K-nearest neighbors; SVM, support vector machine.

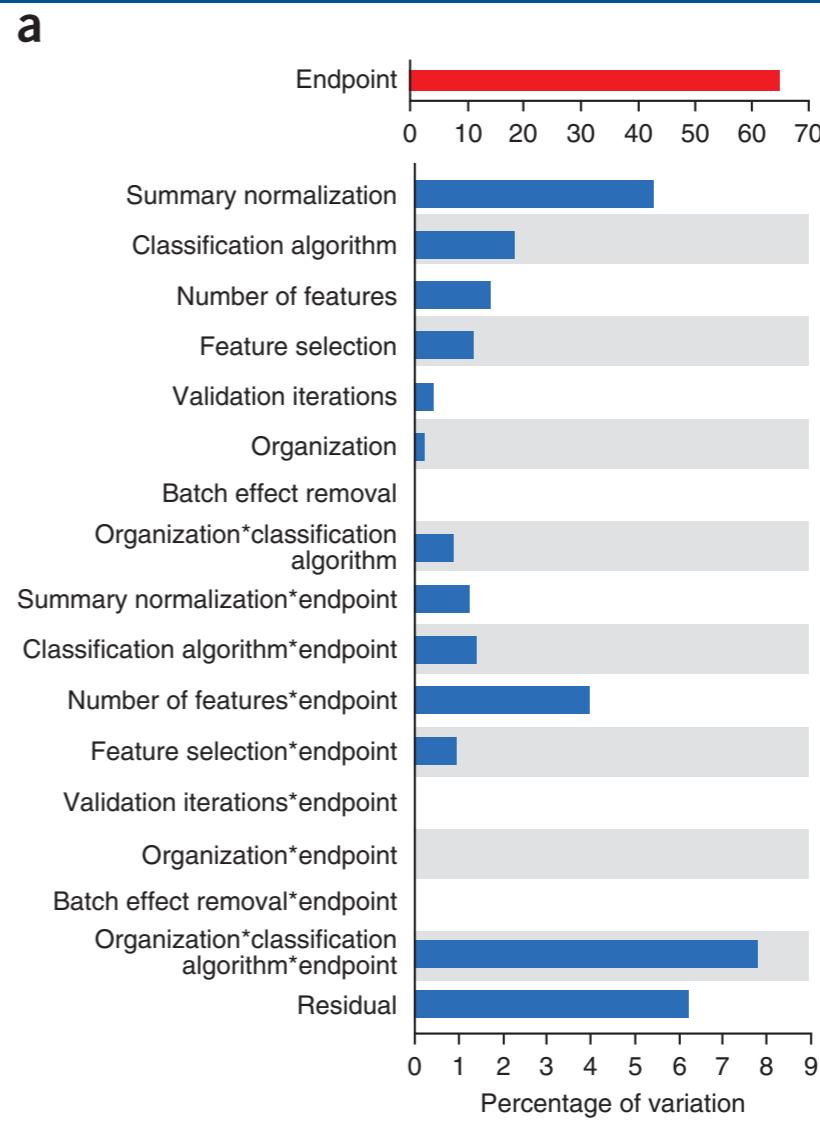


# Summary of the results

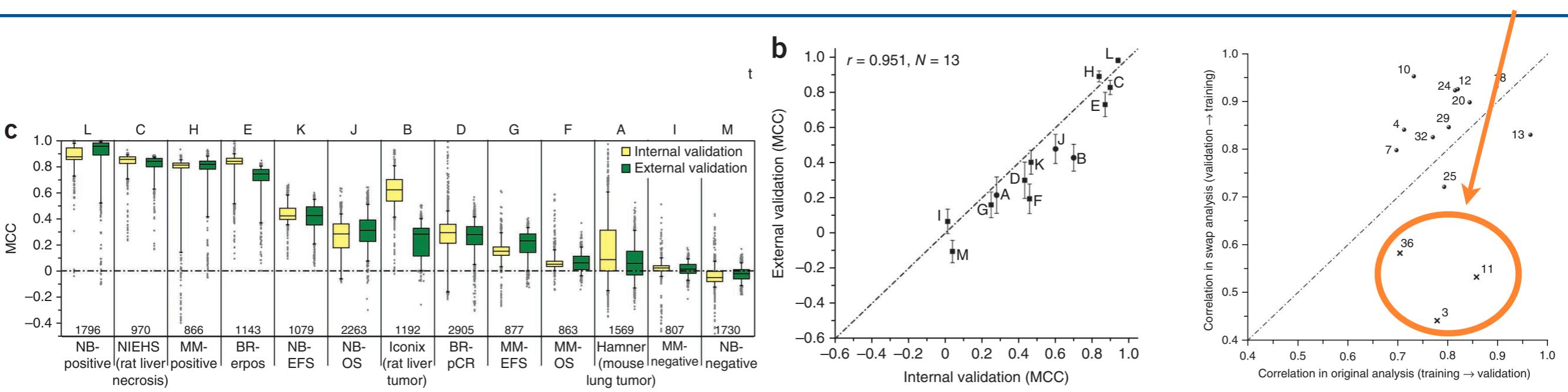
**Table 2 Modeling factor options frequently adopted by MAQC-II data analysis teams**

		Original analysis (training => validation)		
Modeling factor	Option	Number of teams	Number of endpoints	Number of models
Summary and normalization	Loess	12	3	2,563
	RMA	3	7	46
	MAS5	11	7	4,947
Batch-effect removal	None	10	11	2,281
	Mean shift	3	11	7,279
Feature selection	SAM	4	11	3,771
	FC+P	8	11	4,711
	T-Test	5	11	400
	RFE	2	11	647
Number of features	0~9	10	11	393
	10~99	13	11	4,445
	≥1,000	3	11	474
	100~999	10	11	4,298
Classification algorithm	DA	4	11	103
	Tree	5	11	358
	NB	4	11	924
	KNN	8	11	6,904
	SVM	9	11	986

Analytic options used by two or more of the 14 teams that submitted models for all endpoints in both the original and swap experiments. RMA, robust multichip analysis; SAM, significance analysis of microarrays; FC, fold change; RFE, recursive feature elimination; DA, discriminant analysis; Tree, decision tree; NB, naive Bayes; KNN, K-nearest neighbors; SVM, support vector machine.



operator errors

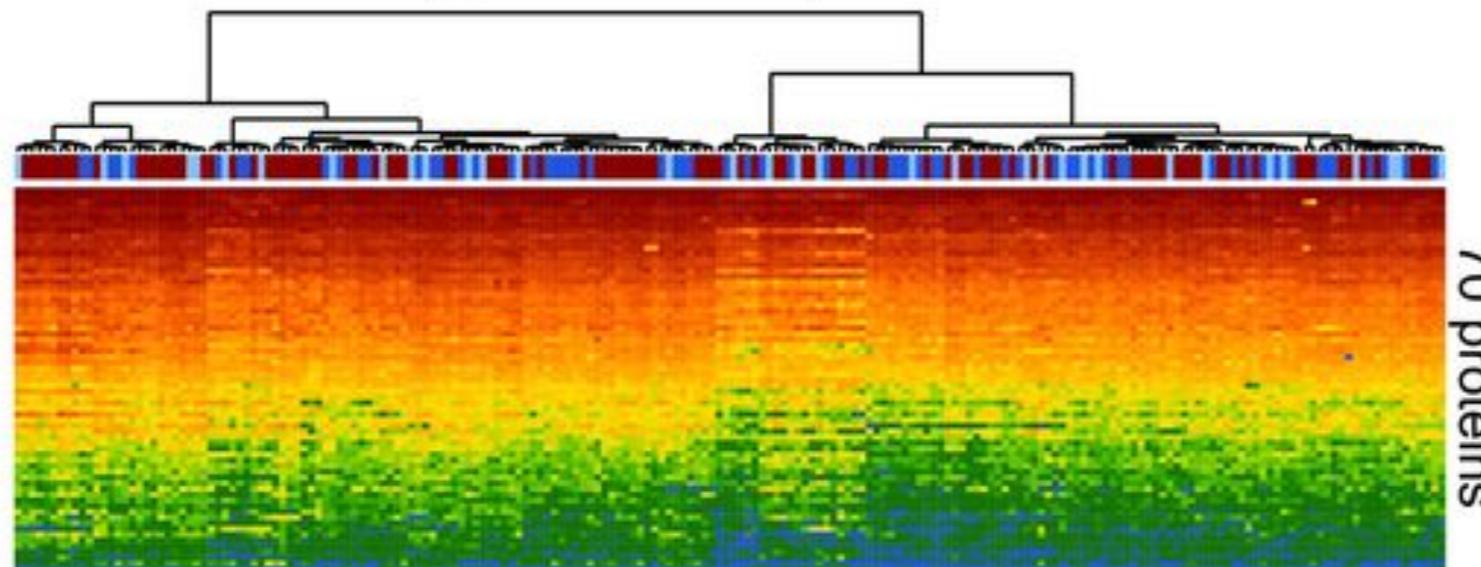


# OUTLINE

- Supervised classification
  - Classifiers, evaluation, issues
- Case studies
  - Alzheimer's diagnosis with plasma signaling proteins
  - Colorectal cancer diagnosis from plasma proteins
  - MACQ II: common practices in assay validation
- Experimental design for prediction
  - Impacts of dimensionality and noise

# DOES THE CLASSIFICATION ACCURACY DEPEND ON THE PROTEIN/SAMPLE SPACE?

*Training cohort in phase 3*



More subjects are better for biomarker discovery?

200 subjects ordered by similarity



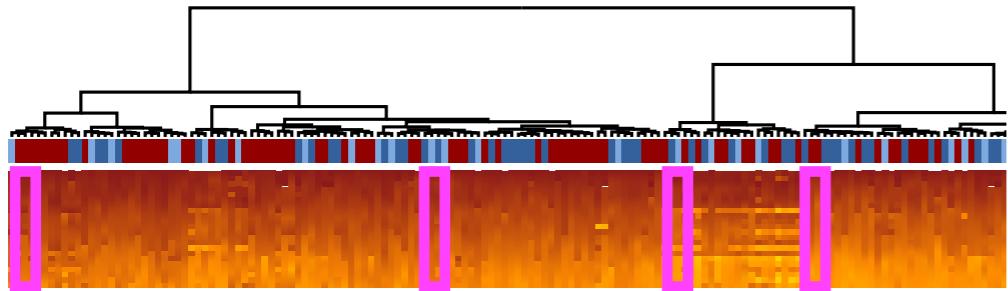
More proteins are better for biomarker discovery?

	DDA	SRM/PRM	DIA
Proteomic coverage	High (1000s~10000)	Low (10s~100s)	Medium (1000s)
Sample throughput	Low (10~100s, even one sample)	High (100s~1000s)	Medium (100s)
Sensitivity	Low	High	Medium

# REPRODUCIBILITY EVALUATION

Resample subjects in training set, while fixing proteins number

## *Training cohort in phase 3*

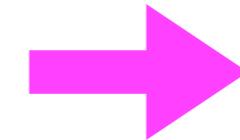


*Protein rank: random forest classifier*

10 predictive proteins

Protein	Mean accuracy decrease
ACTG1	4.29792551
H2AFX	4.17272432
ACTA2	3.64730166
HIST2H2AC	3.26940544
HIST1H2AB	3.10230189
ACTA1	3.07322088
TUBA1C	2.76765106
ACTG2	2.28619635
HSP90AA1	2.02174233
TUBB4A	1.99468791
HIST1H4A	1.86925866
SPTAN1	1.84453388
TUBB2B	1.75602473
COL6A3	1.72005229
TUBB4B	1.70743012
HIST1H2BK	1.68147503
HBB	1.61645207
HIST2H2AB	1.50013107

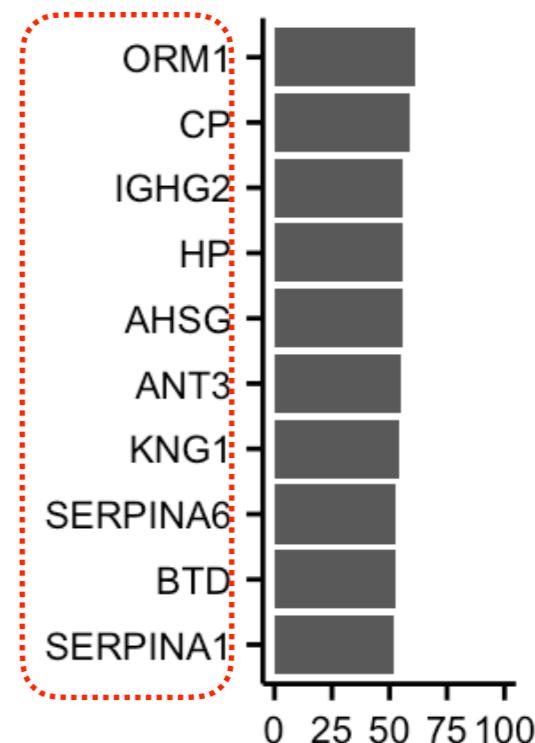
Repeat  
100 times



## *Decision*

	CRC	Control
Truth		
CRC	TP	FN
Control	FP	TN

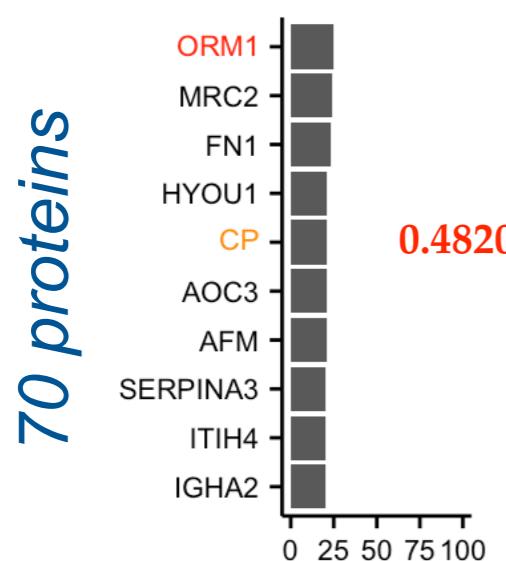
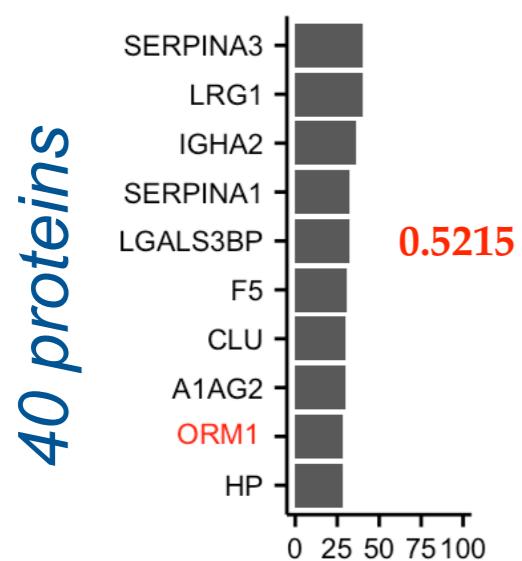
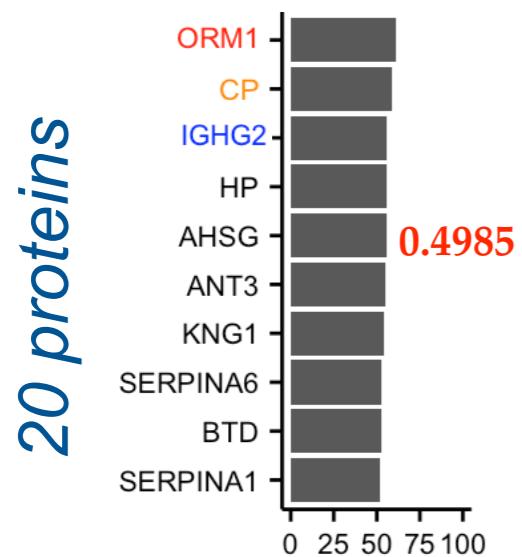
*Top 10 most frequently selected proteins*



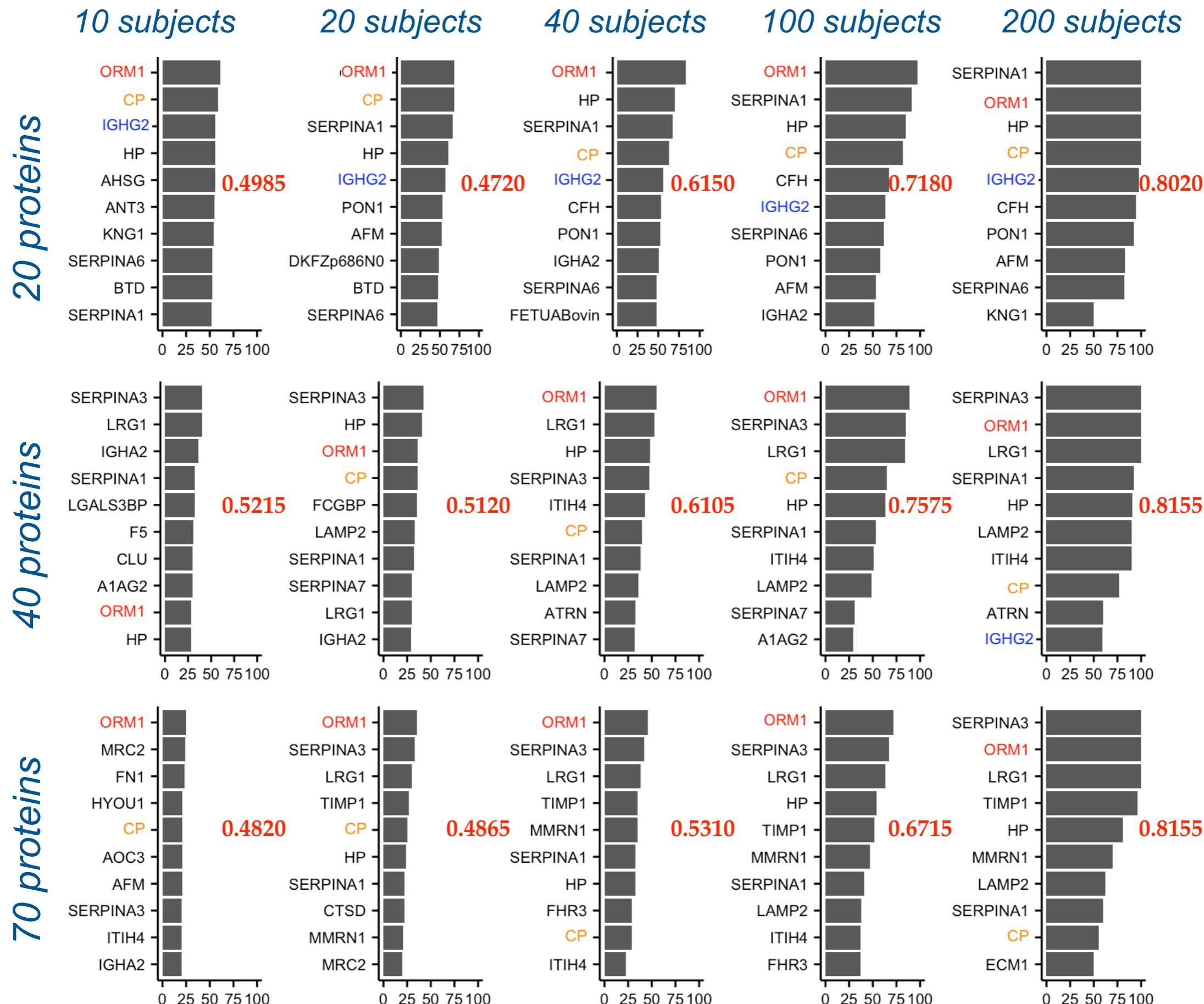
*Number of iterations  
that a protein is identified predictive*

# 200 COLORECTAL CANCER SUBJECTS: SRM

10 subjects

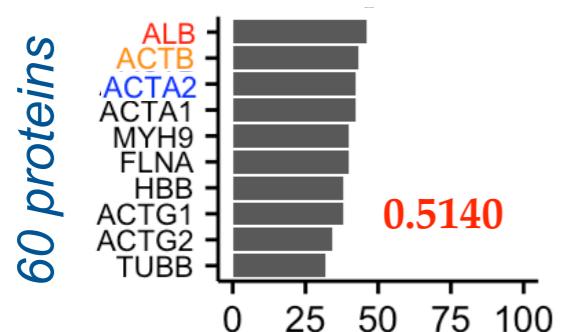


# 200 COLORECTAL CANCER SUBJECTS: SRM

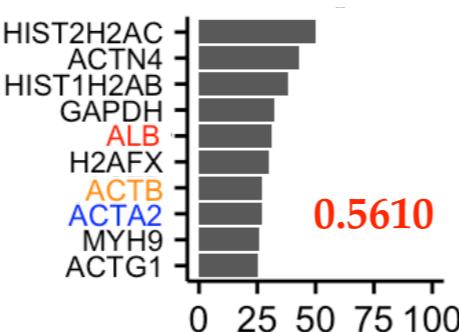


# 75 COLORECTAL CANCER SUBJECTS: DDA

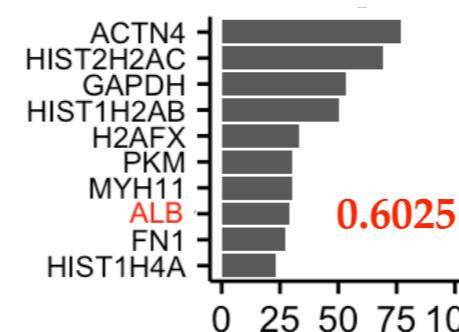
10 subjects



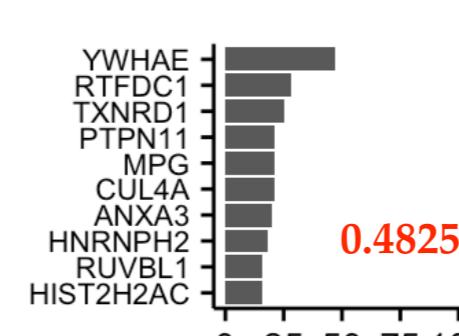
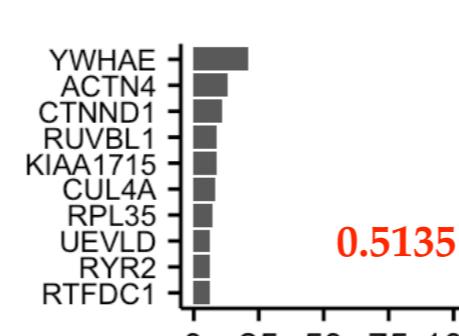
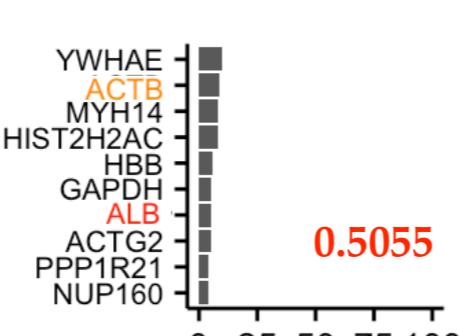
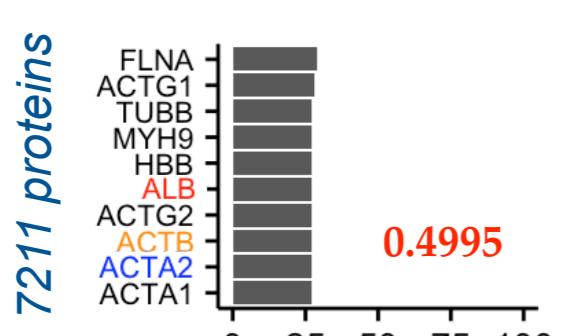
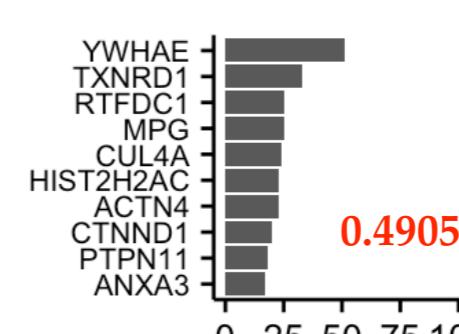
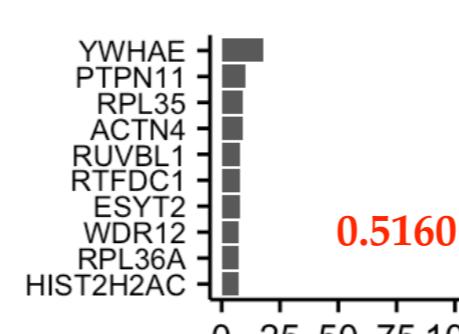
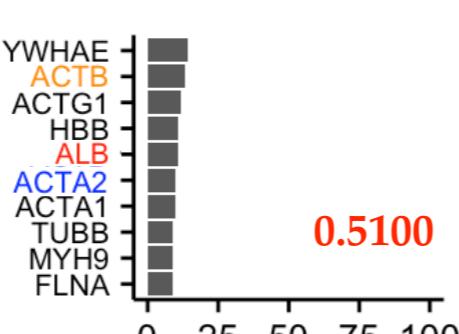
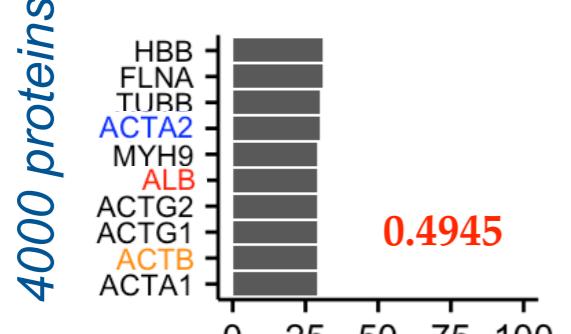
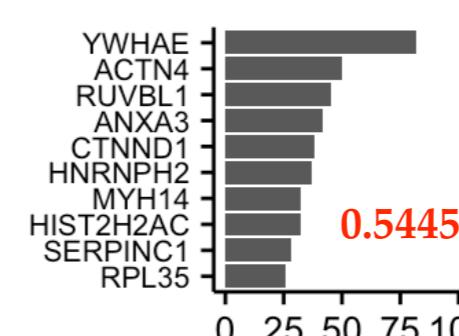
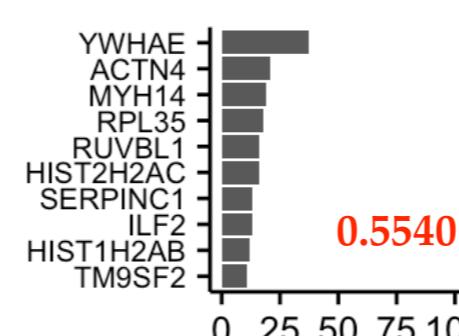
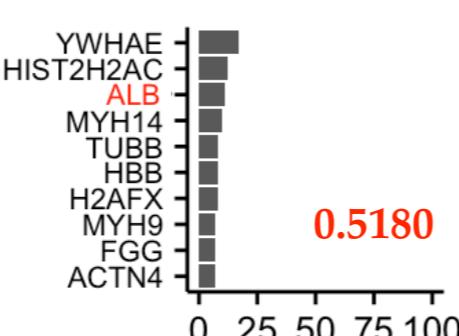
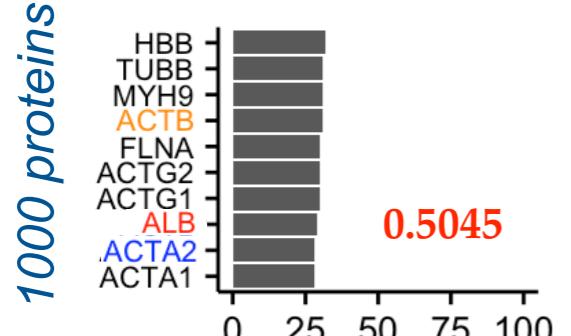
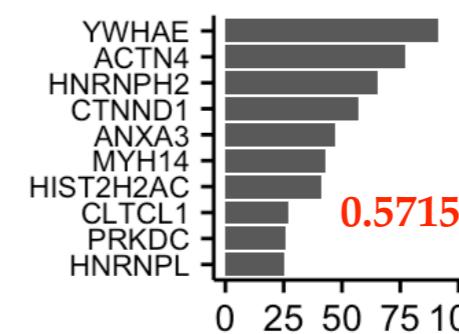
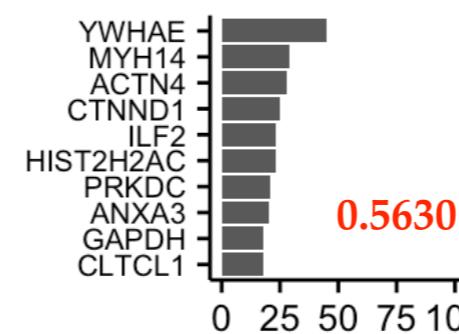
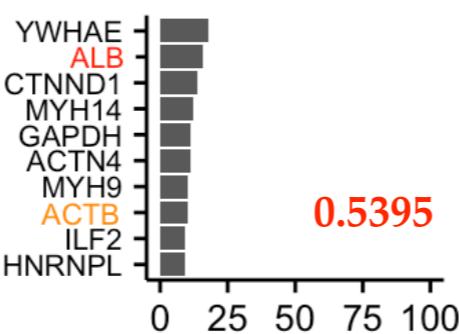
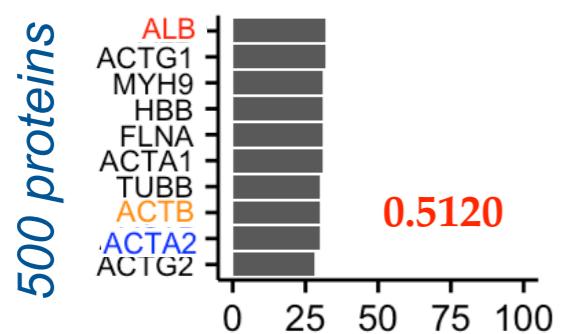
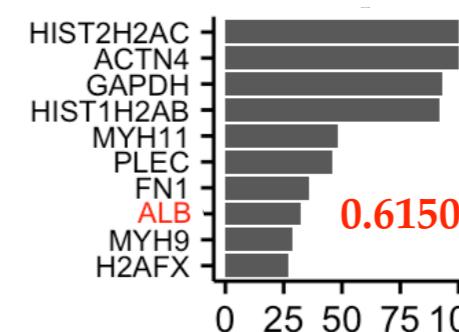
20 subjects



40 subjects

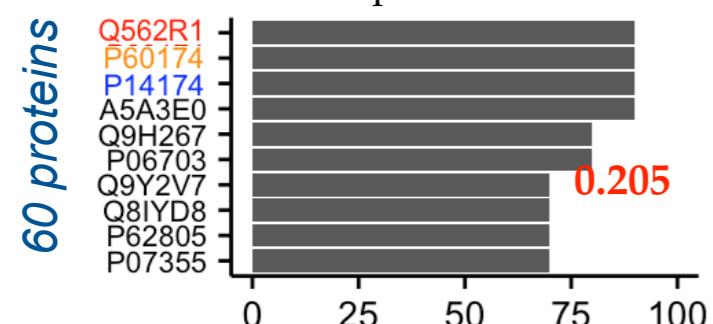


74 subjects

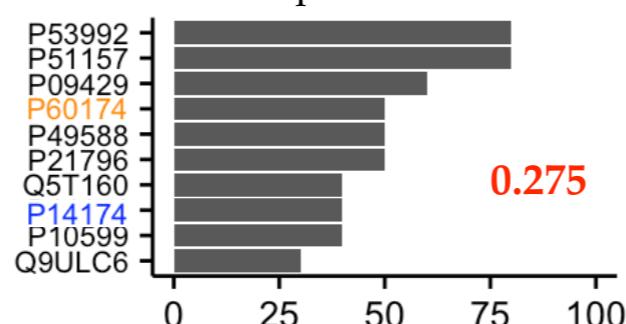


# 60 CELL LINES : DIA

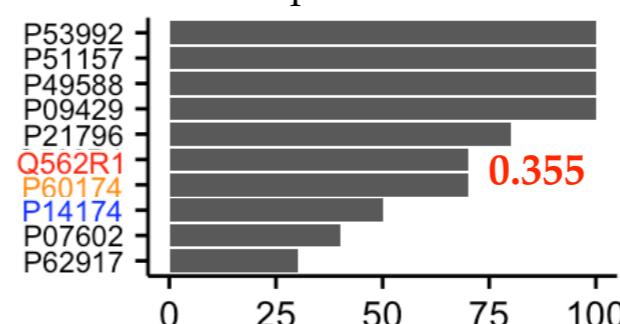
# *10 subjects*



## *20 subjects*



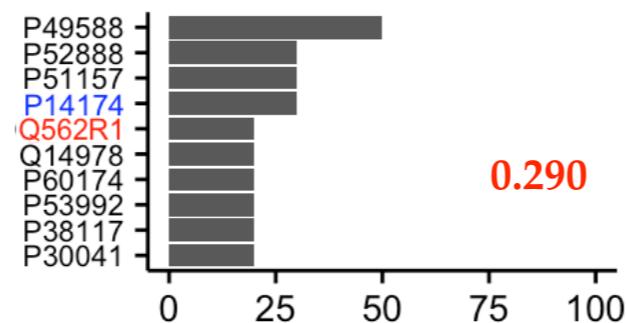
## 40 subjects



*500 proteins*

Protein ID	Abundance (approx.)
P07355	85
Q9Y2V7	75
Q9H267	75
Q8IYD8	75
Q562R1	75
P62805	75
P60174	75
P14174	75
P06703	75
A5A3E0	75

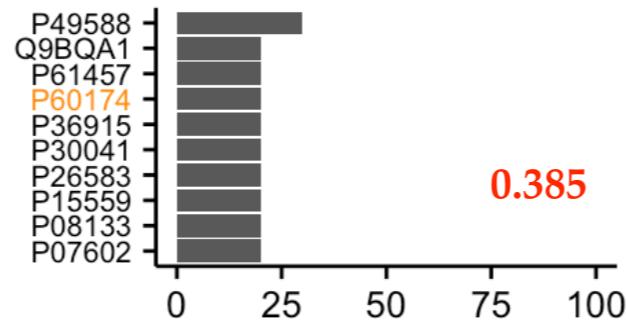
0.190



Protein ID	Percentage
P49588	100
P08133	100
P53992	95
P22234	74
P55011	65
P52888	52
P36915	50
Q9NTK5	45
P51157	42
P28072	30

*1000 proteins*

Protein ID	Abundance (approx.)
P14174	75
A5A3E0	75
Q9Y2V7	75
Q9H267	75
Q8IYD8	75
<b>Q562R1</b>	0.235
P62805	75
<b>P60174</b>	0.235
P07355	75
P06703	75



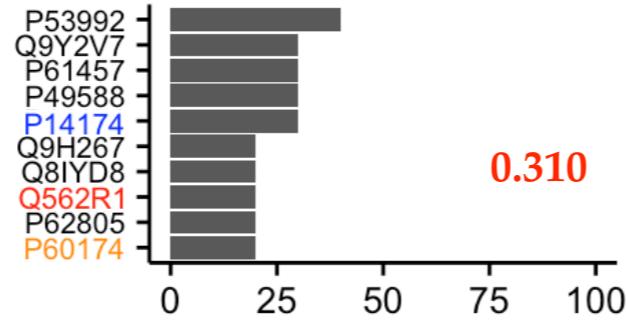
Protein ID	Value
P61457	~90
P49588	~90
P53992	~85
P22234	~85
P08133	~85
P36915	~50
P17174	~50
Q8WUM4	~40
P51157	~40
P31153	~40

0.435

**2000 proteins**

Protein ID	Abundance (approx.)
Q562R1	40
Q9Y2V7	30
Q9H267	30
Q8IYD8	30
P62805	30
P60174	10
P14174	10
P07355	30
P06703	30
A5A3E0	30

0.220



Protein	Value
Q9HAV0	~98
Q9NQW7	~95
P49588	~92
P53992	~74
P61457	~50
Q9H3R2	~42
P52888	~40
P22234	~30
P14550	~30
P08133	~30

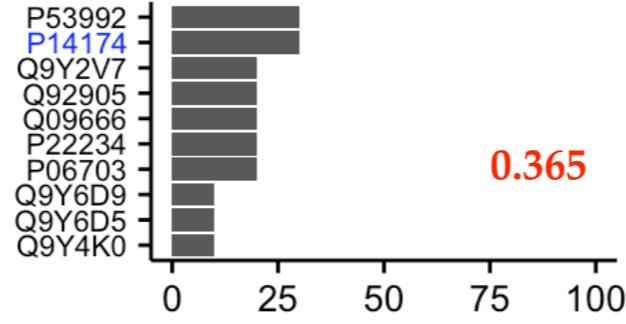
0.420

3717 proteins

Protein ID	Length (approx.)
Q9Y2V7	65
Q9H267	50
Q8IYD8	50
<b>Q562R1</b>	50
P62805	50
<b>P60174</b>	~15
<b>P14174</b>	~15
P07355	50
P06703	50
A5A3E0	50

0 25 50 75 100

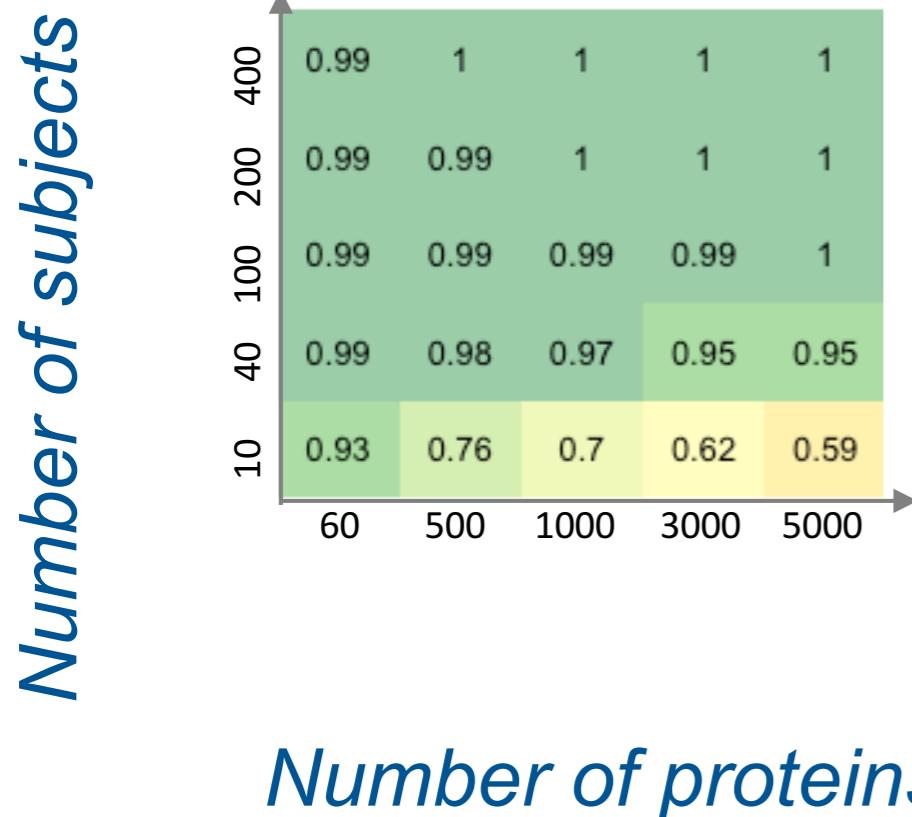
0.160



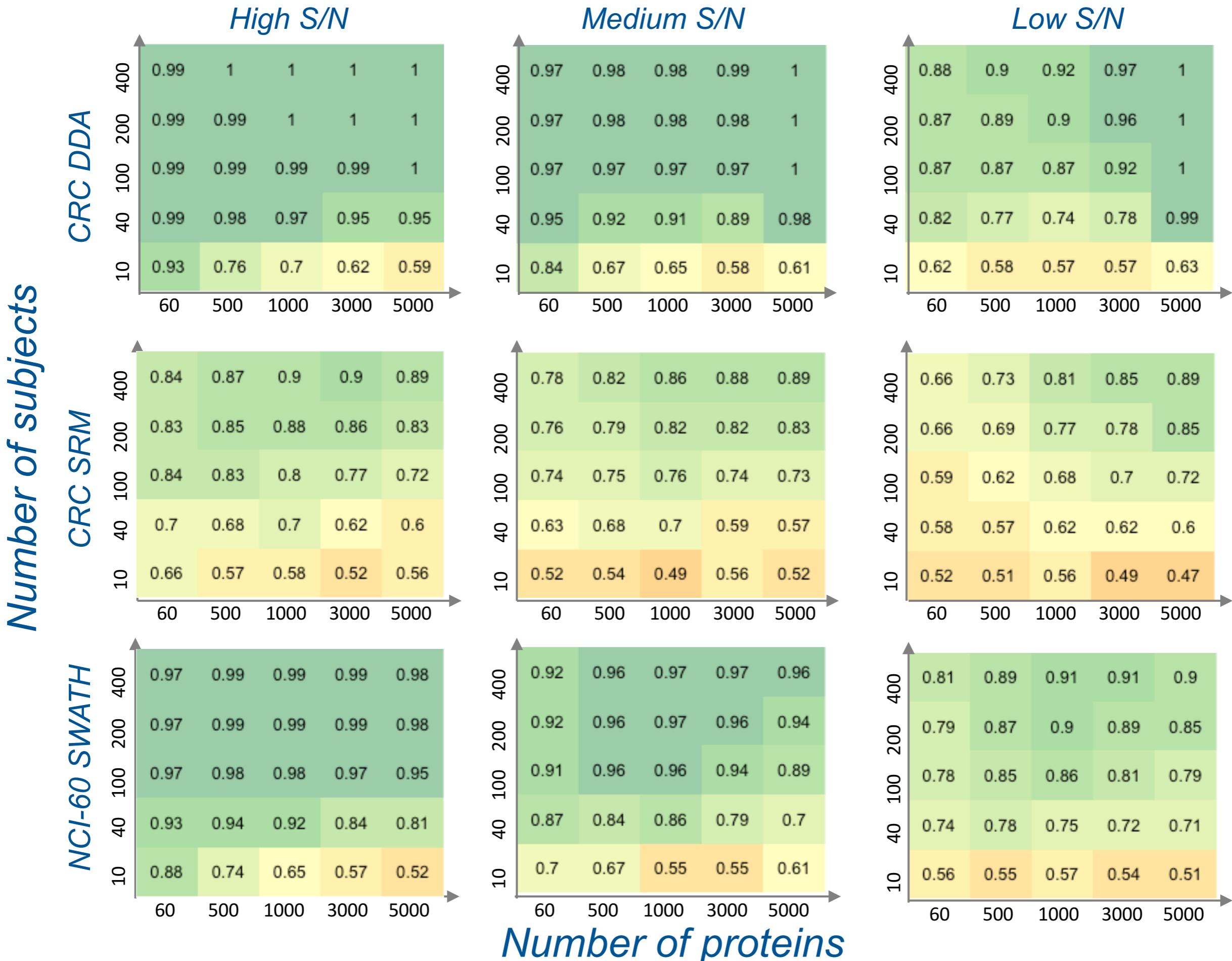
Sample	Expression Level (%)
Q9BRP8	~90
Q9HAV0	~85
P49588	~85
P22059	~85
Q9NQW7	~75
P61457	~40
P53992	~40
P08133	~40
P17174	~30
Q9H0S4	~20

0.400

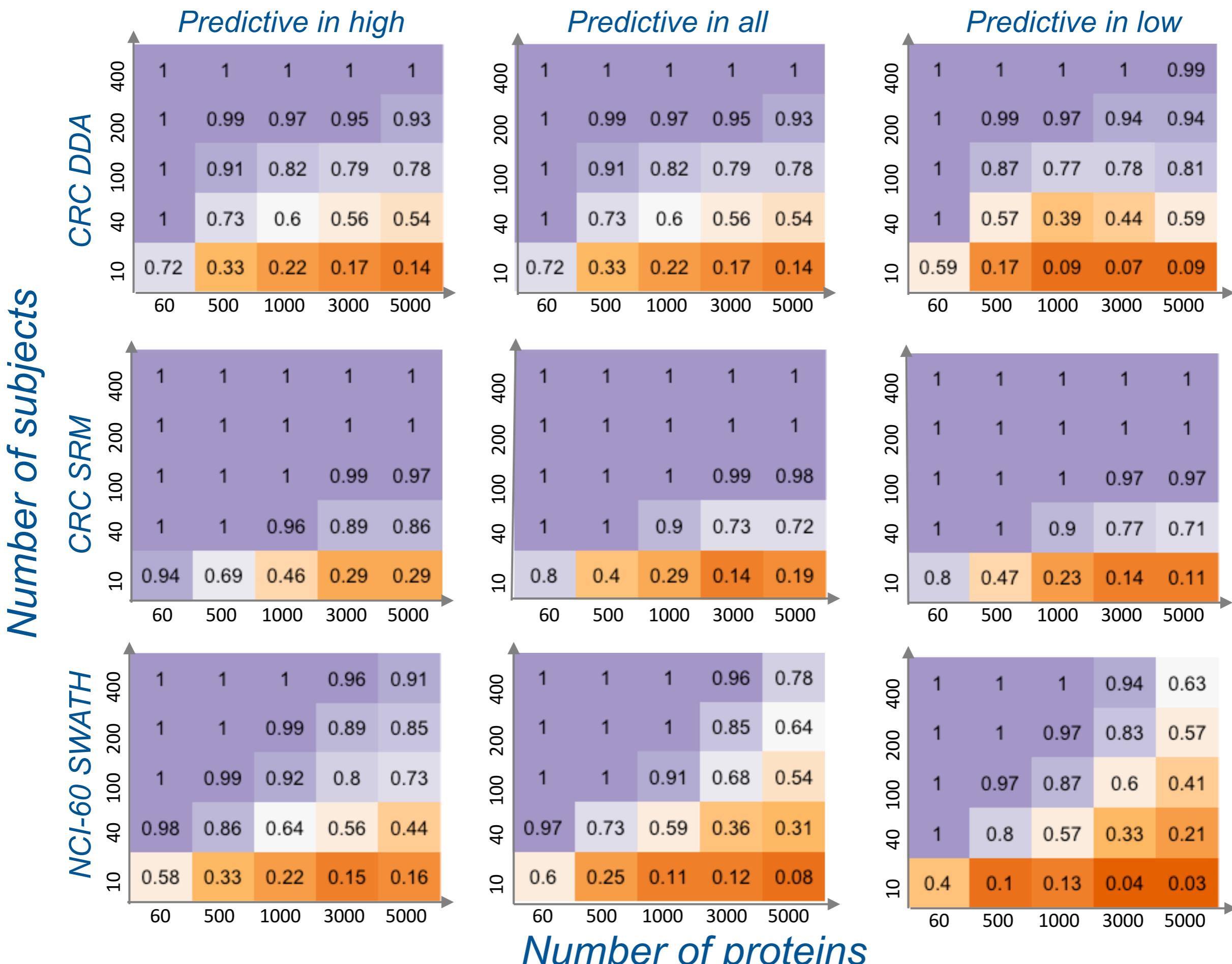
# INCREASING PROTEIN NUMBER DECREASES ACCURACY



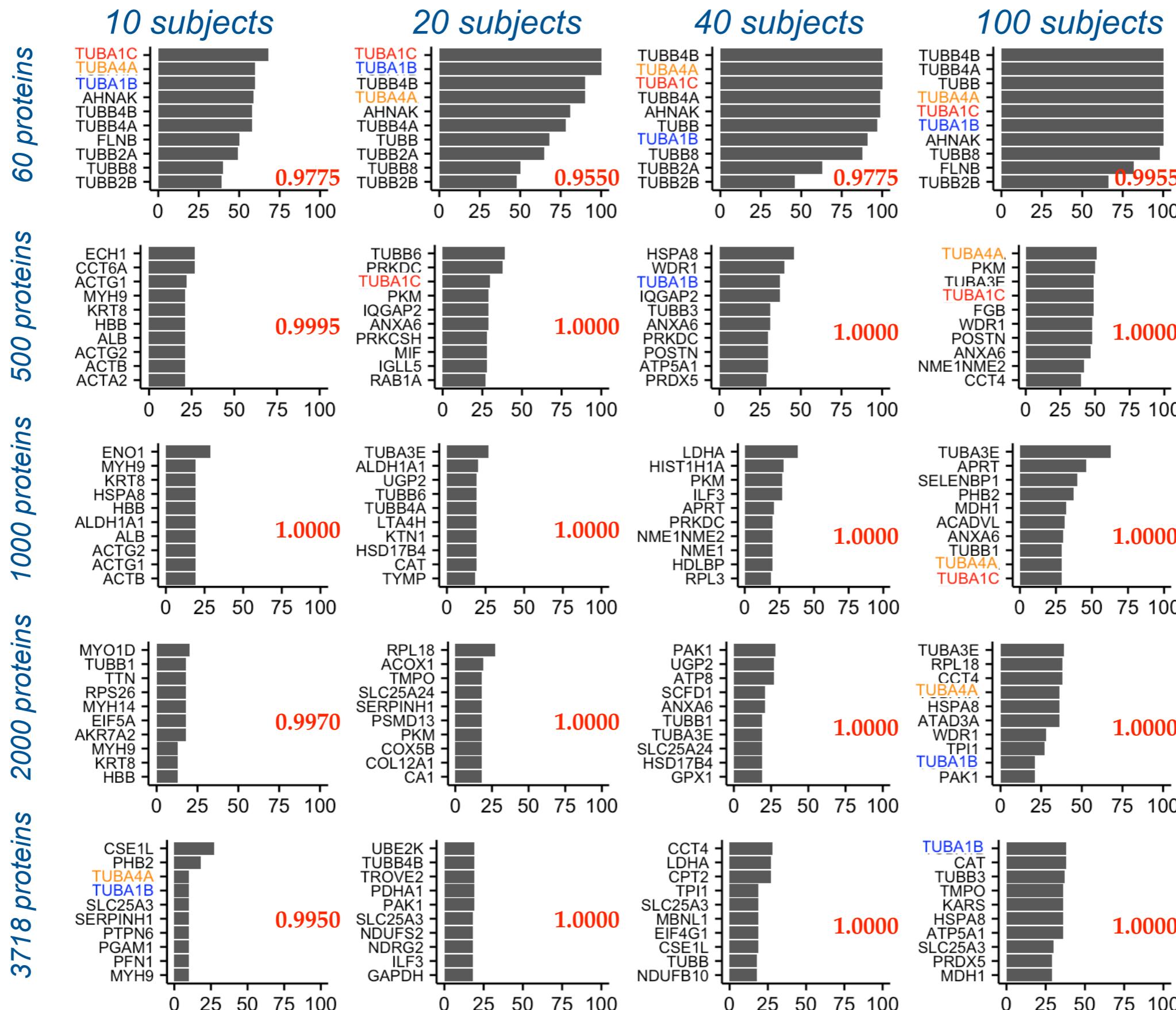
# INCREASING PROTEIN NUMBER DECREASES ACCURACY



# MORE REPLICATES SELECT MORE PREDICTIVE PROTEINS



# PERFECT CLASSIFICATION: BIASED DESIGN



# ACKNOWLEDGEMENTS

## Northeastern University

Kylie Bemis  
Meena Choi  
Eralp Dogu  
Dan Guo  
April Harry  
Ting Huang  
Cyril Galitzine  
Robert Ness  
Sara Taheri  
Tsung-Heng Tsai

## ETH Zurich

Ruedi Aebersold  
Tiannan Guo  
Ruth Huttenhain  
Paola Picotti  
Silvia Surinova  
Bernd Wollscheid

## University of Washington

Michael MacCoss  
Brendan MacLean  
Jarrett Egertson

## Biognosis

Lukas Reiter

## iPRG 2015

Zeynep Eren-Dogu  
Chris Colangelo  
John Cottrell  
Michael Hopman  
Eugene Kapp  
Santa Kim  
Henry Lam  
Tom Neubert  
Magnus Palmblad  
Brett Phinney  
Sue Weintraub,

