



# Recent Advances in Protein Sequence Analysis

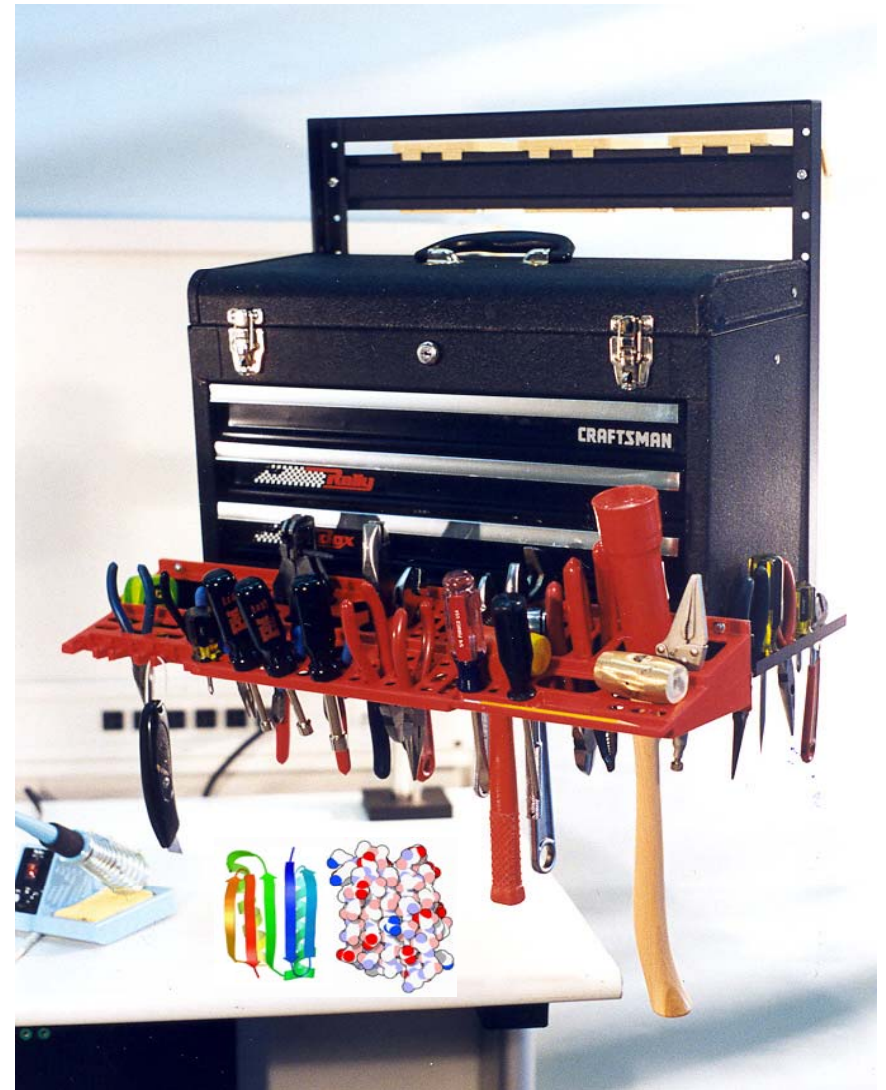
Nick V. Grishin

*Howard Hughes Medical Institute, Department of Biochemistry,*

**University of Texas Southwestern**

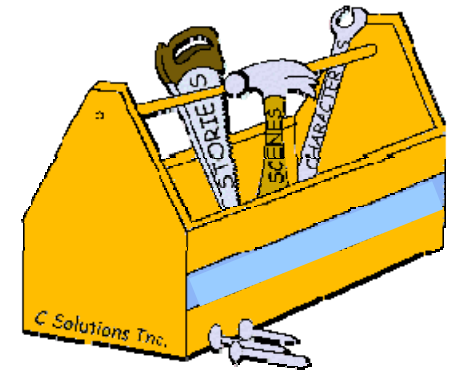
***Medical Center at Dallas***

Assembling  
**a toolbox**  
for analysis of  
protein molecules



# History tour.

## How did it all start?

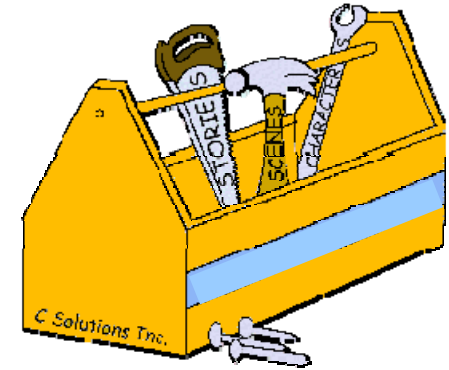


**Question 1:** Why is it that educated people (=experts) can understand biological phenomena so much better than computers?

**Question 2:** Why is it that those experts are so-o-o slow at what they do best?

**Question 3:** Why can't these experts teach computers to do the job right?

# History tour.



## How did it all start?

**Question 1:** Why is it that educated people (=experts) can understand biological phenomena so much better than computers?

Maybe they don't

**Question 2:** Why is it that those experts are so-o-o slow at what they do best?

Lazy?

**Question 3:** Why can't these experts teach computers to do the job right?

Snobbish?

We think we are experts.

We are trying to teach computers to give correct answers –  
and it is hard!

We think we are experts.

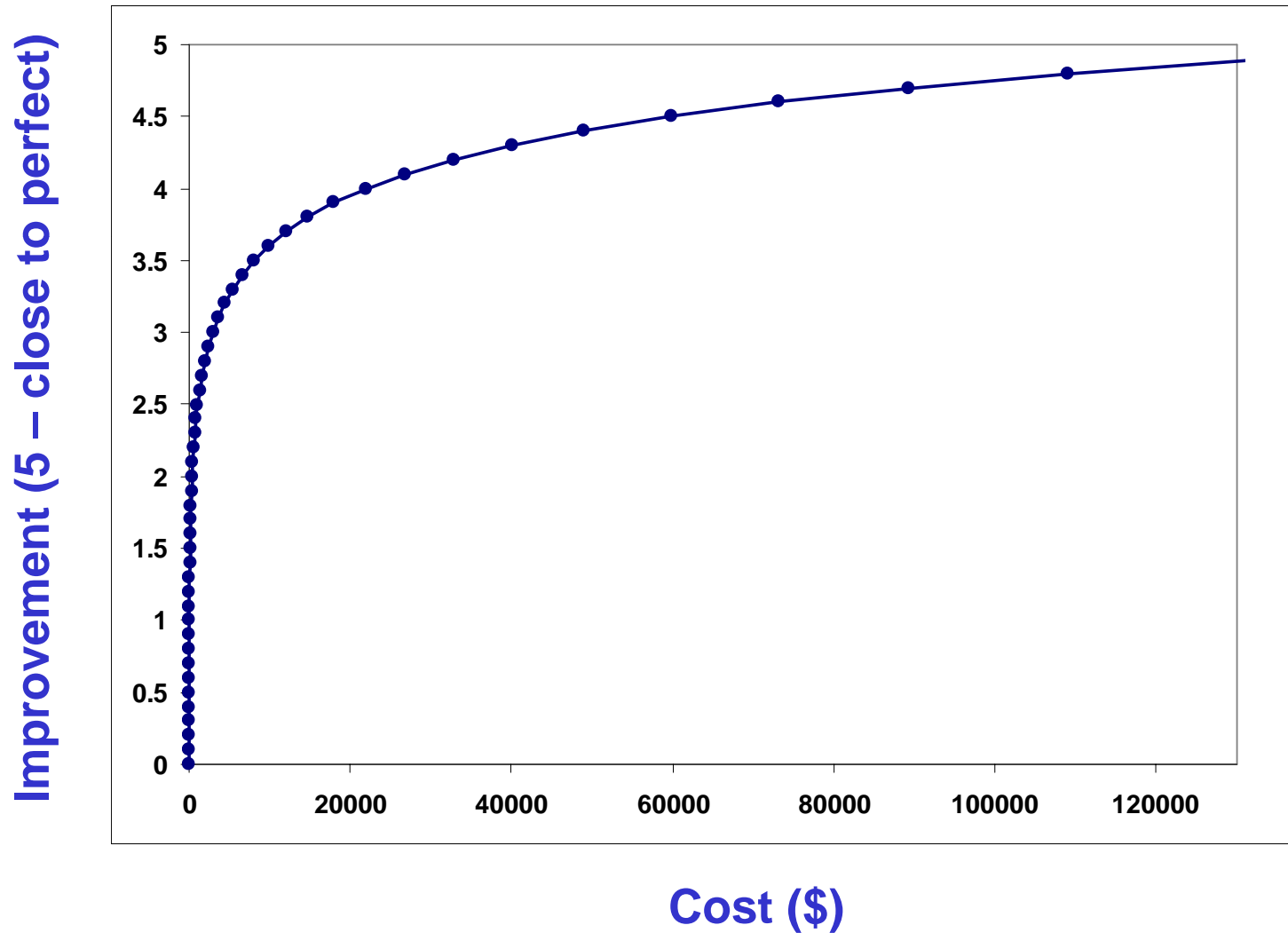
We are trying to teach computers to give correct answers –  
and it is hard!

Answers to what questions?

YOU KNOW ...

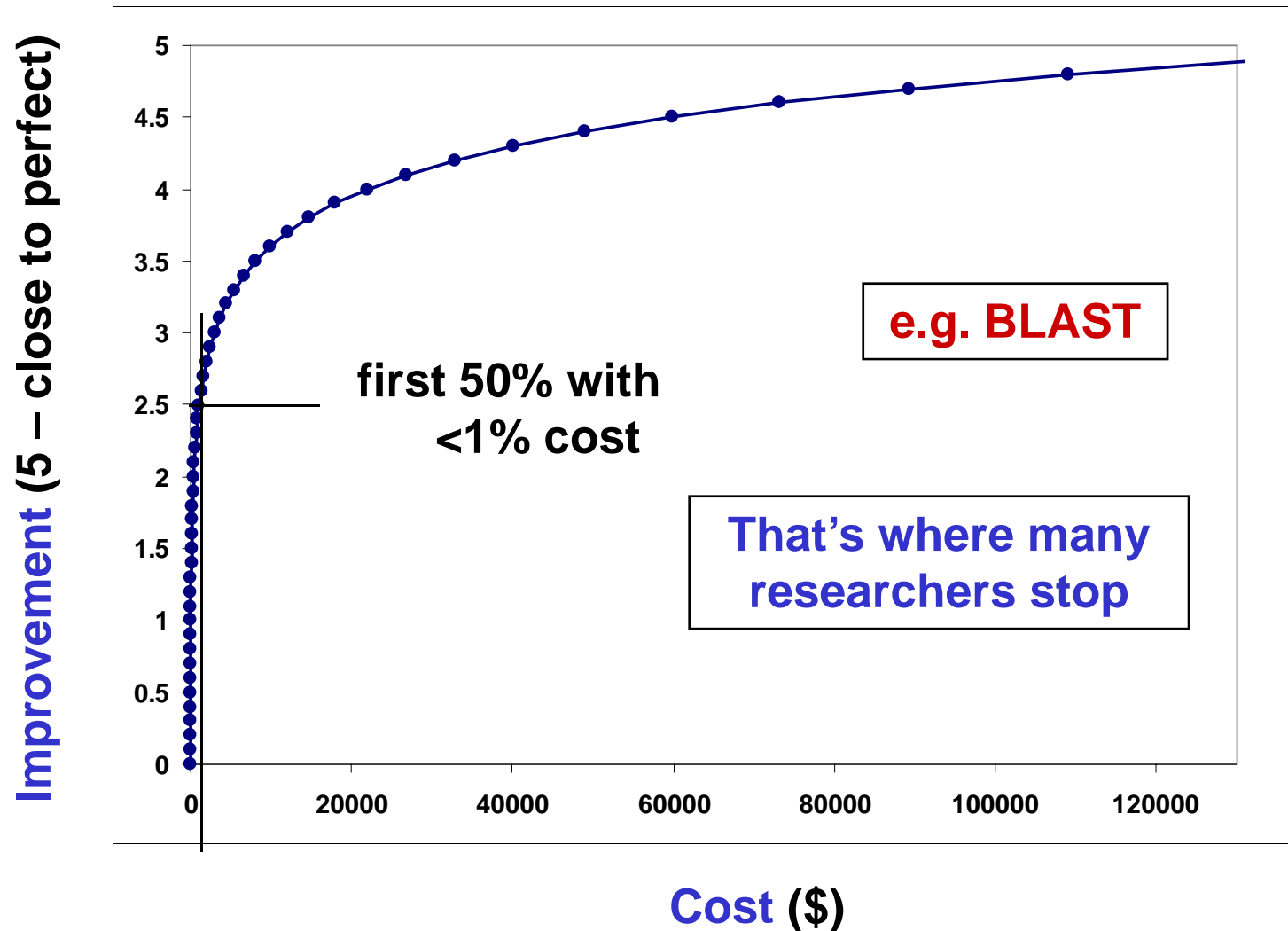
- we have a protein sequence – what is its 3D structure?
  - we have a protein 3D structure – where is the functional site?
  - we have 2 sequences – what is their alignment?
  - we have many related sequences – what is the tree?
- etc. etc. etc.

Universal law of science: cost for an increment in improvement increases exponentially with the improvement

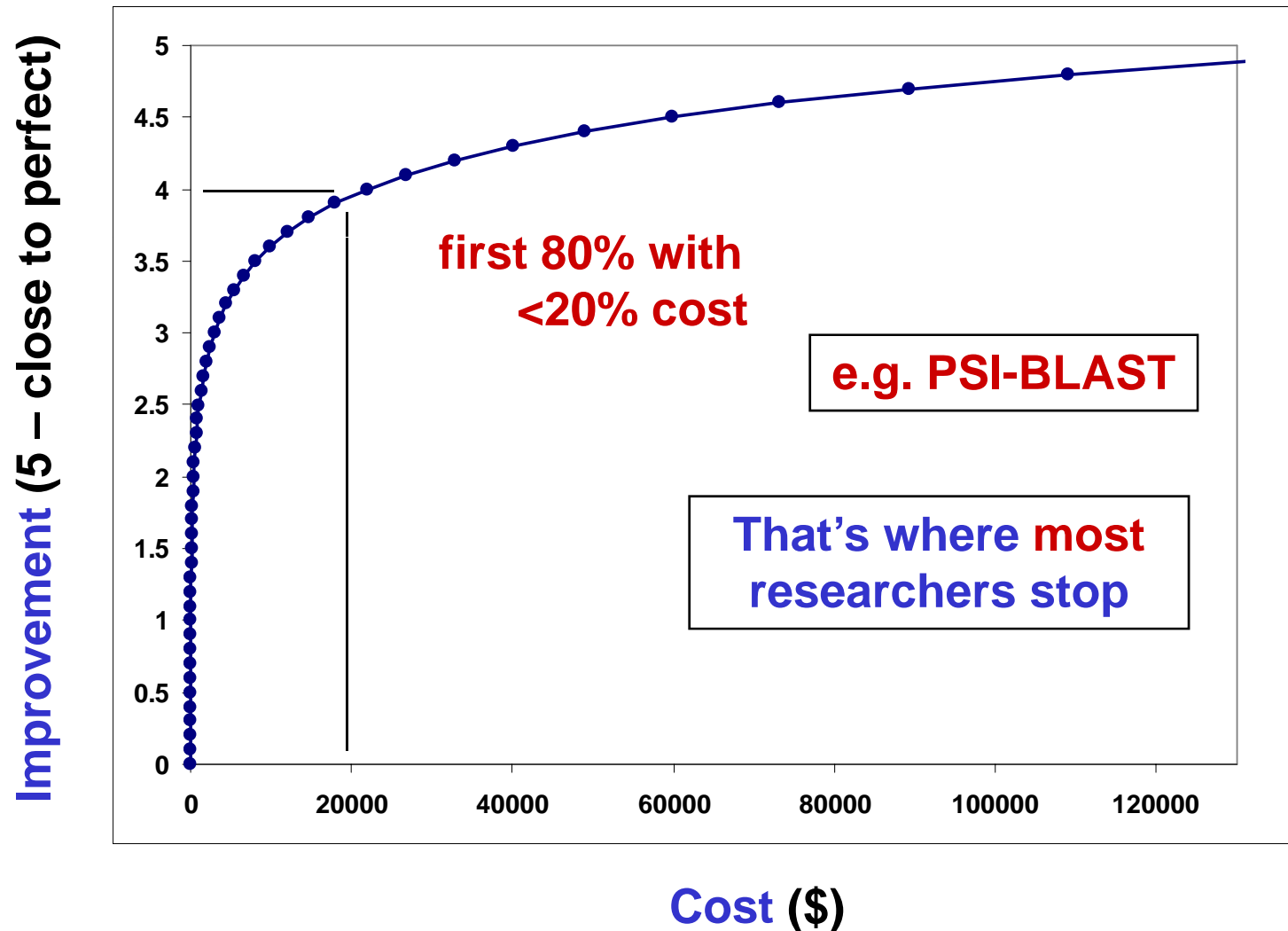




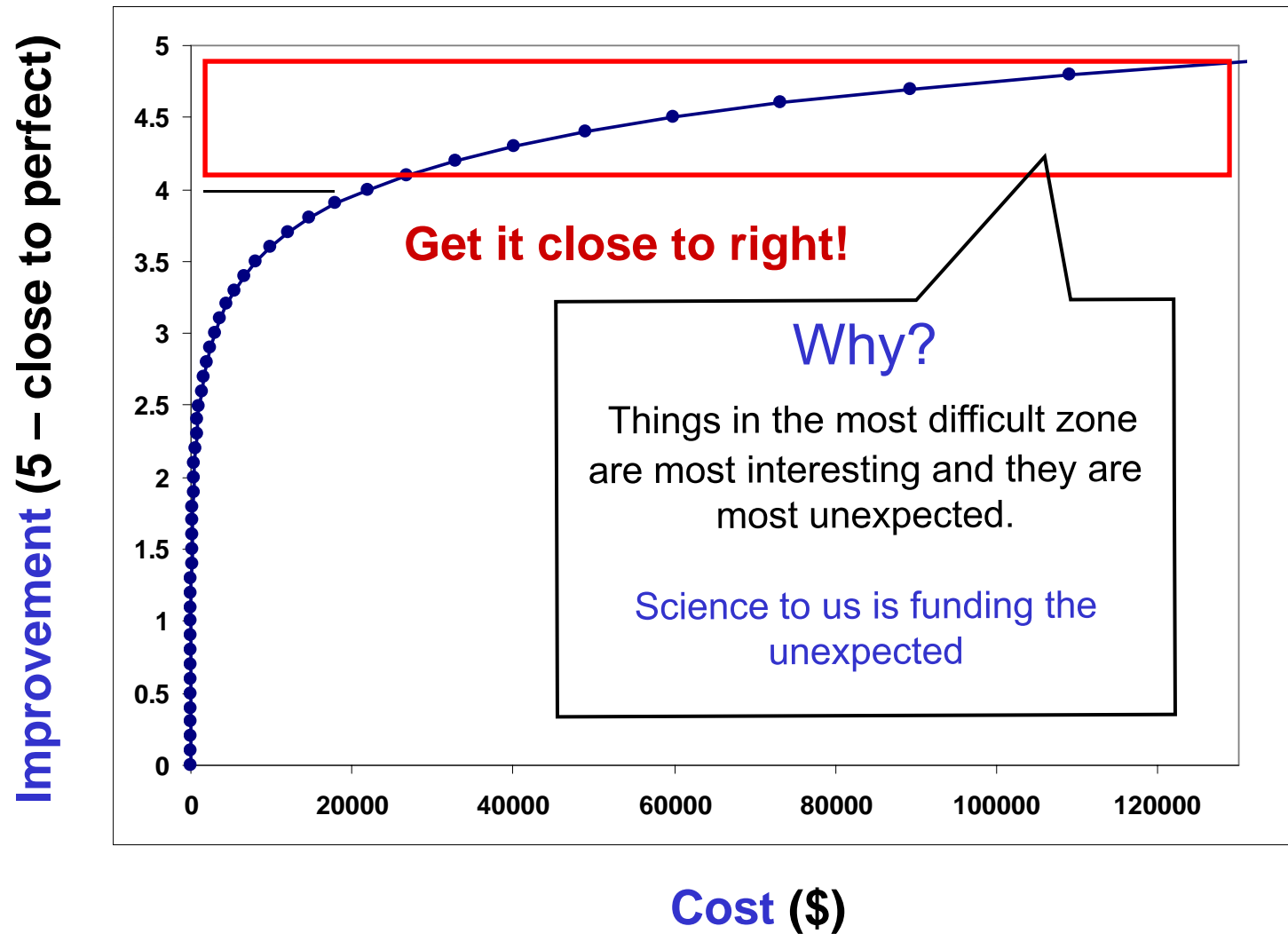
Universal law of science: cost for an increment in improvement increases exponentially with the improvement



Universal law of science: cost for an increment in improvement increases exponentially with the improvement

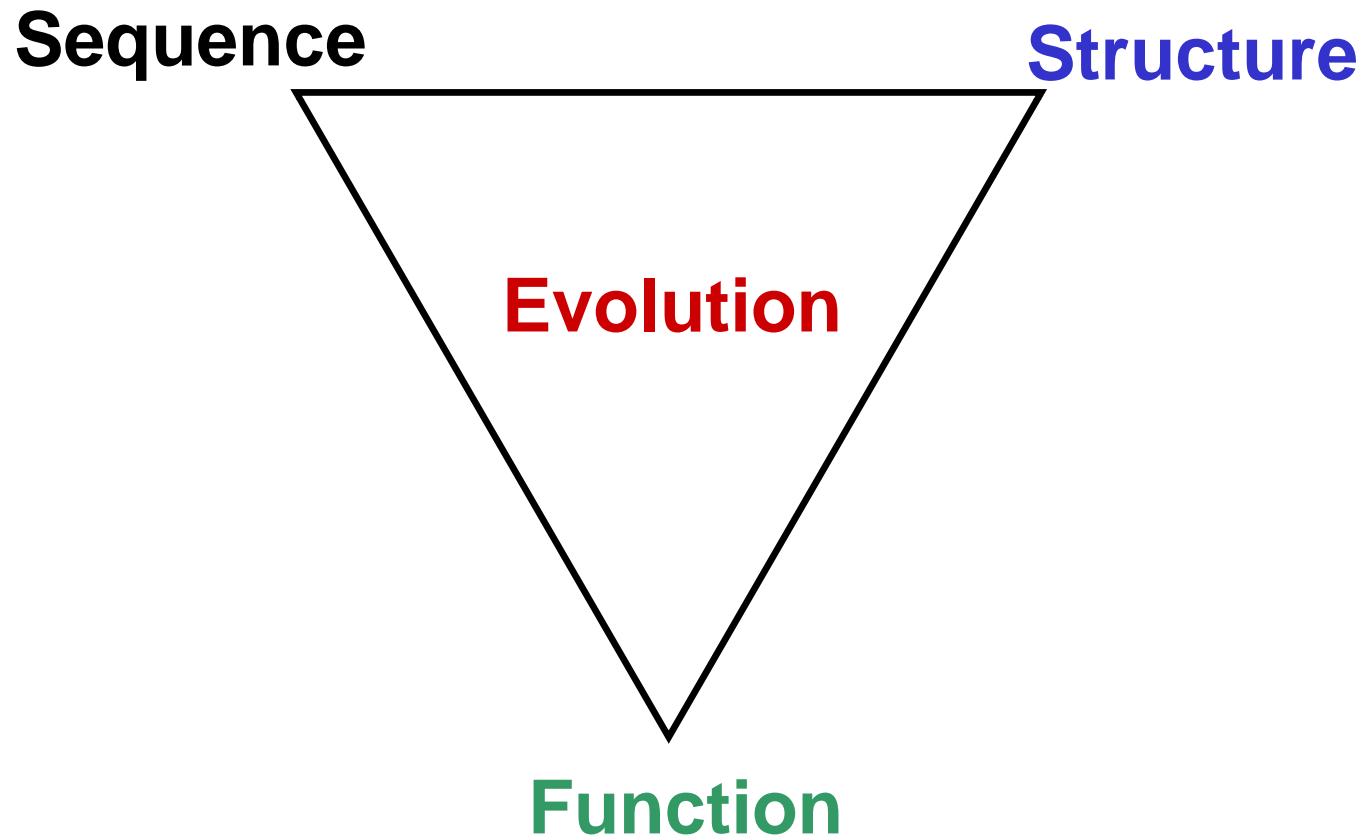


# Our Zone



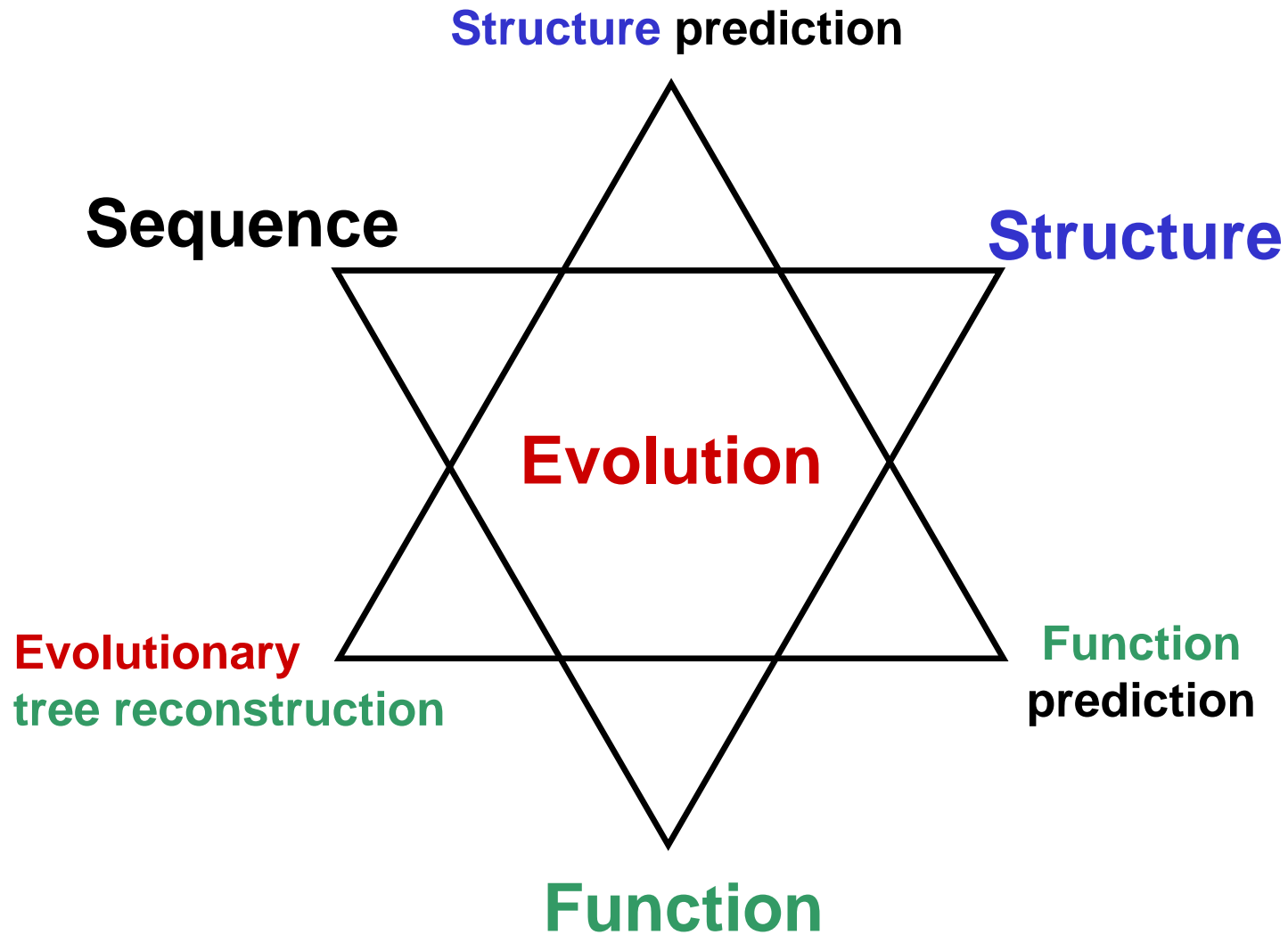
# Why do we need many tools?

---



# Why do we need many tools?

---



# Main tools in the toolbox

---

## Sequence analysis tools:

- Alignment of alignments and alignment similarity search;
  - plain sequence alignments;
  - alignments with predicted sec.str.
- Multiple sequence alignment;
- Sequence space visualization.

## Structure analysis tools:

- Secondary structure delineation;
- Pattern-matching structure similarity search;
- Structure alignment.

## Function prediction tools:

- Prediction of functional sites
  - universally important sites;
  - functional specificity sites.
- Evolutionary tree and ancestral sequence reconstruction.

# Today's agenda

---

1. **COMPASS: Search for similarity between families**
2. **PROMALS: Multiple sequence alignment**

**1. COMPASS:  
Search for similarity between families**

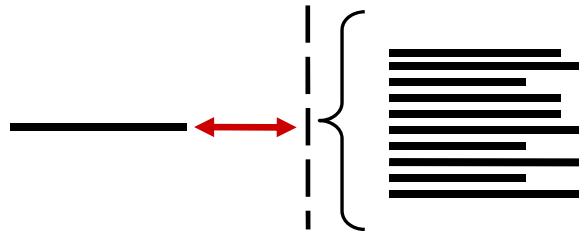
Ruslan Sadreyev



# Comparison of multiple alignments improves similarity detection

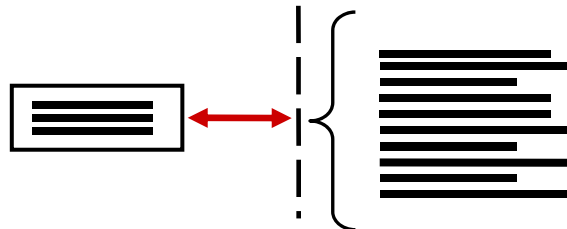
---

## Sequence-sequence (e.g. BLAST)



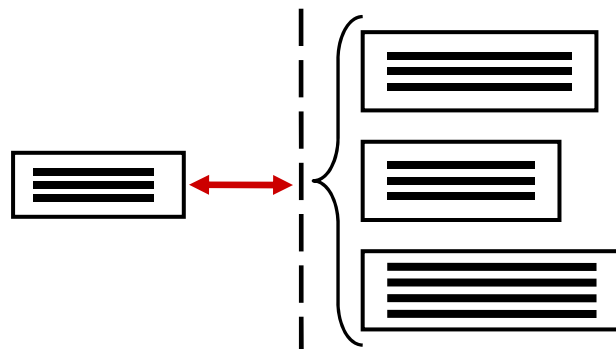
QGV<sup>G</sup>EGPKPAIKLRA  
*VS.*  
RVAGMKPRFVRSVKIVHR

## Alignment-sequence (e.g. PSI-BLAST)



QGV<sup>G</sup>EGPKPAIKLRA  
EG<sup>L</sup>EGPASRF<sup>R</sup>VTV  
KK<sup>V</sup>DGPPV-S<sup>R</sup>MTT  
*VS.*  
RVAGMKPRFVRSVKIVHR

## Alignment-alignment (e.g. COMPASS)



QGV<sup>G</sup>EGPKPAIKLRA  
EG<sup>L</sup>EGPASRF<sup>R</sup>VTV  
KK<sup>V</sup>DGPPV-S<sup>R</sup>MTT  
*VS.*  
RVAG<sup>M</sup>KPR<sup>F</sup>VRSVKIVHR  
IIR<sup>A</sup>SKPK<sup>F</sup>TRS<sup>V</sup>TI-HR  
QLV<sup>G</sup>SKPK<sup>F</sup>TRTLVT-HR

# COMPASS web server

<http://prodata.swmed.edu/compass>

## COMPASS

COmparison of Multiple Protein sequence Alignments with assessment of Statistical Significance.  
Similarity search with alignment (or sequence) as a query against a database of protein families.

**COMPASS:**  
a method for  
COmparison of  
Multiple  
Protein  
Alignments with  
assessment of  
Statistical  
Significance

Enter [query](#) as protein sequence or multiple alignment:

Or file name

Choose [database](#)

Email address to receive the result (optional)

[Job ID](#) (optional)

### Input processing options

Run PSI-BLAST if input is: [a sequence](#)  Yes  No [an alignment](#)  Yes  No

PSI-BLAST options: [Iterations](#)  [Evalue](#)  [Coverage](#)  [Identity](#)

[Gap fraction threshold](#)

### Search options

[Gap opening penalty](#)  [Gap extension penalty](#)  [Matrix](#)

[Effective length of the database](#)  Real Size  User Defined:

### Output formatting options

[Expect](#)  [Significance threshold](#)  [Display up to](#)  hits

[Top sequences to show](#) in query  in subject

[Show consensus sequences](#)  Yes  No [Show fully gapped positions](#)  Yes  No

[Width of alignment segment](#)

[Documentation](#)

[Downloads](#)

Sadreyev and Grishin  
(2003) *JMB*, 326: 317

# Recent changes: 2007

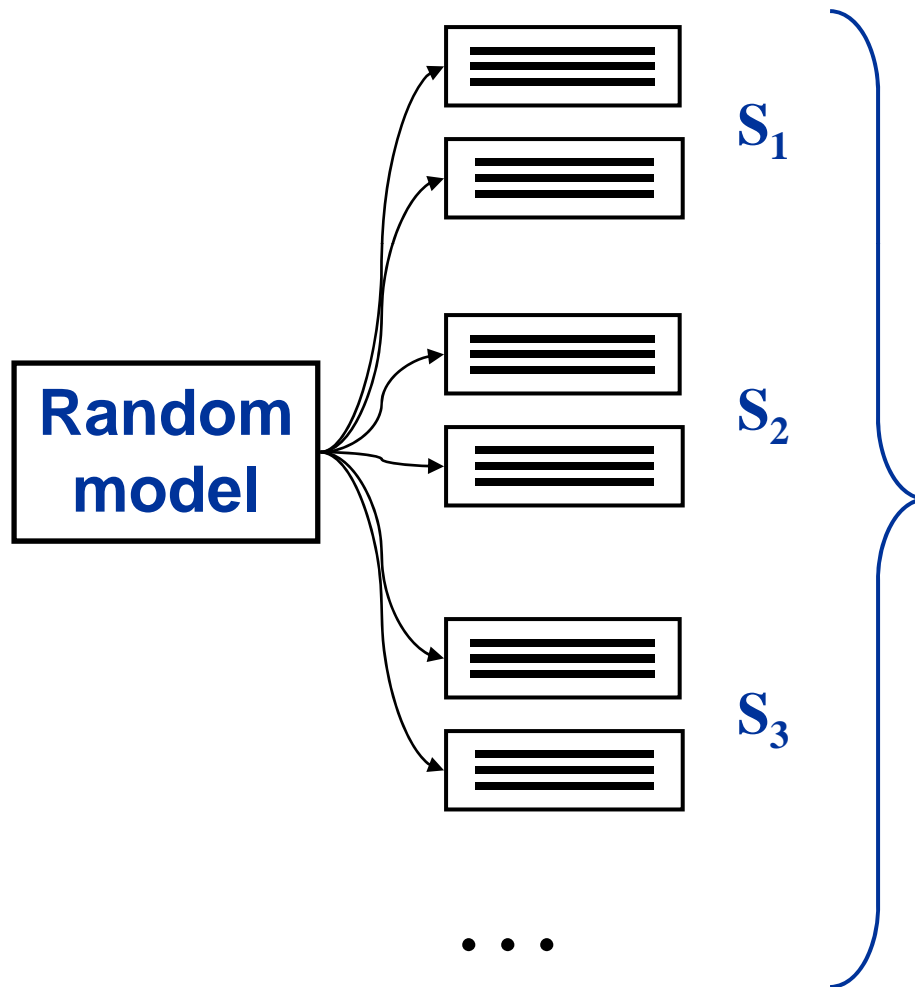
---

**1. New random model for profiles**

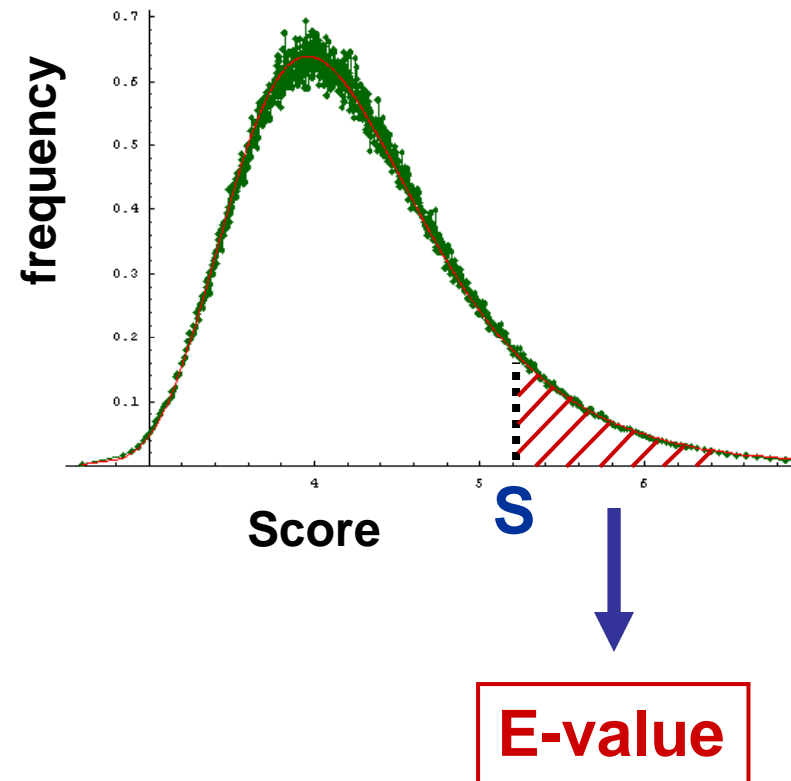
**2. New distribution to describe scores**

# Estimates of statistical significance are based on a random model of alignment comparison

Random decoy profiles



Score distribution



## Old random model

---

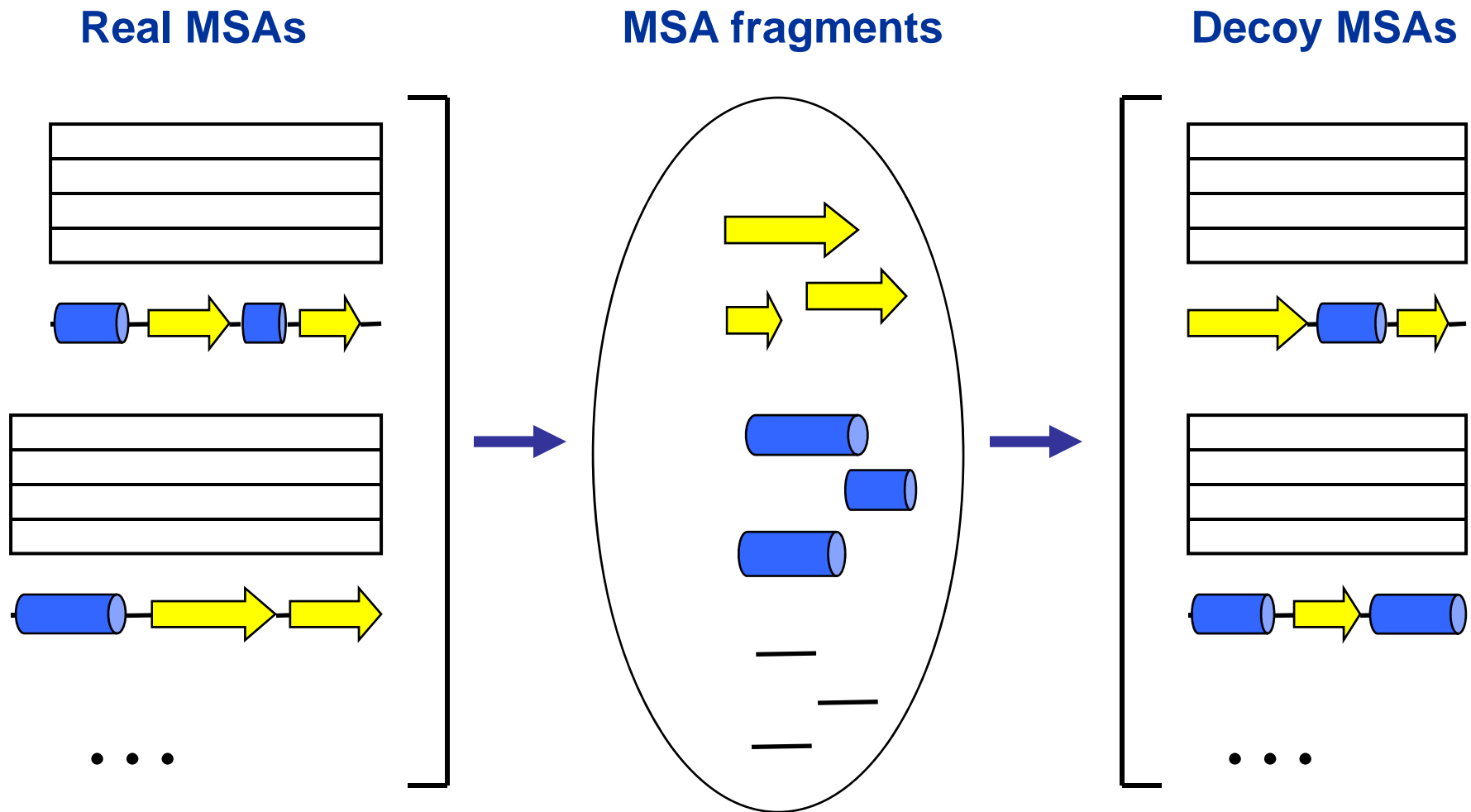
Independent positions: shuffling positions  
makes decoy alignments

This model works very well in  
**BLAST** and **PSI-BLAST**,

however, maybe more realistic models work better

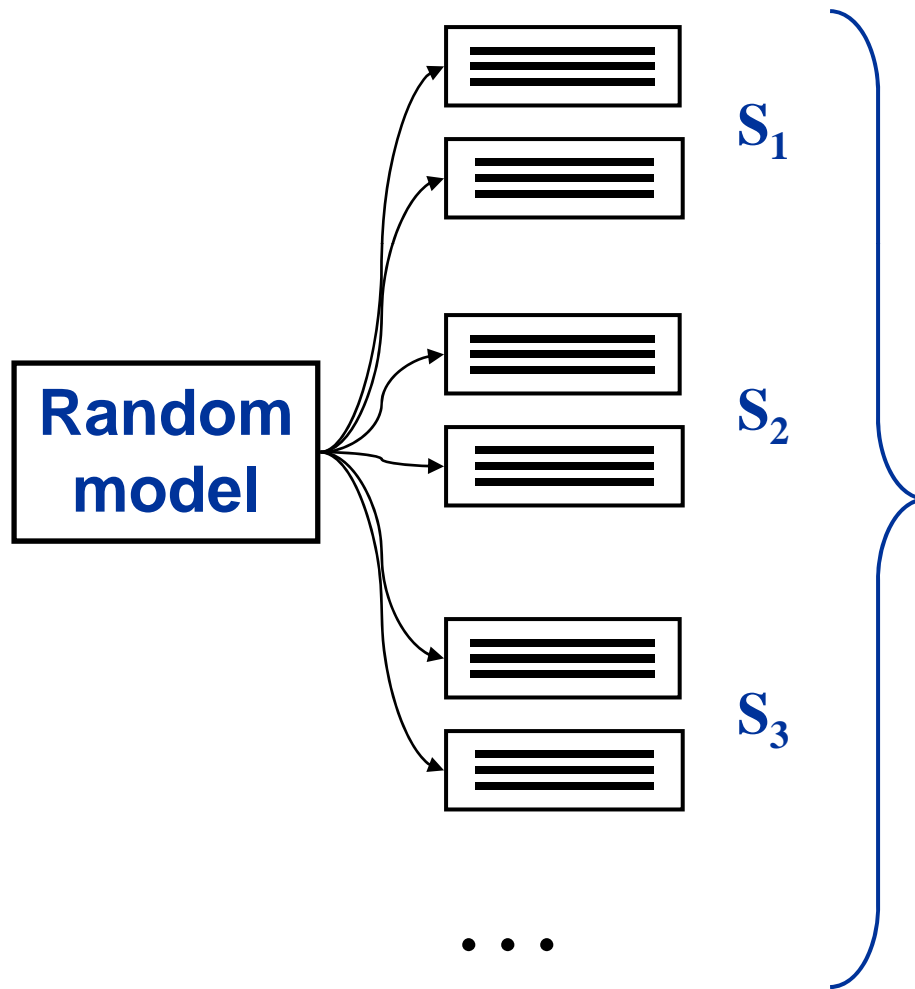
## Reproducing protein features:

Real secondary structure elements are used as building blocks for decoy MSA

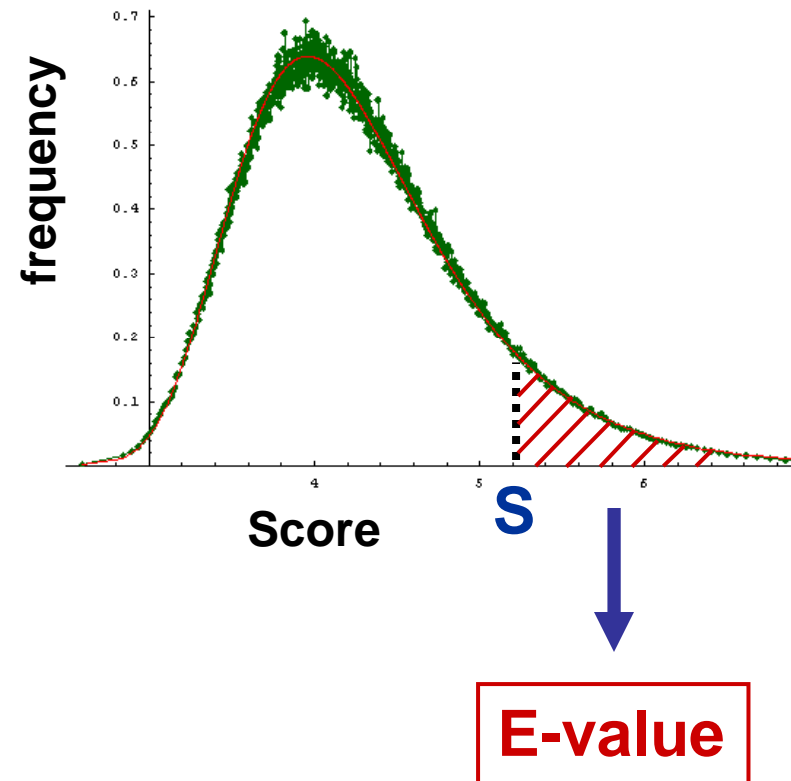


# Estimates of statistical significance are based on a random model of alignment comparison

Random decoy profiles from SS

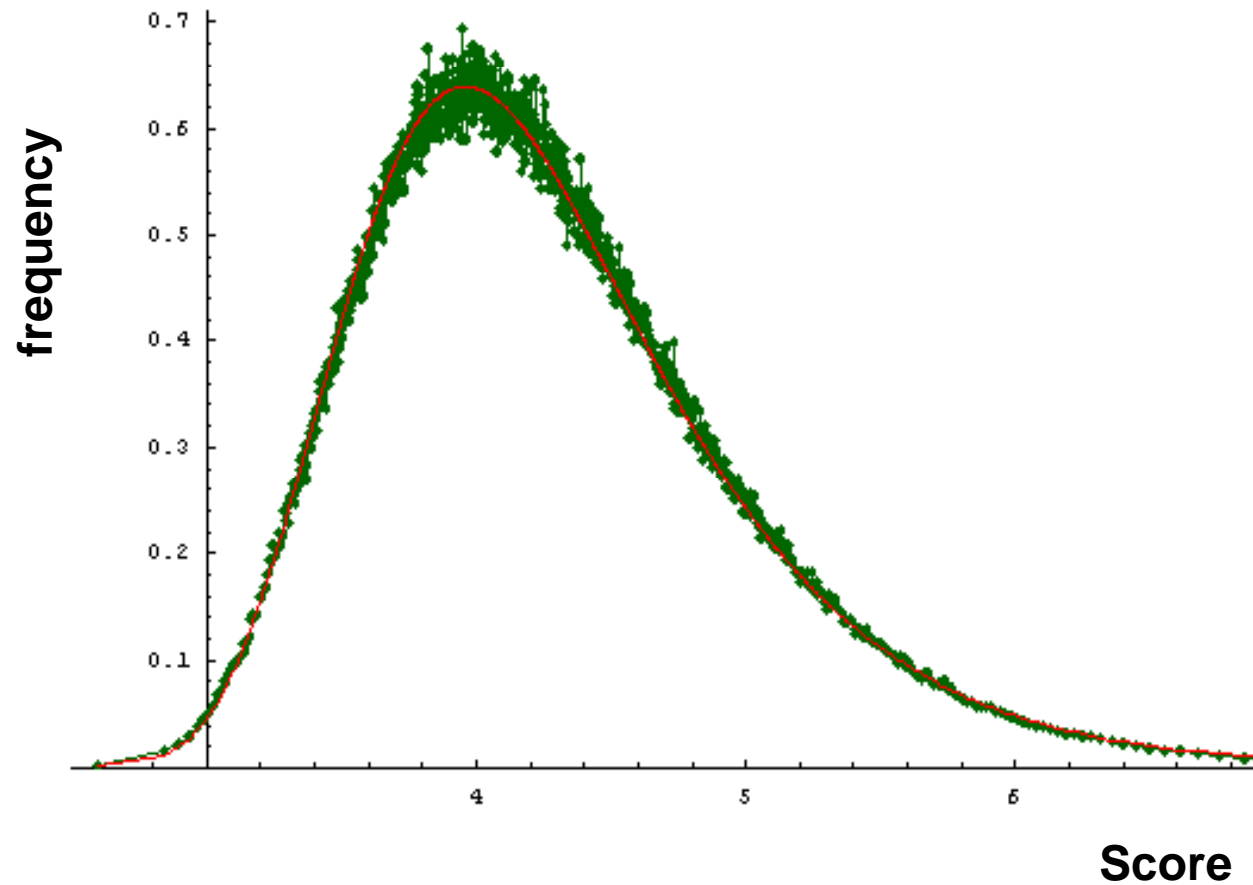


Score distribution



# Distribution of scores for random MSA comparison

---



**Describe empirical distribution with a continuous density function**



# Gumbel Extreme Value Distribution (EVD) is traditionally used to describe similarity cores

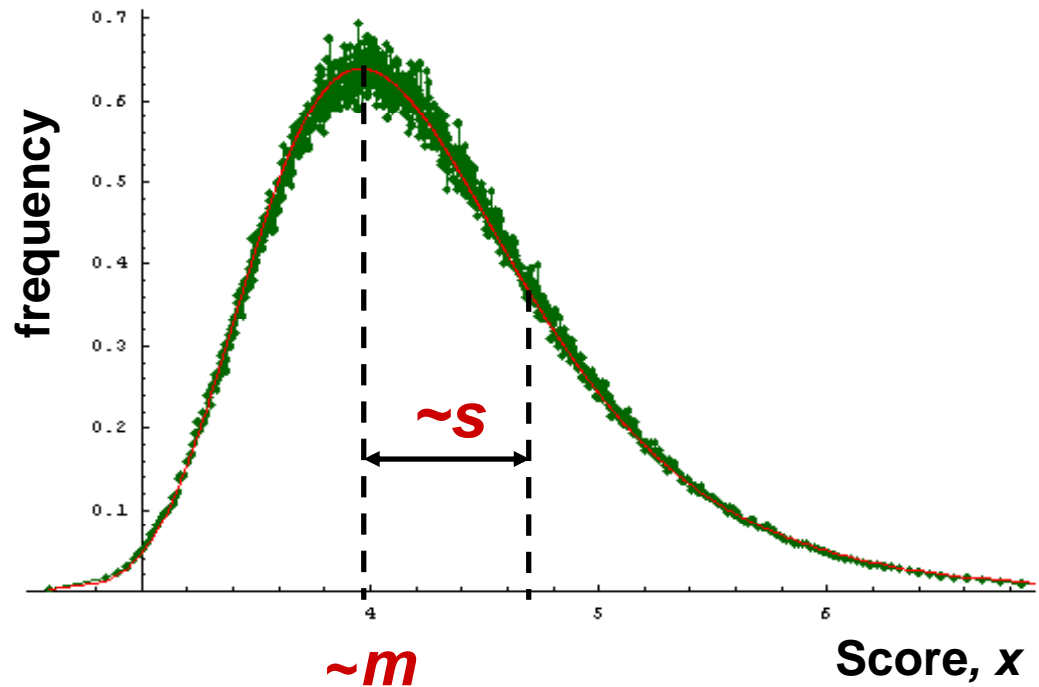
---

EVD pdf:

$$f(x) = C_1 \exp\left(-e^{-\frac{x-m}{s}} - \frac{x-m}{s}\right)$$

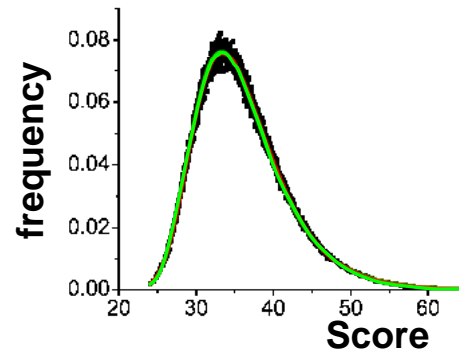
$m$ : location parameter

$s$ : scale parameter



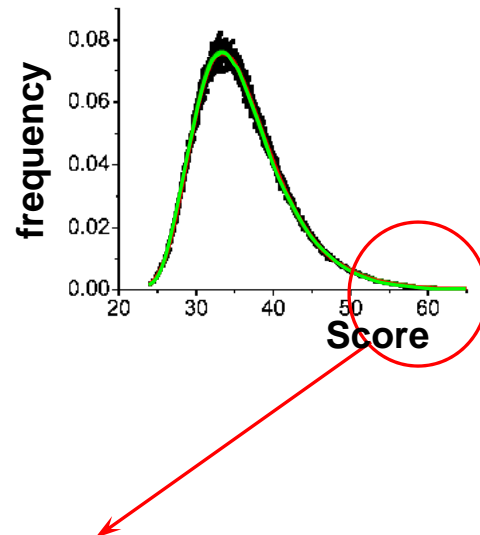
## EVD does not fit empirical score distributions

$$f(x) = C_1 \exp\left(-e^{-\frac{x-m}{s}} - \frac{x-m}{s}\right)$$



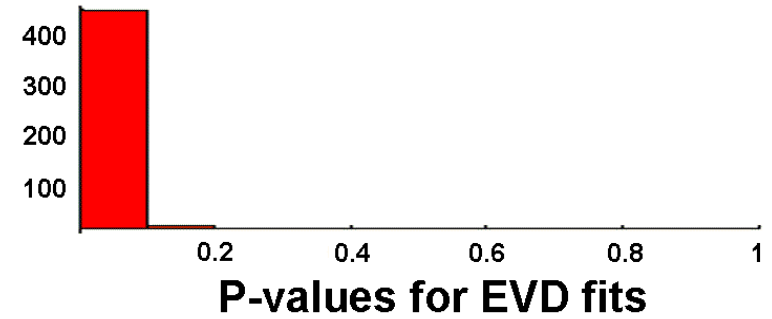
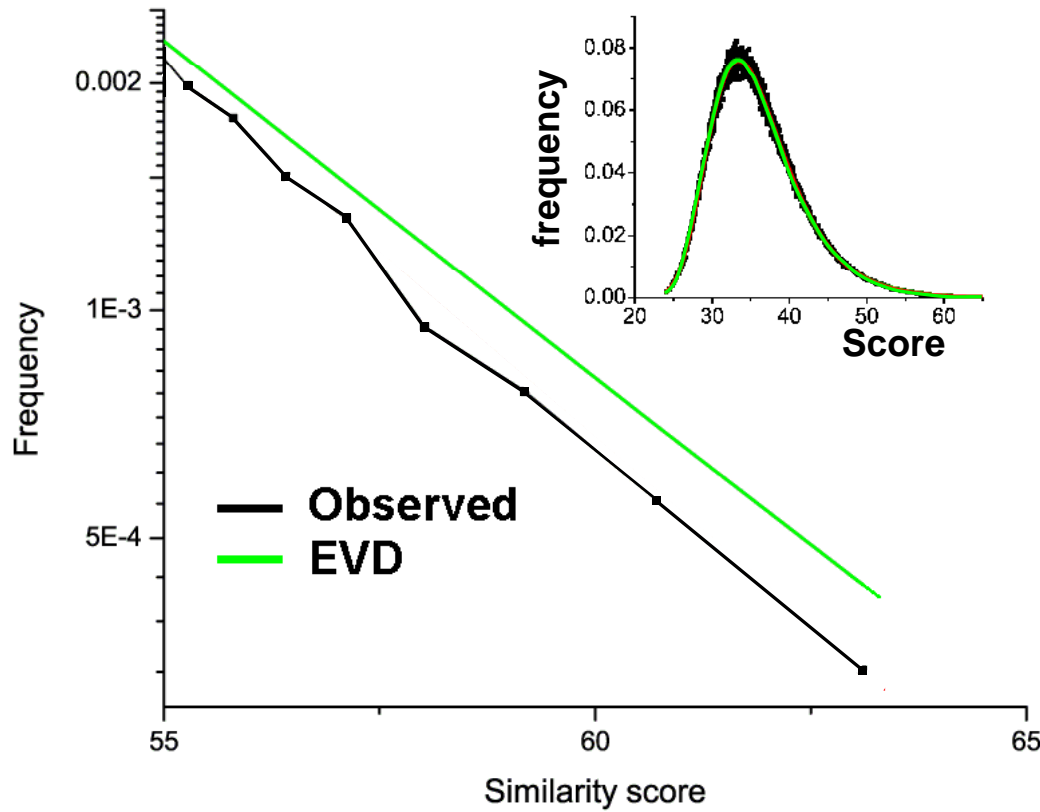
# EVD does not fit empirical score distributions

$$f(x) = C_1 \exp\left(-e^{-\frac{x-m}{s}} - \frac{x-m}{s}\right)$$

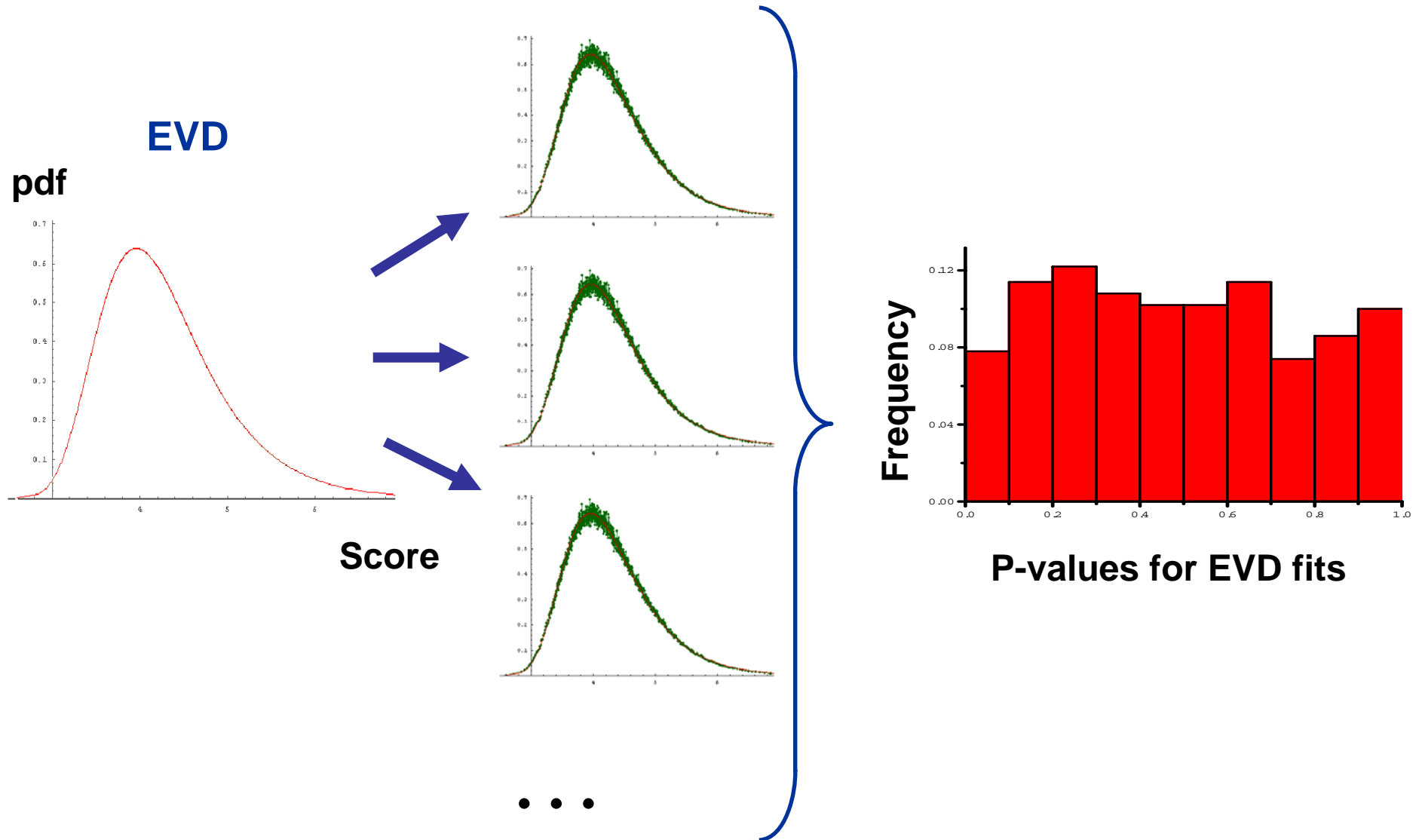


# EVD does not fit empirical score distributions

$$f(x) = C_1 \exp\left(-e^{-\frac{x-m}{s}} - \frac{x-m}{s}\right)$$



For data generated from the same distribution,  
fitting P-values are distributed uniformly



## Scores generated by SS-based model do not obey other standard statistical distributions

---

Distributions of Pearson system

Distributions of Johnson system

Inverse Gaussian (Wald) distribution

Burr

Weibul

Tukey (lambda)

Non-central chi square

Non-central t

$\chi^2$  goodness-of-fit

does not pass

P-values  $< \sim 10^{-5}$

# We had to invent a new distribution

---

How?

Modify EVD!

EVD pdf:  $f(x) = C_1 \exp\left(-e^{-\frac{x-m}{s}} - \frac{x-m}{s}\right)$

Power EVD pdf:  $f(x) = C_2 \exp\left(-e^{-\frac{x-m}{s}} - \beta \frac{x^\alpha - m^\alpha}{s^\alpha}\right)$

WOW!

## A new distribution, power EVD (PEVD), is created by modification of EVD

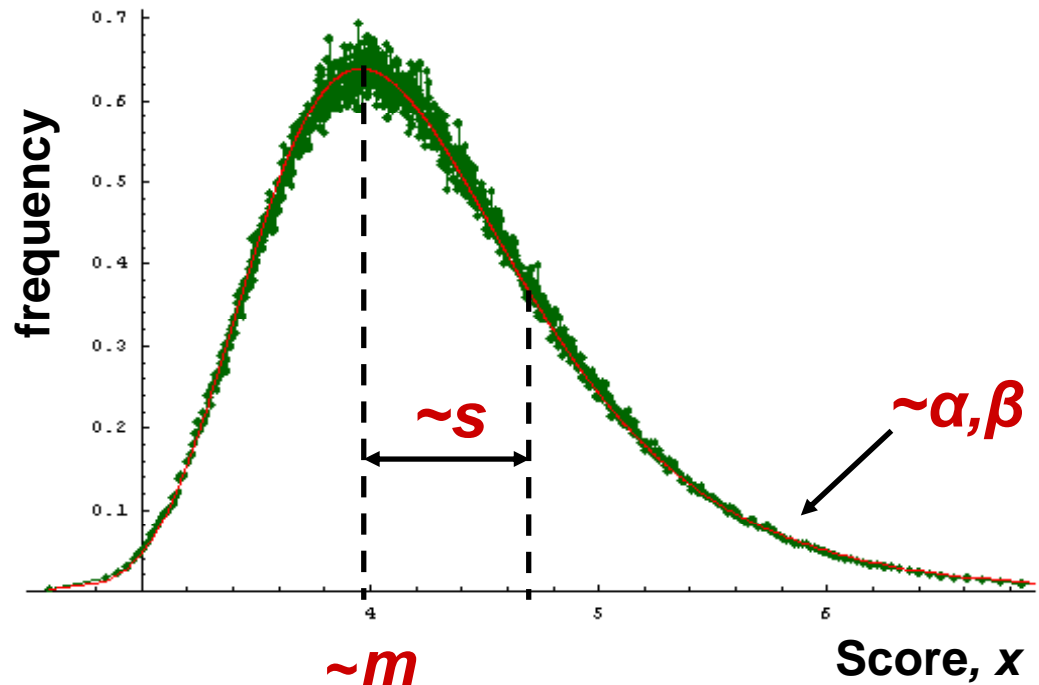
**EVD pdf:** 
$$f(x) = C_1 \exp\left(-e^{-\frac{x-m}{s}} - \frac{x-m}{s}\right)$$

**PEVD pdf:** 
$$f(x) = C_2 \exp\left(-e^{-\frac{x-m}{s}} - \beta \frac{x^\alpha - m^\alpha}{s^\alpha}\right)$$

**$m$** : location parameter

**$s$** : scale parameter

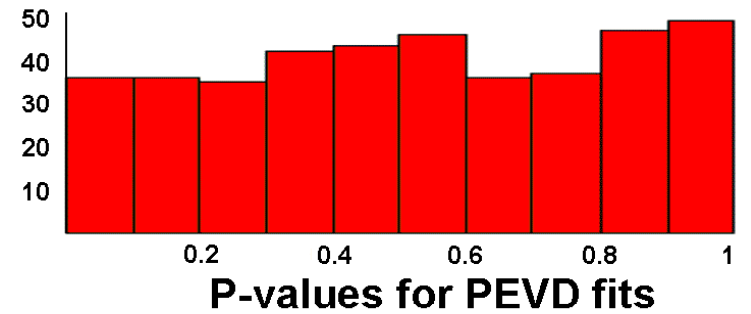
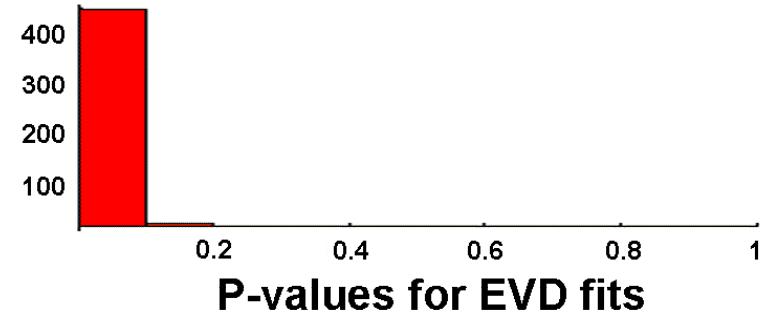
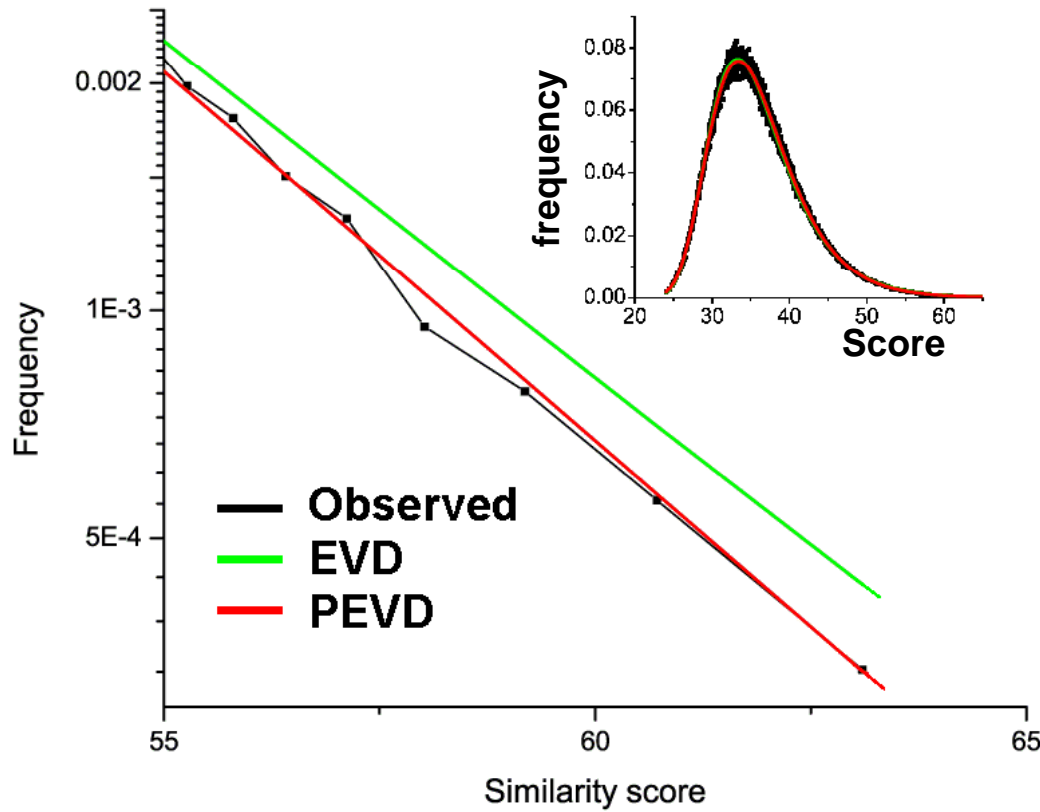
**$\alpha, \beta$** : shape parameters





# Power EVD precisely fits empirical score distributions

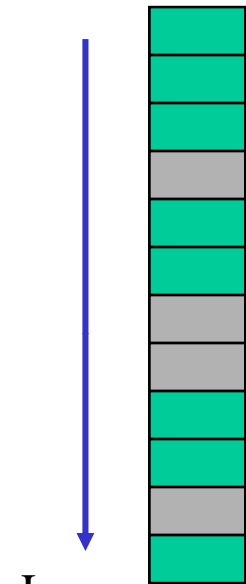
$$f(x) = C_2 \exp\left(-e^{-\frac{x-m}{s}} - \beta \frac{x^\alpha - m^\alpha}{s^\alpha}\right)$$



# The **new random model + new distribution** improve homology detection

Query: 

Database hits:



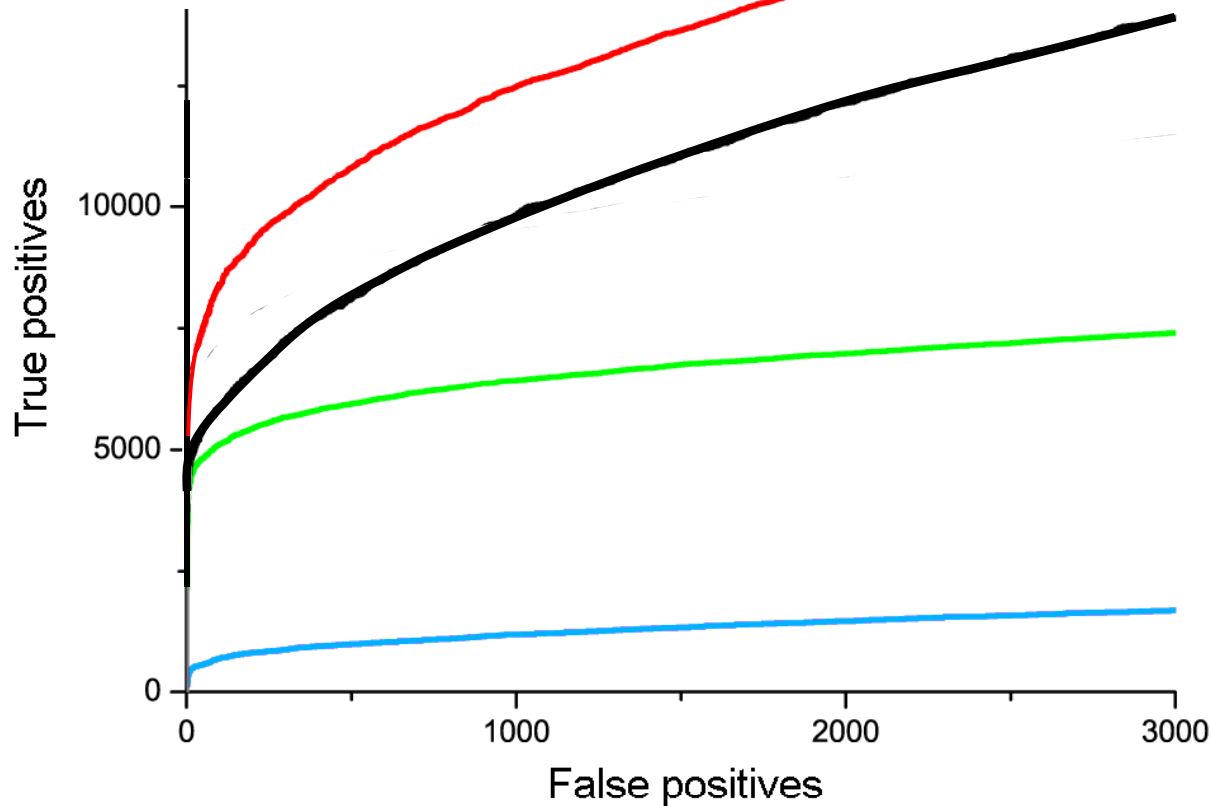
Less  
significant  
E-value

 True Positive  
 False Positive

# The new random model + new distribution improve homology detection

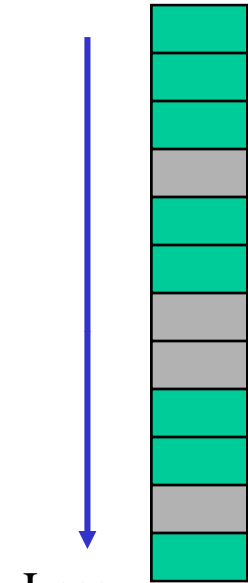
- Random columns
- SS-based model
- SSEARCH
- PSI-BLAST

## ROC curve



Query: 

Database hits:



Less  
significant  
E-value

-  True Positive
-  False Positive

Benchmark: 2900 PSI-BLAST alignments for SCOP domain representatives with known relationships

# Summary

- We developed a **realistic random model** that simulates random MSA comparison by **mimicking** native protein **secondary structure**
- We developed a **precise analytical approximation of the simulated score distributions**, based on a **new distribution function, PEVD**
- Applied to protein similarity searches, the **new model** produces more realistic E-values and (unexpectedly) **improves homology detection**

**2. Towards **accurate**  
multiple sequence alignments  
of **distantly** related proteins**

Jimin Pei

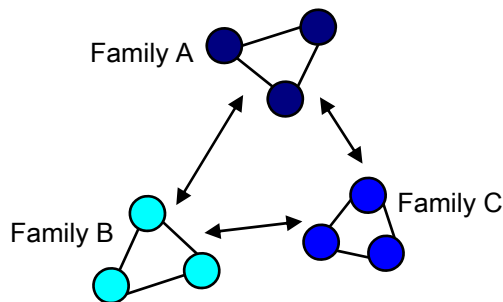
# Multiple sequence alignment

```
BSUB00  RMAHYDSLTDLPNRRHAI SHLTKVLNREHSLHYNTVVFFLDLNRFKVINDAL
ECU738  VMSTRDGMTGVYNRRHWETMLRNEFDNCRRHNRDATLLIIDIDHFKSINDTW
D90790  HEVGMDVLTKLLNRRFLPTIFKREIAHANRTGTPLSVLIIDVDKFKEINDTW
SYCSLL  QISSLDALTQVGNRYLFDSTLEREWQRLQRIREPLALLLCDVDFFKGFNDNY
ECAE00  NIAHRDPLTNI FNRNYFFNEL - - TVQSASAQKTPYCV MIMDIDHFKKVNDTW
AF0348  QAANVDSL TGLANRAAYNAHM - ERLTAADAPS - - IGLLLIDVDRLKQVNDIL
D90796  IRSNMDVLTGLPGRRVLDESFDHQLRNAEPLN - - LYLMLLDIDRFKLVNDTY
Y4LL_R  HMARHDALTGLPNRQFLREEF - ERLSDHIAPSTR LAILCLDLDGFKAINDAY
Y07I_M  YLADHDDL TGLHNRRALLQHLDQRLAPGQPGP - - VAALFLDL DRLKAINDYL
.....
```

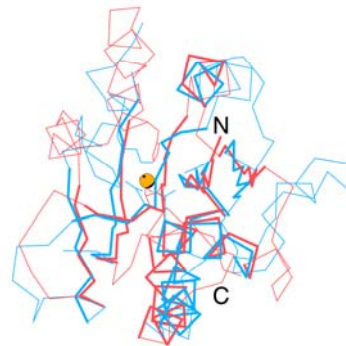
# Multiple sequence alignment

```
BSUB00  RMAHYDSLTDLPNRRHAISHLTKVLNREHSLHYNTVVFFLDLNRFKVINDAL
ECU738  VMSTRDGMTGVYNRRHWETMLRNEFDNCRRHNRDATLLIIDIDHFKSINDTW
D90790  HEVGMDVLTKLLNRRFLPTIFKREIAHANRTGTPLSVLIIDVDKFKEINDTW
SYCSLL  QISSLDALTQVGNRYLFSTLEREWQRLQRIREPLALLLCDVDFFKGFNDNY
ECAE00  NIAHRDPLTNIFNRNYFFNEL--TVQSASAQKTPYCVMIMDIDHFKKVNDTW
AF0348  QAANVDSLTGLANRAAAYNAHM-ERLTAADAPS--IGLLLIDVDRLKQVNDIL
D90796  IRSNMDVLTGLPGRVLDESFDHQLRNAEPLN--LYLMLLDIDRFKLVNDTY
Y4LL_R  HMARHDALTGLPNRQFLREEF-ERLSDHIAPSTRLAILCLDLDGFKAINDAY
Y07I_M  YLADHDDLTGLHNRRALLQHLDQRLAPGQPGP--VAALFLDLDRLKAINDYL
.....
```

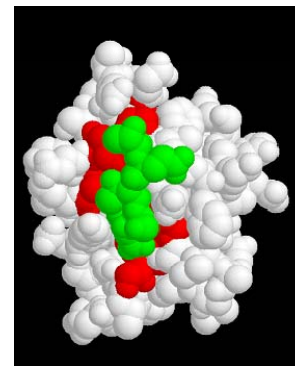
Protein similarity search  
and classification



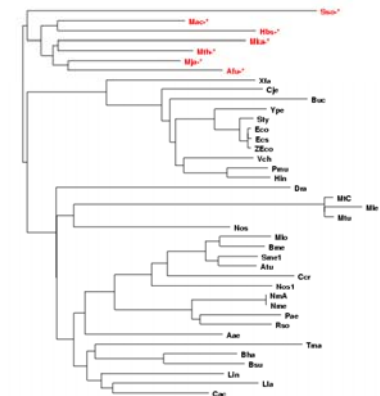
Structure modeling



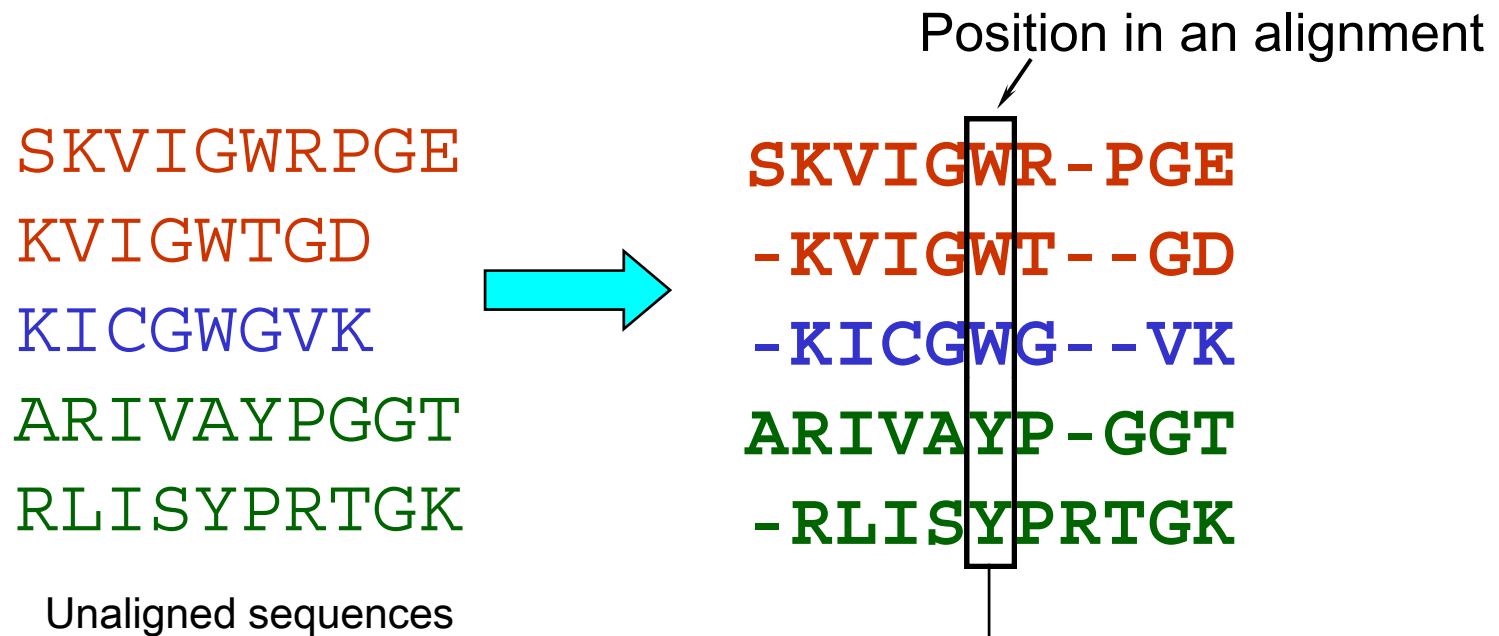
Active site prediction  
experimental design



Phylogenetic analysis



# Meaning of alignments



- Homologous
- Structurally equivalent
- Similar function



**How is the alignment made?**

**ClustalW** –

**the most widely used alignment program**

# ClustalW – the most widely used program

**+** [ch.EMBnet.org](http://www.ch.embnet.org)

Home Services Courses Links Contacts

## ClustalW

Valid format for input is: FASTA(Pearson)  
max number of sequences = 30  
max total length of sequences = 10000  
[Help page](#)  
For more than 30 sequences please use [ClustalW-XXL](#)

---

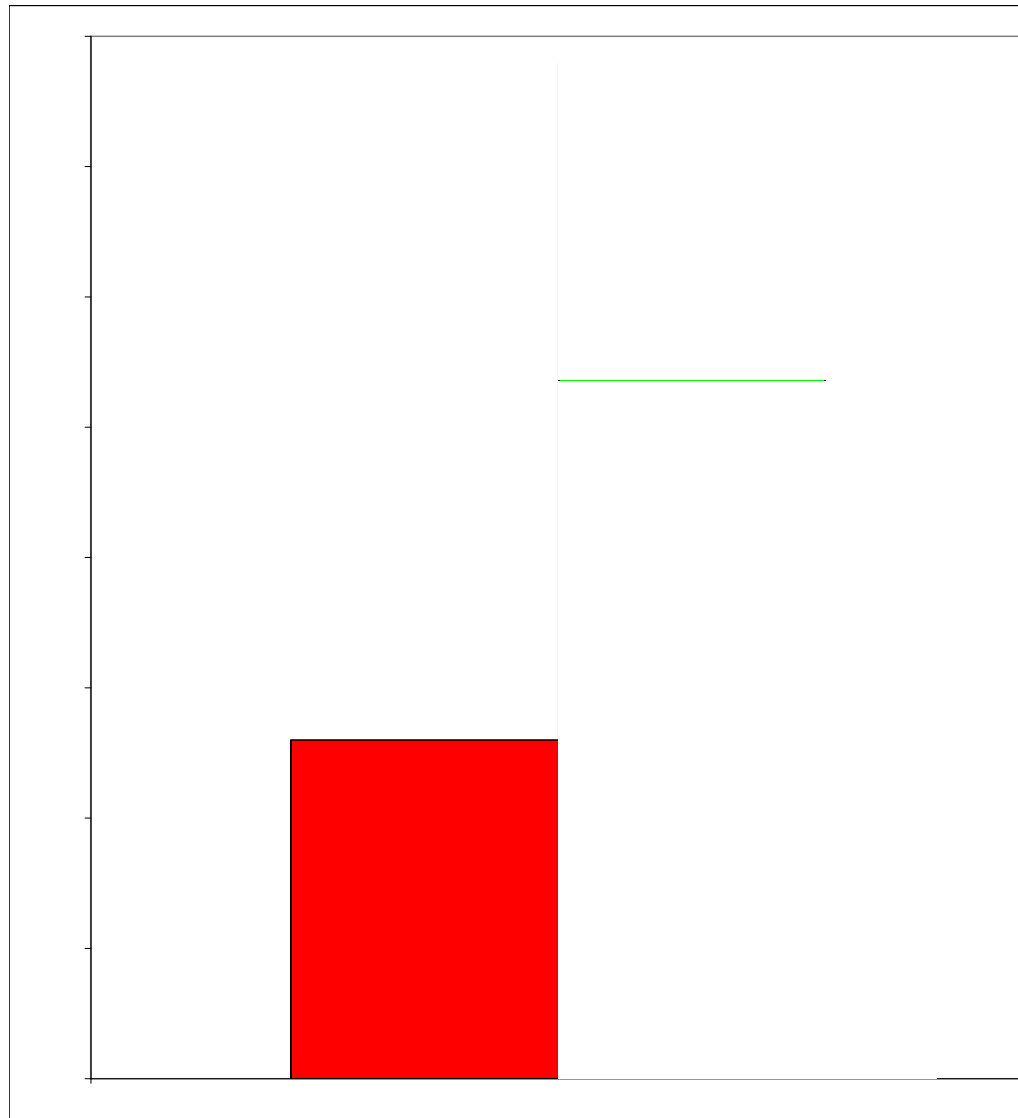
Scoring matrix :	Blosum ▾		
Opening gap penalty :	10	Extending gap penalty :	0.05
End gap penalty :	10	Separation gap penalty :	0.05
Output format :	Clustal ▾	Output order :	Input ▾

Input sequences:  
(see above for valid formats)

Run ClustalW Clear Input

Thompson et al. (1994). <http://www.ch.embnet.org/software/ClustalW.html>

# How accurate are these alignments?



ClustalW accuracy

# How accurate are these alignments?



ClustalW accuracy

PROMALS accuracy

About 3 times better than ClustalW

# PROMALS:

(**PRO**file **M**ultiple **A**lignment with  
predicted **L**ocal **S**tructure)

# http://prodata.swmed.edu/promals

## The PROMALS multiple sequence alignment server

PROMALS constructs multiple protein sequence alignments using information from database searches and secondary structure prediction. [\[Documentation\]](#)

Enter your sequences in [FASTA](#) format:

Or upload a local file containing your sequences:

Enter your [email](#) address to receive the result ([recommended](#)):

Alignment options:

- [Weight for amino acid scores](#):
- [Weight for predicted secondary structure scores](#):
- [Identity threshold above which fast alignment is applied](#):

Enter a name for your job ([recommended](#)):

[PROMALS Documentation](#)

[Reference](#): Pei, J. and Grishin, N. V. (submitted). Towards accurate multiple sequence alignments of distantly related proteins.

Comments, suggestions and bug reports to: [jpei@chop.swmed.edu](mailto:jpei@chop.swmed.edu)

# What did we do to achieve this?



ClustalW accuracy

PROMALS accuracy

About **3 times** better  
than ClustalW

**First of all,**

**ClustalW is not that bad ...**



## ... for similar sequences

```
Q2BMK3      MLAKTVREQIQRPADLVARYGGEEFIVVLPDTDEEGAMAVAGQICVAVAS
Q3A8D4      QVARMLQSVVARPGDLVARYGGEEFALILPQTD-HGAKFLGESCRAAVAG
Q36PG9      ALASILSDEVQRSGDLVARYGGEEFAILLPTTDVAGAQQVAERMRLSVAR
Q2BQL8      TVAQTIKHSIQRAQDMVCRYGGEEFVVILPETDL DGAQMIAERIRKAIK
Q3XUK3      ALAHTISL-HLRPGDIAARYGGEEFAVVLPDTDAVSGRMIAERLRTAVEA
Q9HXT9      QVAGAIREGCSRSSDLAARYGGEEFAMVLPGTSPGGARLLAEKVRRTVES
P73713      TIGRILQSNIRGS-DIACRYGGEEMTIVLPQTSLEDTLVKAESLRQAIAS
Q36SI5      MVGDVLATCFRGS-DTVCRYGGEEFSVLMPGASLDEARQRAEQLRAAISA
Q747B7      EAAAVFRGCIRTS-DIAARYGGEEFVVIMPETTRELALLAAEKLRRAVEE
Q2DK38      KTADI I KASLRDM-DIVARYGGEEFCAILPGTSKKESIVVAERIRVGIEK
```

**ClustalW good alignment**

## ... for similar sequences

```
Q2BMK3    MLAKTVREQIQRPADLVARYGGEEFIVVLPDTDEEGAMAVAGQICVAVAS
Q3A8D4    QVARMLQSVVARPGDLVARYGGEEFALILPQTD-HGAKFLGESCRAAVAG
Q36PG9    ALASILSDEVQRSGDLVARYGGEEFAILLPTTDVAGAQQOVAERMRLSVAR
Q2BQL8    TVAQTIKHSIQRAQDMVCRYGGEEFVVILPETDLDG AQMIAERIRKAI AK
Q3XUK3    ALAHTISL-HLRPGDIAARYGGEEFAVVLPDTDAVSGRMI AERLRTA VEA
Q9HXT9    QVAGAIREGCSRSSDLAARYGGEEFAMVLPGTSPGGARLLAEKVRRTVES
P73713    TI GRILQSNIRGS-DIACRYGGEE MTIVLPQTSLED TLVKAESLRQAIAS
Q36SI5    MVGDV L ATCFRGS-DTVCRYGGEEFSVLMPGASLDEARQRAEQLRAAISA
Q747B7    EAAAVFRGCI RTS-DIAARYGGEEFVVIMPETTRELALLAEKLRRAVEE
Q2DK38    KTADI IKASLRDM-DIVARYGGEEFCAILPGTSKKESIVVAERIRVGI EK
```

**ClustalW good alignment**

# Here are distantly related sequences: diguanylate cyclase and adenylate cyclase

## ClustalW alignment

```
1w25 -----NRRYMTGQLDSLVRKRALGGDPVSALL-----IDIDFFKKINDTFGHDIGDEV-----LREFALRLAS
1wc4 -PEPRLITILFSDIVGFTRMSNALQSQGVALLNEYLGE MTRAVFENQGTVDK FVGD AIMALYGAPEEMSPSEQVRRAIATARQ

1w25 NVRAI-DLPCRYGEE-----FVIMPDTALADALRI-AERIRMHVSGSPFTVAHGREML--NVTISIGVSATAGEGD
1wc4 MLVALEKLNQGWQERGLVGRNEVPPVRFRCGIHQGMVVGLFGSQERSDFTAIGPSVNIAARLQEATAPNSIMVSAMVAQYVPD

1w25 TPEALLKRADEGVYQAKASGRNAVVGKAA--
1wc4 E-----EIIKREFLELKGIDEPVMTCVINPN
```

sequence identity = 12%

## DALI alignment based on structural comparison

```
1w25 NRRYMTGQLDSLVRKRALGGDPVSALLI DIDFFKKINDTFGHDIGDEVLRREFALRLASNVRA-IDLP-CRYGGEFVIMPDT-
1wc4 -----PEPR----LITILFSDIVGFTRMSNALQSQGVALLNEYLGE MTRAVFENQGTVDK FVGD AIMALYGAPE

1w25 -----ALADALRIAERIRMHVSG-SPFTVAHGREMLN-----VTISIGVSATAGEGDT-----PEALLKRADEGVYQ
1wc4 EMSPSEQVRRAIATARQMLVALEKLNQGW-QERGLVGRNEVPPVRFRCGIHQGMVVGLFGSQERSDFTAIGPSVNIAARLQEATAPNSIMVSAMVAQYVPD

1w25 AKASGRNAVVGKAA-----
1wc4 TA---PNSIMVSAMVAQYVPDEEIIKREFLELKGIDEPVMTCVINPN
```

sequence identity = 12%

1. Pei and Grishin 2001
2. Steegborn et al. 2005
3. Holm and Sander 1998

Red: alpha-helix blue: beta-strand

# 1. ClustalW alignment

```
1w25 -----NRRYMTGQLDSLVRKRALGGDPVSALL-----IDIDFFKKINDTFGHDIGDEV-----LREFALRLAS
1wc4 -PEPRLITILFSDIVGFTRMSNALQSQGVAELLNEYLGEMTRAVFENQGTVDKFVGDAIMALYGAPEEMSPSEQVRRAIATARQ

1w25 NVRAI-DLPCRYGGEE-----FVVIMPDTALADALRI-AERIRMHVSGSPFTVAHGREML--NVTISIGVSATAGEGD
1wc4 MLVALEKLNQGWQERGLVGRNEVPPVRFRCGIHQMAVVGLFGSQERSDFTAIGPSVNIAARLQEATAPNSIMVSAMVAQYVPD

1w25 TPEALLKRADEGVYQAKASGRNAVVGKAA--
1wc4 E-----EIKREFLEELKGIDEPVMTCVINPN
```

 :  $\alpha$ -helix aligned to  $\beta$ -strand!

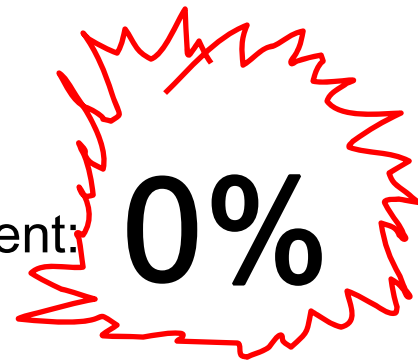
# 2. DALI alignment based on structural comparison

```
1w25 NRRYMTGQLDSLVRKRALGGDPVSALLIDIDFFKKINDTFGHDIGDEVLREFALRLASNVRA-IDLP-CRYGGEFVVIMPPDT-
1wc4 -----PEPR-----LITILFSDIVGFTRMSNALQSQGVAELLNEYLGEMTRAVFENQGTVDKFVGD-DAIMALYGAPE

1w25 -----ALADALRIAERIRMHVSG-SPFTVAHGREMLN-----VTISIGVSATAGEGDT-----PEALLKRADEGVYQ
1wc4 EMSPSEQVRRAIATARQMLVALEKLNQGW-QERGLVGRNEVPPVRFRCGIHQMAVVGLFGSQERSDFTAIGPSVNIAARLQEA

1w25 AKASGRNAVVGKAA-----
1wc4 TA---PNSIMVSAMVAQYVPDEEEIKREFLEELKGIDEPVMTCVINPN
```

Accuracy of the above ClustalW alignment:



# Alignment-based structural superposition

## 1. ClustalW alignment

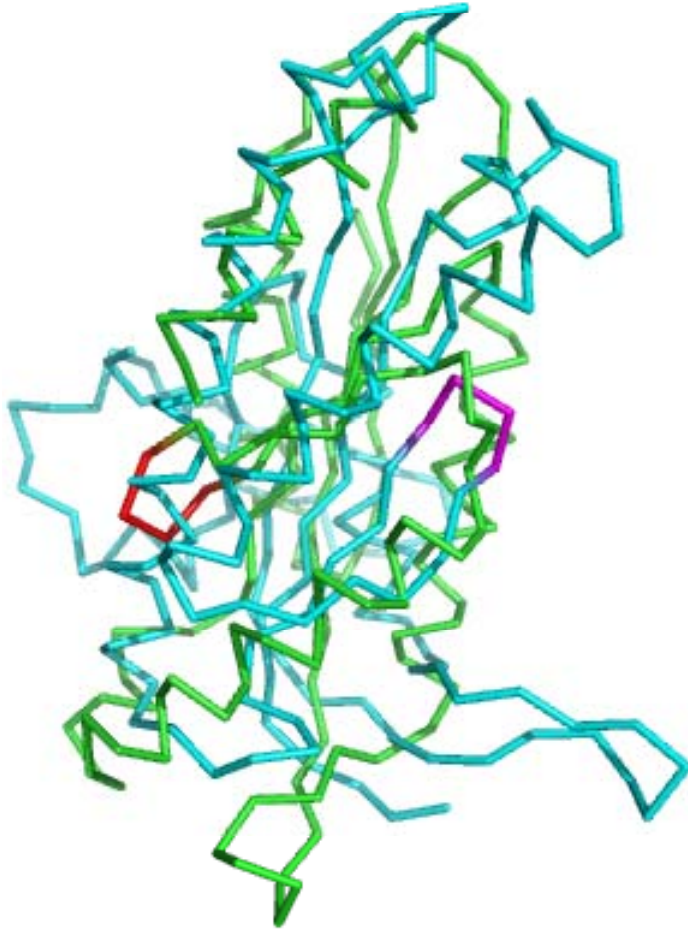
```
1w25  -----NRRYMTGQLDSLVKRATLGGDPVSALL-----IDIDFFKKINDTFGHDIGDEV-----LREFALRLAS
1wc4  -PEPRLITILFSDIVGFTRMSNALQSQGVAELLNEYLGEMTRAVFENQGTVDKFVGDAIMALYGAPEEMSPSEQVRRAIATARQ

1w25  NVRAI-DLPCRYGGEE-----FVVIMPDTALLADALRI-AERIRMHVSGSPFTVAHGREML--NVTISIGVSATAGEGD
1wc4  MLVALEKLNQGWQERGLVGRNEVPPVRFRCGIHQGMAVVGLFGSQERSDFTAIGPSVNIAARLQEATAPNSIMVSAMVAQYVPD

1w25  TPEALLKRADEGVYQAKASGRNAVVGKAA--
1wc4  E-----EIKREFLELKGIDEPVMTCVINPN
```

ClustalW superposition

# Alignment-based structural superposition



ClustalW superposition

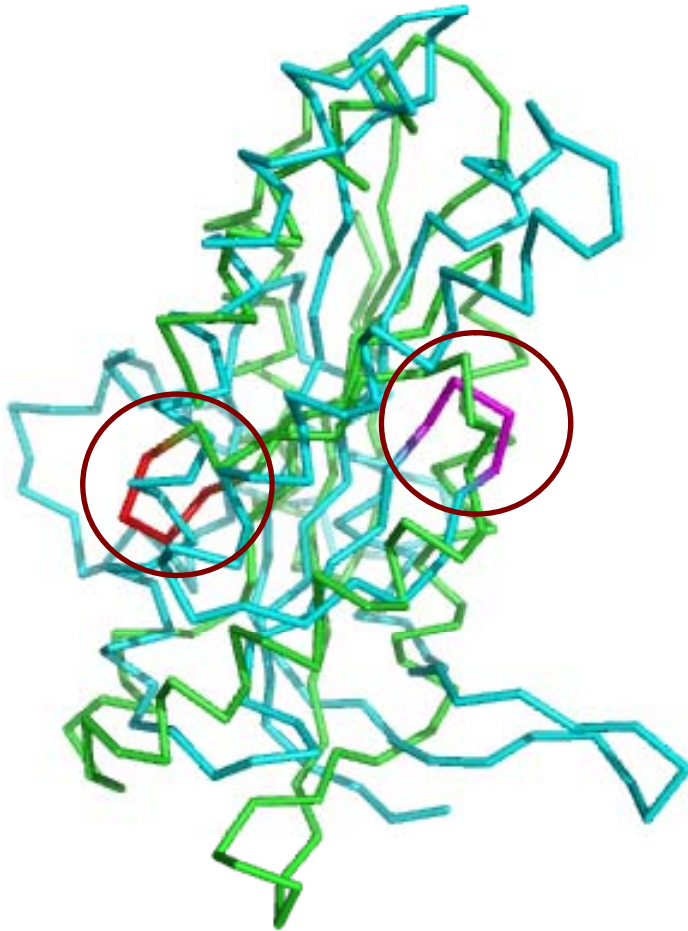
# Alignment-based structural superposition

```
1w25 NRRYMTGQLDSLVKRATLGGGDPVSALLIDIDFFKKINDTFGHDIGDEVLREFALRLASNVRA-IDLP-CRYGGEEFVVIMPDT-  
1wc4 -----PEPR----LITILFSDIVGFTRMSNALQSQGVAELLNEYLGEMTRAVFENQGTVDKFVG-DAIMALYGAPE  
  
1w25 -----ALADALRIAERIRMHVSG-SPFTVAHGREMLN-----VTISIGVSATAGEGDT-----PEALLKRADEGVYQ  
1wc4 EMSPSEQVRRAIIATARQMLVALEEKLNQGW-QERGLVGRNEVPPVRFRCGIHQGMAVVGLFGSQERSDFTAIGPSVNIAARLQEA  
  
1w25 AKASGRNAVVGKAA-----  
1wc4 TA---PNSIMVSAMVAQYVPDEEIKREFLELKGIDEPVMTCVINPN
```

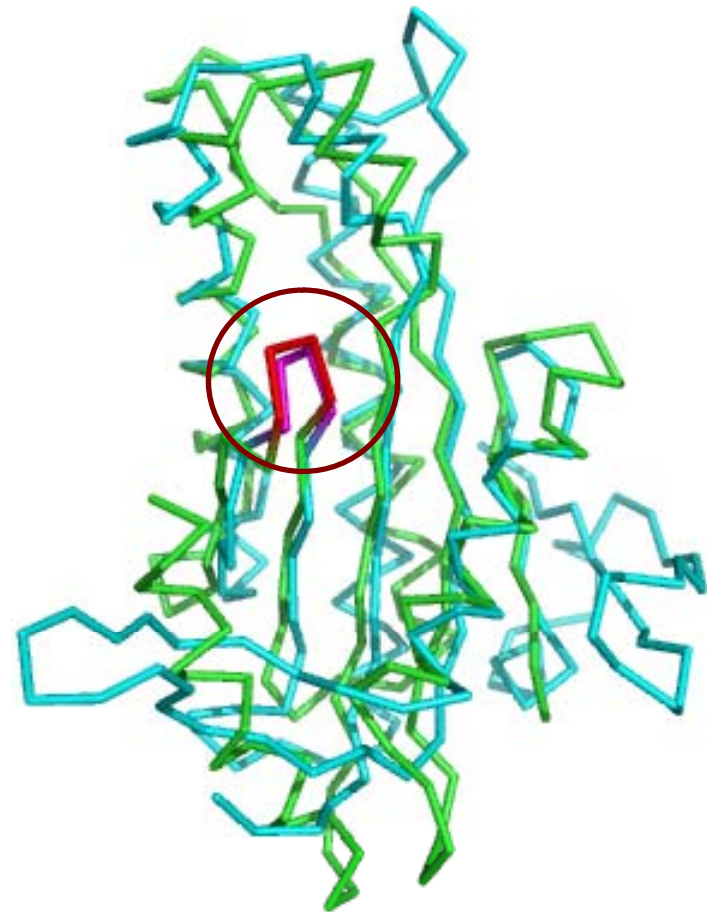
ClustalW superposition

DALI superposition

## Alignment-based structural superposition



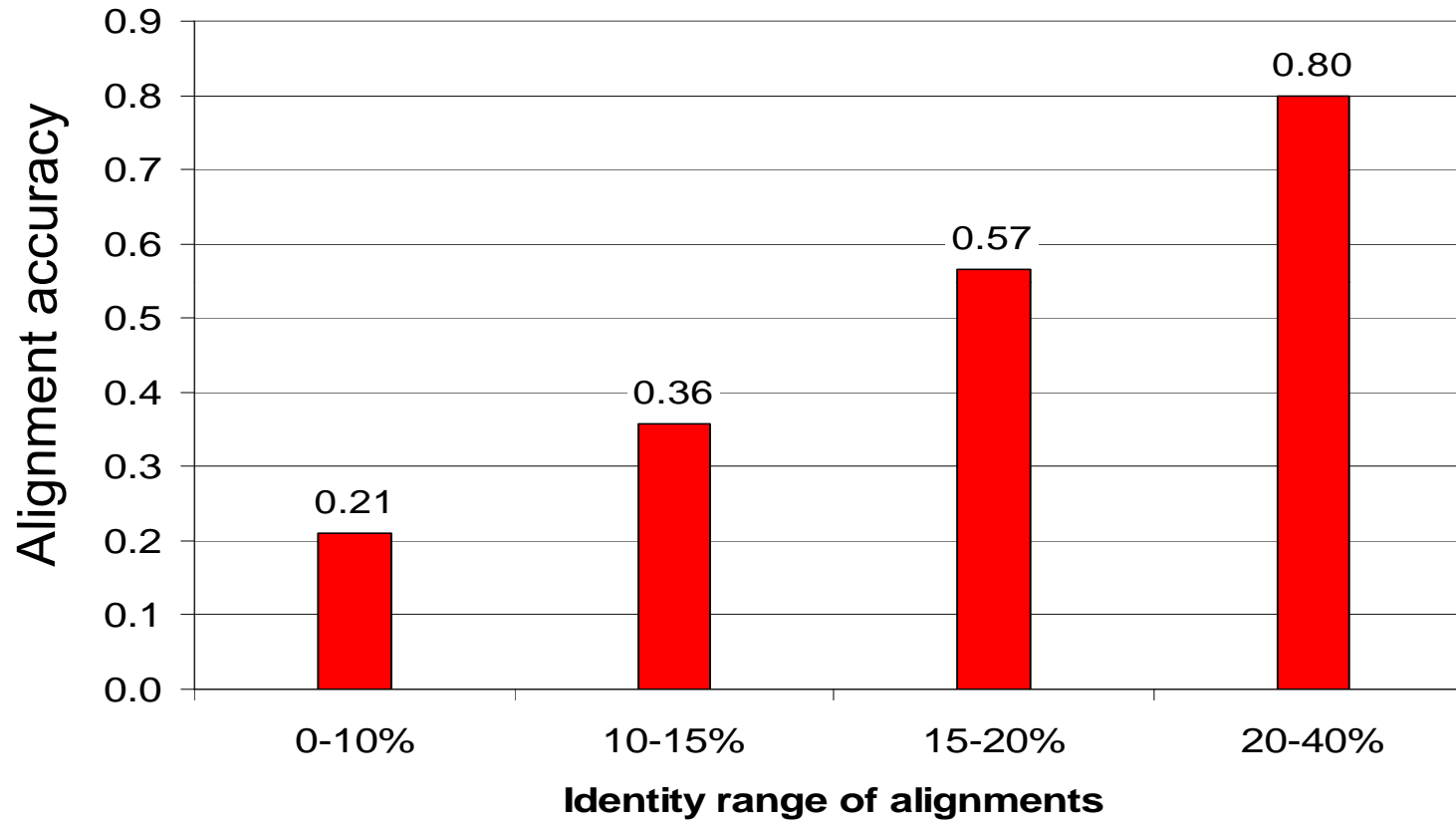
ClustalW superposition



DALI superposition

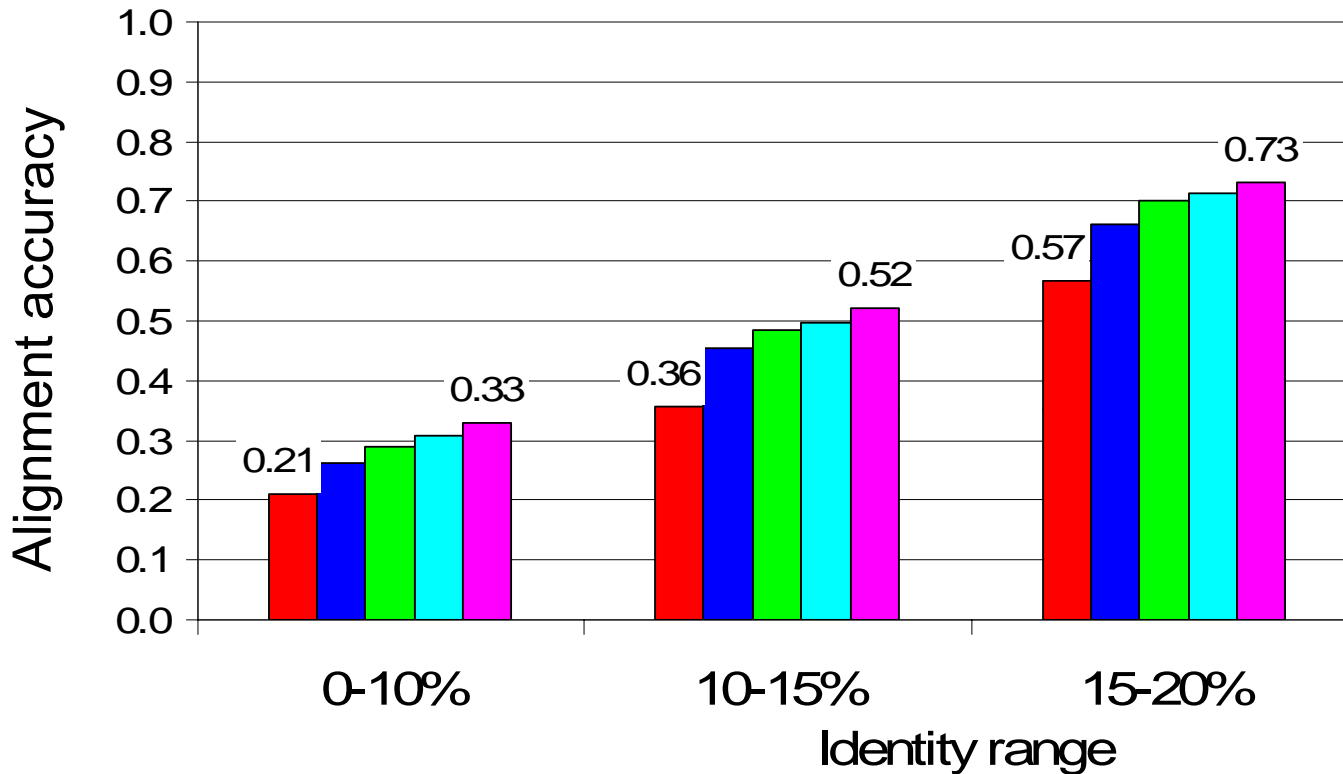


# ClustalW alignment accuracy



Tests on 1785 domain pairs from SCOP (Murzin A. et al. 1995) database.

## What about other methods?



■ ClustalW (Thompson J. et al. 1994)

■ MUSCLE (Edgar R. 2004)

■ ProbCons (Do C. et al. 2005)

■ MAFFT (Kotoh K. et al. 2005)

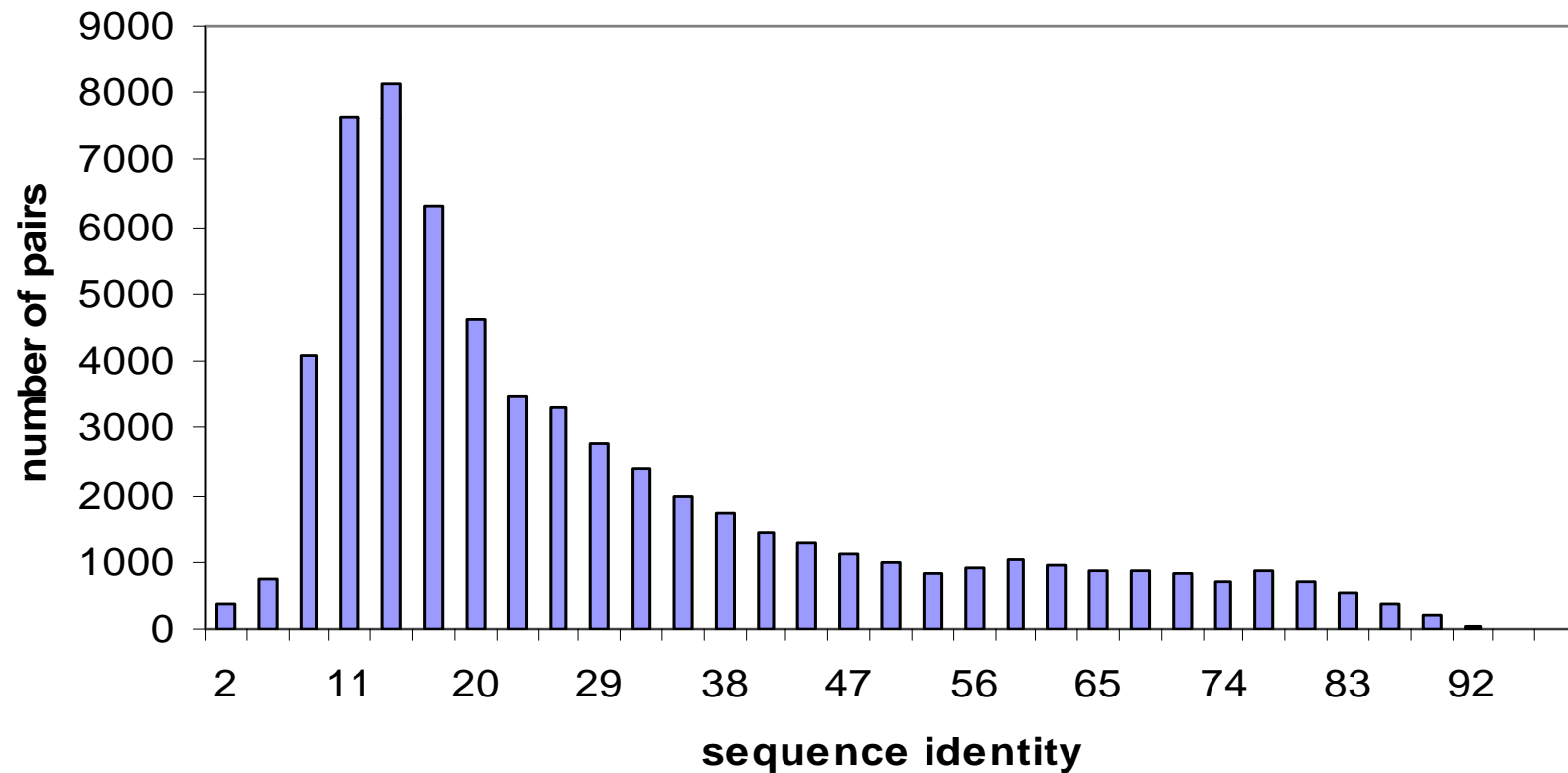
■ MUMMALS (Pei and Grishin 2006)

**Why do we care about remote homologs,**

**i.e. alignments of sequence pairs with  
identity less than 20% ?**

## Why do we care about remote homologs? Reason 1

Sequence identity distribution for proteins with **significant structural similarity** (Dali Z-score >7.0) in FSSP<sup>1</sup> database



1. Holm and Sander, 1996

## Why do we care about remote homologs? Reason 2

### Distant homologs help prediction of functional residues

	Motif 1	Motif 2
<i>H. sapiens</i>	VA <b>H</b> F <b>HH</b> I	LC <b>H</b> SF <b>C</b>
<i>M. musculus</i>	VA <b>H</b> F <b>HH</b> I	LC <b>H</b> SF <b>C</b>
<i>D. melanogaster</i>	VA <b>H</b> L <b>HH</b> I	LV <b>H</b> A <b>F</b> <b>C</b>
<i>S. cerevisiae</i>	LA <b>H</b> A <b>HH</b> A	IL <b>H</b> A <b>L</b> <b>C</b>
<i>S. pombe</i>	MA <b>H</b> I <b>HH</b> T	LV <b>H</b> A <b>F</b> <b>C</b>
	▲ ▲ ▲	▲ ▲
<i>B. halodurans</i>	LV <b>H</b> F <b>RY</b> L	FA <b>H</b> F <b>C</b> <b>I</b>
<i>B. subtilis</i>	AL <b>H</b> F <b>RY</b> L	TA <b>H</b> F <b>I</b> <b>I</b>
<i>A. aeolicus</i>	SA <b>H</b> L <b>AY</b> W	FA <b>H</b> F <b>S</b> <b>A</b>
<i>L. plantarum</i>	LA <b>H</b> L <b>VN</b> I	ML <b>H</b> F <b>L</b> <b>D</b>
<i>L. plantarum</i>	AM <b>H</b> L <b>VN</b> L	SV <b>H</b> W <b>L</b> <b>I</b>
<i>B. anthracis</i>	LF <b>H</b> T <b>SQ</b> -	AI <b>H</b> V <b>L</b> <b>N</b>
<i>P. aeruginosa</i>	AL <b>H</b> L <b>LV</b> N	LL <b>H</b> A <b>S</b> <b>I</b>
<i>V. cholerae</i>	MA <b>H</b> F <b>AG</b> G	GV <b>H</b> F <b>L</b> <b>F</b>

- RCEs are CAAX prenyl proteases identified in **eukaryotes**. (Dolence *et al.* 2000)

- Computational methods identified distant homologs of RCEs in many **bacteria**. (Pei and Grishin, 2001 )

- Recent mutagenesis studies confirmed our predictions. (Plummer *et al.* 2005)

▲ mutations result in complete loss of activity

▲ mutations do not affect enzyme activity

**Our goal (for a few years) has been**

**to improve alignment quality of distantly  
related sequences**

# PROMALS – PROfile Multiple Alignment with predicted Local Structure

PROMALS **input**: unaligned protein sequences

PROMALS **output**: multiple sequence alignment

PROMALS **algorithm** builds alignment of distantly related sequences by utilizing three main sources:

1. Predicted secondary structure
2. Homologous sequences from database searches
3. Complex but reasonable probabilistic models

Source #1: secondary structure

## Secondary structure is more conserved than sequence

DALI structural alignment colored by **real secondary structures**

```
1w25 NRRYMTGQLDSLVKRATLGGDPVSALLIDIDFFKINDTFGHDIGDEVLREFALRLASNVRA-IDLP-CRYGGEEFVVIMPDT-  
1wc4 -----PEPR----LITILFSDIVGFTRMSNALQSQGVAELLNEYLGEMTRAVFENQGTVDKFVG-DAIMALYGAPE  
  
1w25 -----ALADALRIAERIRMHVSG-SPFTVAHGREMLN-----VTISIGVSATAGEGDT-----PEALLKRADEGVYQ  
1wc4 EMSPSEQVRRAIATARQMLVALEKLNQGW-QERGLVGRNEVPPVRFRCGIHQGMAVVGLFGSQERSDFTAIGPSVNIAARLQEA  
  
1w25 AKASGRNAVVGKAA-----  
1wc4 TA---PNSIMVSAMVAQYVPDEEIIKREFLELKGIDEPVMTCVINPN
```

sequence identity = 12%

## Secondary structure prediction is about 80% accurate

DALI structural alignment colored by **PSIPRED<sup>1</sup> predicted secondary structures**

```
1w25 NRRYMTGQLDSLVKRATLGGDPVSALLIDIDFFKINDTFGHDIGDEVLREFALRLASNVRA-IDLP-CRYGGEEFVVIMPDT-  
1wc4 -----PEPR----LITILFSDIVGFTRMSNALQSQGVAELLNEYLGEMTRAVFENQGTVDKFVG-DAIMALYGAPE  
  
1w25 -----ALADALRIAERIRMHVSG-SPFTVAHGREMLN-----VTISIGVSATAGEGDT-----PEALLKRADEGVYQ  
1wc4 EMSPSEQVRRAIATARQMLVALEKLNQGW-QERGLVGRNEVPPVRFRCGIHQGMAVVGLFGSQERSDFTAIGPSVNIAARLQEA  
  
1w25 AKASGRNAVVGKAA-----  
1wc4 TA---PNSIMVSAMVAQYVPDEEIIKREFLELKGIDEPVMTCVINPN
```



Source #2: homologous sequences

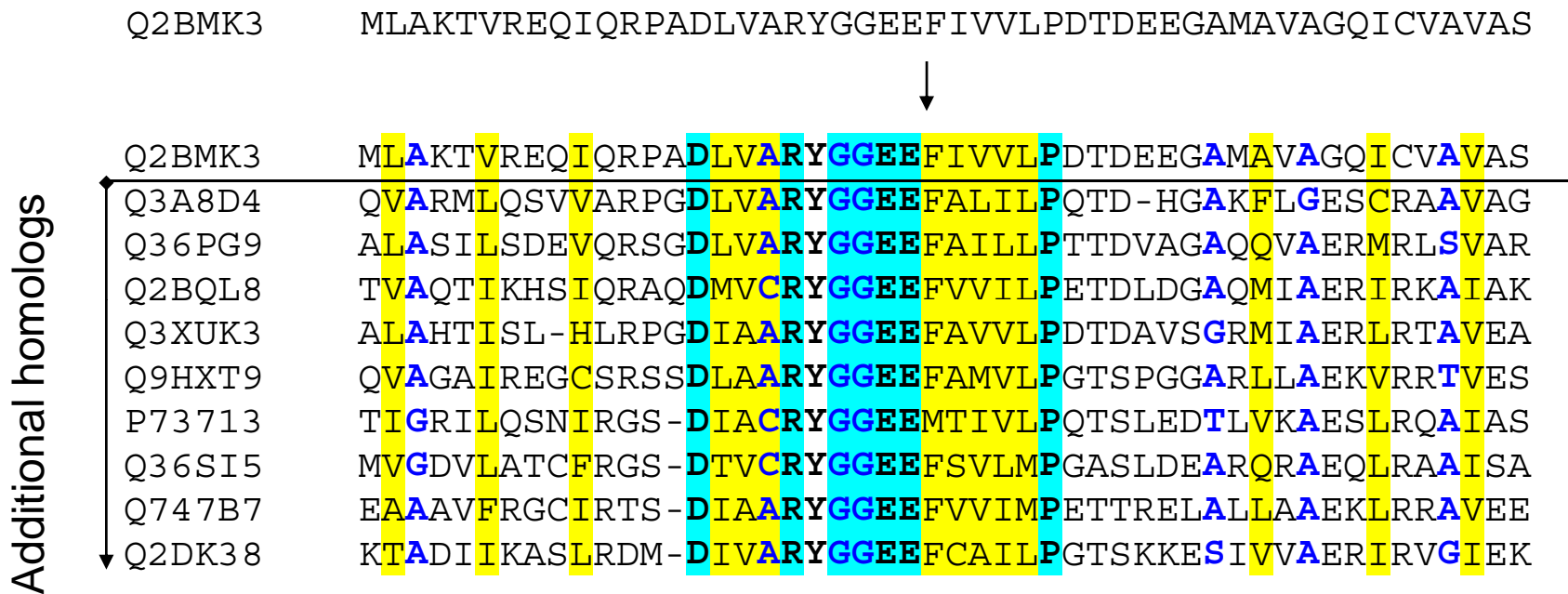
## **More homologs bring up important sequence features through averaging**

Q2BMK3

MLAKTVREQIQRPADLVARYGGEEFIVVLPDTDEEGAMAVAGQICVAVAS

Source #2: homologous sequences

## More homologs bring up important sequence features through averaging



A profile derived from multiple sequence alignment contains position-specific information about:

(1) amino acid usage

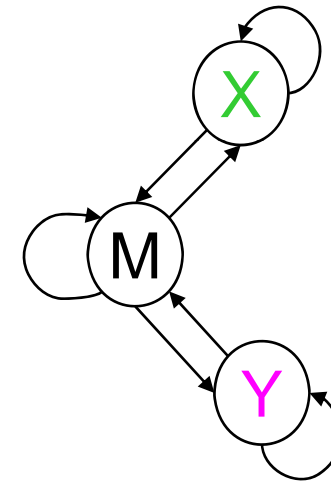
(2) amino acid conservation

**Cyan** : invariant position  
**Yellow** : hydrophobic position  
**Blue** : small residues

# Statistical models of profile-profile alignment

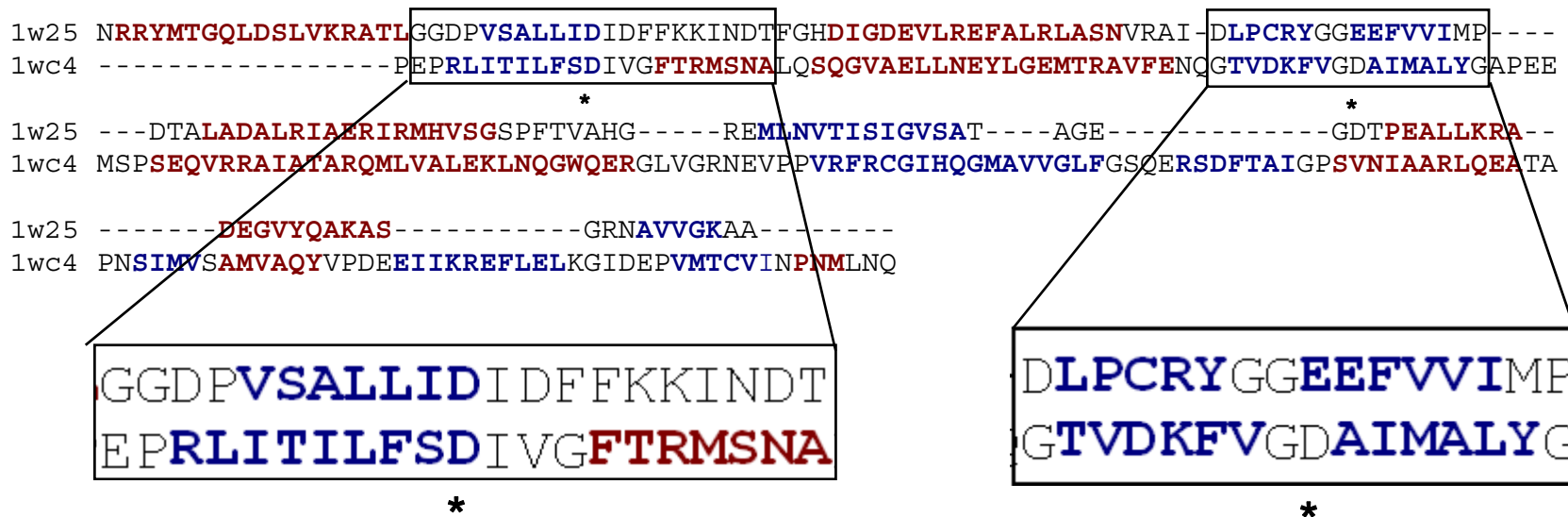
Predicted SS: **hhhhhhhhhhhhc** **cccccccccccccccc**  
 ...  
 Added homologs {  
 LKVISNRL LALVHP-EDAVCRLGGDEFALILNHT  
 LVEIAGR IRSIAKD-DYVLSRSGGDEFVVVVPDC  
 LVEVSERLQRALRQ-TDTVARLGGDEFLLIILDQV  
 LLYIGERVQAAVGE-QGQTFRRGGNEFVVLLPAV  
 LRHVTERLRNFLKQ-SDILCRLSGDEFVVLRVGI  
 LKYVASEI IKNIRK-TDCAVRFGGDEILVAFPDT  
 LKDIARI IRESIRG-TDIAVRIGGDEFLLIILPNS  
 Seq1: **LVRISAATRDVRS-RDIVVRYGGEEFLVLLTHV**  
 Hidden states: MMMMMMMMMMMMMMMYMMMMMMMMMMMMMMMMXX  
 Seq2: **LNEFFRVVDTVGRHGGFVNKFQGDAAALAFG--**  
 Added homologs {  
 LDNHDTIVCHEIQRFGGREVNTAGDGFVATFT--  
 LNELFARFDKLAAE NHCLRILKILGDCYYCVSG--  
 LNSMYSKFDRLT SVHDVYKVETIGDAYMVVVG--  
 LNIYFGKMADVITHGGTIDEFMGDGILVLF--  
 IKTHNDIMRRQLRIYGGYE V KTEGDAFMVAFP--  
 LNEYMSCMVDCEIQTGGVVDKFIGDAIMAIWG--  
 ...  
 Predicted SS: **hhhhhhhhhhhhhhc** **cccccccccccccccc**

## Hidden Markov model

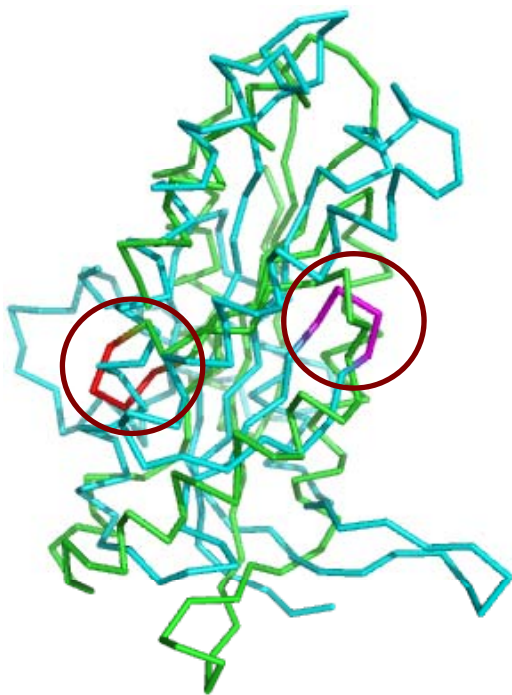


M: emit an aligned position pair  
 X: emit a position in first profile  
 Y: emit a position in second profile

## PROMALS alignment example: diguanylate cyclase and adenylate cyclase



Red: predicted alpha-helix  
 Blue: predicted beta-strand  
 \*: metal-binding residues



ClustalW superposition

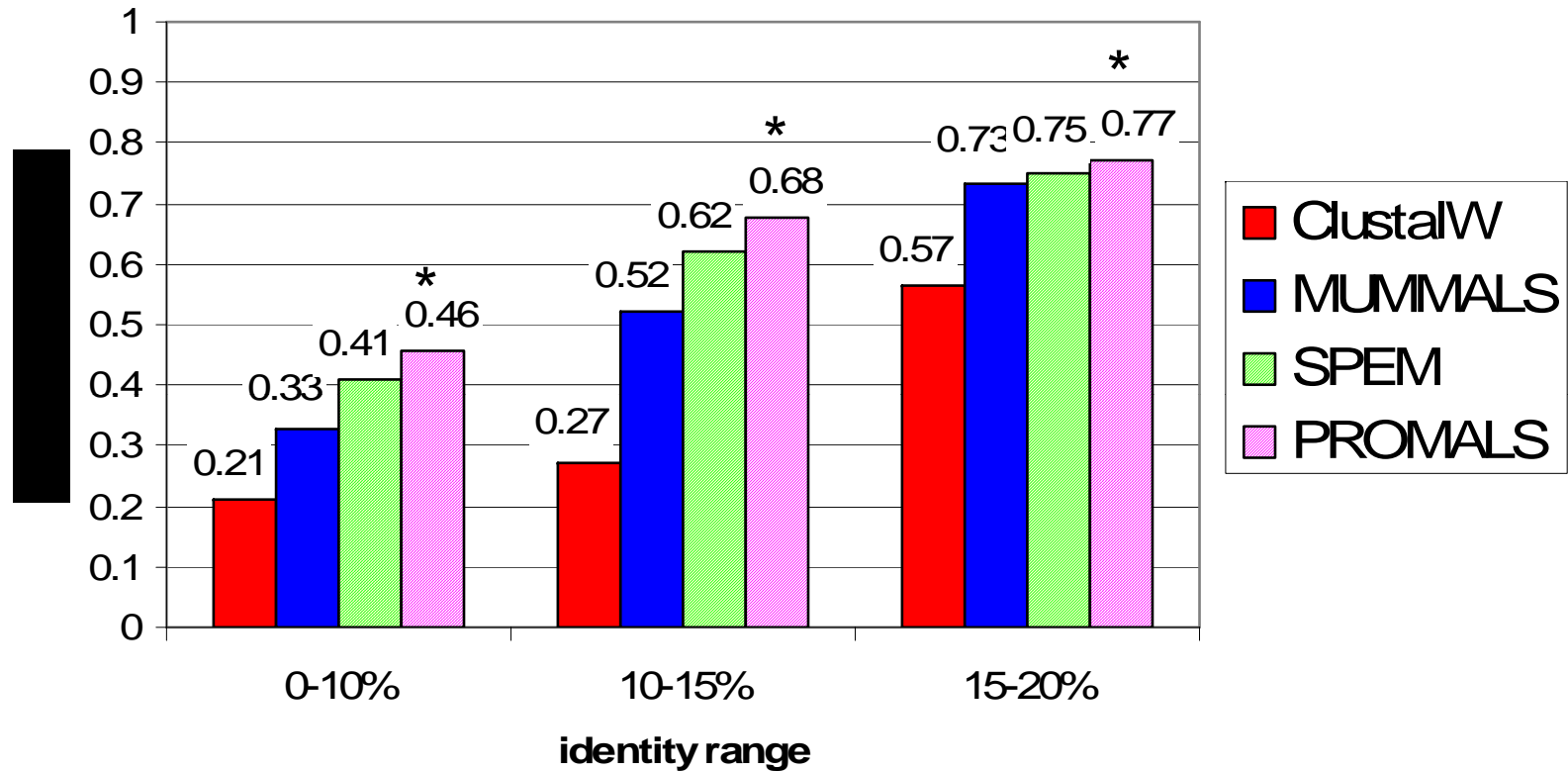


DALI structural  
superposition



**PROMALS** superposition

# Tests on SCOP domain pairs binned by sequence identity

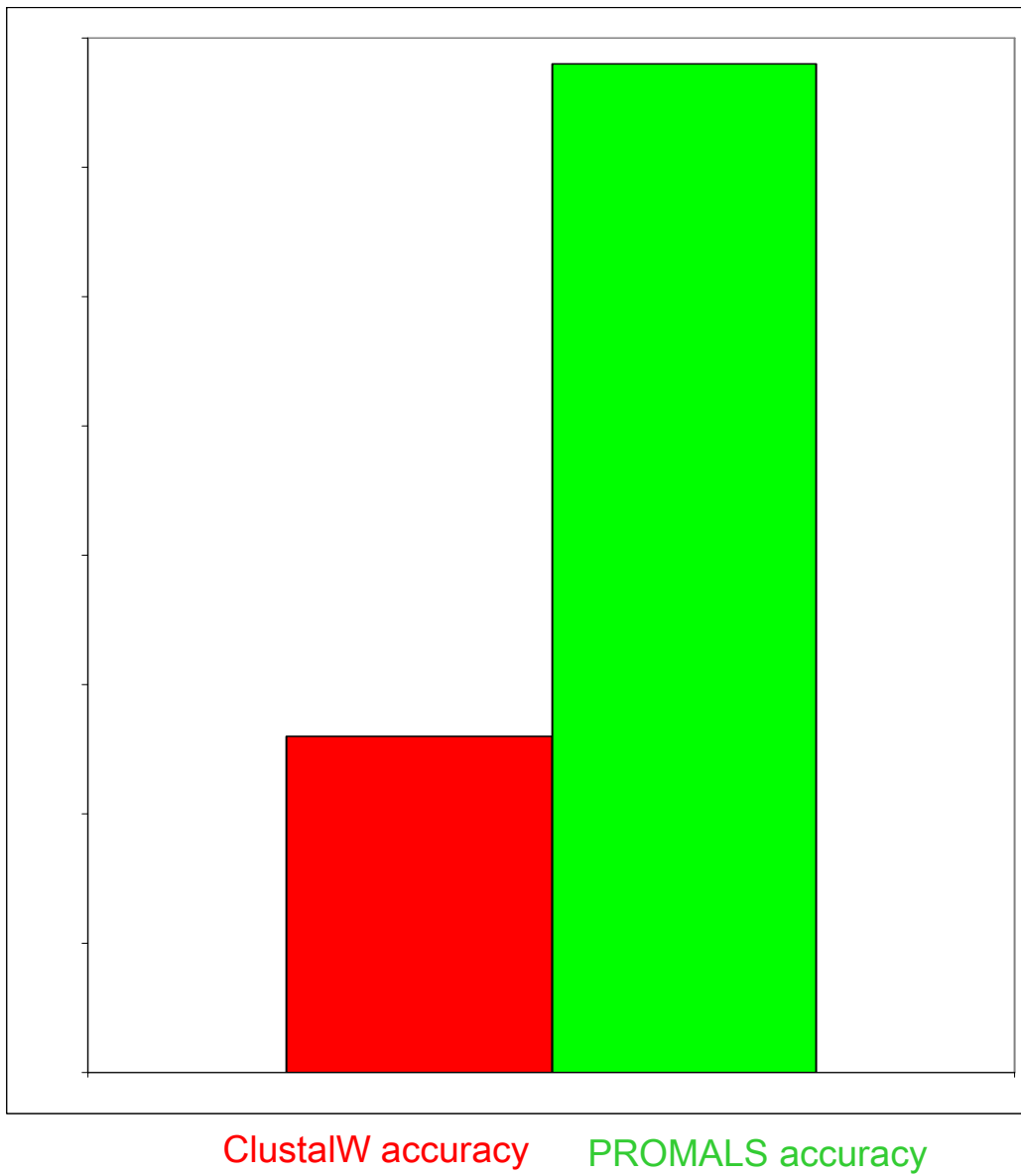


**ClustalW** and **MUMMALS**: methods that do not use additional homologs and predicted secondary structures.

**SPEM** and **PROMALS**: methods that use additional homologs and predicted secondary structures.

\* PROMALS is statistically better than other methods ( $P < 0.0001$ )

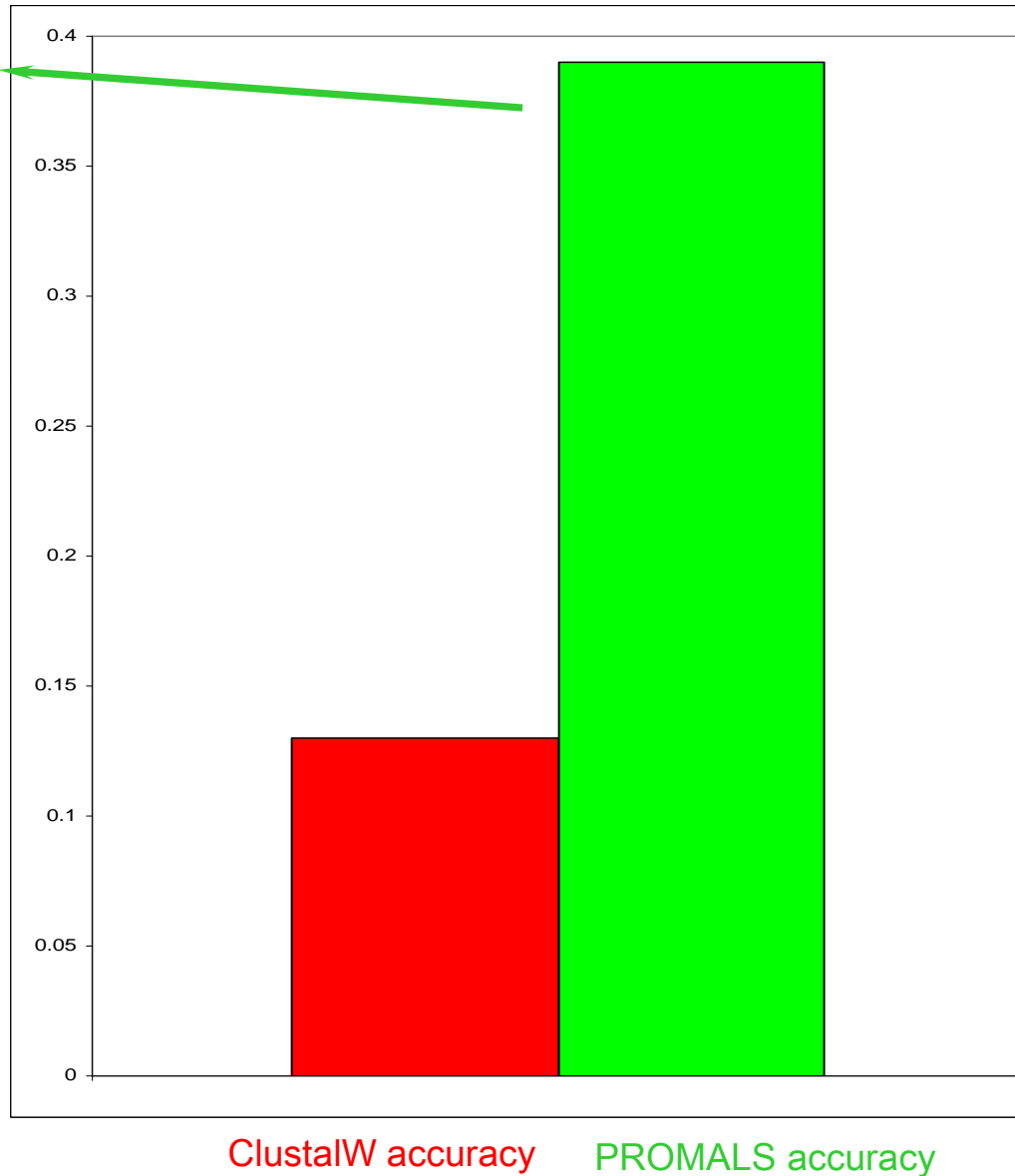
# How accurate are PROMALS alignments?



# How accurate are PROMALS alignments?

40%

Accuracy for  
sequence pairs  
with  
~7% identity

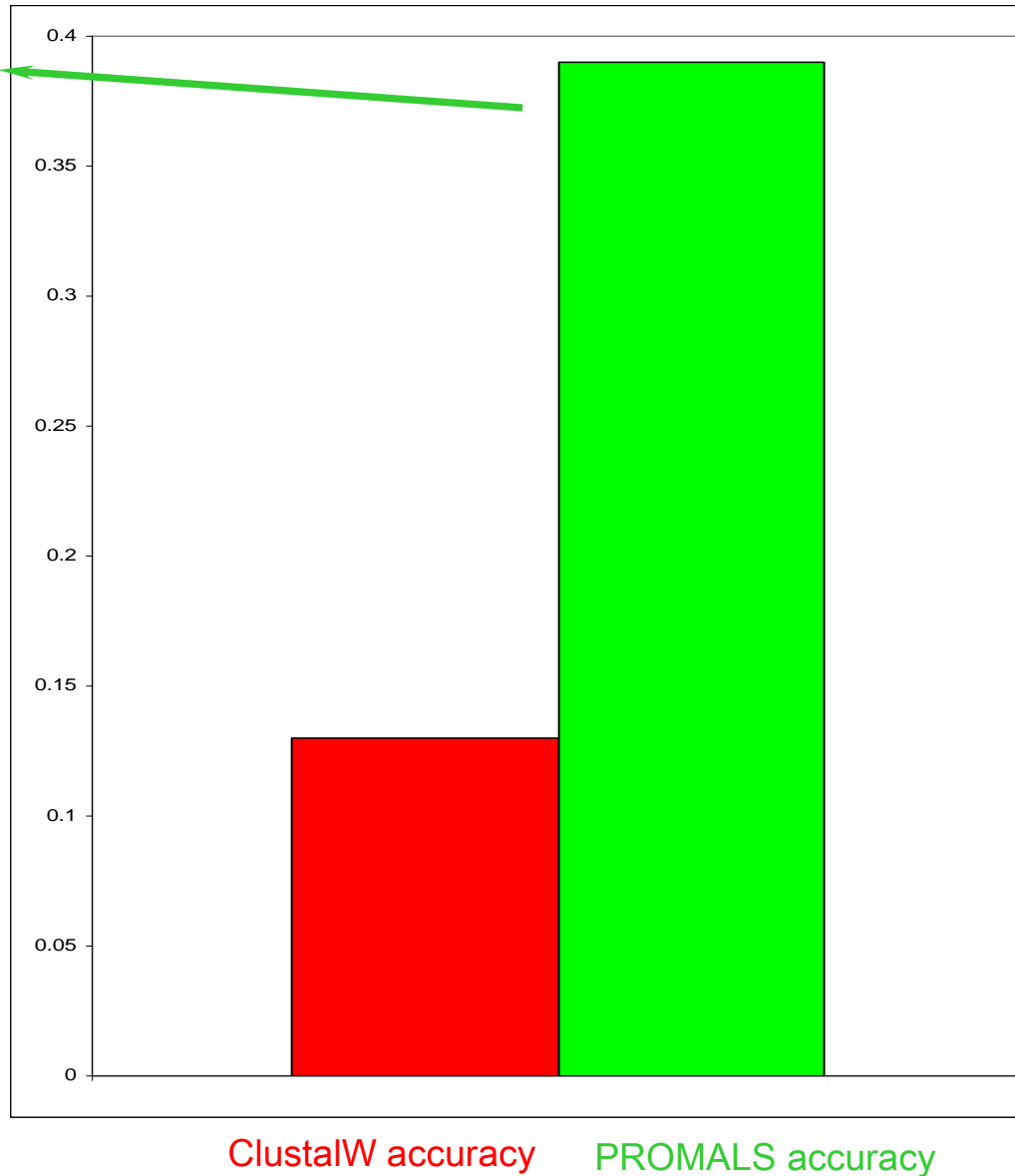




# How accurate are PROMALS alignments?

40%

Accuracy for  
sequence pairs  
with  
~7% identity



<http://prodata.swmed.edu/promals>

## The PROMALS multiple sequence alignment server

PROMALS constructs multiple protein sequence alignments using information from database searches and secondary structure prediction. [\[Documentation\]](#)

Enter your sequences in [FASTA](#) format:

Or upload a local file containing your sequences:

Enter your [email](#) address to receive the result ([recommended](#)):

Alignment options:

- [Weight for amino acid scores](#):
- [Weight for predicted secondary structure scores](#):
- [Identity threshold above which fast alignment is applied](#):

Enter a name for your job ([recommended](#)):

[PROMALS Documentation](#)

[Reference](#): Pei, J. and Grishin, N. V. (submitted). Towards accurate multiple sequence alignments of distantly related proteins.

*Comments, suggestions and bug reports to:* [jpei@chop.swmed.edu](mailto:jpei@chop.swmed.edu)

## What PROMALS does not do

---

It does not use

explicit **3D structure modeling** techniques

It uses only input sequences and internal sequence database

It predicts secondary structure from sequences, but **does not build 3D models**

We have a decent alignment program,

where is the catch?

---

**SPEED (or the lack of it) !**

**ClustalW** takes seconds to minutes per alignment

**PROMALS** takes minutes to hours per alignment:

average is about **30 min per family**,  
some large families take much, much longer

# We have a decent alignment program, what **NOT** to do with it?

---

## **GI-GO** effects:

non-homologous proteins should not be an input

Low complexity proteins should not be an input:

**NQQQQQNNNSSSQQQQQQQQQQSSSTTTTQQQQQQQQQNN**

since the concept of an alignable position that can be traced to a common ancestor does not apply to them

**Membrane proteins** should be used with caution, since their amino acid composition is different, and we still have too few structures of them to test our algorithms thoroughly

# Acknowledgement

## Our group

Lisa Kinch	Erik Nelson
Jimin Pei	Ming Tang
Sara Cheek	Yuan Qi
Shuoyong Shi	Jamie Wrabl
Indraneel M.	Ruslan Sadreyev
Yong Wang	Hua Cheng
Yi Zhong	Bong-Hyun Kim
Wei Cai	Dorothee Staber

## Collaborators

Eugene Koonin	NCBI, NIH
Yuri Wolf	NCBI, NIH
Eugene Shakhnovich	Harvard
Andrei Osterman	Burnham
Leszek Rychlewski	Bioinfobank, Poland

**HHMI, NIH, UTSW,  
The Welch Foundation**

