



# *Winning the Rat Race*

## Sub-challenge 2: Inter Species Prediction of Protein Phosphorylation

**Gyan Bhanot**

Rutgers University

**Michael Biehl**

University of Groningen

**Adel Dayarian**

Kavli Institute of Theoretical Physics

**Sahand Hormoz**

UC Santa Barbara



## *The team*

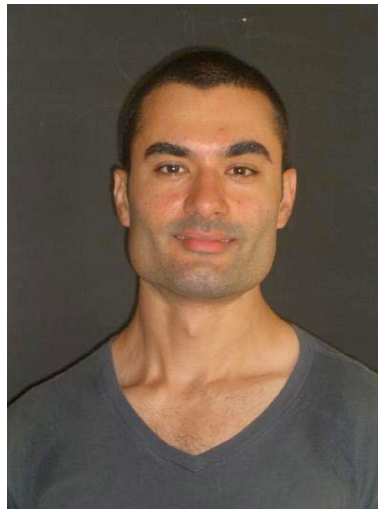
**Gyan Bhanot**

**Rutgers University and  
Inst. for Adv. Study, Princeton**



**Sahand Hormoz**

**Kavli Institute for Theoretical Physics,  
University of California, Santa Barbara**



**Adel Dayarian**

**Kavli Institute for Theoretical Physics,  
University of California, Santa Barbara**





**set-up:** phosphorylation data only

**naïve prediction:**

assume “human  $\approx$  rat” (in terms of phosphorylation)

**machine learning approach:**

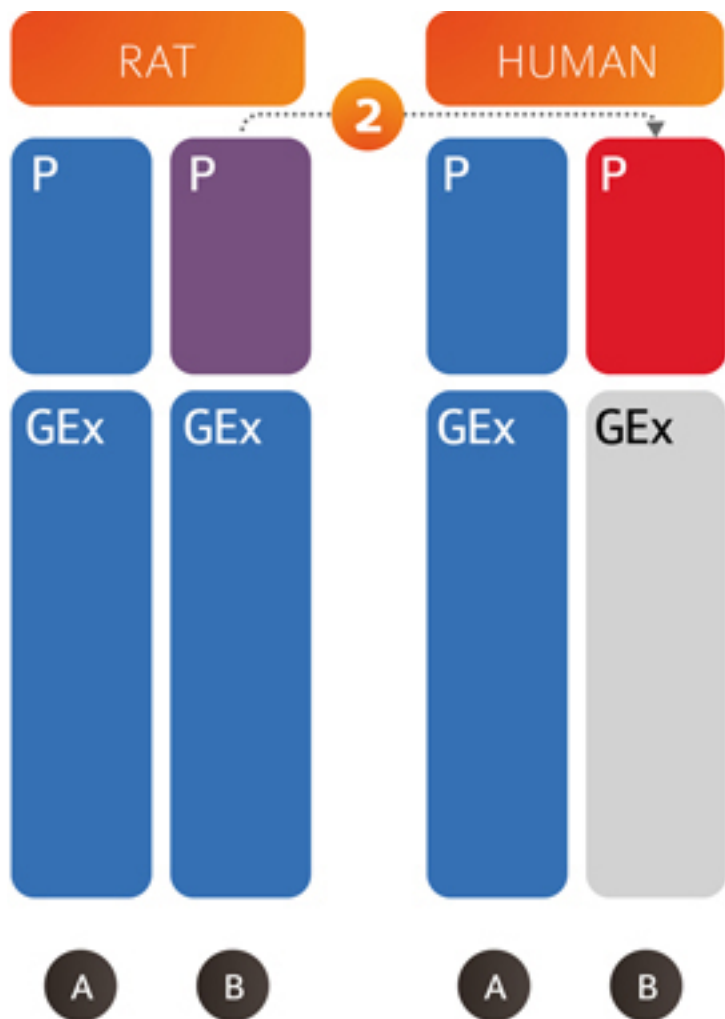
Learning Vector Quantization (LVQ) classifier

**combination of predictions:**

based on validation procedures to  
estimate performances



[www.sbvimprover.com](http://www.sbvimprover.com)



### Legend:

P Phosphorylation

GEx Gene expression

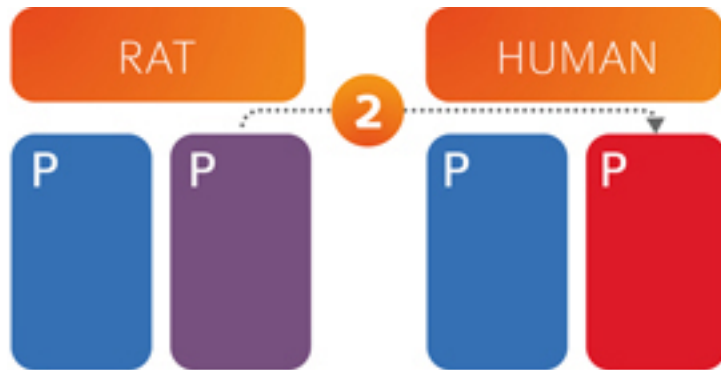
A B Stimulus subset

Not provided

Provided data

Predicted data

Provided after 1 July



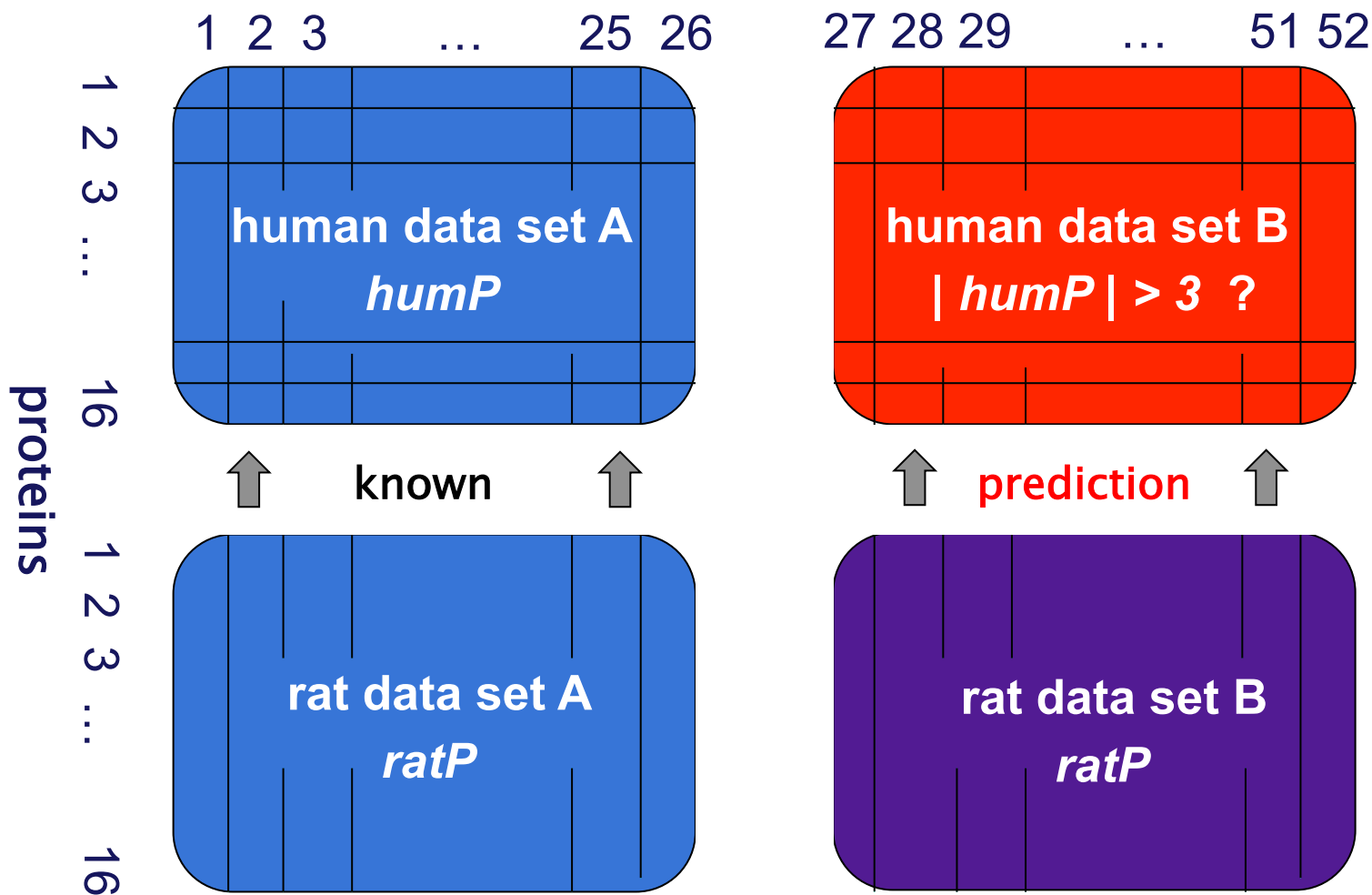
[www.sbvimprover.com](http://www.sbvimprover.com)

**preprocessing:** median of 3 replicates  
maximum absolute level of (@5 min, @25 min)

**target:** recommended threshold of 3.0  
for phosphorylation



stimuli

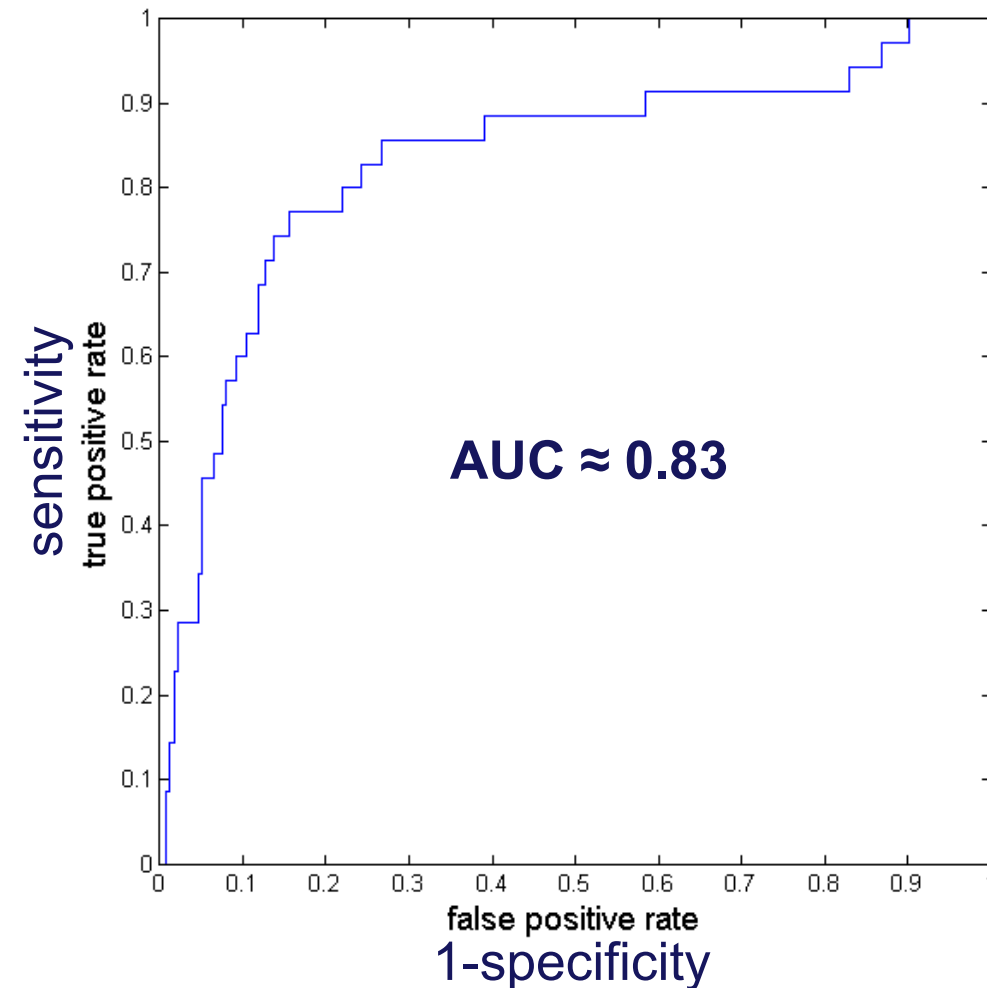


assume similar activation in both species: “human  $\approx$  rat”

prediction *score*, corresponding to **threshold 3** for activation

$$c_{naive} = \frac{1}{2} \left[ 1 + \tanh \left( \frac{|ratP| - 3}{5} \right) \right] \in [0, 1]$$

- precise (monotonic!) form is irrelevant for ROC, PR etc.
- threshold 0.5 for crisp classification
- here: scaling factor yields values well-spread in [0,1]



## ROC

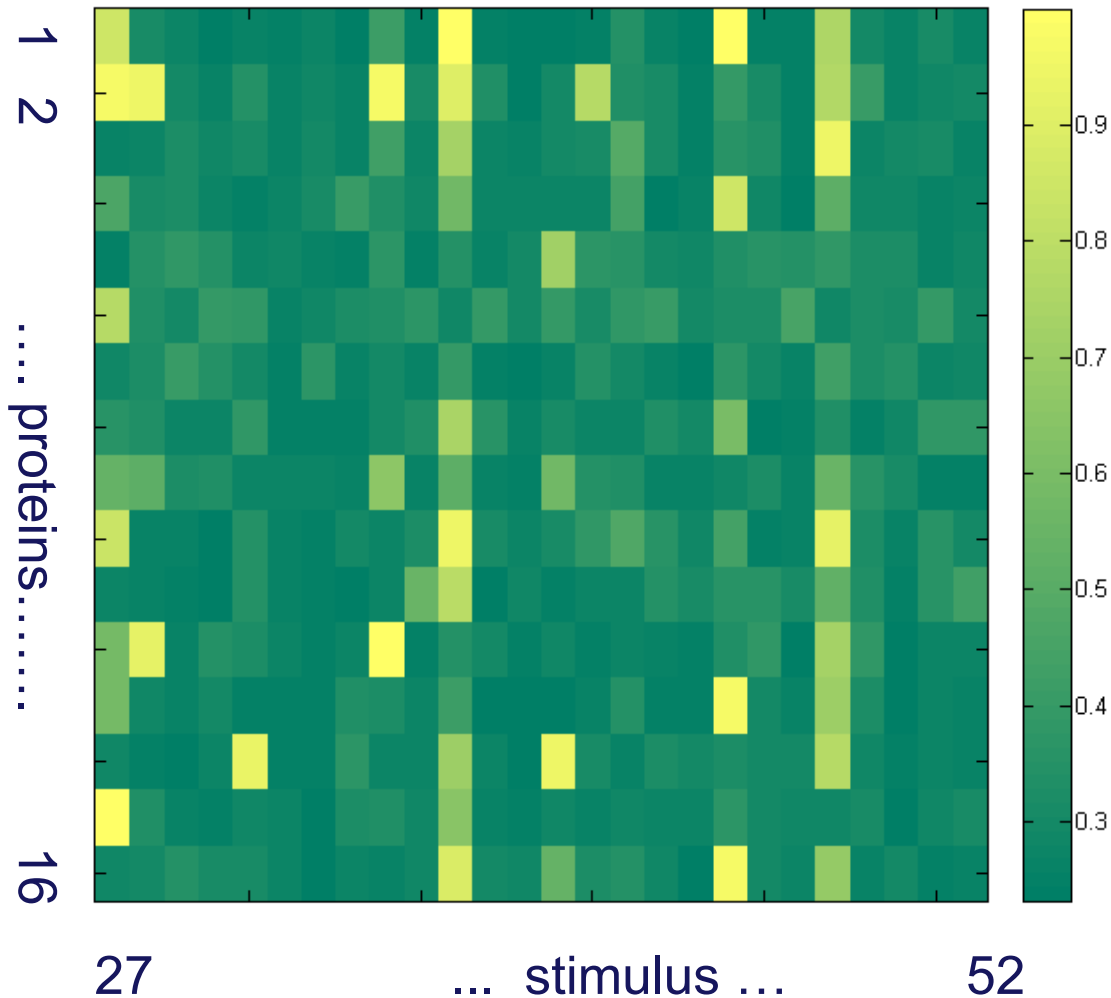
with respect to the full panel  
(416 predictions) of  
 $|humP| > 3$

AUC: probability for  $C_{naive}(+) > C_{naive}(-)$

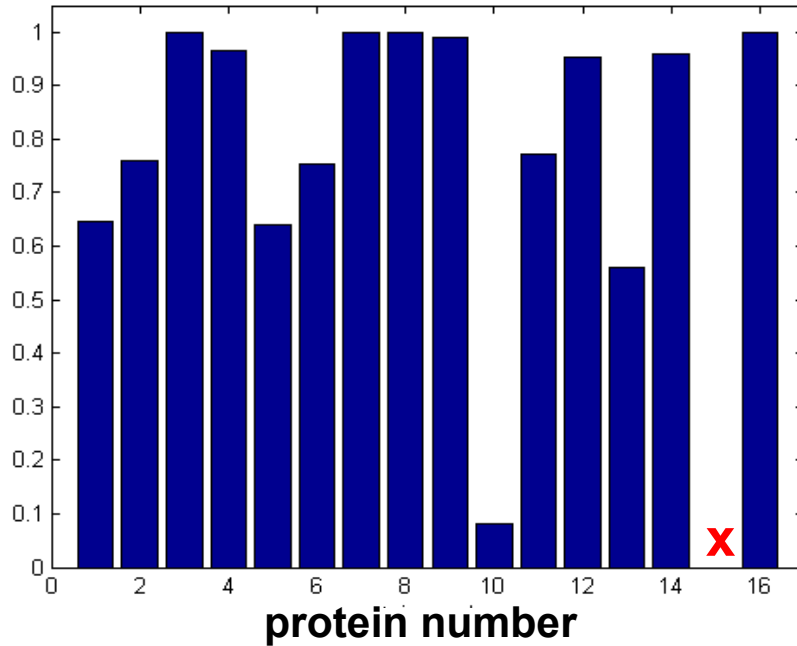
for a random pair of 1 positive and 1 negative sample



$C_{naive}$



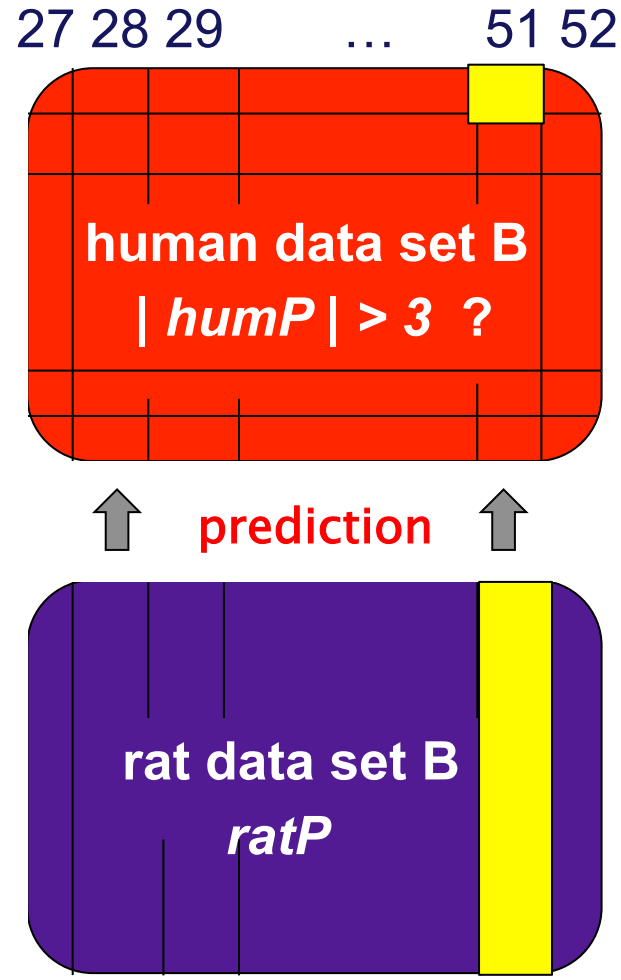
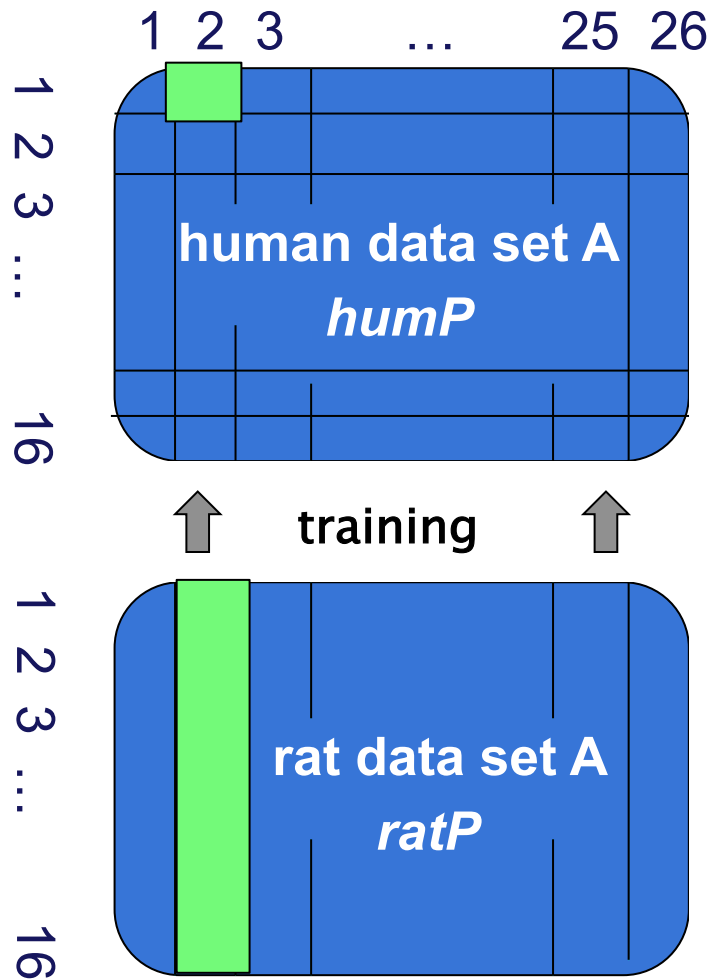
color-coded certainty  
for  $|humP| > 3$   
in data set B



protein specific AUC (ROC)  
for prediction of data set A

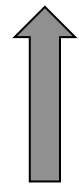
x (only negative samples)

## stimuli



$$y \in \{0, 1\}$$

16 separate  
binary  
classification  
problems



16-dim.  
vectors

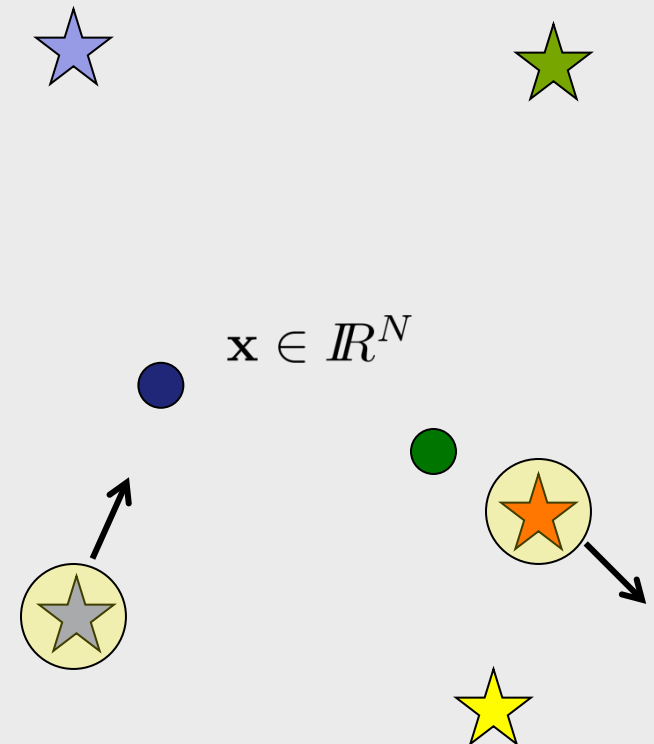
$$\mathbf{x} \in \mathbb{R}^{16}$$

N-dimensional data, feature vectors

- identification of **prototype vectors** from labeled example data
- distance based **classification** (e.g. Euclidean)

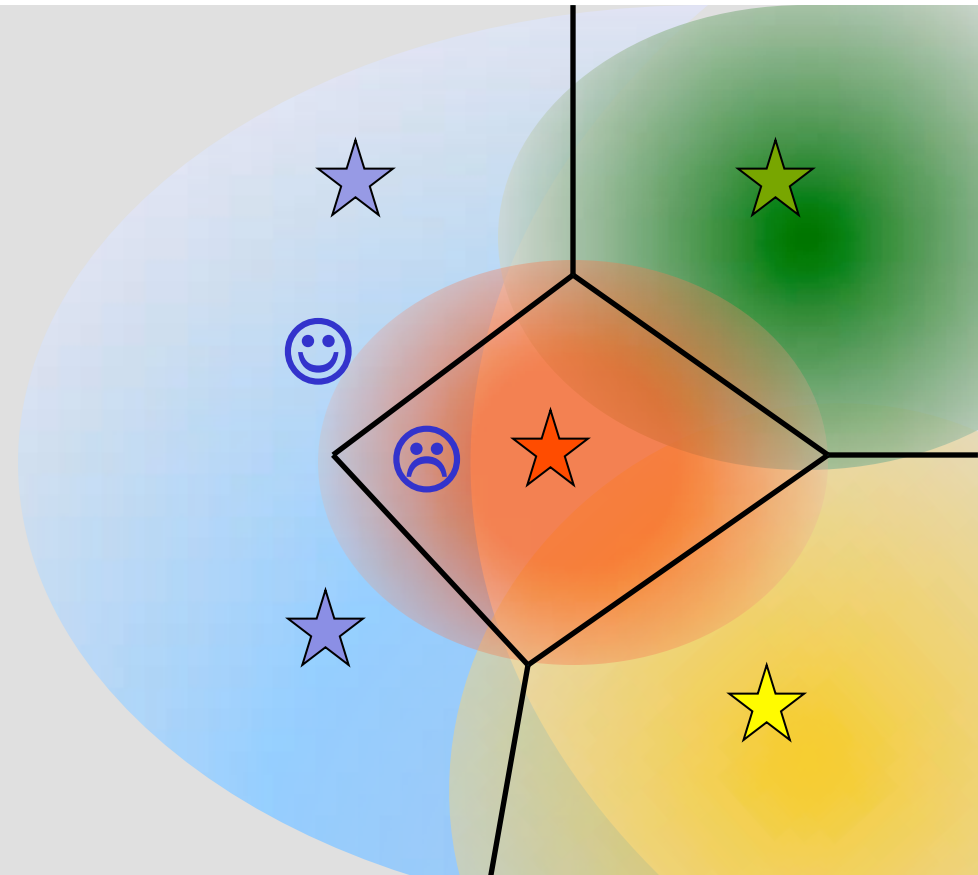
heuristic scheme: **LVQ1** [Kohonen, 1990]

- **initialize** prototype vectors for different classes
- present a **single example**
- identify the **winner** (closest prototype)
- move the winner
  - closer **towards the data** (same class)
  - **away from the data** (different class)
- present training data many times



- N-dimensional data, feature vectors
- identification of **prototype vectors** from labeled example data
- distance based **classification** (e.g. Euclidean)

- **distance-based classification**  
[here: Euclidean distances]
- **tessellation of feature space**  
[piece-wise linear]
- aim: **discrimination of classes**  
(  $\neq$  vector quantization  
or density estimation )
- **generalization ability**  
correct classification of *new* data



here: 16-dim. data  $\mathbf{x} \in \mathbb{R}^{16}$   $y(\mathbf{x}) \in \{0, 1\}$

one prototype per class  $\mathbf{w}^0, \mathbf{w}^1 \in \mathbb{R}^{16}$

**Nearest prototype classification:**

$$d^0 = \sum_i (w_i^0 - x_i)^2$$

$$d^1 = \sum_i (w_i^1 - x_i)^2$$

$$y(\mathbf{x}) = \begin{cases} 1 & \text{if } d^1(\mathbf{x}) < d^0(\mathbf{x}) \\ 0 & \text{else.} \end{cases}$$

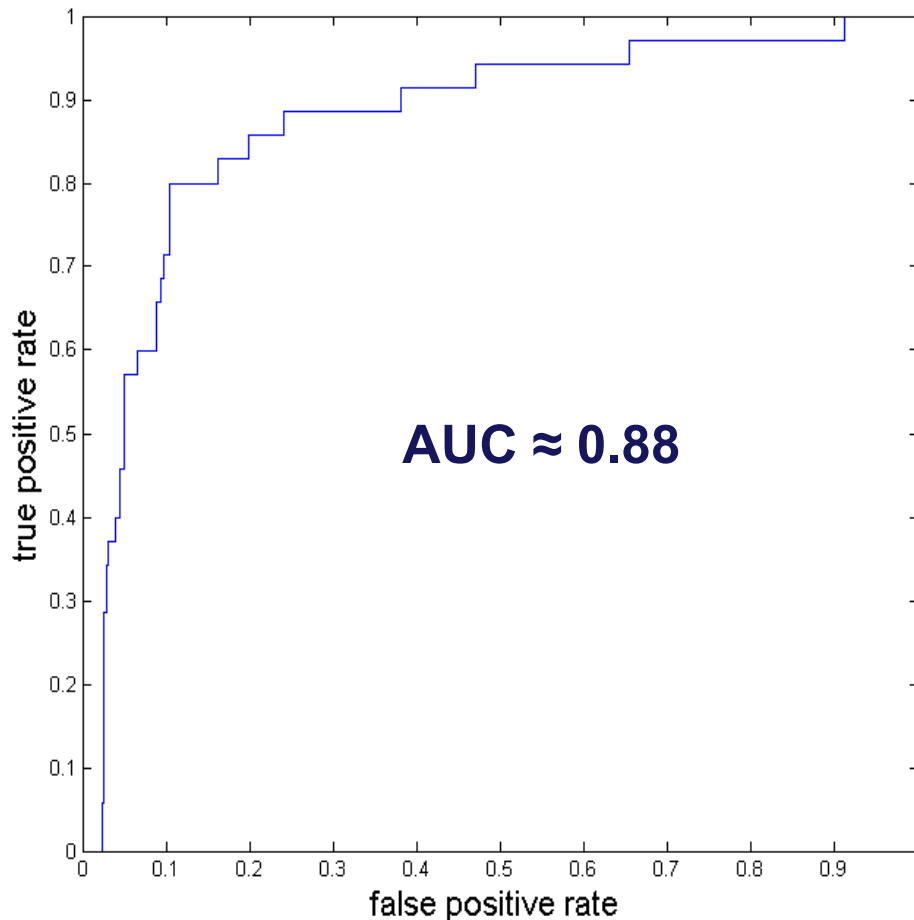
prediction score / certainty for activation

$$c_{LVQ} = \frac{1}{2} \left[ 1 + \tanh \left( \frac{d^0 - d^1}{200} \right) \right]$$

- precise (monotonic!) form is irrelevant for ROC, PR etc.
- crisp classification for threshold 0.5
- here: scaling factor yields range of values similar to naïve prediction

**validation:** 26 Leave-One-Out training processes:  
split data set A (complete) in 25 training / 1 test sample  
(if training set is all negative: accept naïve prediction)

**prediction:** *ensemble average* of certainties over the 26 LVQ systems



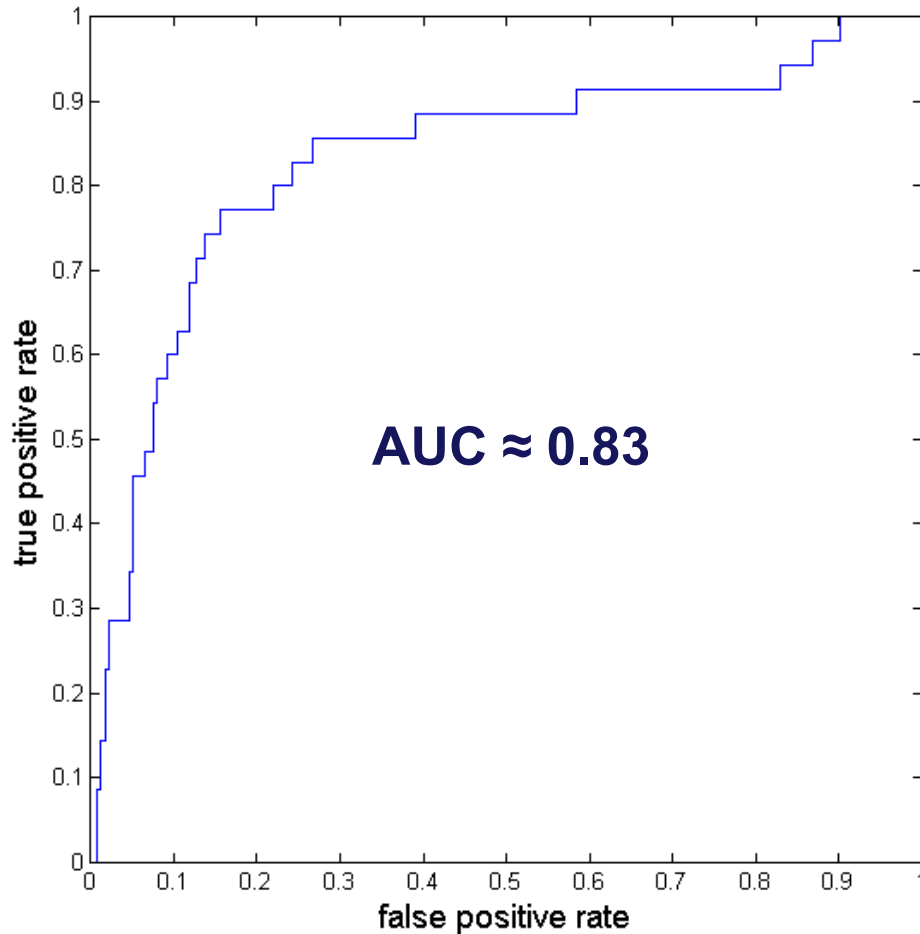
## ROC

with respect to the full panel  
(416 predictions) of

$$| \text{hum}P | > 3$$

obtained in the Leave-One-Out  
validation scheme





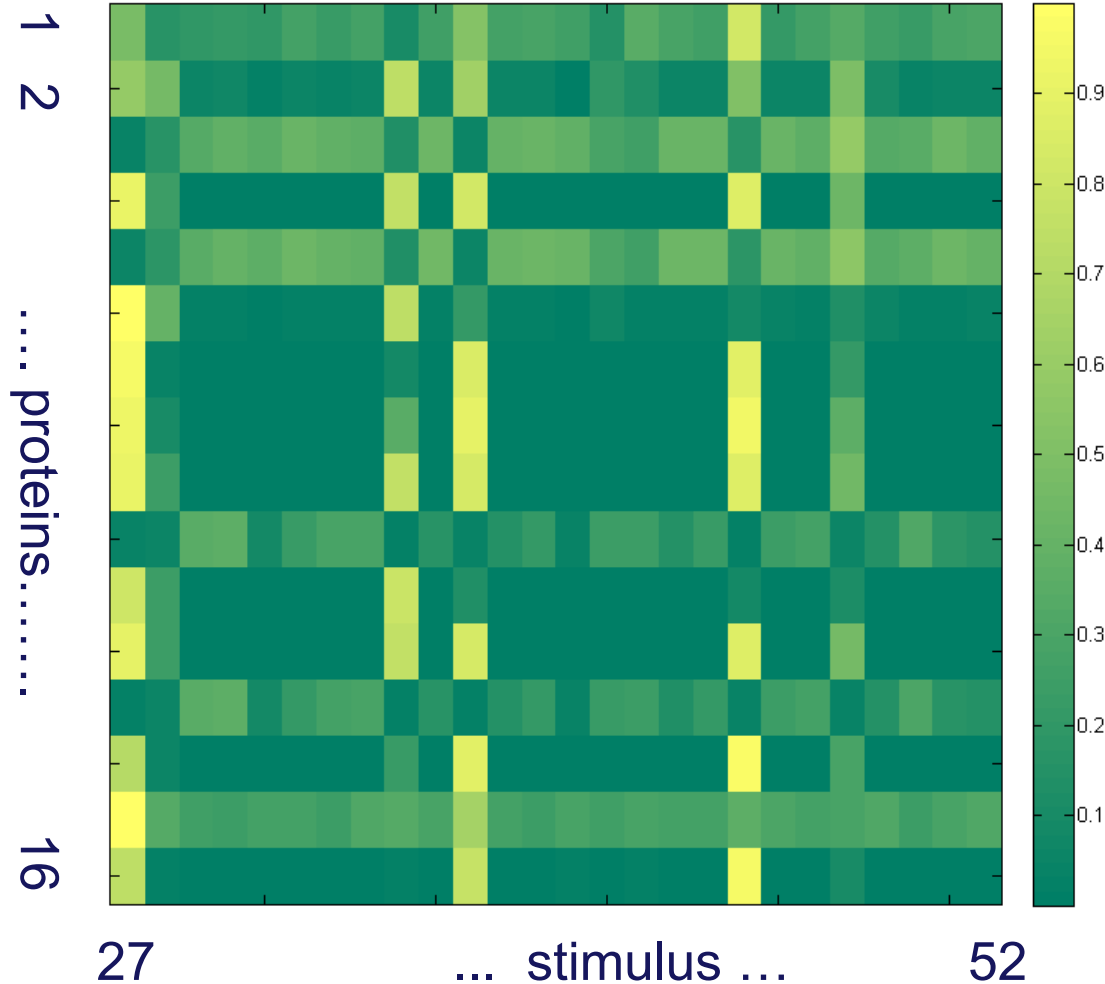
## ROC

with respect to the full panel  
(416 predictions) of

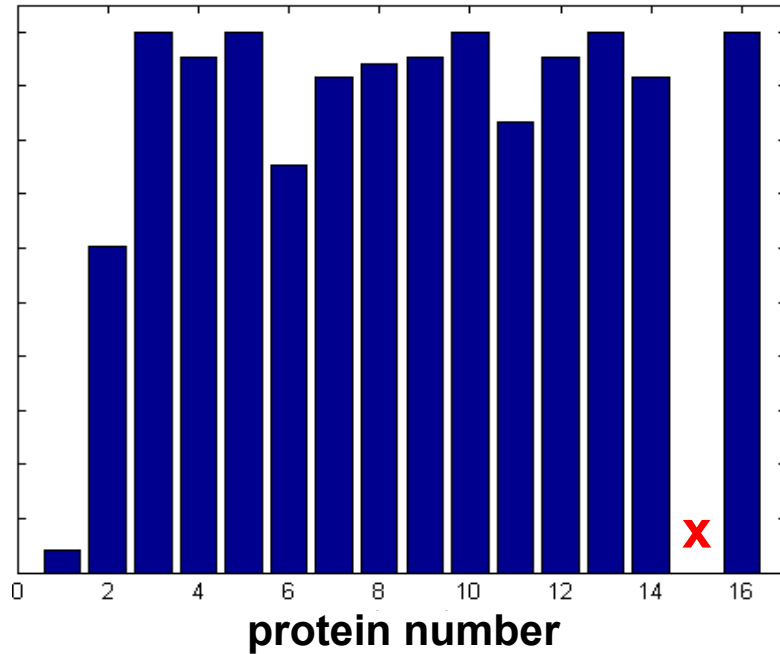
$$| \text{humP} | > 3$$

**AUC  $\approx$  0.83**

$C_{LVQ}$



color-coded certainty  
for  $|humP| > 3$   
in **data set B**

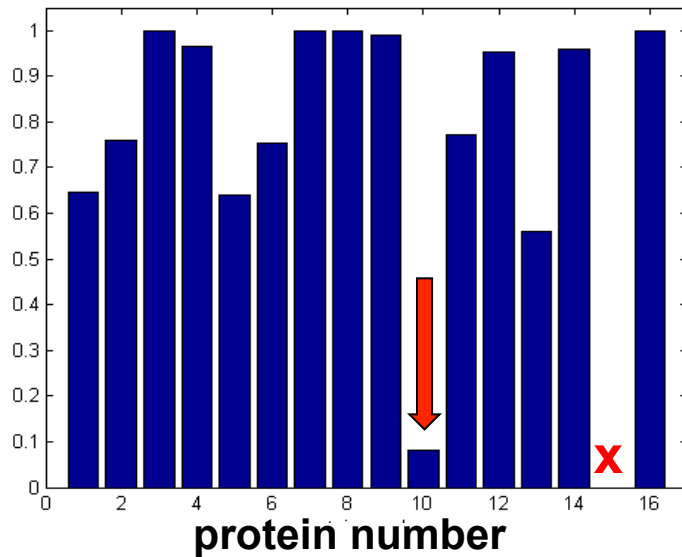


**protein specific AUC (ROC)**  
for prediction of **data set A**

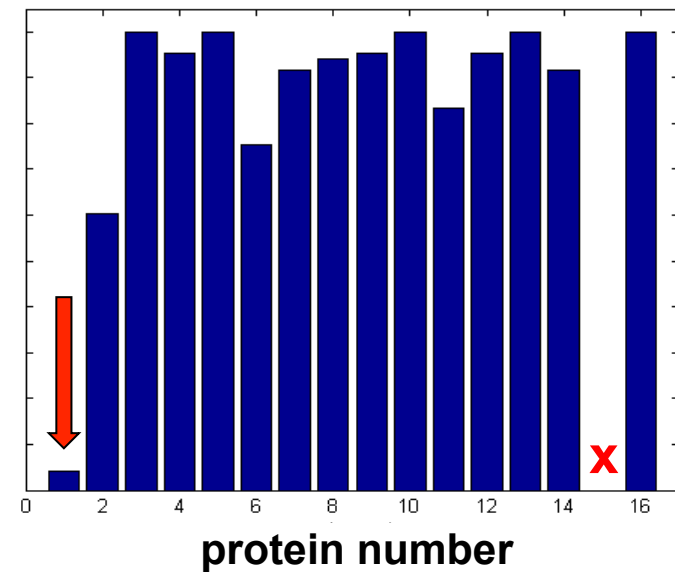
**x** (only negative samples)

## protein specific AUC (ROC) in data set A

$C_{naive}$



$C_{LVQ}$

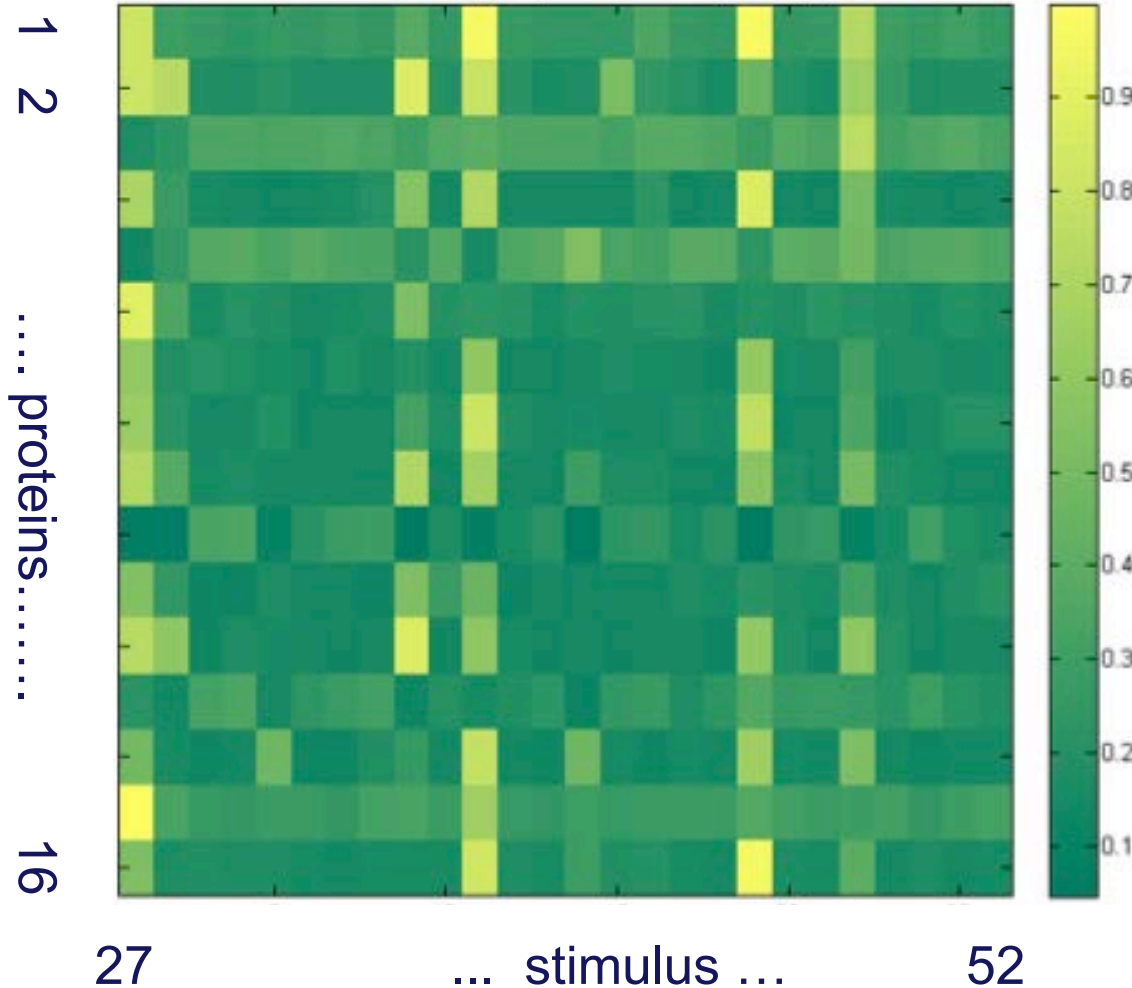


weighting  
scheme:

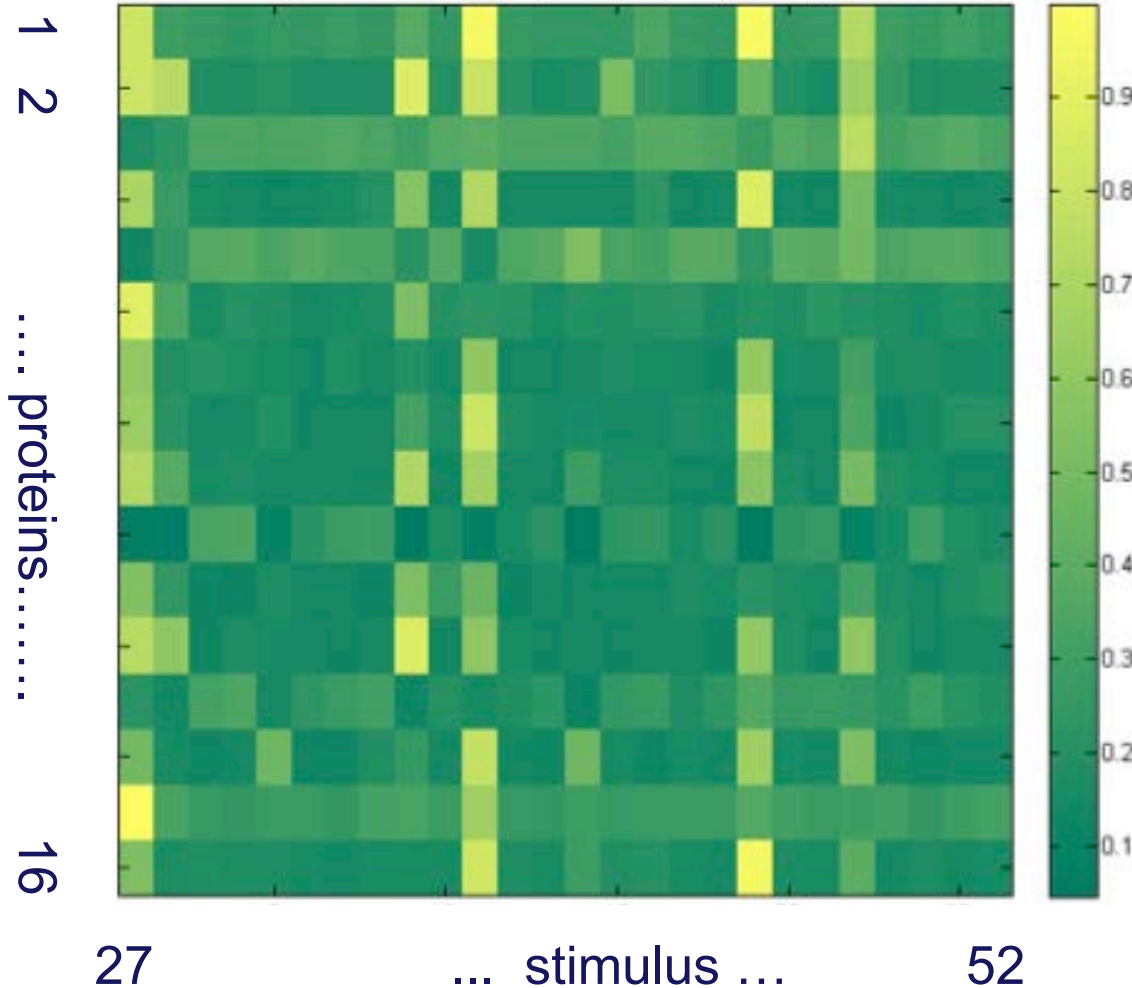
$$C_{final} = \frac{AUC_{naive} \cdot C_{naive} + AUC_{LVQ} \cdot C_{LVQ}}{AUC_{naive} + AUC_{LVQ}}$$

x (naïve prediction for protein 15)

$C_{final}$



color-coded certainty  
for  $|humP| > 3$   
in **data set B**

$C_{final}$ 

color-coded certainty  
for  $|humP| > 3$   
in **data set B**

**AUC (PR) = 0.54**

**Pearson = 0.75**

**BAC = 0.77**

## Analysis based on Phosphorylation data only

- 1) naïve prediction: “human = rat”
- 2) machine learning approach: LVQ classifier

## Leave-One-Out procedure

combination of classifiers based on  
protein-specific validation performance

## Result:

(slight) improvement over naïve prediction is possible

## Extend machine learning analysis

feature weighting schemes, e.g. Matrix Relevance LVQ (\*)  
should give further insight and help to understand,  
for instance, protein-specific difficulties

(\*) [www.cs.rug.nl/biehl](http://www.cs.rug.nl/biehl)

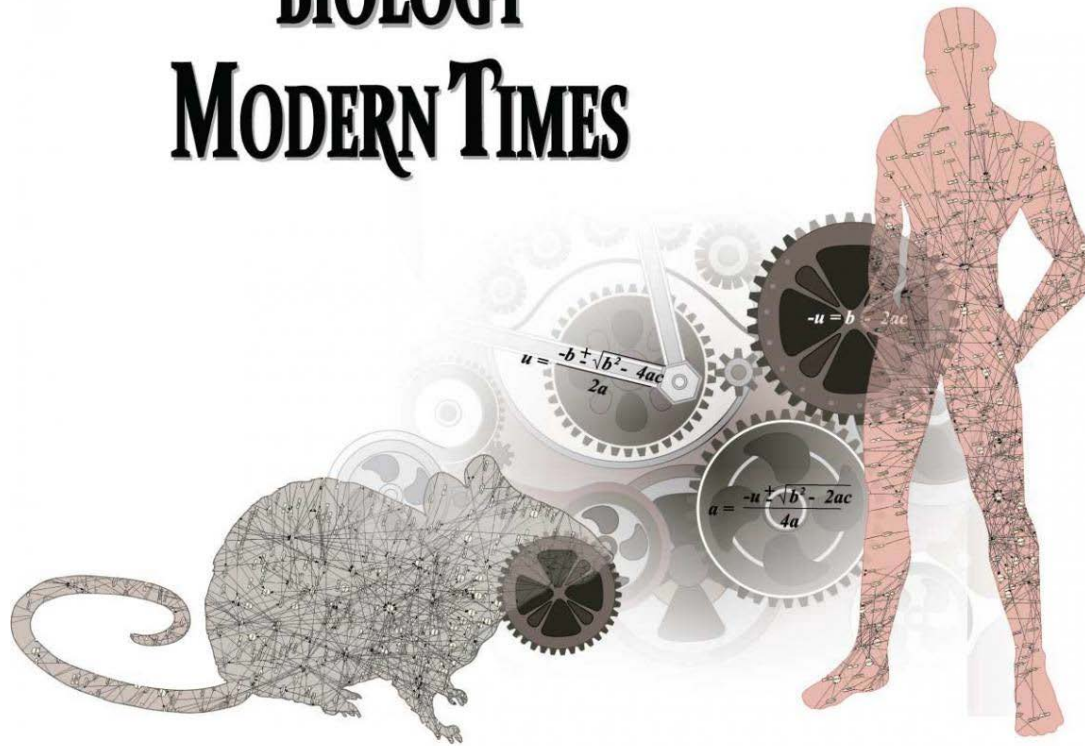
## Include gene expression data

try to understand the underlying mechanisms  
and differences / common properties of rat/human  
samples



## TRANSLATIONAL SYSTEMS BIOLOGY MODERN TIMES

[www.sbvimprover.com](http://www.sbvimprover.com)



**Rats and humans may be closer than we think!**



Thanks for  
 a great challenge!

Thanks for  
 real teamwork!

Thank you for your  
 attention!