

## Key-sites (1)

### 1) **National Center for Biotechnology Information (NCBI):**

[www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)

- GQuery
- Pubmed
- Gene

### 2) **European Bioinformatics Institute (EBI):** [www.ebi.ac.uk](http://www.ebi.ac.uk)

- Services

## Key-sites (2)

4) **Nucleic Acid Research** (NAR): [nar.oxfordjournals.org](http://nar.oxfordjournals.org)

Database Issue → Compilation paper (Tables)

Web Server Issue → Editorial ([bioinformatics.ca/links\\_directory](http://bioinformatics.ca/links_directory))

5) **Wikipedia**: [www.wikipedia.org](http://www.wikipedia.org)

6) **Google**: [www.google.com](http://www.google.com)

# Sequence-structure-function relationships

... ACG TCA GTA CAT CCG TAA ...

**Sequenza gene**



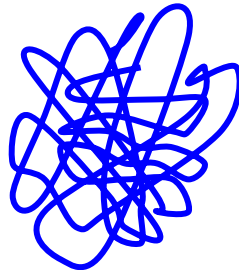
**Codice genetico**

... T S V H P K ...

**Sequenza aminoacidica**



**Protein folding**



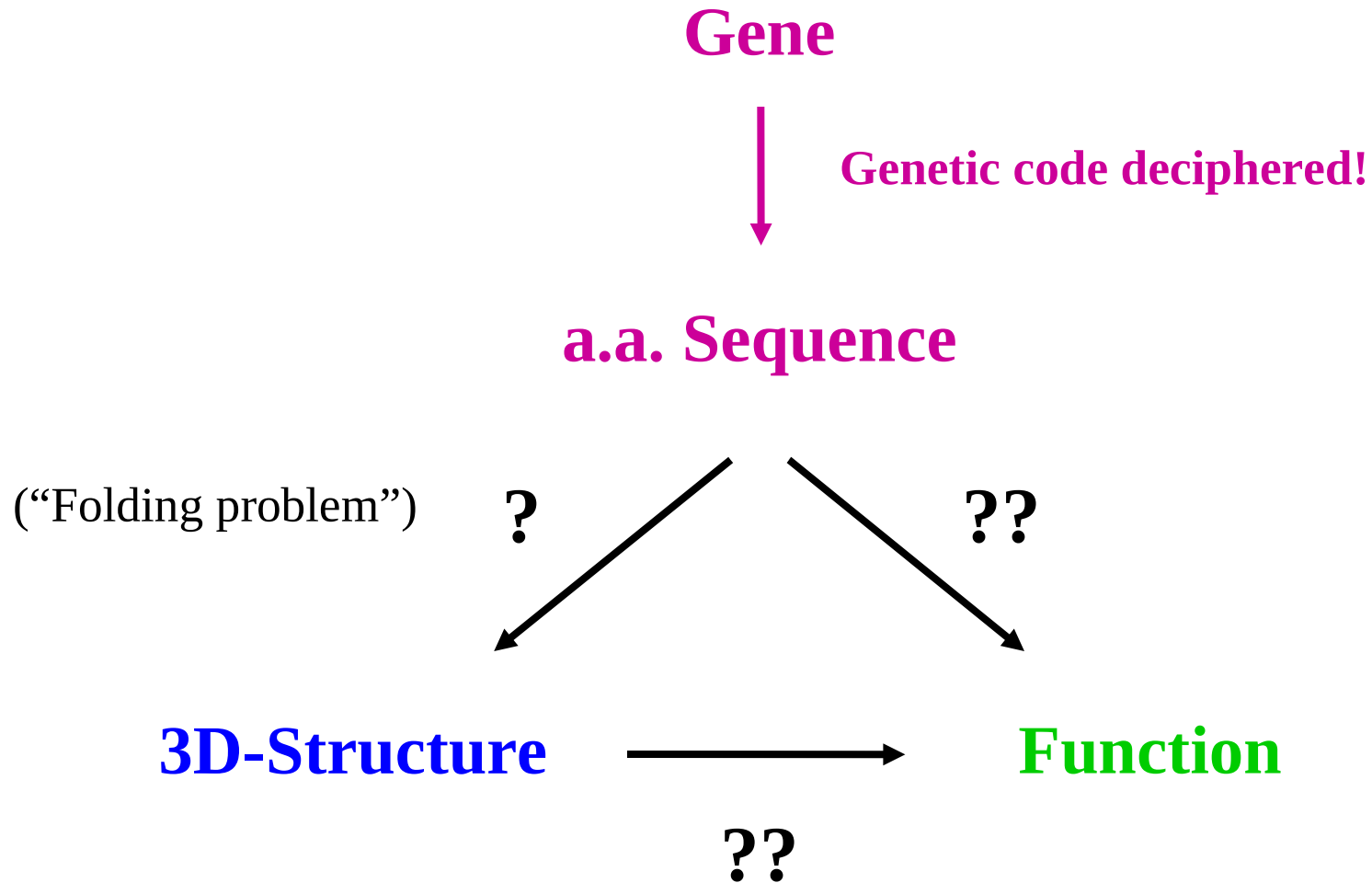
**Struttura proteica**



**ATP → ADP+P**

**Funzione biologica**

# Sequence-structure-function relationships



# Gene-protein sequence

Gene  
(nt sequence)

Genetic code deciphered!

a.a. Sequence  
(primary structure)

Algorithm:  
nt triplet – aa  
correspondence

...  
ACG  
TCA  
GTA  
CAT  
CCG  
TAA  
...

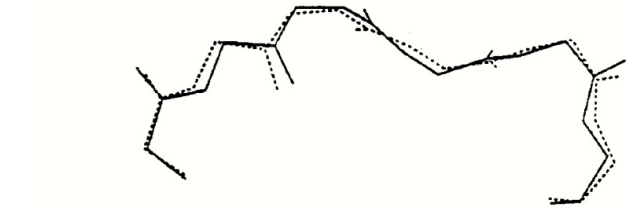
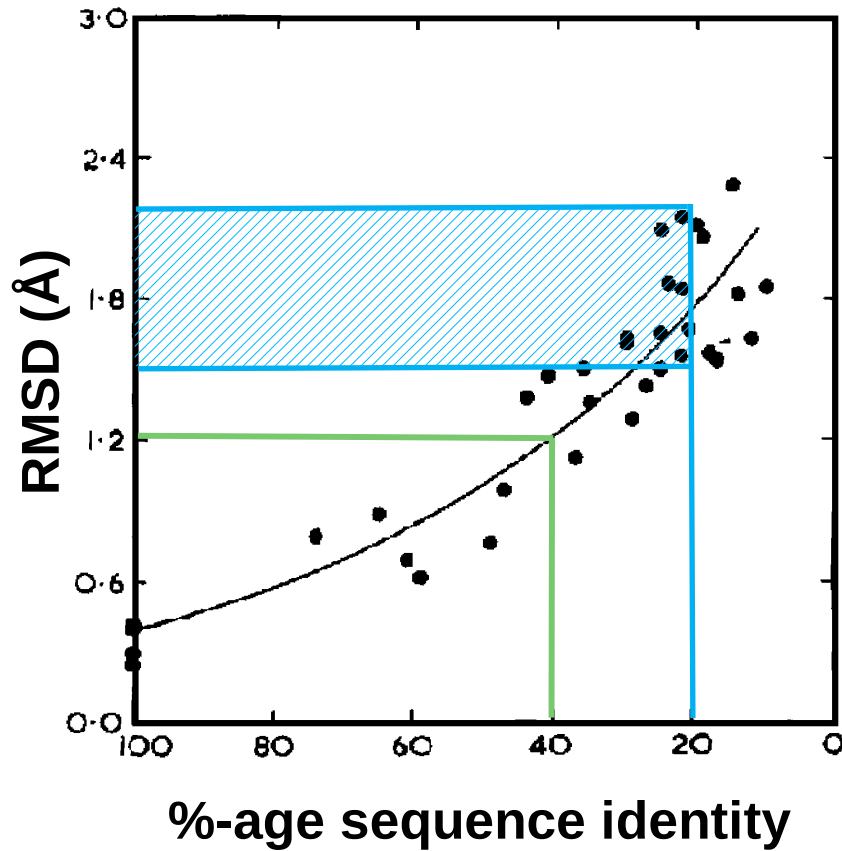
GGG	CGG	AGG	TGG
GGA	CGA	AGA	TGA
GGC	CGC	AGC	TGC
GGT	CGT	AGT	TGT
GCG	CCG	ACG	TCG
GCA	CCA	ACA	TCA
GCC	CCC	ACC	TCC
GCT	CCT	ACT	TCT
GAG	CAG	AAG	TAG
GAA	CAA	AAA	TAA
GAC	CAC	AAC	TAC
GAT	CAT	AAT	TAT
GTG	CTG	ATG	TTG
GTA	CTA	ATA	TTA
GTC	CTC	ATC	TTC
GTT	CTT	ATT	TTT

Gly	Arg	***	Trp
		Ser	Cys
Ala	Pro	Thr	Ser
Glu	Gln	Lys	***
Asp	His	Asn	Tyr
Val	Leu	Met	Leu
		Ile	Phe

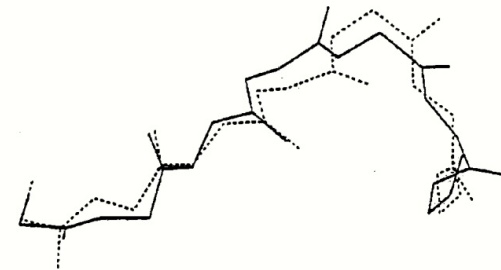
...  
T  
S  
V  
H  
P  
K  
...

# Sequence-structure relationships

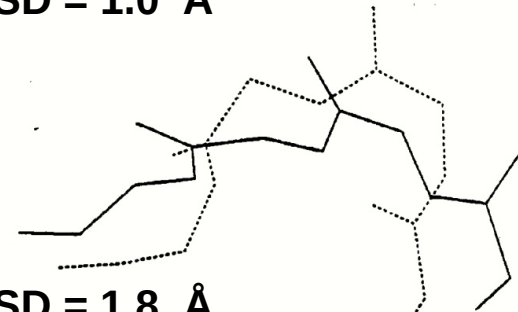
## Structure similarity



RMSD = 0.5 Å



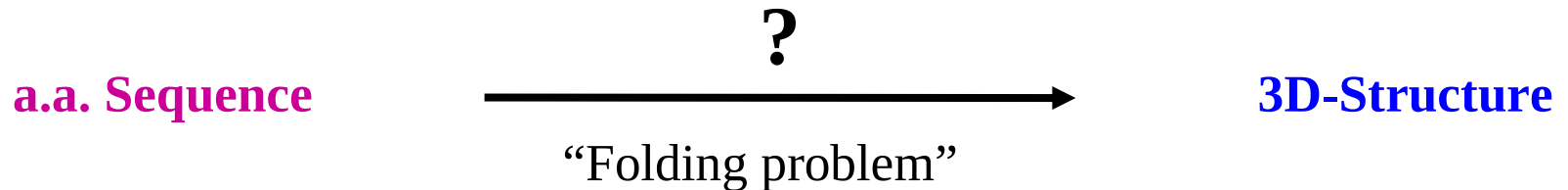
RMSD = 1.0 Å



RMSD = 1.8 Å

Similar sequences  $\longrightarrow$  Similar structures

# Sequence-structure relationships



- i) Structure of protein A is known
- ii) Structure of protein B is unknown
- iii) Proteins A and B have similar sequences => similar structures  
=> Structure of protein B can be modelled on the structure of protein A  
(Homology / Comparative Modelling)

**evolutionary** Relationships between sequences and structures in DBs **+++**

**ab initio** Physical-Chemical Principles **-**

# Homology

**Homology** = Evolutionary relationship = Common ancestor

‘*High*’ or ‘*Low*’ Homology

**No!!!**

‘**Close**’ or ‘**Distant**’ Homology

**Yes**

‘*Measured*’ Homology

**No**

Homology is inferred from measurable parameters:

- %-age of sequence identity
- structure similarity (RMSD)





# Protein sequence analysis

- Name (full name or symbol)
- Sequence identifier (s)
- Name or identifier of homologue in different species
- ...

*e.g., granulocyte colony-stimulating factor*

**National Center for Biotechnology Information (NCBI):**

[www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)

- GQuery
- Pubmed
- Gene

# Protein sequence analysis

Fasta format:

>1-line

AA sequence (1 or more lines)

```
>gi|4503079|ref|NP_000750.1| granulocyte colony-stimulating  
factor isoform a precursor [Homo sapiens]  
MAGPATQSPMKLMALQLLLWHSALWTVQEATPLGPASSLPQSFLLKCLEQVRKIQGDGAALQ  
EKLVSSECATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLSQLHSGLFLYQGLLQA  
LEGISPELGPTLDTLQLDVADFATTIWQQMEELGMAPALQPTQGAMPAFASAFQRRAGGVLV  
ASHLQSFLEVSYRVLRHLAQP
```

3D structure? Biological function?

Residues essential for 3D structure / function?

Save Fasta sequence



# Protein sequence analysis

- 1) NCBI gene annotation
- 2) UniProt annotation
- 3) BLAST

# Sequence Comparison Methods (SCM)

## Blast

1.) Provide the **Query**  
sequence  
in Fasta format

```
>gi|4503079|ref|NP_000750.1| granulocyte  
colony-stimulating factor isoform a precursor  
[Homo sapiens]  
MAGPATQSPMKLMALQLLLWHSALWTVQEATPLGPASSLPQSFLLKC  
LEQVRKIQGDGAALQEKLVSECATYKLCHPEELVLLGHSLGIPWAPL  
SSCPSQALQLAGCLSQLHSGLFLYQGILLQALEGISPELGPTLDTLQL  
DVADFATTIWQQMEELGMAPALQPTQGAMPAFASAFQRRAGGVLVAS  
HLQSFLEVSYRVLRLHLAQP
```

2.) Select the sequence **Database**

```
NR; PDB; RefSeq; SwissProt; Organism; Exclude;  
Specialized
```

# Sequence Comparison Methods (SCM)

## BLAST

### 3.) Algorithm **parameters**:

- Max target sequences
- Expect threshold
- Filter Low complexity regions

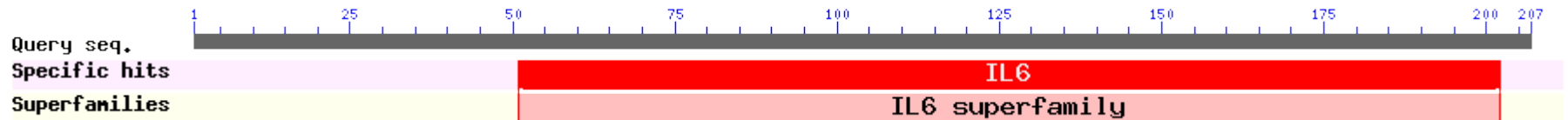
# Sequence Comparison Methods (SCM)

## Program output

```
>gi|4503079|ref|NP_000750.1| granulocyte colony-stimulating factor isoform a precursor  
[Homo sapiens]  
MAGPATQSPMKLMALQLLLWHSALWTVQEATPLGPASSLPQSFLKCLEQVRKIQGDGAALQEKLVS  
ECATYKLCHPEELVLLGHSL  
GIPWAPLSSCPSQALQLAGCLSQLHSGFLYQGGLLQALEGISPELGPTLDTLQLDVADFATTIWQ  
QMEELGMAPALQPTQGAMP  
AFA  
SAFQRRAGGVLVASHLQSFLEVS  
YRVLRL  
HLAQP
```

### 1.) Conserved domains (CDD)

Putative conserved domains have been detected, click on the image below for detailed results.



**Domain-based DBs:**

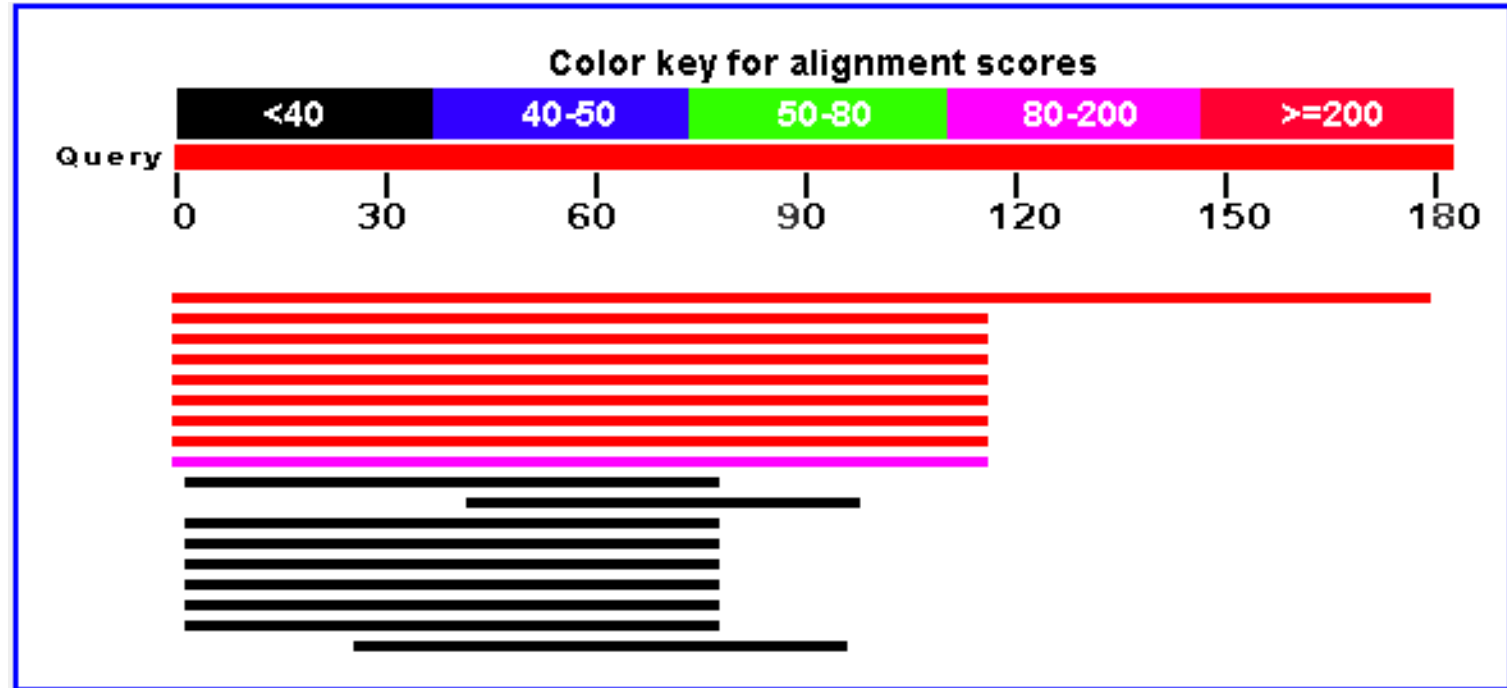
**CDD; Pfam; SMART; Superfamily; etc.**



# Sequence Comparison Methods (SCM)

## 2.) Graphic view of matched sequences

### Distribution of Blast Hits on the Query Sequence



# Sequence Comparison Methods (SCM)

## 3.) List of matched sequences

### ▼ Descriptions

Sequences producing significant alignments:

	Score (Bits)	E Value	
<a href="#">gb ACF41164.1 </a> granulocyte colony stimulating factor [Homo sa...	<a href="#">351</a>	8e-97	
<a href="#">pdb 1CD9 A</a> Chain A, 2:2 Complex Of G-Csf With Its Receptor >p...	<a href="#">206</a>	3e-53	<b>S</b>
<a href="#">pdb 1RHG A</a> Chain A, The Structure Of Granulocyte-Colony-Stimu...	<a href="#">206</a>	4e-53	<b>S</b>
<a href="#">ref NP_757373.1 </a> colony stimulating factor 3 isoform b precu...	<a href="#">205</a>	5e-53	<b>UG</b>
<a href="#">ref NP_757374.1 </a> colony stimulating factor 3 isoform c [Homo ...	<a href="#">205</a>	7e-53	<b>UG</b>
<a href="#">gb AAA03056.1 </a> human granulocyte-colony stimulating factor	<a href="#">204</a>	1e-52	<b>G</b>
<a href="#">pdb 1GNC A</a> Chain A, Structure And Dynamics Of The Human Granu...	<a href="#">201</a>	1e-51	<b>S</b>
<a href="#">ref NP_000750.1 </a> colony stimulating factor 3 isoform a precu...	<a href="#">200</a>	2e-51	<b>UG</b>
<a href="#">dbj BAG60399.1 </a> unnamed protein product [Homo sapiens]	<a href="#">114</a>	2e-25	<b>G</b>
<a href="#">dbj BAG62733.1 </a> unnamed protein product [Homo sapiens]	<a href="#">33.5</a>	0.40	<b>G</b>
<a href="#">dbj BAG62358.1 </a> unnamed protein product [Homo sapiens]	<a href="#">33.1</a>	0.45	<b>G</b>
<a href="#">ref NP_000591.1 </a> interleukin 6 precursor [Homo sapiens] >sp P...	<a href="#">33.1</a>	0.51	<b>UG</b>
<a href="#">pdb 1IL6 A</a> Chain A, Human Interleukin-6, Nmr, Minimized Avera...	<a href="#">33.1</a>	0.52	<b>S</b>
<a href="#">emb CAA27991.1 </a> 26 kDa protein (aa 32-212) [Homo sapiens]	<a href="#">33.1</a>	0.52	<b>G</b>
<a href="#">pdb 1P9M B</a> Chain B, Crystal Structure Of The Hexameric Human ...	<a href="#">32.7</a>	0.55	<b>S</b>
<a href="#">pdb 1ALU A</a> Chain A, Human Interleukin-6	<a href="#">32.7</a>	0.55	<b>S</b>
<a href="#">dbj BAG60317.1 </a> unnamed protein product [Homo sapiens]	<a href="#">32.7</a>	0.67	<b>G</b>
<a href="#">dbj BAG62348.1 </a> unnamed protein product [Homo sapiens]	<a href="#">32.3</a>	0.91	<b>G</b>

**Related  
Structures**

**Gene Info**

**UniGene Info**

# Sequence Comparison Methods (SCM)

## 3.) List of matched sequences

**Description** → pair-wise alignment

**Query cover** → %age of input sequence matched

**E-value** → probability that the matched sequence is not homologous

**Max ident** → % of sequence identity of the longest fragment

**Accession** → page with protein description

# Sequence Comparison Methods (SCM)

## 4.) Alignments of matched sequences to Query (download)

### 4.1) 'Easy' Results: clear homology

#### ▼ Alignments

☐ Select All

[Get selected sequences](#)

[Distance tree of results](#)

[Multiple alignment](#) NEW

> ☐ [gb|ACF41164.1](#) granulocyte colony stimulating factor [Homo sapiens]  
Length=182

Score = 351 bits (900), Expect = 8e-97, Method: Compositional matrix adjust.  
Identities = 179/179 (100%), Positives = 179/179 (100%), Gaps = 0/179 (0%)

Query	1	CLEQVRKIQDDGAALQEKLCATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLS	60
		CLEQVRKIQDDGAALQEKLCATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLS	
Sbjct	1	CLEQVRKIQDDGAALQEKLCATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLS	60
Query	61	QLHSGFLFYQGLLQALEGISPELGPTLDTLQLDVADFATTNWQKNWERNFGGNWAPFVPL	120
		QLHSGFLFYQGLLQALEGISPELGPTLDTLQLDVADFATTNWQKNWERNFGGNWAPFVPL	
Sbjct	61	QLHSGFLFYQGLLQALEGISPELGPTLDTLQLDVADFATTNWQKNWERNFGGNWAPFVPL	120
Query	121	XKAPNHQGGLAQLRPGRHFRXISLVFTQERARGEERXGVXVMSYVCSQKYLKKXIXTL	179
		XKAPNHQGGLAQLRPGRHFRXISLVFTQERARGEERXGVXVMSYVCSQKYLKKXIXTL	
Sbjct	121	XKAPNHQGGLAQLRPGRHFRXISLVFTQERARGEERXGVXVMSYVCSQKYLKKXIXTL	179

# Sequence Comparison Methods (SCM)

## 4.) Alignments of matched sequences to Query

### 4.1) 'Easy' Results: clear homology

```
>[ ]pdb|1RHG|A [S] Chain A, The Structure Of Granulocyte-Colony-Stimulating Factor
And Its Relationship To Those Of Other Growth Factors
pdb|1RHG|B [S] Chain B, The Structure Of Granulocyte-Colony-Stimulating Factor
And Its Relationship To Those Of Other Growth Factors
pdb|1RHG|C [S] Chain C, The Structure Of Granulocyte-Colony-Stimulating Factor
And Its Relationship To Those Of Other Growth Factors
pdb|2D9Q|A [S] Chain A, Crystal Structure Of The Human Gcsf-Receptor Signaling
Complex
Length=174
```

```
Score = 206 bits (523), Expect = 4e-53, Method: Compositional matrix adjust.
Identities = 103/116 (88%), Positives = 104/116 (89%), Gaps = 0/116 (0%)
```

Query	1	CLEQVRKIQDDGAALQEKLCATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLS	60
		CLEQVRKIQ DGAALQEKLCATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLS	
Sbjct	17	CLEQVRKIQGDGAALQEKLCATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLS	76
Query	61	QLHSGFLFYQGLLQALEGISPELGPTLDTLQLDVADFATTNWQKNWERNFGGNWAP	116
		QLHSGFLFYQGLLQALEGISPELGPTLDTLQLDVADFATT WQ+ E P	
Sbjct	77	QLHSGFLFYQGLLQALEGISPELGPTLDTLQLDVADFATTIWQQMEELGMAPALQP	132

# Sequence Comparison Methods (SCM)

## 4.) Alignments of matched sequences to Query

### 4.1) 'Easy' Results: clear homology

```
>[dbj|BAG60399.1|] [G] unnamed protein product [Homo sapiens]  
Length=164
```

```
GENE ID: 1440 CSF3 | colony stimulating factor 3 (granulocyte) [Homo sapiens]  
(Over 100 PubMed links)
```

```
Score = 114 bits (284), Expect = 2e-25 Method: Compositional matrix adjust.  
Identities = 67/116 (57%), Positives = 68/116 (58%), Gaps = 36/116 (31%)
```

```
Query 1 CLEQVRKIQDDGAALQEKLCATYKLCHPEELVLLGHSLGIPWAPLSSCPSQALQLAGCLS 60  
      CLEQVRKIQ DGAALQEKL AGCLS  
Sbjct 43 CLEQVRKIQGDGAALQEKL-----AGCLS 66  
  
Query 61 QLHSGFLFLYQGLLQALEGISPELGPTLDTLQLDVADFATTNWQKNWERNFGGNWAP 116  
      QLHSGFLFLYQGLLQALEGISPELGPTLDTLQLDVADFATT WQ+ E P  
Sbjct 67 QLHSGFLFLYQGLLQALEGISPELGPTLDTLQLDVADFATTIWQQMEELGMAPALQP 122
```

# Sequence Comparison Methods (SCM)

## 4.) Alignments of matched sequences to Query

### 4.2) 'Difficult' Results: homology?

```
>[dbj|BAG62733.1] [G] unnamed protein product [Homo sapiens]  
Length=198
```

```
GENE ID: 3569 IL6 | interleukin 6 (interferon, beta 2) [Homo sapiens]  
(Over 100 PubMed links)
```



```
Score = 33.5 bits (75), Expect = 0.40, Method: Compositional matrix adjust.  
Identities = 19/77 (24%), Positives = 37/77 (48%), Gaps = 1/77 (1%)
```

```
Query   3      EQVRKIQDDGAALQEKLCATYKLCHPPEELVLLGHSLGIP-WAPLSSCPSQALQLAGCLSQ   61  
          +Q+R I D  +AL+++ C      +C   +  L  ++L +P  A      C              CL +  
Sbjct  55      KQIRYILDGISALRKETCNKSNMCESSKEALAENNLNLPKMAEKDGCFCQSGFNEETCLVK  114  
  
Query   62      LHSGLEFLYQGLLQALEG      78  
          + +GL  ++   L+ L+  
Sbjct  115     IITGLLEFEVYLEYLQN      131
```

# Sequence Comparison Methods (SCM)

## 4.) Alignments of matched sequences to Query

### 4.2) 'Difficult' Results: homology?

```
>  pdb|1ALU|A  Chain A, Human Interleukin-6  
Length=186
```

```
Score = 32.7 bits (73), Expect = 0.55, Method: Compositional matrix adjust.  
Identities = 19/77 (24%), Positives = 37/77 (48%), Gaps = 1/77 (1%)
```

```
Query   3      EQVRKIQDDGAALQEKLCATYKLCHPEELVLLGHSLGIP-WAPLSSCPSQALQLAGCLSQ 61  
          +Q+R I D  +AL+++ C      +C   +  L  ++L +P  A      C          CL +  
Sbjct  29      KQIRYILDGISALRKETCNKSNMCESSKEALAENNLNLPKMAEKDGCFQSGFNEETCLVK 88  
  
Query   62      LHSGLFLYQGLLQALEG 78  
          + +GL  ++  L+ L+  
Sbjct  89      IITGLLEFEVYLEYLQN 105
```



# Sequence Comparison Methods (SCM)

## Homologous or Not-homologous?

### ▼ Descriptions

Sequences producing significant alignments:		Score (Bits)	E Value
<a href="#">gb ACF41164.1 </a>	granulocyte colony stimulating factor [Homo sa...	<a href="#">351</a>	<a href="#">8e-97</a>
<a href="#">pdb 1CD9 A</a>	Chain A, 2:2 Complex Of G-Csf With Its Receptor >p...	<a href="#">206</a>	<a href="#">3e-53</a> <b>S</b>
<a href="#">pdb 1RHG A</a>	Chain A, The Structure Of Granulocyte-Colony-Stimu...	<a href="#">206</a>	<a href="#">4e-53</a> <b>S</b>
<a href="#">ref NP_757373.1 </a>	colony stimulating factor 3 isoform b precur...	<a href="#">205</a>	<a href="#">5e-53</a> <b>UG</b>
<a href="#">ref NP_757374.1 </a>	colony stimulating factor 3 isoform c [Homo ...	<a href="#">205</a>	<a href="#">7e-53</a> <b>UG</b>
<a href="#">gb AAA03056.1 </a>	human granulocyte-colony stimulating factor	<a href="#">204</a>	<a href="#">1e-52</a> <b>G</b>
<a href="#">pdb 1GNC A</a>	Chain A, Structure And Dynamics Of The Human Granu...	<a href="#">201</a>	<a href="#">1e-51</a> <b>S</b>
<a href="#">ref NP_000750.1 </a>	colony stimulating factor 3 isoform a precur...	<a href="#">200</a>	<a href="#">2e-51</a> <b>UG</b>
<a href="#">dbj BAG60399.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">114</a>	<a href="#">2e-25</a> <b>G</b>
<a href="#">dbj BAG62733.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">33.5</a>	<a href="#">0.40</a> <b>G</b>
<a href="#">dbj BAG62358.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">33.1</a>	<a href="#">0.45</a> <b>G</b>
<a href="#">ref NP_000591.1 </a>	interleukin 6 precursor [Homo sapiens] >sp P...	<a href="#">33.1</a>	<a href="#">0.51</a> <b>UG</b>
<a href="#">pdb 1IL6 A</a>	Chain A, Human Interleukin-6, Nmr, Minimized Avera...	<a href="#">33.1</a>	<a href="#">0.52</a> <b>S</b>
<a href="#">emb CAA27991.1 </a>	26 kDa protein (aa 32-212) [Homo sapiens]	<a href="#">33.1</a>	<a href="#">0.52</a> <b>G</b>
<a href="#">pdb 1P9M B</a>	Chain B, Crystal Structure Of The Hexameric Human ...	<a href="#">32.7</a>	<a href="#">0.55</a> <b>S</b>
<a href="#">pdb 1ALU A</a>	Chain A, Human Interleukin-6	<a href="#">32.7</a>	<a href="#">0.55</a> <b>S</b>
<a href="#">dbj BAG60317.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">32.7</a>	<a href="#">0.67</a> <b>G</b>
<a href="#">dbj BAG62348.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">32.3</a>	<a href="#">0.91</a> <b>G</b>

# Sequence Comparison Methods (SCM)

## Homologous or Not-homologous?

- 1.) **% Sequence Identity (%\_ID)**
- 2.) **Expect value (E-value)**
- 3.) **Conservation of key-residues**

# Sequence Comparison Methods (SCM)

## Homologous or Not-homologous?

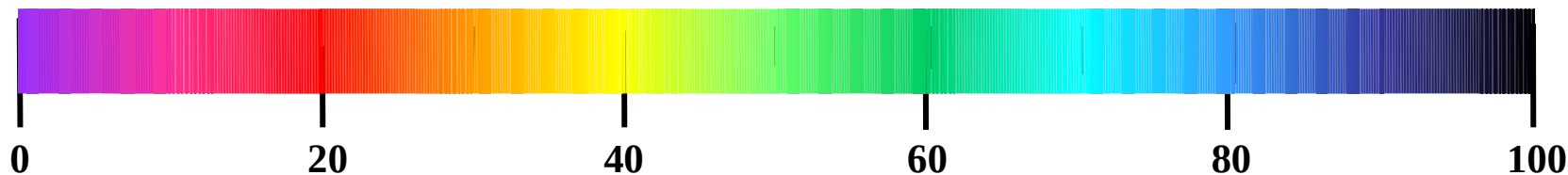
### 1.) % Sequence Identity (%\_ID)

▼ Descriptions			
Sequences producing significant alignments:			
	Score (Bits)	E Value	
→ <a href="#">gb ACF41164.1</a> granulocyte colony stimulating factor [Homo sa...	<a href="#">351</a>	<a href="#">8e-97</a>	100
<a href="#">pdb 1CD9 A</a> Chain A, 2:2 Complex Of G-Csf With Its Receptor >p...	<a href="#">206</a>	<a href="#">3e-53</a> <b>S</b>	
→ <a href="#">pdb 1RHG A</a> Chain A, The Structure Of Granulocyte-Colony-Stimu...	<a href="#">206</a>	<a href="#">4e-53</a> <b>S</b>	88
<a href="#">ref NP_757373.1 </a> colony stimulating factor 3 isoform b precur...	<a href="#">205</a>	<a href="#">5e-53</a> <b>UG</b>	
<a href="#">ref NP_757374.1 </a> colony stimulating factor 3 isoform c [Homo ...	<a href="#">205</a>	<a href="#">7e-53</a> <b>UG</b>	
<a href="#">gb AAA03056.1 </a> human granulocyte-colony stimulating factor	<a href="#">204</a>	<a href="#">1e-52</a> <b>G</b>	
<a href="#">pdb 1GNC A</a> Chain A, Structure And Dynamics Of The Human Granu...	<a href="#">201</a>	<a href="#">1e-51</a> <b>S</b>	
<a href="#">ref NP_000750.1 </a> colony stimulating factor 3 isoform a precur...	<a href="#">200</a>	<a href="#">2e-51</a> <b>UG</b>	
→ <a href="#">dbj BAG60399.1 </a> unnamed protein product [Homo sapiens]	<a href="#">114</a>	<a href="#">2e-25</a> <b>G</b>	57
→ <a href="#">dbj BAG62733.1 </a> unnamed protein product [Homo sapiens]	<a href="#">33.5</a>	<a href="#">0.40</a> <b>G</b>	24
<a href="#">dbj BAG62358.1 </a> unnamed protein product [Homo sapiens]	<a href="#">33.1</a>	<a href="#">0.45</a> <b>G</b>	
<a href="#">ref NP_000591.1 </a> interleukin 6 precursor [Homo sapiens] >sp P...	<a href="#">33.1</a>	<a href="#">0.51</a> <b>UG</b>	
<a href="#">pdb 1IL6 A</a> Chain A, Human Interleukin-6, Nmr, Minimized Avera...	<a href="#">33.1</a>	<a href="#">0.52</a> <b>S</b>	
<a href="#">emb CAA27991.1 </a> 26 kDa protein (aa 32-212) [Homo sapiens]	<a href="#">33.1</a>	<a href="#">0.52</a> <b>G</b>	
<a href="#">pdb 1P9M B</a> Chain B, Crystal Structure Of The Hexameric Human ...	<a href="#">32.7</a>	<a href="#">0.55</a> <b>S</b>	
→ <a href="#">pdb 1ALU A</a> Chain A, Human Interleukin-6	<a href="#">32.7</a>	<a href="#">0.55</a> <b>S</b>	24
<a href="#">dbj BAG60317.1 </a> unnamed protein product [Homo sapiens]	<a href="#">32.7</a>	<a href="#">0.67</a> <b>G</b>	
<a href="#">dbj BAG62348.1 </a> unnamed protein product [Homo sapiens]	<a href="#">32.3</a>	<a href="#">0.91</a> <b>G</b>	

# Sequence Comparison Methods (SCM)

## Homologous or Not-homologous?

### 1.) % Sequence Identity (%\_ID)



**Random**  
(~ 20%)

**‘Twilight’**

**(Close) Homology**  
(> 40%)

Length ~ 100 a.a.

**Homology:**

< Length, > %\_ID (e.g., 10 a.a.)

> Length, < %\_ID

Not ‘low-complexity’

‘low complexity’ can be

‘masked’

# Sequence Comparison Methods (SCM)

## Homologous or Not-homologous?

2.) Expect value (E-value):      Number of matches (with a certain score)

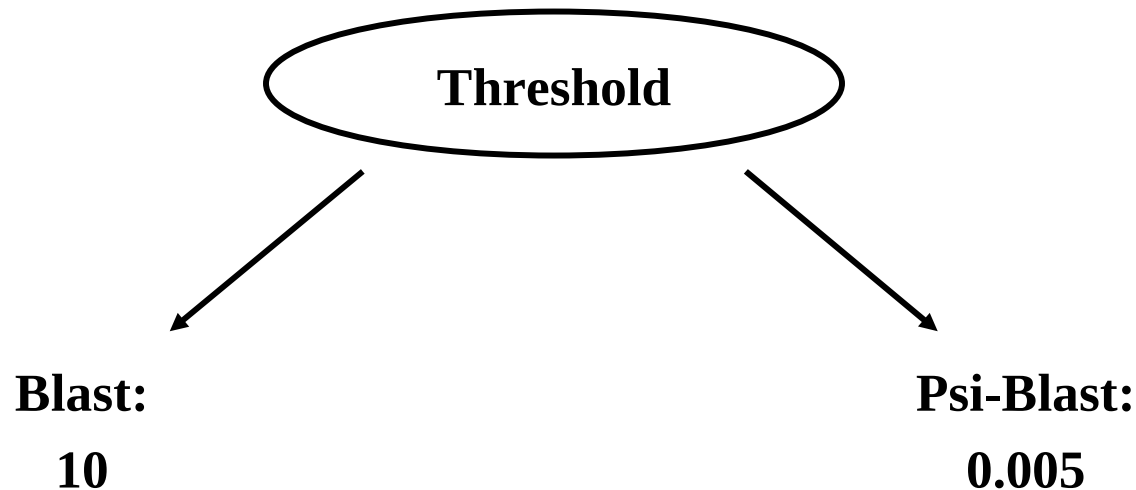
“expected to be found merely by chance”

▼ Descriptions		Score (Bits)	E Value	
Sequences producing significant alignments:				
→ <a href="#">gb ACF41164.1 </a>	granulocyte colony stimulating factor [Homo sa...	<a href="#">351</a>	<a href="#">8e-97</a>	
<a href="#">pdb 1CD9 A</a>	Chain A, 2:2 Complex Of G-Csf With Its Receptor >p...	<a href="#">206</a>	<a href="#">3e-53</a>	<a href="#">S</a>
→ <a href="#">pdb 1RHG A</a>	Chain A, The Structure Of Granulocyte-Colony-Stimu...	<a href="#">206</a>	<a href="#">4e-53</a>	<a href="#">S</a>
<a href="#">ref NP_757373.1 </a>	colony stimulating factor 3 isoform b precur...	<a href="#">205</a>	<a href="#">5e-53</a>	<a href="#">UG</a>
<a href="#">ref NP_757374.1 </a>	colony stimulating factor 3 isoform c [Homo ...	<a href="#">205</a>	<a href="#">7e-53</a>	<a href="#">UG</a>
<a href="#">gb AAA03056.1 </a>	human granulocyte-colony stimulating factor	<a href="#">204</a>	<a href="#">1e-52</a>	<a href="#">G</a>
<a href="#">pdb 1GNC A</a>	Chain A, Structure And Dynamics Of The Human Granu...	<a href="#">201</a>	<a href="#">1e-51</a>	<a href="#">S</a>
<a href="#">ref NP_000750.1 </a>	colony stimulating factor 3 isoform a precur...	<a href="#">200</a>	<a href="#">2e-51</a>	<a href="#">UG</a>
→ <a href="#">dbj BAG60399.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">114</a>	<a href="#">2e-25</a>	<a href="#">G</a>
→ <a href="#">dbj BAG62733.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">33.5</a>	<a href="#">0.40</a>	<a href="#">G</a>
<a href="#">dbj BAG62358.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">33.1</a>	<a href="#">0.45</a>	<a href="#">G</a>
<a href="#">ref NP_000591.1 </a>	interleukin 6 precursor [Homo sapiens] >sp P...	<a href="#">33.1</a>	<a href="#">0.51</a>	<a href="#">UG</a>
<a href="#">pdb 1IL6 A</a>	Chain A, Human Interleukin-6, Nmr, Minimized Avera...	<a href="#">33.1</a>	<a href="#">0.52</a>	<a href="#">S</a>
<a href="#">emb CAA27991.1 </a>	26 kDa protein (aa 32-212) [Homo sapiens]	<a href="#">33.1</a>	<a href="#">0.52</a>	<a href="#">G</a>
<a href="#">pdb 1P9M B</a>	Chain B, Crystal Structure Of The Hexameric Human ...	<a href="#">32.7</a>	<a href="#">0.55</a>	<a href="#">S</a>
→ <a href="#">pdb 1ALU A</a>	Chain A, Human Interleukin-6	<a href="#">32.7</a>	<a href="#">0.55</a>	<a href="#">S</a>
<a href="#">dbj BAG60317.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">32.7</a>	<a href="#">0.67</a>	<a href="#">G</a>
<a href="#">dbj BAG62348.1 </a>	unnamed protein product [Homo sapiens]	<a href="#">32.3</a>	<a href="#">0.91</a>	<a href="#">G</a>

# Sequence Comparison Methods (SCM)

**Homologous or Not-homologous?**

**2.) Expect value (E-value):**      **Number of matches (with a certain score)**  
**“expected to be found merely by chance”**



# Sequence Comparison Methods (SCM)

**Homologous** or **Not-homologous**?

**Threshold problem:**

**For any given  
E-value threshold:**

**Matches above  
(better than)  
threshold**

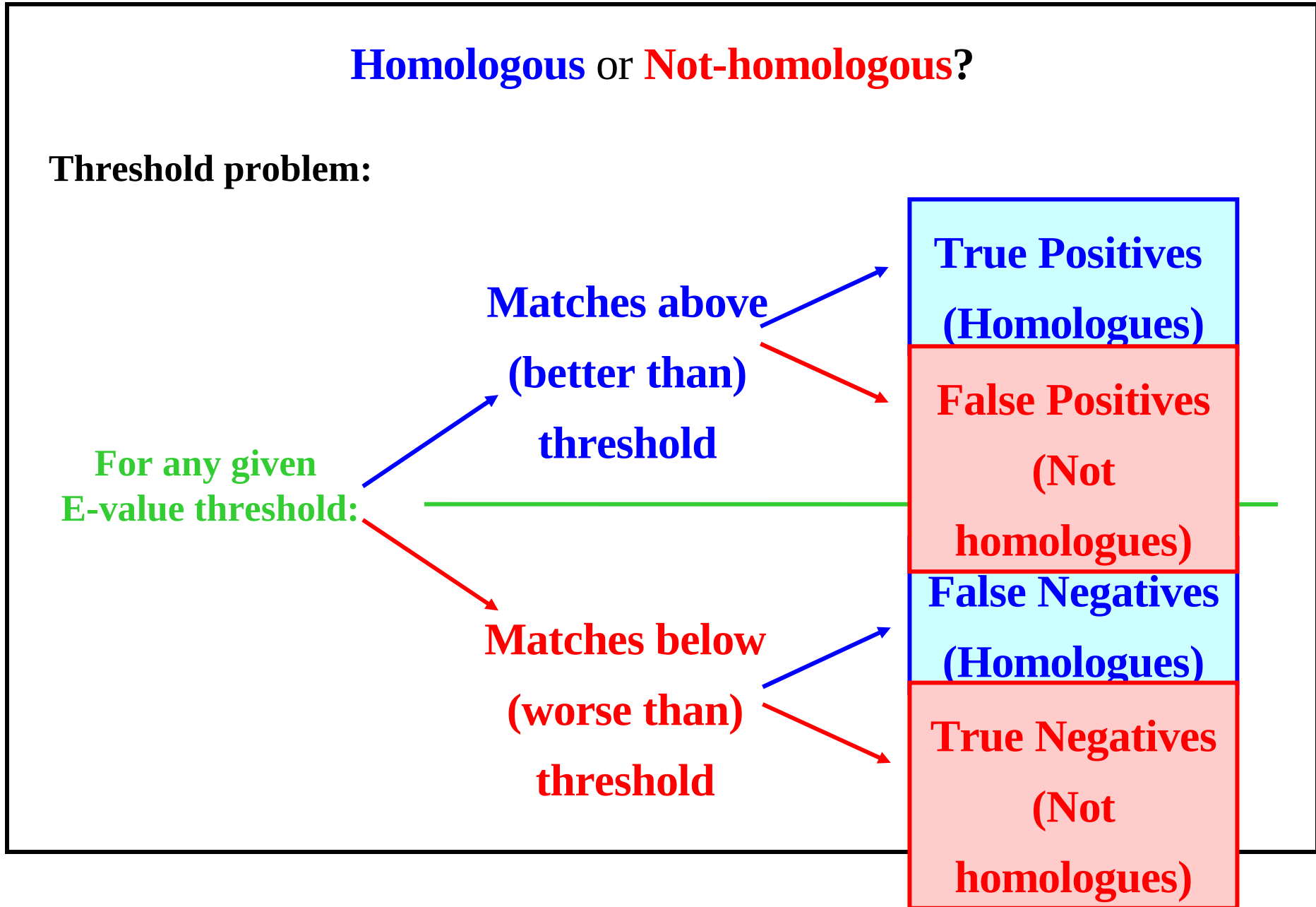
**Matches below  
(worse than)  
threshold**

**True Positives  
(Homologues)**

**False Positives  
(Not  
homologues)**

**False Negatives  
(Homologues)**

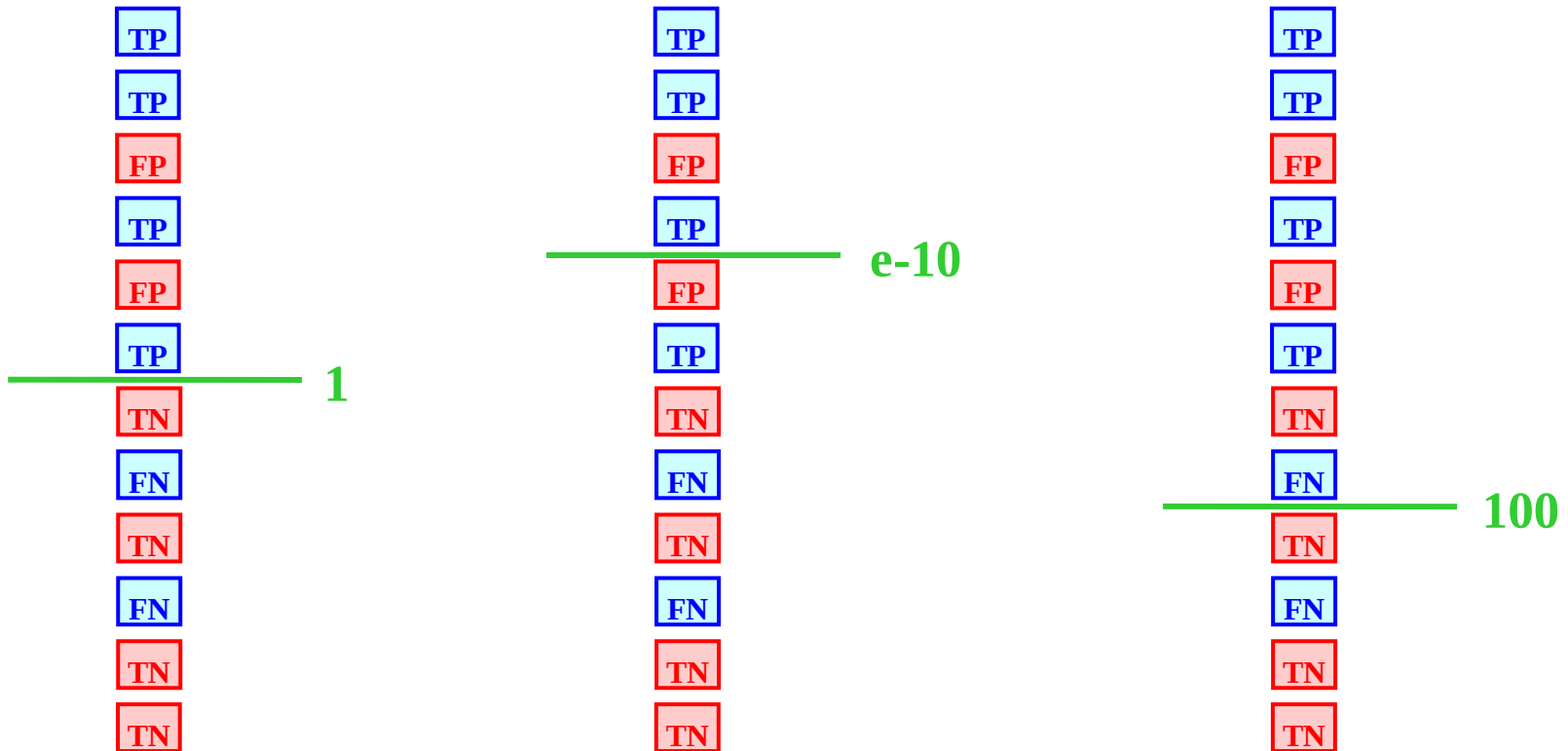
**True Negatives  
(Not  
homologues)**



# Sequence Comparison Methods (SCM)

**Homologous** or **Not-homologous**?

2.) Expect value (E-value):





# Sequence Comparison Methods (SCM)

## Homologous or Not-homologous?

- 1.) % Sequence Identity (%\_ID)
- 2.) Expect value (E-value)
- 3.) **Conservation of key-residues**

# Protein sequence-structure-function relationships

3D-Structures are more conserved than a.a sequences



1alu & 1bgc

1.1 Å RMSD

14/71 = 20 % Seq. ID



1cnt & 1ax8

0.9 Å RMSD

8/71 = 11 % Seq. ID



1hgu & 1lki

1.7 Å RMSD

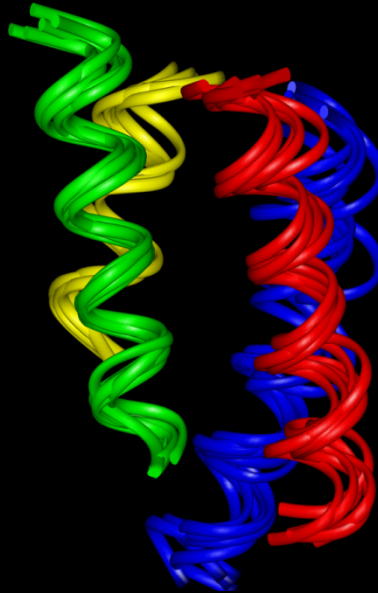
9/71 = 13 % Seq. ID

# Protein sequence-structure-function relationships

Just a small number of residues (“key-residues”) are required to maintain the protein fold



4-alpha-helical cytokines family



common core: 71 a.a.

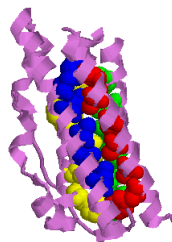


11 highly conserved “key-residues” surrounded by 10 a.a. conserving a generic hydrophobic-neutral character

# Protein sequence-structure-function relationships

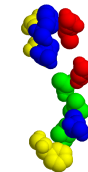
## 4-helical cytokines family

		----- $\alpha$ A-----				----- $\beta$ 1-----					----- $\alpha$ B-----				
		1*	4	11		1	6				1*	3		13	
1rcb (1)	-----hkcdit-	LQE	IIKTLNSL	teqktlcte	-----	LTVTDI	faasknt	-----	TE	KETFCRAATVL	rqfyshhekdt	rc			
leer (1)	----apprlicdsrv-	LER	YLLEAKEA	ekittgcaehcsln	-----	EKITVP	dtkvnfyawkrmevgqqav	-----	EV	WQGLALLSEAV	lrgqallvkssq	p			
lete (1)	----tqdcsfqhspi-	SSD	FAVKIREL	sdyllqd	-----	YPVTVA	snlqddel	-----	c	GG-LWRLVLAQRWM	erlktva	-----			
2gmf (4)	rspspstqpwehvna-	IQE	ARRLLNLS	rdtaaemn	-----	ETVEVI	semfdlqept	-----	CL	QTRLELYKQGL	rgs	-----			
1hul (5)	-----iptsal-	VKE	TLALLSTH	rtllianet	-----	LRIPVP	vhknh	-----	QL	CTEEIFQGI	gtlesqtvqg	---			
3ink (6)	-----stkktqlq-	LEH	LLLDLQMI	lginnyknpkltrmlt	-----	FKFYMP	kkat	-----	E	LKHLQCLEEEL	kpleevlnlaqsk				
1scf (11)	-----	NVKDVTKL	vanlpkd	-----	-----	YMITLK	yvpqmdvlpshc	-----	WI	SEMVVQLSDSL	tdlldkfs-nise				
1rcb (1)	-----997592-	066	00900630	392818009	-----	360211	4549992	-----	86	94300000300	9900692691990				
leer (1)	----99399036595-	059	22930960	49508406970509	-----	590500	3293999409959706000	-----	40	29004502900	8906443998996				
lete (1)	----99904199540-	498	05541950	3993999	-----	860421	50059294	-----	02	00800100930	6904992	-----			
2gmf (4)	995677981990370-	095	08921992	93886449	-----	680700	2992759915	-----	00	43109203901	634	-----			
1hul (5)	-----984531-	096	03620692	696057299	-----	290310	71910	-----	11	005502900	3507891976	---			
3ink (6)	-----98992392-	079	01900640	29007499095499069	-----	890920	9918	-----	8	49106003700	9506910961459				
1scf (11)	-----	79720970	3780799	-----	-----	770606	309129967150	-----	00	02004101710	67029959-9999				



Common Core

Structure similarity





'Key'-residues

Structure & sequence similarity






# Sequence Comparison Methods (SCM)

**Pair-wise sequence comparison methods do not recognize “key-residues” for protein structure/function**

**All positions of the alignment are the same and have the same weight on the computed parameters (i.e., %\_ID, E-value, etc.)**

```
>  pdb|1ALU|A  Chain A, Human Interleukin-6  
Length=186
```

```
Score = 32.7 bits (73), Expect = 0.55, Method: Compositional matrix adjust.  
Identities = 19/77 (24%) Positives = 37/77 (48%), Gaps = 1/77 (1%)
```

```
Query   3      EQVRKIQDDGAALQEKLCATYKLCHPEELVLLGHSLGIP-WAPLSSCPSQALQLAGCLSQ 61  
          +Q+R I D +AL+++ C   +C   + L ++L +P A   C           CL +  
Sbjct  29      KQIRYILDGISALRKETCNKSNMCESSKEALAENNLNLPKMAEKDGCFQSGFNEETCLVK 88  
                
Query   62      LHSGLFLYQGLLQALEG 78  
          + +GL ++ L+ L+  
Sbjct   89      IITGLLEFEVYLEYLQN 105
```

# Sequence Comparison Methods (SCM)

## Multiple sequence alignments (MSA)

- **More informative than pair-wise alignments**
- **Different positions have different conservation**
- **May allow to recognize “key-residues” for protein structure/function**
- **Input sequences:**
  - **Relatively high number**
  - **Similar enough to produce correct alignments (eliminate ‘outliers’, i.e., < 20 %\_ID)**
  - **Different enough to distinguish between conserved and variable positions (make non-redundant, i.e., eliminate > 80 %\_ID)**

# Sequence Comparison Methods (SCM)

## ‘High-quality’ MSA

### Dps proteins

H.pylori	B	1J14	Q	A	D	A	I	V	L	F	M	K	V	H	N	F	H	W	N	V	K	G	T	D	F	F	N	V	H	K	A	T	E	E	I	Y	E	E	F	A	D	M	F	D	D	L	A	E	R	I	V	Q	I	L	E	D	Y	K	Y	L	L	A	K	-	L	Q	K	S	I	W	
H.hepaticus	B		Q	A	D	A	A	V	F	Y	V	K	V	H	N	F	H	W	N	V	K	G	M	D	F	Y	P	T	H	K	A	T	E	E	I	Y	E	K	Y	A	D	V	F	D	D	V	A	E	R	V	L	Q	I	L	S	D	Y	E	Y	F	V	G	E	-	L	Q	K	A	I	W	
V.cholerae	B	3IQ1	L	A	N	Y	Q	V	F	Y	M	N	T	R	G	Y	H	W	N	I	Q	G	K	E	F	F	E	L	H	A	K	F	E	E	I	Y	T	D	L	Q	L	K	I	D	E	L	A	E	R	I	L	T	L	V	D	G	F	S	I	L	I	R	E	-	Q	E	K	L	V	W	
S.degradans	B		L	A	D	S	Y	V	L	Y	L	K	T	H	N	F	H	W	N	V	T	G	P	M	F	Q	T	L	H	N	M	F	M	D	Q	Y	T	E	A	W	T	A	L	D	T	I	A	E	R	I	R	T	L	L	E	G	Q	E	T	L	I	E	V	-	H	E	K	N	A	W	
L.pneumophila	B		L	A	D	T	Y	A	L	Y	L	K	T	Q	N	Y	H	W	H	V	T	G	P	Q	F	K	S	L	H	E	L	F	F	E	M	Q	Y	K	E	L	A	E	A	V	D	Q	I	A	E	R	I	R	I	L	A	K	D	N	M	M	I	V	A	A	-	H	E	K	A	H	W
B.anthraxis	B	1JIG	V	A	N	W	N	V	L	Y	V	K	L	H	N	Y	H	W	Y	V	T	G	P	H	F	F	T	L	H	E	K	F	E	E	F	Y	N	E	A	G	T	Y	I	D	E	L	A	E	R	I	L	A	L	V	N	D	Y	S	A	L	H	T	T	-	L	E	Q	H	V	W	
B.anthraxis	B	1J15	V	A	D	W	S	V	L	F	T	K	L	H	N	F	H	W	Y	V	K	G	P	Q	F	F	T	L	H	E	K	F	E	E	L	Y	T	E	S	A	H	I	D	E	I	A	E	R	I	L	A	I	M	K	D	Y	E	M	M	Y	T	E	-	L	E	K	H	A	W		
S.aureus	B	2D5K	V	A	N	W	T	V	A	Y	T	K	L	H	N	F	H	W	Y	V	K	G	P	N	F	F	S	L	H	V	K	F	E	E	L	Y	N	E	A	S	Q	Y	V	D	E	L	A	E	R	I	L	A	L	S	Q	D	F	T	N	I	Q	T	S	-	V	D	K	H	N	W	
S.epidermidis	B		V	A	N	W	T	V	A	Y	T	K	L	H	N	F	H	W	Y	V	K	G	P	N	F	F	S	L	H	T	K	F	E	E	L	Y	N	E	A	S	Q	Y	V	D	D	L	A	E	R	I	L	A	L	S	K	D	F	S	K	I	Q	T	S	-	V	D	K	H	N	W	
B.subtilis	B	2CHP	L	S	N	W	F	L	L	Y	S	K	L	H	R	F	H	W	Y	V	K	G	P	H	F	F	T	L	H	E	K	F	E	E	L	Y	D	H	A	A	E	T	V	D	T	I	A	E	R	L	L	A	L	V	N	D	Y	K	Q	I	I	E	E	-	V	E	K	Q	V	W	
S.pyogenes	B	2WLA	V	A	D	L	S	V	A	A	S	I	V	H	Q	V	H	W	Y	M	R	G	P	G	F	L	Y	L	H	P	K	M	D	E	L	L	D	S	L	N	A	N	L	D	E	M	S	E	R	L	I	T	L	V	E	V	Y	L	Y	L	K	T	E	-	A	E	K	T	I	W	
L.monocytogenes	B	2IY4	V	A	N	L	N	V	F	T	V	K	I	H	Q	I	H	W	Y	M	R	G	H	N	F	F	T	L	H	E	K	M	D	D	L	Y	S	E	F	G	E	Q	M	D	E	V	A	E	R	L	L	A	L	V	G	T	L	E	L	L	K	A	S	-	I	D	K	H	I	W	
O.oeni	B		I	A	D	I	S	Q	L	K	V	N	V	Q	T	H	W	Y	M	R	G	E	N	F	F	R	L	H	P	L	M	D	E	Y	G	D	Q	L	S	E	Q	L	D	Q	I	A	E	R	L	I	A	L	V	D	Q	F	K	Y	L	K	D	E	-	T	D	K	N	I	W		
E.coli	B	1F33	V	I	Q	F	I	D	L	S	L	I	T	K	Q	A	H	W	N	M	R	G	A	N	F	I	A	V	H	E	M	L	D	G	F	R	T	A	L	I	D	H	L	D	T	M	A	E	R	A	V	Q	L	A	D	R	Y	A	I	V	S	R	D	-	L	D	K	F	L	W	
S.enterica	B		V	I	Q	F	I	D	L	S	L	I	T	K	Q	A	H	W	N	M	R	G	A	N	F	I	A	V	H	E	M	L	D	G	F	R	T	A	L	T	D	H	L	D	T	M	A	E	R	A	V	Q	L	A	D	R	Y	A	V	V	S	R	D	-	L	D	K	F	L	W	
B.melitensis	B	3GE4	L	A	A	T	I	D	L	A	L	I	T	K	Q	A	H	W	N	L	K	G	P	Q	F	I	A	V	H	E	M	L	D	G	F	R	A	E	L	D	H	V	D	T	I	A	E	R	A	V	Q	L	I	E	R	Y	G	D	V	S	R	S	-	L	D	K	A	L	W		

### Bacterioferritins

S. enterica	B		L	G	N	E	L	V	A	I	N	Q	Y	F	L	H	A	R	M	F	K	N	W	G	L	T	R	L	N	D	V	E	Y	H	E	S	I	D	E	M	K	H	A	D	K	Y	I	E	R	I	L	F	D	L	R	L	E	L	-	E	L	A	D	-	E	E	G	H	I	D		
E. coli	B	2HTN	L	G	N	E	L	V	A	I	N	Q	Y	F	L	H	A	R	M	F	K	N	W	G	L	K	R	L	N	D	V	E	Y	H	E	S	I	D	E	M	K	H	A	D	R	Y	I	E	R	I	L	F	D	L	A	L	E	L	-	D	L	R	D	-	E	E	G	H	I	D		
Y. pestis	B		L	G	N	E	L	V	A	I	N	Q	Y	F	L	H	A	R	M	F	K	N	W	G	L	M	R	L	N	D	K	E	Y	H	E	S	I	D	E	M	K	H	A	D	K	Y	I	E	R	I	L	F	D	L	A	L	E	L	-	S	L	V	D	-	E	E	E	H	I	D		
C.B.pennsylvanicus	B		L	S	D	E	L	V	A	V	N	Q	Y	F	L	H	S	K	I	F	N	N	W	G	L	E	R	L	N	K	I	E	Y	Q	E	C	V	D	E	L	D	H	A	D	L	Y	A	K	R	I	L	F	D	L	S	L	E	F	-	H	L	K	D	-	E	E	K	H	I	D		
A. vinelandii	B	1SOF	L	G	N	E	L	I	A	I	N	Q	Y	F	L	H	A	R	M	F	E	D	W	G	L	E	K	L	G	K	H	E	Y	H	E	S	I	D	E	M	K	H	A	D	K	L	I	K	R	I	L	F	D	L	K	L	E	Q	-	A	L	E	S	-	E	E	D	H	I	D		
M. capsulatus	B		L	T	N	E	L	T	A	I	N	Q	Y	F	L	H	A	R	M	F	K	N	W	G	F	G	K	L	N	D	H	E	H	E	Y	K	E	S	I	E	E	M	K	H	A	D	R	L	I	E	R	I	L	F	D	L	Q	L	E	Q	-	Q	L	E	S	-	E	E	E	H	V	D
S. salaskensis	B		L	K	N	E	L	T	A	I	N	Q	Y	W	L	H	Y	R	M	L	D	N	W	G	V	A	R	L	A	H	F	E	R	E	E	S	I	E	E	M	K	H	A	D	K	L	A	D	R	I	L	F	D	L	A	L	E	E	-	E	L	E	S	-	E	E	H	V	D			
H. baltica	B		L	K	N	E	L	T	A	I	N	Q	Y	F	L	H	S	R	M	L	K	D	W	G	V	S	V	L	A	E	K	E	Y	K	E	S	I	E	E	M	Q	H	A	D	W	L	I	D	R	I	L	F	D	L	K	L	E	H	-	D	L	E	N	-	E	E	E	H	V	D		
B. melitensis	B		L	F	L	E	L	G	A	V	N	Q	Y	W	L	H	Y	R	L	L	N	D	W	G	Y	T	R	L	A	K	K	E	R	E	E	S	I	E	E	M	H	A	D	K	L	I	D	R	I	L	F	D	L	K	G	E	Y	-	D	L	A	D	-	E	E	G	H	I	D			
Bradyrhizobium sp.	B		L	R	S	E	L	T	A	I	N	Q	Y	W	L	H	Y	R	L	L	N	N	W	G	L	L	E	M	A	K	V	W	R	K	E	S	I	E	E	M	H	A	D	K	F	T	D	R	I	L	F	D	L	A	A	E	I	-	G	M	K	D	-	E	E	H	I	D				
P. aeruginosa	B		L	T	G	E	L	A	A	R	D	Q	Y	F	I	H	S	R	M	Y	E	D	W	G	F	S	K	L	Y	E	R	L	N	H	E	M	E	E	E	T	Q	H	A	D	A	L	L	R	R	I	L	L	D	L	K	L	E	R	-	H	L	A	D	T	E	E	D	H	A	Y		
R. palustris	B		L	R	G	E	L	T	A	I	S	Q	Y	W	L	H	Y	R	L	L	A	N	W	G	L	K	D	M	A	K	V	W	R	K	E	S	I	E	E	M	E	H	A	D	L	L	T	D	R	I	L	F	D	L	A	A	E	M	-	G	M	K	D	-	E	E	H	I	D			
P. fluorescens	B		L	T	G	E	L	A	A	R	D	Q	Y	F	V	H	S	R	M	Y	E	D	W	G	F	T	K	L	Y	E	R	I	N	H	E	M	E	E	E	A	A	H	A	D	A	L	M	R	R	I	L	M	D	L	R	L	E	Y	-	K	L	H	D	T	E	E	D	H	T	Y		
M. capsulatus	B		L	A	G	E	L	A	A	I	D	Q	Y	F	I	H	A	M	M	Y	R	D	W	G	F	H	V	L	Y	E	H	T	A	H	E	M	Q	E	Q	A	H	A	S	A	L	I	R	R	I	L	F	D	L	G	V	E	H	-	A	L	D	D	T	E	E	D	H	C	L			
I. loihensis	B		L	A	F	E	L	T	S	I	D	Q	Y	T	S	H	S	R	Q	Y	E	D	M	G	L	M	K	L	Y	E	R	I	N	H	E	I	D	D	E	R	G	H	A	D	L	L	I	R	R	I	L	F	D	L	K	L	E	H	-	N	L	K	D	T	E	E	D	H	A	Y		
M. bovis	B		L	T	S	E	L	T	A	I	N	Q	Y	F	L	H	S	K	M	Q	D	N	W	G	F	T	E	L	A	A	H	T	R	A	E	S	F	D	E	M	R	H	A	E	E	I	T	D	R	I	L	L	D	L	A	I	E	Y	-	D	V	A	D	-	E	E	E	H	I	D		

# Sequence Comparison Methods (SCM)

## Multiple sequence alignments

Cobalt alignment; or:

### 1) Get sequences to align:

- putative homologs detected from a Blast search (saved as text)

### 2) Align all sequences to one another

- ClustalW, T-Coffee: [www.ebi.ac.uk](http://www.ebi.ac.uk) -> tools -> sequence analysis



# Sequence Comparison Methods (SCM)

## Clustal programs

### ClustalW2:

- Input sequences
- Multiple Sequence Alignment Options: Aligned vs. Input
- Output (%-age sequence identity)
- Alignment:
  1. Sequences with very different length
  2. Outliers (<20% sequence identity)
  3. Redundant (>80 % sequence identity)

# Sequence Comparison Methods (SCM)

## Multiple sequence alignment methods

### Edit/Visualize MSA:

- ClustalX (<ftp://ftp-igbmc.u-strasbg.fr/pub/ClustalX/>)
- JalView (<http://www.jalview.org/download.html>)
- BioEdit (<http://www.mbio.ncsu.edu/BioEdit/bioedit.html>)
- WebLogo (<http://weblogo.berkeley.edu/logo.cgi>)

# Sequence Comparison Methods (SCM)

## **‘High-quality’ MSA**

**At the basis of a number of structure/function prediction methods:**

- Homology modelling
- domains
- natively unfolded regions
- TM regions
- solvent accessibility
- secondary structures
- ...



# Protein sequence-structure-function relationships

## Most proteins are modular

**Domains: structural, functional, folding and evolutionary units**  
**(30-700 a.a.; 100 a.a. on average)**

**Analysis and prediction: domain – not whole protein – level**

**Tools for prediction of:**

- **Domains**
- **Disordered regions**
- **Trans-membrane regions**
- **Secondary structure elements**

# Protein sequence-structure-evolutionary relationships

## Structural Classification Of Proteins (SCOP) DB ([scop.mrc-lmb.cam.ac.uk/scop/](http://scop.mrc-lmb.cam.ac.uk/scop/))

**DOMAINS:** Structural, functional, folding and evolutionary protein units

**FAMILY:** Close evolutionary relationship (structure and sequence similarity)

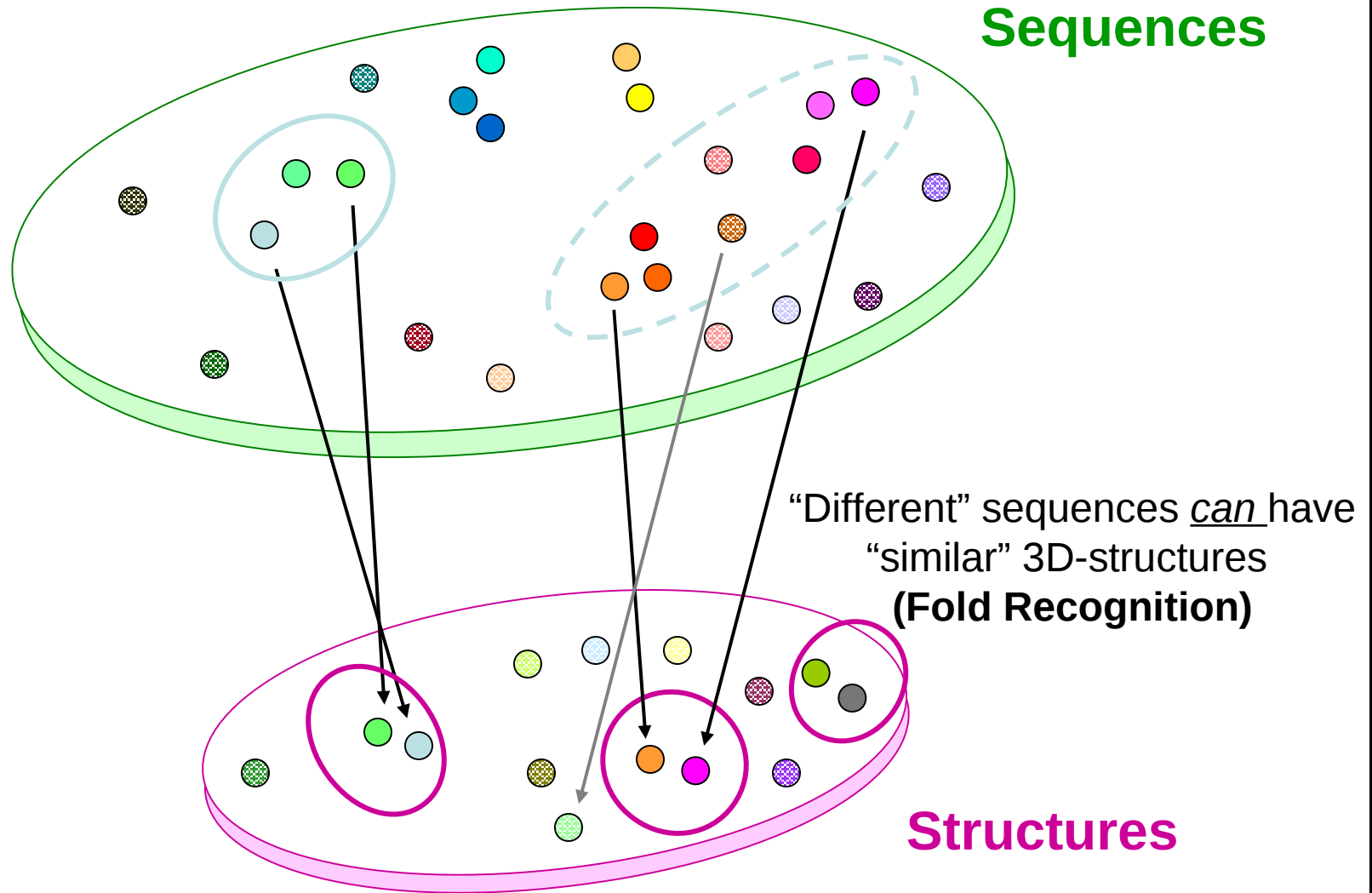
**SUPERFAMILY:** Distant evolutionary relationship (structure similarity)

---

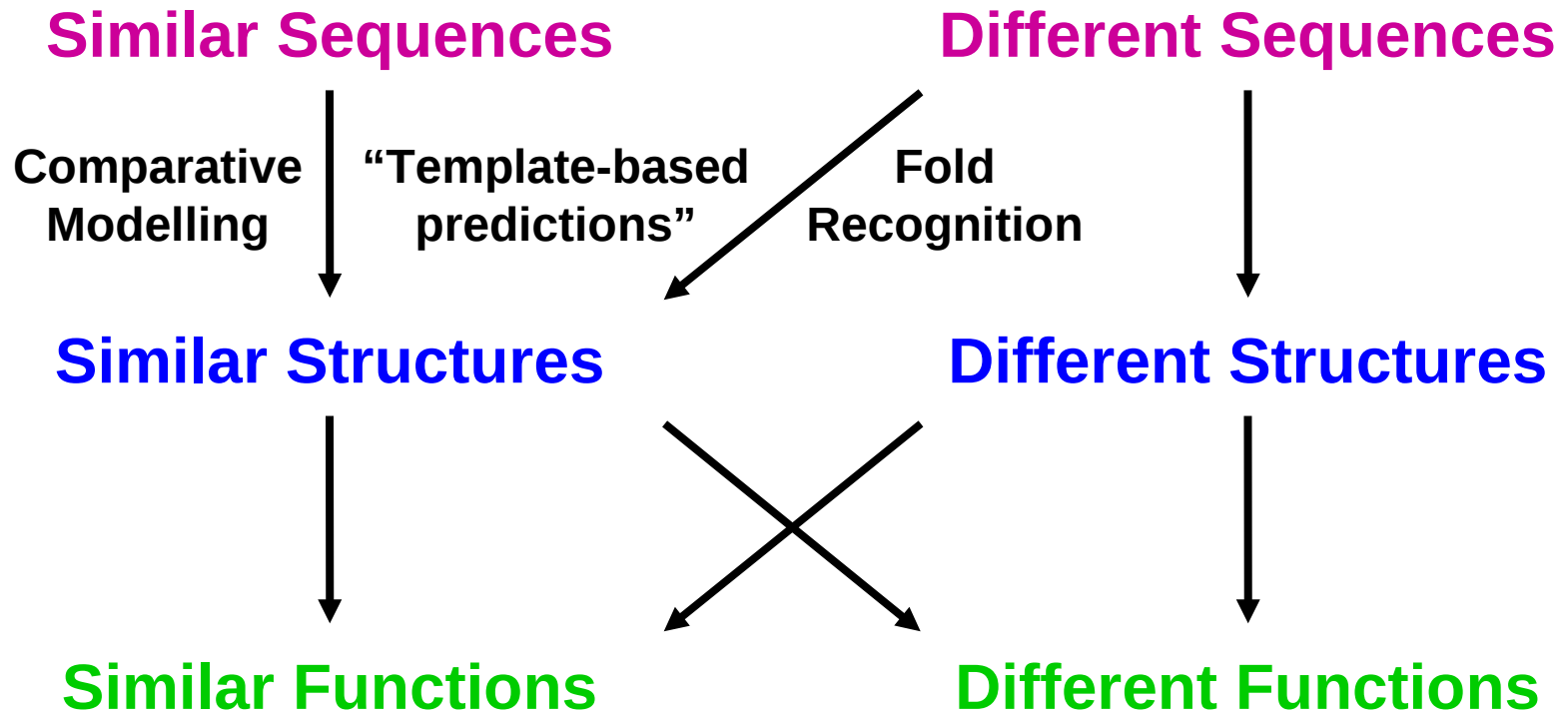
**FOLD:** Uncertain evolutionary relationship (structural similarity)

**CLASS:** Secondary structure content ( $\alpha$ ;  $\beta$ ;  $\alpha/\beta$ ;  $\alpha+\beta$ ; ...)

# Protein sequence-structure-function relationships

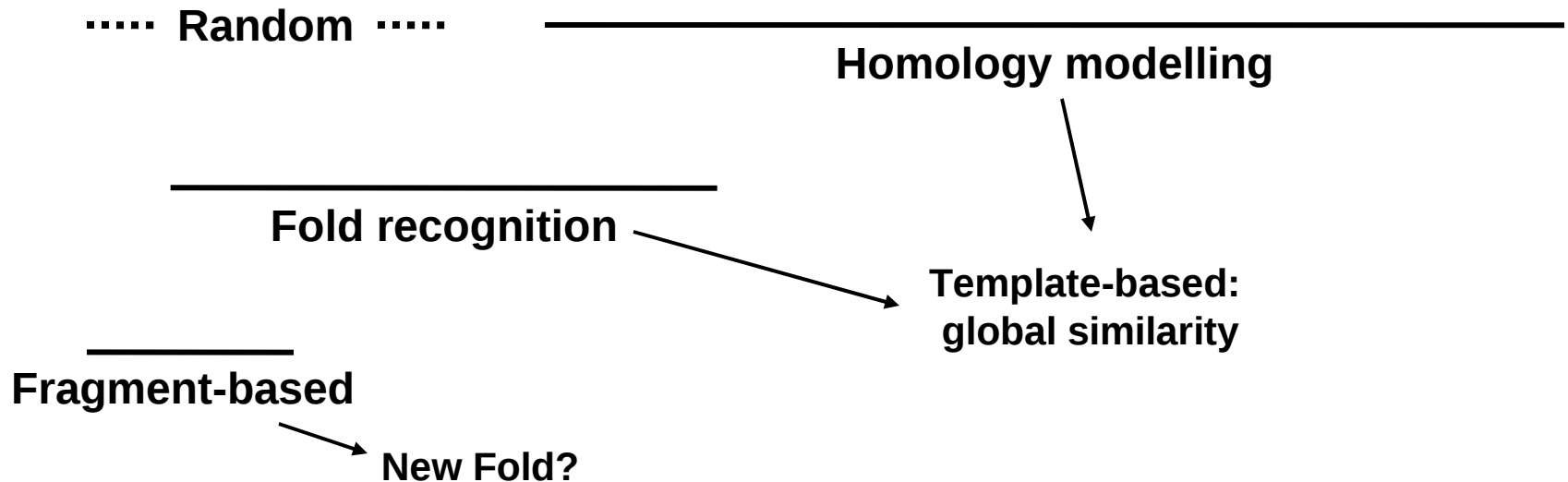
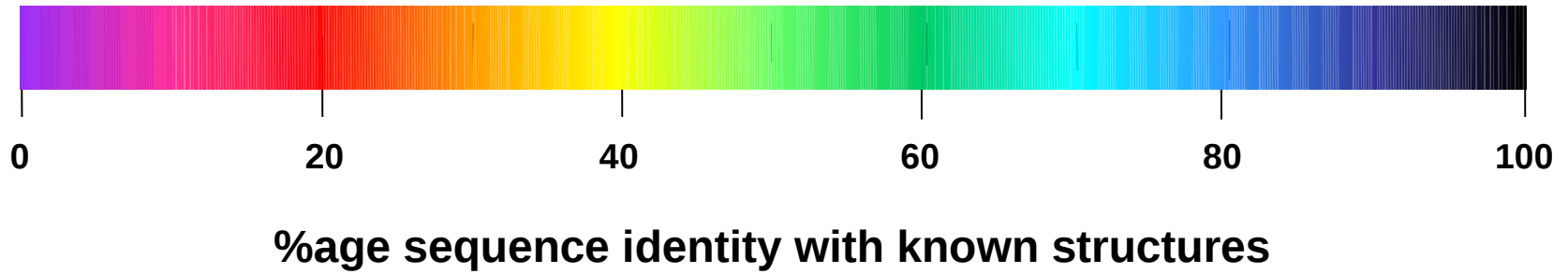


# Protein sequence-structure-function relationships

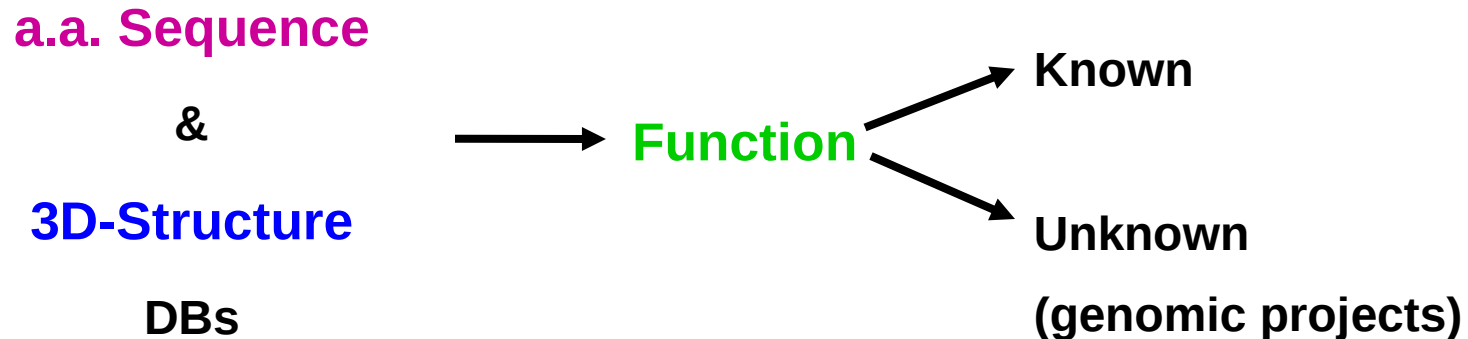
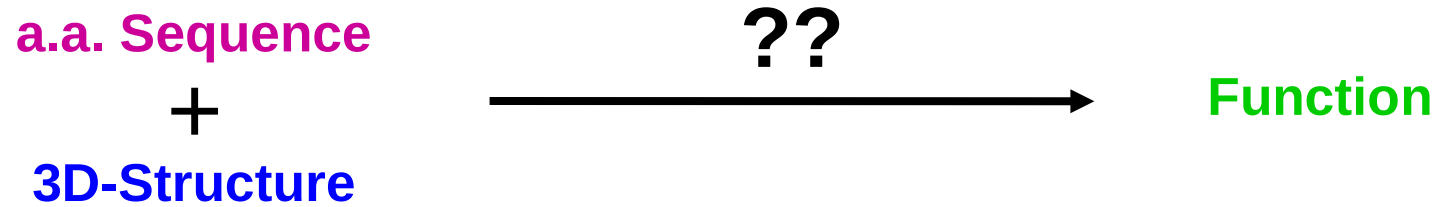




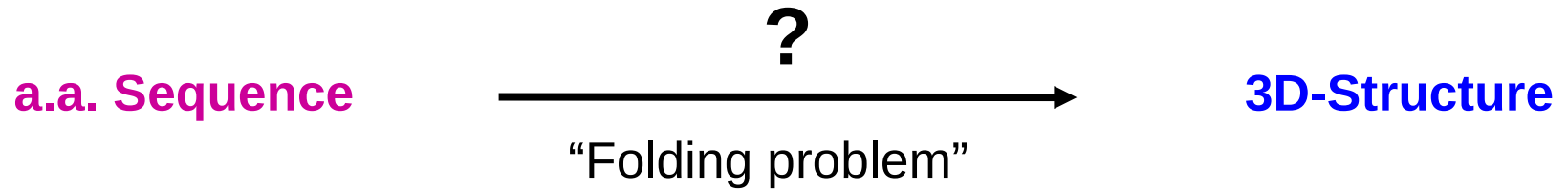
# Protein structure prediction methods



# Structure-function relationships



# Sequence-structure relationships



*ab initio*  
(the "Holy grail")

Physical-Chemical Principles  
(*"nature folds proteins without  
searching the DBs..."*)

—

evolutionary

Relationships between  
known sequences and  
structures (DBs!!!)

+++

# Sequence-structure-function relationships

**Similar sequences**



**Similar Structures**



**Similar Functions**

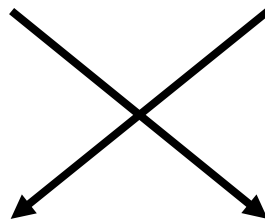
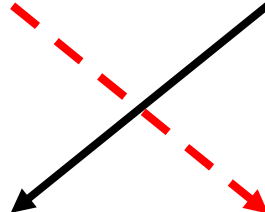
**Different Sequences**



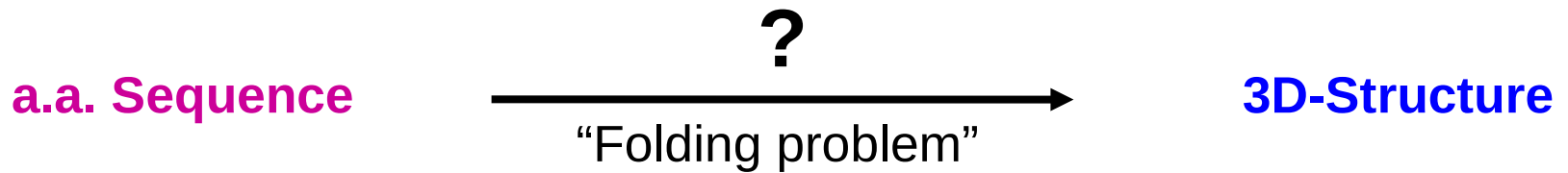
**Different Structures**



**Different Functions**



# Sequence-structure relationships



*ab initio*

Physical-Chemical Principles

—

**evolutionary**

Relationships between sequences  
and structures in DBs

**+++**

## **Global: "Template-based predictions"**

- High sequence similarity: Homology modelling **++/+++**
- Low sequence similarity: Fold recognition **+/++**

## **Local: "New-Fold predictions"**

- No sequence similarity: Fragment-based methods **-/±**

# Protein sequence analysis

1) NCBI (UniProt) annotation

2) BLAST

- Conserved domains (CDD) → annotate residues (51-202)
- Similar sequences → download Fasta complete
- Pair-wise alignments → download (text format)
- Multiple alignment (COBALT) → download (text format)