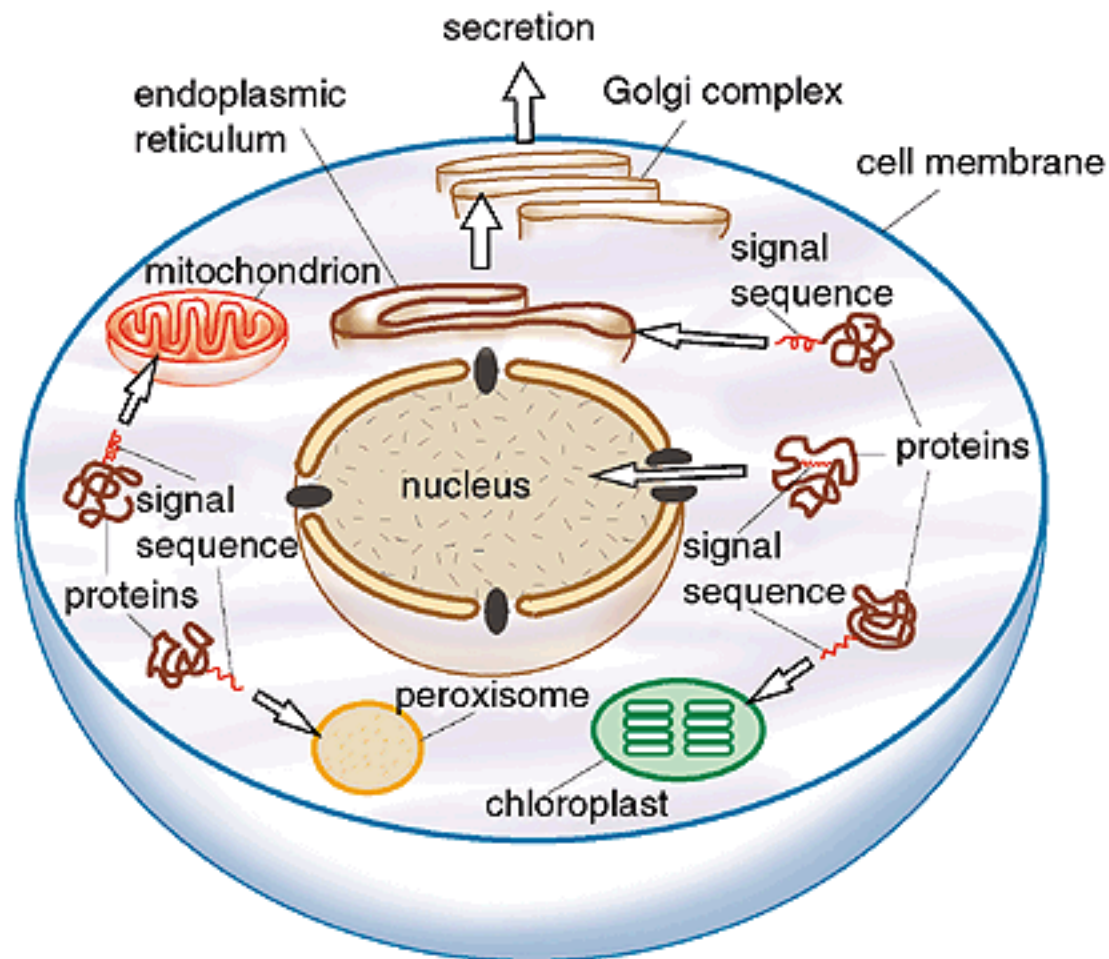# Membrane Protein Bioinformatics

## Gunnar von Heijne

Center for Biomembrane Research

Department of Biochemistry and Biophysics
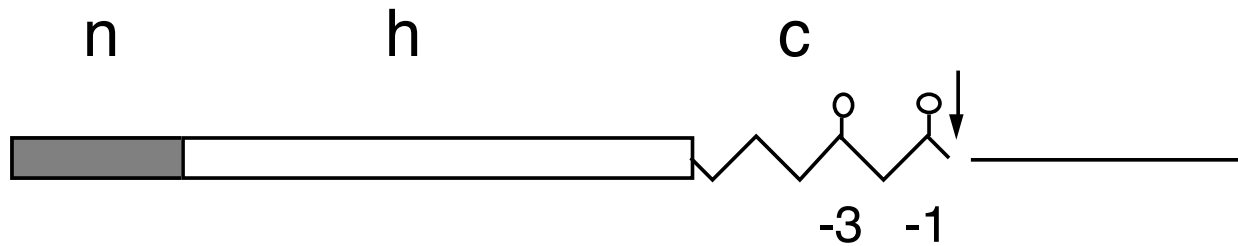
Stockholm University

Center for Biomembrane Research

# Protein sorting in a eukaryotic cell

# The signal peptide



n-region: positively charged

h-region: hydrophobic

c-region: more polar, small residues in -1, -3

# An early signal peptide predictor

**A new method for predicting signal sequence cleavage sites**

Gunnar von Heijne

Research Group for Theoretical Biophysics, Department of Theoretical Physics, Royal Institute of Technology, S-100 44 Stockholm, Sweden

ABSTRACT
A new method for identifying secretory signal sequences and for predicting the site of cleavage between a signal sequence and the mature exported protein is described. The predictive accuracy is estimated to be around 75-80% for both prokaryotic and eukaryotic proteins.

# An early signal peptide predictor

A new method for predicting signal sequence c[...]

Gunnar von Heijne

Research Group for Theoretical Biophysics, Depa[...]
Technology, S-100 44 Stockholm, Sweden

ABSTRACT
A new method for identifying secr[...]
the site of cleavage between a [...]
protein is described. The predict[...]
75-80% for both prokaryotic and eukar[...]

**Table 1   Amino acid counts for eukaryotic signal sequences**
The average composition (last column) is from Ref.(10)

| | -13 | -12 | -11 | -10 | -9 | -8 | -7 | -6 | -5 | -4 | -3 | -2 | -1 | +1 | +2 | Expected |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 16 | 13 | 14 | 15 | 20 | 18 | 18 | 17 | 25 | 15 | 47 | 6 | 80 | 18 | 6 | 14.5 |
| C | 3 | 6 | 9 | 7 | 9 | 14 | 6 | 8 | 5 | 6 | 19 | 3 | 9 | 8 | 3 | 4.5 |
| D | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 3 | 0 | 5 | 0 | 10 | 11 | 8.9 |
| E | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 3 | 7 | 0 | 7 | 0 | 13 | 14 | 10.0 |
| F | 13 | 9 | 11 | 11 | 6 | 7 | 18 | 13 | 4 | 5 | 0 | 13 | 0 | 6 | 4 | 5.6 |
| G | 4 | 4 | 3 | 6 | 3 | 13 | 3 | 2 | 19 | 34 | 5 | 7 | 39 | 10 | 7 | 12.1 |
| H | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 5 | 0 | 0 | 6 | 0 | 4 | 2 | 3.4 |
| I | 15 | 15 | 8 | 6 | 11 | 5 | 4 | 8 | 5 | 1 | 10 | 5 | 0 | 8 | 7 | 7.4 |
| K | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 4 | 0 | 2 | 0 | 11 | .9 | 11.3 |
| L | 71 | 68 | 72 | 79 | 78 | 45 | 64 | 49 | 10 | 23 | 8 | 20 | 1 | 8 | 4 | 12.1 |
| M | 0 | 3 | 7 | 4 | 1 | 6 | 2 | 2 | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 2.7 |
| N | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 3 | 3 | 0 | 10 | 0 | 4 | 7 | 7.1 |
| P | 2 | 0 | 2 | 0 | 0 | 4 | 1 | 8 | 20 | 14 | 0 | 1 | 3 | 0 | 22 | 7.4 |
| Q | 0 | 0 | 0 | 1 | 0 | 6 | 1 | 0 | 10 | 8 | 0 | 18 | 3 | 19 | 10 | 6.3 |
| R | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 7 | 4 | 0 | 15 | 0 | 12 | 9 | 7.6 |
| S | 9 | 3 | 8 | 6 | 13 | 10 | 15 | 16 | 26 | 11 | 23 | 17 | 20 | 15 | 10 | 11.4 |
| T | 2 | 10 | 5 | 4 | 5 | 13 | 7 | 7 | 12 | 6 | 17 | 8 | 6 | 3 | 10 | 9.7 |
| V | 20 | 25 | 15 | 18 | 13 | 15 | 11 | 27 | 0 | 12 | 32 | 3 | 0 | 8 | 17 | 11.1 |
| W | 4 | 3 | 3 | 1 | 1 | 2 | 6 | 3 | 1 | 3 | 0 | 9 | 0 | 2 | 0 | 1.8 |
| Y | 0 | 1 | 4 | 0 | 0 | 1 | 3 | 1 | 1 | 2 | 0 | 5 | 0 | 1 | 7 | 5.6 |

# An early signal peptide predictor

A new m

Gunnar v

Research
Technolog

Received

ABSTRAC
A new
the si
protein
75-80%

| | Table 1 | | Amino acid counts for eukaryotic signal sequences | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | The average composition (last column) is from Ref.(10) | | | | | | | | | | | | |
| | -13 | -12 | -11 | -10 | -9 | -8 | -7 | -6 | -5 | -4 | -3 | -2 | -1 | +1 | +2 | Expected |
| A | 16 | 13 | 14 | 15 | 20 | 18 | 18 | 17 | 25 | 15 | 47 | 6 | 80 | 18 | 6 | 14.5 |
| C | 3 | 6 | 9 | 7 | 9 | 14 | 6 | 8 | 5 | 6 | 19 | 3 | 9 | 8 | 3 | 4.5 |
| D | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 3 | 0 | 5 | 0 | 10 | 11 | 8.9 |
| E | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 3 | 7 | 0 | 7 | 0 | 13 | 14 | 10.0 |
| F | 13 | 9 | 11 | 11 | 6 | 7 | 18 | 13 | 4 | 5 | 0 | 13 | 0 | 6 | 4 | 5.6 |
| G | 4 | 4 | 3 | 6 | 3 | 13 | 3 | 2 | 19 | 34 | 5 | 7 | 39 | 10 | 7 | 12.1 |
| H | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 5 | 0 | 0 | 6 | 0 | 4 | 2 | 3.4 |
| I | 15 | 15 | 8 | 6 | 11 | 5 | 4 | 8 | 5 | 1 | 10 | 5 | 0 | 8 | 7 | 7.4 |
| K | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 4 | 0 | 2 | 0 | 11 | 9 | 11.3 |
| L | 71 | 68 | 72 | 79 | 78 | 45 | 64 | 49 | 10 | 23 | 8 | 20 | 1 | 8 | 4 | 12.1 |
| M | 0 | 3 | 7 | 4 | 1 | 6 | 2 | 2 | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 2.7 |
| N | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 3 | 3 | 0 | 10 | 0 | 4 | 7 | 7.1 |
| P | 2 | 0 | 2 | 0 | 0 | 4 | 1 | 8 | 20 | 14 | 0 | 1 | 3 | 0 | 22 | 7.4 |
| Q | 0 | 0 | 0 | 1 | 0 | 6 | 1 | 0 | 10 | 8 | 0 | 18 | 3 | 19 | 10 | 6.3 |
| R | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 7 | 4 | 0 | 15 | 0 | 12 | 9 | 7.6 |
| S | 9 | 3 | 8 | 6 | 13 | 10 | 15 | 16 | 26 | 11 | 23 | 17 | 20 | 15 | 10 | 11.4 |
| T | 2 | 10 | 5 | 4 | 5 | 13 | 7 | 7 | 12 | 6 | 17 | 8 | 6 | 3 | 10 | 9.7 |
| V | 20 | 25 | 15 | 18 | 13 | 15 | 11 | 27 | 0 | 12 | 32 | 3 | 0 | 8 | 17 | 11.1 |
| W | 4 | 3 | 3 | 1 | 1 | 2 | 6 | 3 | 1 | 3 | 0 | 9 | 0 | 2 | 0 | 1.8 |
| Y | 0 | 1 | 4 | 0 | 0 | 1 | 3 | 1 | 1 | 2 | 0 | 5 | 0 | 1 | 7 | 5.6 |

Convert to a weight-matrix:

$$W(a,i) = \ln(N(a,i)/\langle N(a) \rangle)$$

Scan W along the sequence. For each position of W, sum the weights $W(a,i)$ corresponding to the sequence. The position with the maximum score is the predicted cleavage site.

# An early signal peptide predictor

A new m...

Gunnar v...

Research

Technolog...

Received...

**ABSTRACT**
A new ...
the si...
protein...
75-80%...

Table 1  Amino acid counts for eukaryotic sig...
The average composition (last column)...

| | -13 | -12 | -11 | -10 | -9 | -8 | -7 | -6 | -5 | -4 | -3 | -2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 16 | 13 | 14 | 15 | 20 | 18 | 18 | 17 | 25 | 15 | 47 | 6 |
| C | 3 | 6 | 9 | 7 | 9 | 14 | 6 | 8 | 5 | 6 | 19 | 3 |
| D | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 3 | 0 | 5 |
| E | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 3 | 7 | 0 | 7 |
| F | 13 | 9 | 11 | 11 | 6 | 7 | 18 | 13 | 4 | 5 | 0 | 13 |
| G | 4 | 4 | 3 | 6 | 3 | 13 | 3 | 2 | 19 | 34 | 5 | 7 |
| H | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 5 | 0 | 0 | 6 |
| I | 15 | 15 | 8 | 6 | 11 | 5 | 4 | 8 | 5 | 1 | 10 | 5 |
| K | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 4 | 0 | 2 |
| L | 71 | 68 | 72 | 79 | 78 | 45 | 64 | 49 | 10 | 23 | 8 | 20 |
| M | 0 | 3 | 7 | 4 | 1 | 6 | 2 | 2 | 0 | 0 | 0 | 1 |
| N | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 3 | 3 | 0 | 10 |
| P | 2 | 0 | 2 | 0 | 0 | 4 | 1 | 8 | 20 | 14 | 0 | 1 |
| Q | 0 | 0 | 0 | 1 | 0 | 6 | 1 | 0 | 10 | 8 | 0 | 18 |
| R | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 7 | 4 | 0 | 15 |
| S | 9 | 3 | 8 | 6 | 13 | 10 | 15 | 16 | 26 | 11 | 23 | 17 |
| T | 2 | 10 | 5 | 4 | 5 | 13 | 7 | 7 | 12 | 6 | 17 | 8 |
| V | 20 | 25 | 15 | 18 | 13 | 15 | 11 | 27 | 0 | 12 | 32 | 3 |
| W | 4 | 3 | 3 | 1 | 1 | 2 | 6 | 3 | 1 | 3 | 0 | 9 |
| Y | 0 | 1 | 4 | 0 | 0 | 1 | 3 | 1 | 1 | 2 | 0 | 5 |



**Distribution of maximum scores for signal sequences and cytosolic proteins.** Open squares: cytosolic proteins; solid squares: signal sequences.

# A modern predictor: SignalP



www.cbs.dtu.dk

# A modern predictor: SignalP



SignalP-4.0 prediction (euk networks): ERP44_HUMAN
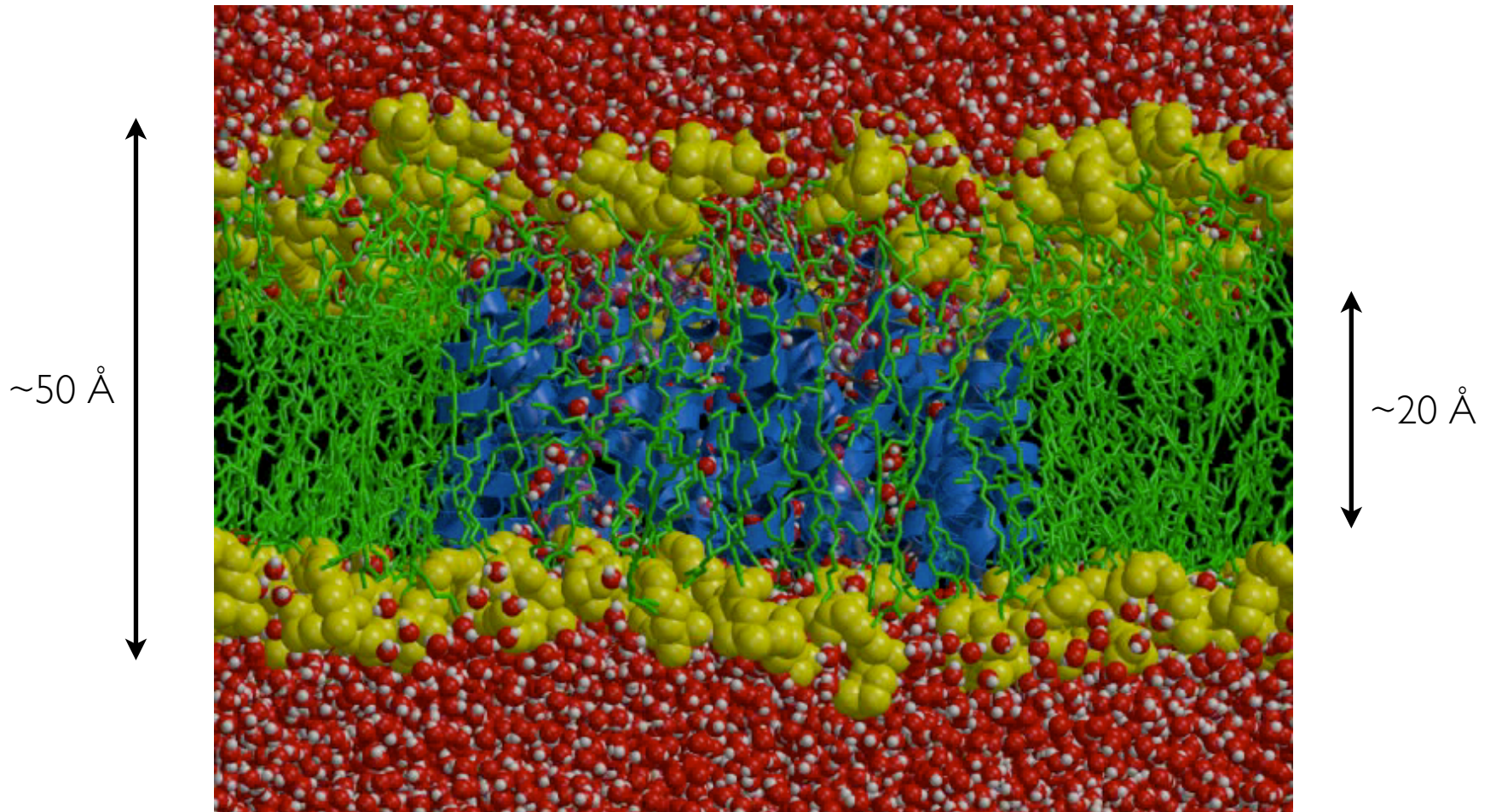
```
# Measure    Position    Value      Cutoff     signal peptide?
  max. C        30        0.427
  max. Y        30        0.586
  max. S         9        0.950
  mean S       1-29       0.821
       D       1-29       0.713     0.450      YES
Name=sp_Q9BS26_ERP44_HUMAN            SP='YES' Cleavage site between pos. 29 and 30: VTT-EI
```
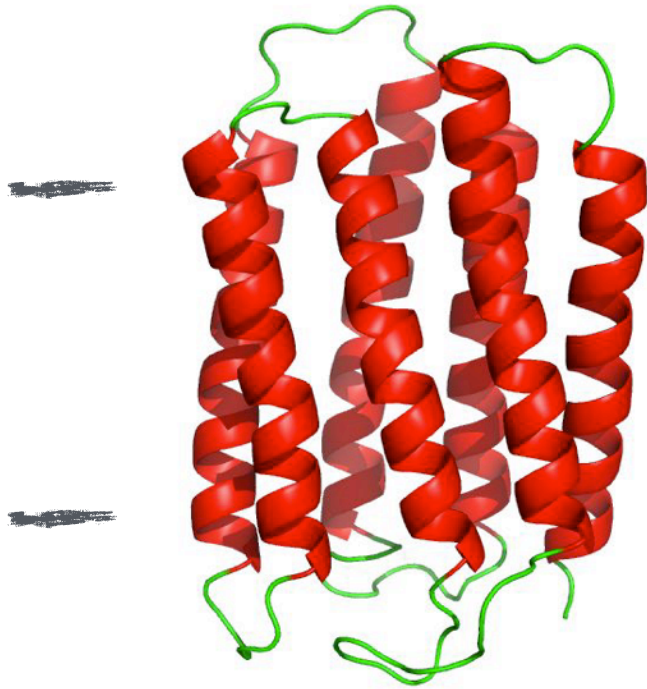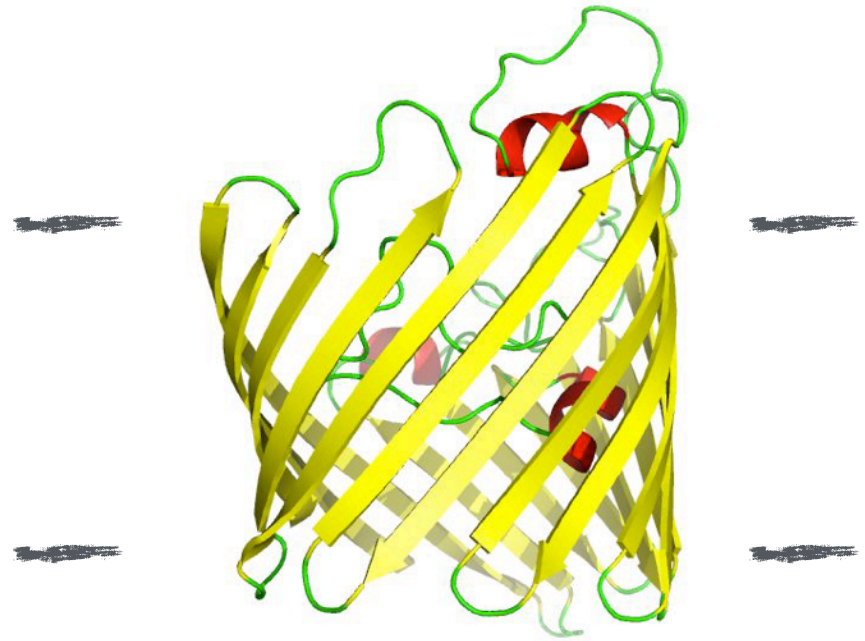
# A lipid membrane



~50 Å

~20 Å

# Two architectures
## (Quart.Rev.Biophys. 32:285)



Helix bundle

ß-barrel

# A helix-bundle membrane protein

# Three important characteristics of the helix-bundle membrane proteins



periplasm

Trp & Tyr

~20 hydrophobic residues

Lys & Arg

cytoplasm

# TM helices are typically 20-30 residues long



Helix Length (Number of Residues)

# Loops connecting the TM helices tend to be short

# The positive-inside rule applies to all (?) organisms



# of genomes with a significant bias in the distribution of amino acid X in intra- vs. extra-cellular loops

# Topology prediction

MDSQRNLLVIALLFVSFMIWQAWE....



out

in

# Popular topology predictors

TMHMM (HMM)
HMMTOP (HMM)
Prodiv-TMHMM (MSA, HMM)
Phobius (HMM)
MEMSAT (MSA, dynamic programming)
TOPCONS (consensus method)
....
SCAMPI (h-plot, PI-rule)
PHD (MSA, NN, PI-rule)
TopPred (h-plot, PI-rule)
....
Kyte & Doolittle (h-plot)
SOAP (h-plot)

# TopPred

Step 1: Make a hydrophobicity plot



*E. coli* LacY

# TopPred



Step 2:
- construct all possible topologies

- rank based on $\Delta+$

# TMHMM



## A hidden Markov model (HMM)

# Helix and loop models in TMHMM
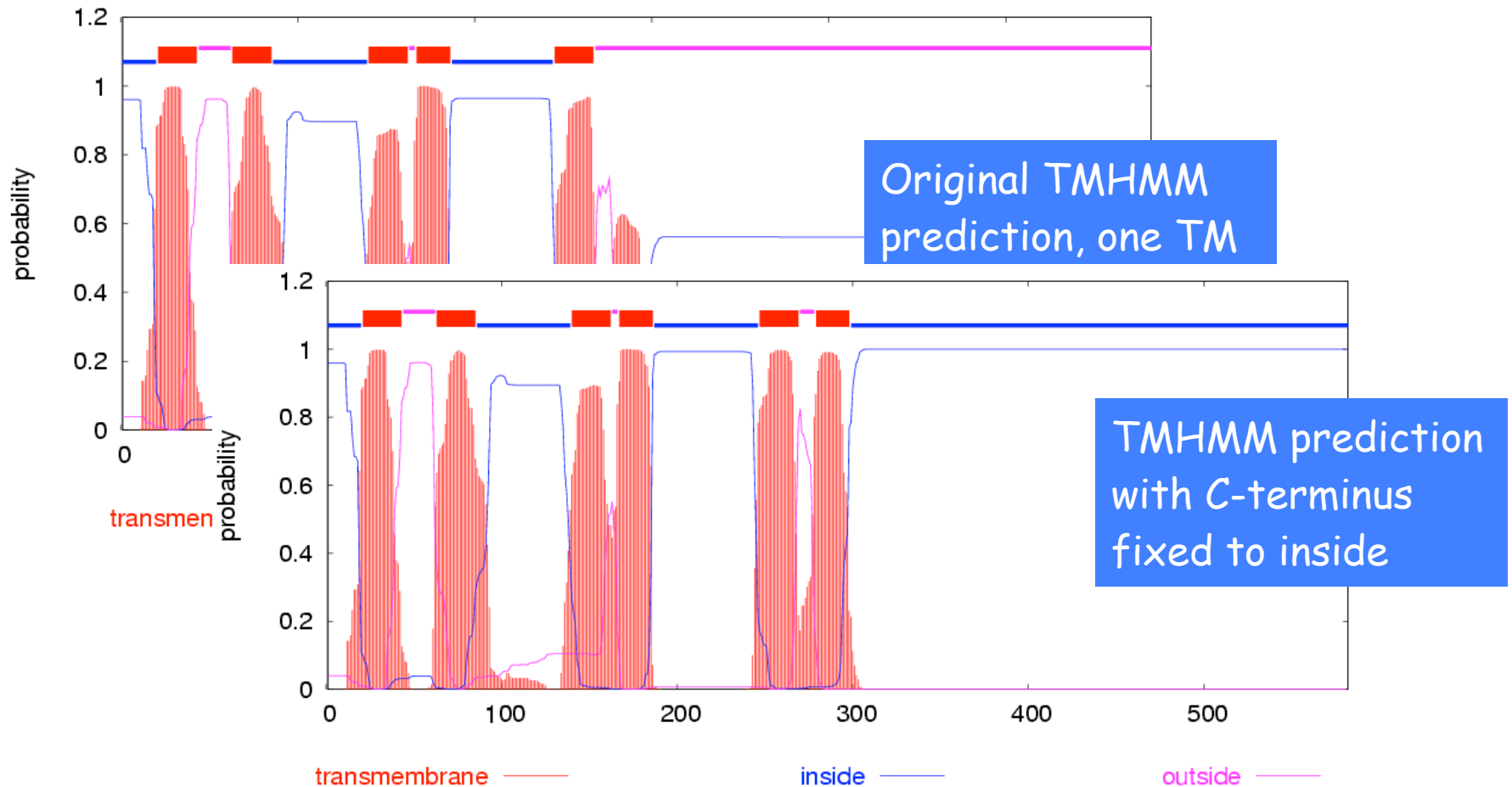
# TOPCONS
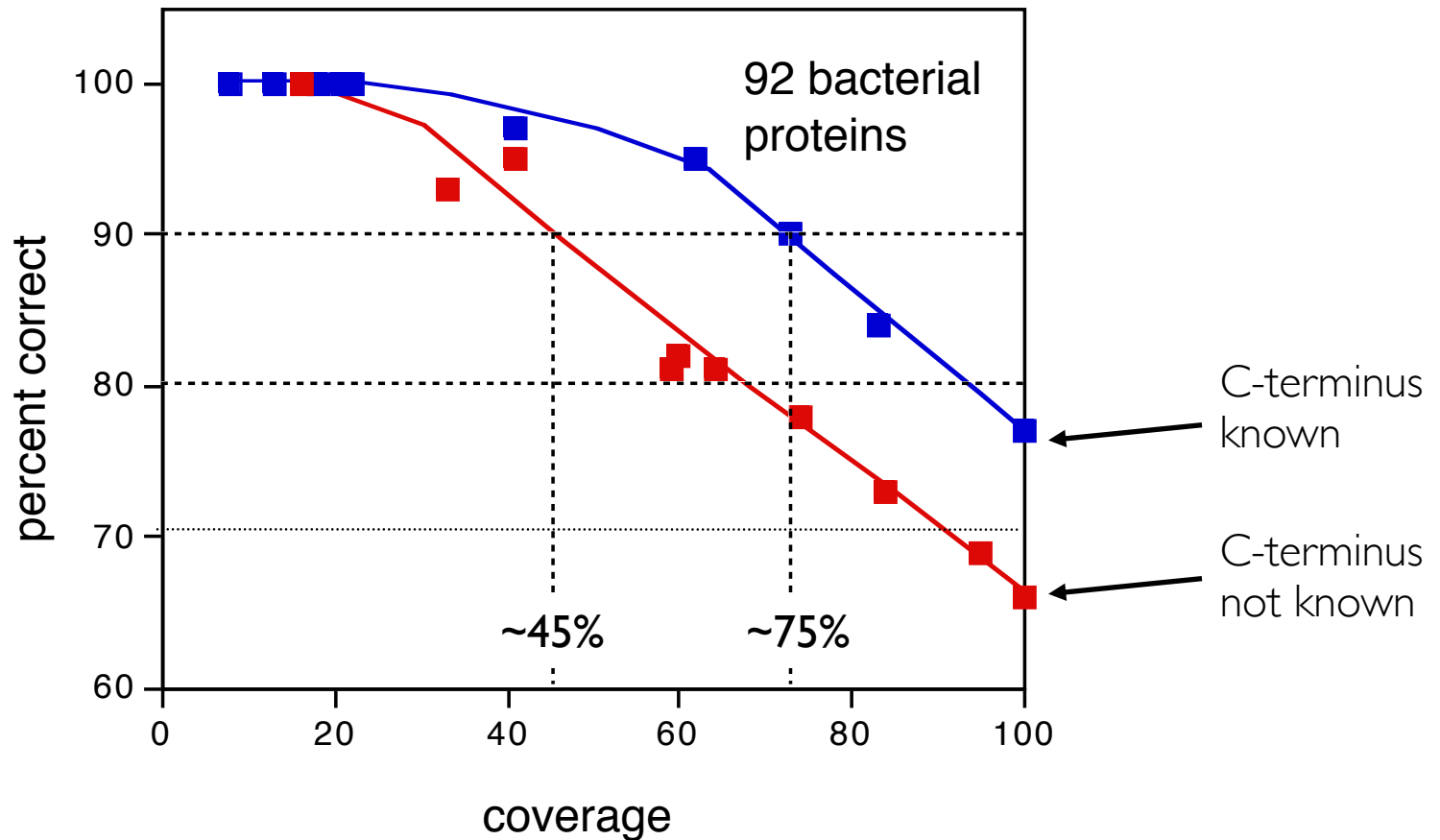
# TOPCONS

# How good are topology predictors?

Discrimination membrane/soluble: sensitivity 97%; specificity 95%

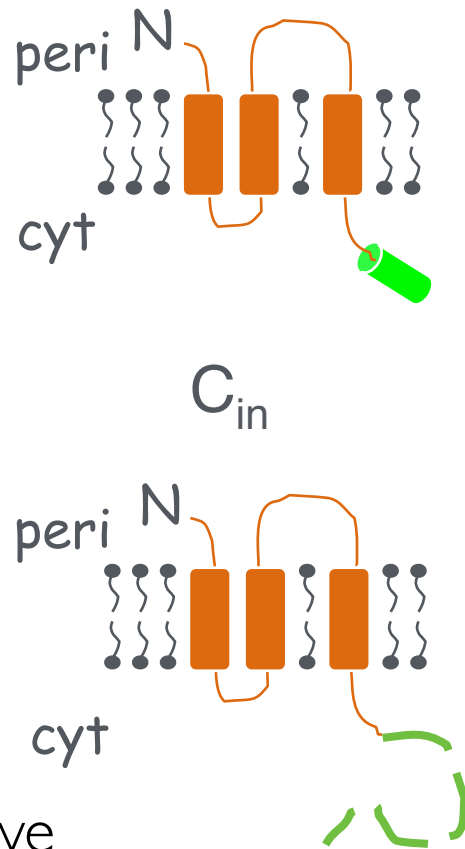Topology: single sequence 70-75%; multi-sequence 80-85%
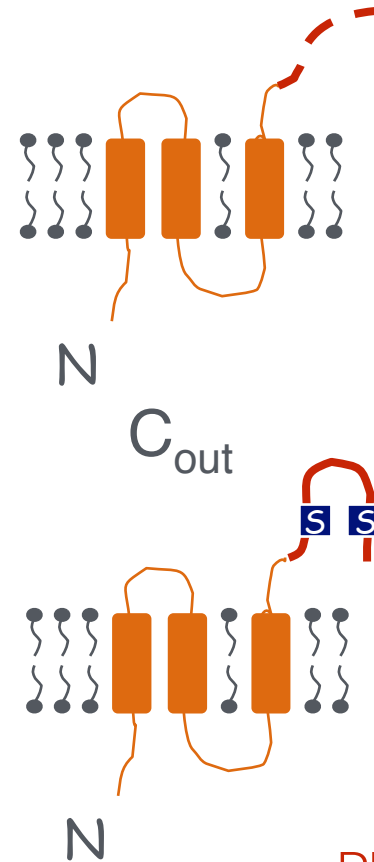
# Experimental constraints help



Original TMHMM prediction, one TM

TMHMM prediction with C-terminus fixed to inside

transmembrane ——    inside ——    outside ——

# Experimental constraints help

# Experimental constraints help



peri N

cyt

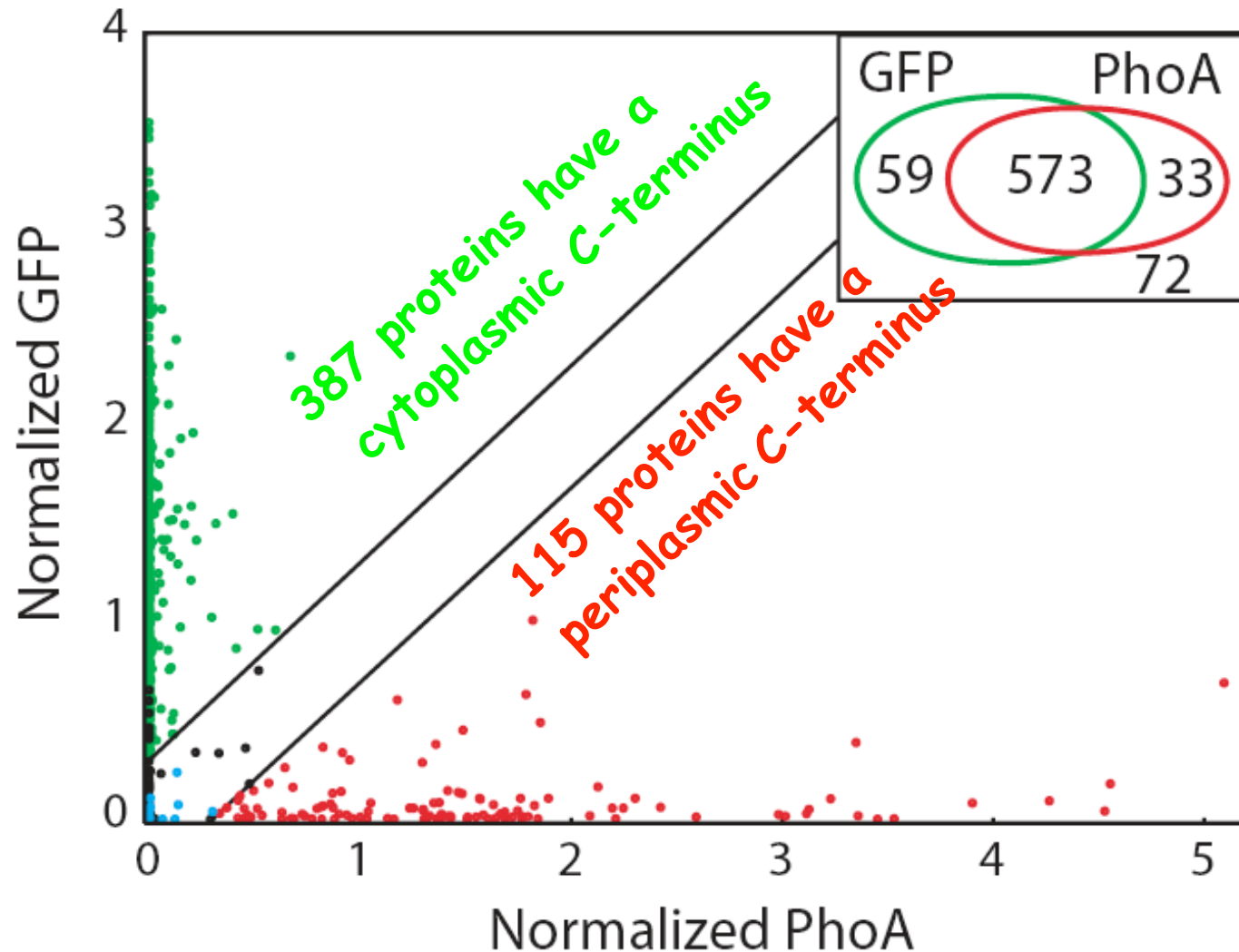$C_{in}$

peri N

cyt

N

$C_{out}$
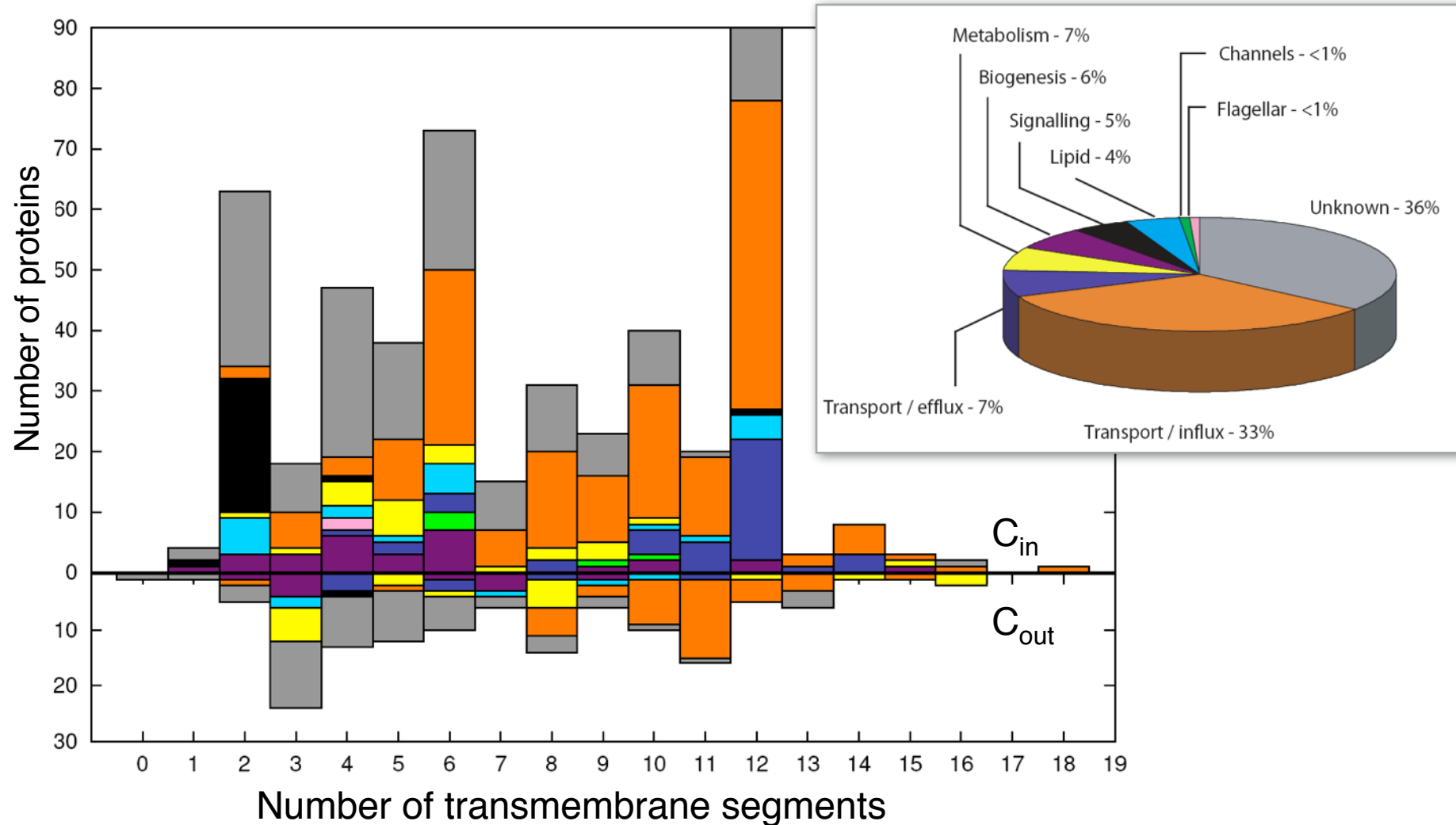
S S

N

GFP is only active
in the cytoplasm

PhoA is only active
in the periplasm

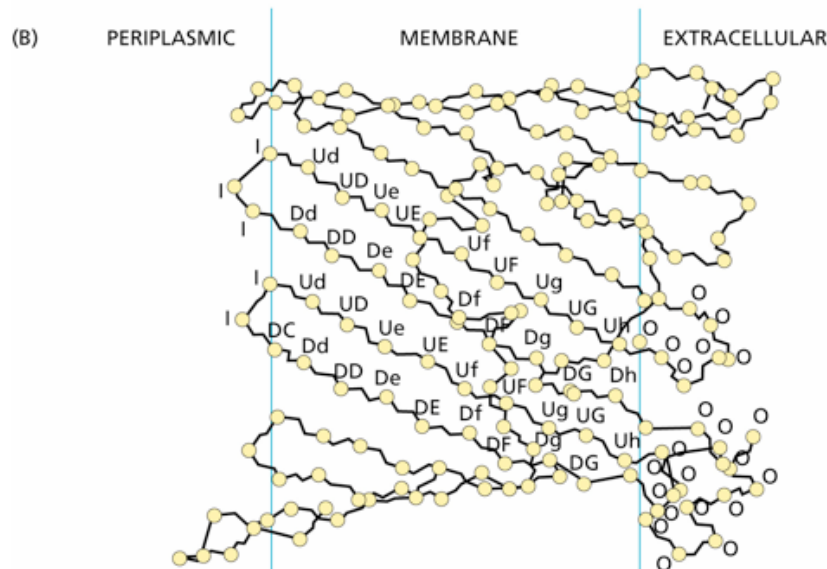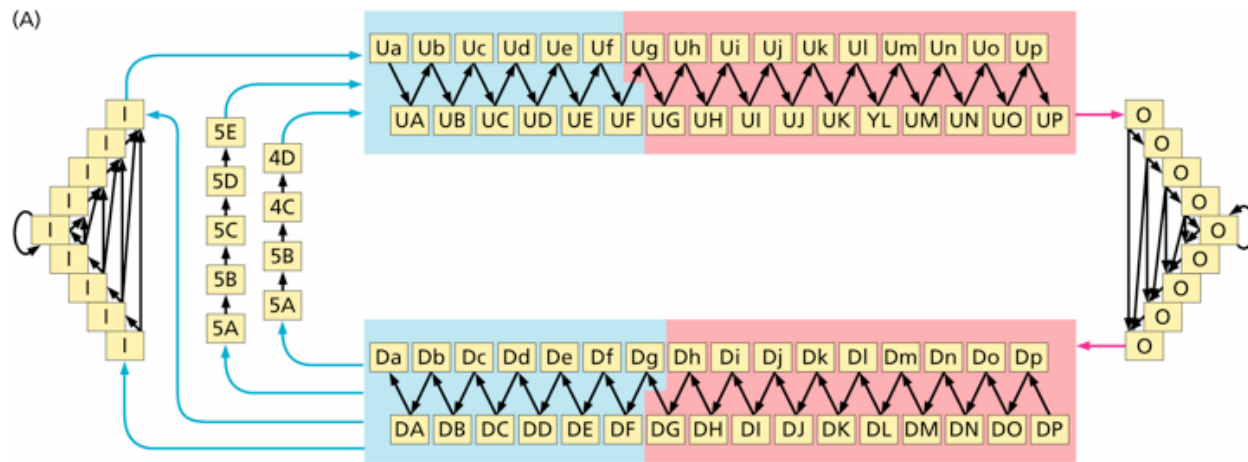# Experimental constraints help

# Experimental constraints help

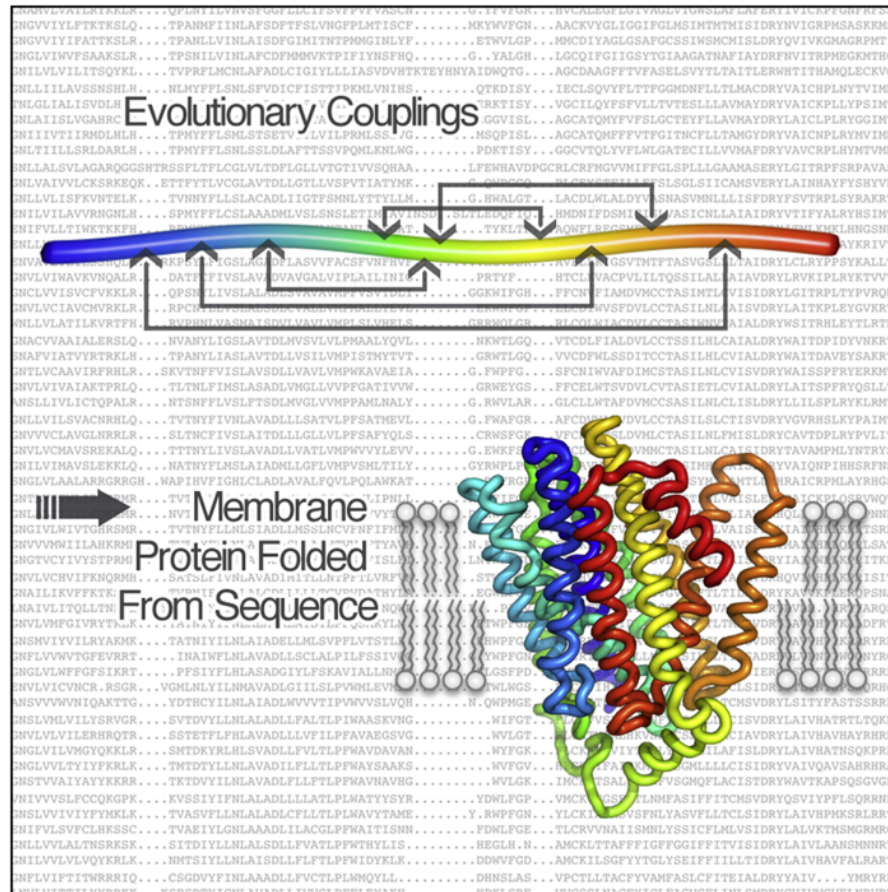# 'Exporting' experimental constraints using alignments

BLAST alignment

E. coli query

hit

C

C

225/38 bacterial/eukaryotic genomes

158 k/139 k predicted membrane proteins

51.000/15.000 assignments

# A HMM for β-barrel membrane proteins



(A)

(B) PERIPLASMIC    MEMBRANE    EXTRACELLULAR

# 3D structure prediction by co-evolving residues
## (Hopf et al., Cell 149: 1607)



Only works if there are hundreds of sequences in the multiple alignment

# Don't forget...

Signal peptide prediction: SignalP

Two architectures: helix bundle and β-barrel
Hydrophobic transmembrane helices
Alternating (Hyf-X) β-strands
The positive-inside rule
TopPred
TMHMM
TOPCONS
Experimental constraints help
3D structure prediction possible for large
protein families