

Chemoinformatika a bioinformatika

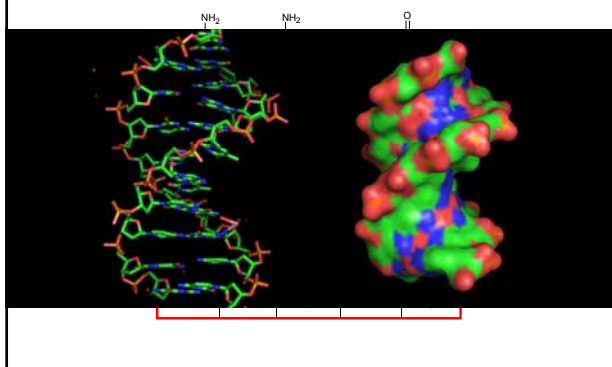
Sequence alignment



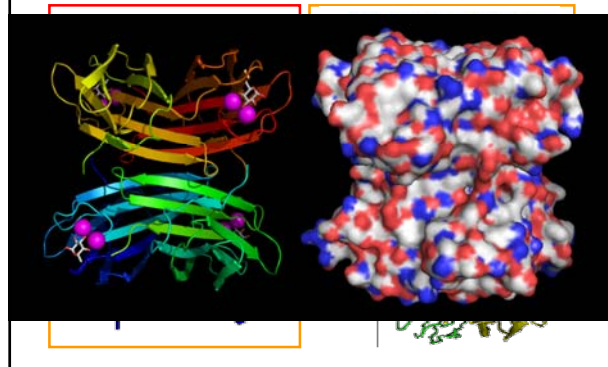
Osnova

1. Struktura biomakromolekul – sekvence
2. Alignment a jeho typy
3. Užívané algoritmy
4. Multiple sequence alignment
5. Programové balíky
6. Benchmark – porovnávání alignmentů

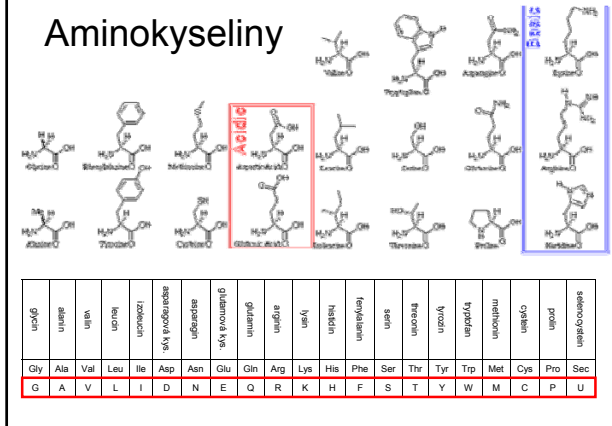
Nukleové kyseliny a báze



Struktura proteinů

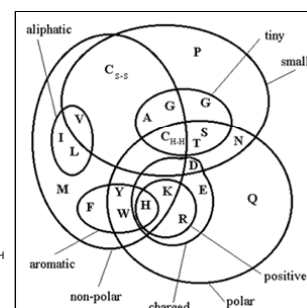
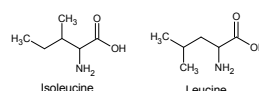


Aminokyseliny



Třídění aminokyselin

Aminokyseliny s podobnými vlastnostmi mohou plnit v proteinu stejné funkce – bývají vzájemně zastupitelné



Alignment

Srovnání (přiložení) dvou či více sekvencí (aminokyselinových, nukleotidových) na základě jejich vzájemné podobnosti.

Pairwise alignment – dvě sekvence

Multiple sequence alignment – více sekvencí

Vstupní data

Sekvence AK (nt) v určitém formátu – dnes desítky formátů, mnohé obsahují krom sekvence i doplňující data

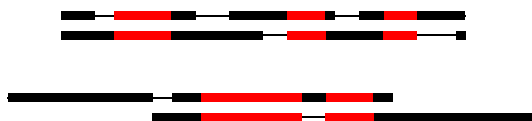
Bližší např.
<http://emboss.sourceforge.net/docs/themes/SequenceFormats.html>

- **FASTA formát**

```
>název_popis dle vlastní volby..|
SEKVENCESEKVENCESEKVENCESEKVENEC
ESEKVENCESEKVENCE
```

Pair-wise alignment

- Srovnání dvou sekvencí
- Sekvence mohou být seřazeny v celé své délce (global alignment) nebo jen v určitém regionu (local alignment).



Local alignment

Hledá úseky dvou sekvencí, které si podle zvolených kritérií dobře odpovídají. Nesnaží se zahrnout celé sekvence, pokud si jejich některé části neodpovídají.



Global alignment

Vychází z předpokladu, že obě srovnávané sekvence jsou víceméně shodné v celé své délce. Alignment k sobě přikládá celé sekvence (od počátku do konce) a to včetně částí, které si nepříliš odpovídají.



Algoritmy

- Témeř výhradně se užívají **heuristické algoritmy** – nalezení výsledku v dostatečně krátkém čase
- Vývoj algoritmů je prováděn v návaznosti na srovnávání výsledků s tzv. zlatým standardem – alignment na základě známé 3D struktur

DNA matice

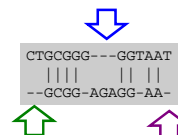
A	1			
T	-10000	1		
G	-10000	-10000	1	
C	-10000	-10000	-10000	1
	A	T	G	C

Jako pozitivní je uvažována pouze shoda, jakákoliv substituce je vysoce penalizována; jsou však povoleny mezery.

Mezery (Gaps)

Příčiny vzniku mezer:

- **Bodová mutace** (velmi častá příčina)
 - Nepřesný crossover při meiose (inzerce nebo delece řetězce bází)
 - DNA slippage během replikace (vzniká repetice – opakující se sekvence v řetězci)
 - Inzerce retroviru
 - Translokace DNA mezi chromozomy
- Mezery nacházíme na **začátku** řetězce, **uprostřed** nebo na jeho **konci**.



Mezery umožňují alignment sekvencí, kdy v jedné z nich došlo k delecí. Zvyšují však také možnost alignmentu náhodných sekvencí. Jejich přítomnost je proto vždy „**penalizována**“, a to více než substituce.

Čím nižší je penalizace mezer, tím lepší (dokonalejší) bude alignment, ovšem z biologického hlediska může jít o nesmysl.

Jednotlivé programy obvykle penalizují **přítomnost mezer** (gap open) a také zvyšují penalizaci s **délkou mezer** (gap ext).

Krátká mezera:

```
ATCTTCAGTGTTCCTGTTTGGCC...ATTAGTTCGCTC
|||||
ATCTTCAGTGTTCCTGTTTGGCCGATTAGTTCGCTC
```

Dlouhá mezera:

```
ATCTTCAGTGTTCCTGTTTGGCC.....ATTAGTTCGCTC
|||||
ATCTTCAGTGTTCCTGTTTGGCCGCCCCCCCCCCCCCCCCCATTTAGTTCGCTC
```

Skóre

Každé dvojici sekvencí je ve výsledku přiřazeno číslo – skóre, které určuje míru jejich podobnosti

Čím vyšší je skóre, tím vyšší je podobnost.

Podle použité matice může být skóre i záporné.

Příklad výpočtu

AABBCCDDEF
AADKKKEFGG

Ve chvíli, kdy zafixujeme pozici dvou sekvencí, pak můžeme snadno vypočítat skóre pro dané přiložení (příklad BLOSUM 62):

skóre pro dané přiložení = skóre na bázi jednotlivých aa + celková penalizace metod

Například, celkové pozitivní skóre na úrovni jednotlivých aa

```
A A B B C C D D - - E E F
A A - - - - D D K K E F G G
4+4      +6+6  +1+5+6      = 32
```

Naopak, pro každou mezeru (-) je dána penalizace: první výskyt zleva -10, každá následující -1.

```
A A B B C C D D - - E E F
A A - - - - D D K K E F G G
-10-1-1-1  -10-1      = -24
```

Celkové skóre 32 – 24 = 8

Multiple sequence alignment - MSA

(mnohonásobné srovnání)

Multiple alignment slouží k:

- Nalezení „diagnostického vzoru“ (diagnostic patterns) na jehož základě jsou **charakterizovány proteinové rodiny**
- Odhalení či dokázání **homologie** mezi novou sekvencí a sekvencemi v databázích
- Určení vzájemné příbuznosti sekvencí v rámci skupiny – tvorba **fylogenetických stromů**
- **Predikci** sekundární a terciární **struktury** nových proteinů
- Navržení primerů (oligonukleotidů) pro PCR

DIST = percentage divergence (/100)
 Length = number of sites used in comparison
 1 vs. 2 DIST = 0.6491; length = 114
 1 vs. 3 DIST = 0.6842; length = 114
 1 vs. 4 DIST = 0.9298; length = 114
 1 vs. 5 DIST = 0.9035; length = 114
 1 vs. 6 DIST = 0.9396; length = 114
 1 vs. 7 DIST = 0.9825; length = 114
 2 vs. 3 DIST = 0.3772; length = 114
 2 vs. 4 DIST = 0.9123; length = 114
 2 vs. 5 DIST = 0.8947; length = 114
 2 vs. 6 DIST = 0.9123; length = 114
 2 vs. 7 DIST = 0.9386; length = 114
 3 vs. 4 DIST = 0.9123; length = 114
 3 vs. 5 DIST = 0.9386; length = 114
 3 vs. 6 DIST = 0.9298; length = 114
 3 vs. 7 DIST = 0.9474; length = 114
 4 vs. 5 DIST = 0.9211; length = 114
 4 vs. 6 DIST = 0.9035; length = 114
 4 vs. 7 DIST = 0.9649; length = 114
 5 vs. 6 DIST = 0.9561; length = 114
 5 vs. 7 DIST = 0.9211; length = 114
 6 vs. 7 DIST = 0.9649; length = 114

Neighbor-joining Method
 Saitou, N. and Nei, M. (1987) The Neighbor-joining Method:
 A New Method for Reconstructing Phylogenetic Trees.
 Mol. Biol. Evol., 4(4), 406-425
 This is an UNROOTED tree
 Numbers in parentheses are branch lengths
 Cycle 1 = SEQ: 2 (0.17807) joins SEQ: 3 (0.19912)
 Cycle 2 = SEQ: 1 (0.34101) joins Node: 2 (0.13706)
 Cycle 3 = SEQ: 5 (0.44298) joins SEQ: 7 (0.47807)
 Cycle 4 = SEQ: 4 (0.44518) joins SEQ: 6 (0.45833)
 Cycle 5 (Last cycle, trichotomy):
 Node: 1 (0.12171) joins
 Node: 4 (0.01864) joins
 Node: 5 (0.02083)

.nj soubor

.ph soubor

```
(
(PAIL:0.34101,
(RSIL:0.17807,
(CVIL:0.19912)
):0.13706)
:0.12171,
(
(BCLA:0.44518,
(BCLC:0.45833)
):0.01864,
(
(BCLB:0.44298,
(BCLD:0.47807)
):0.02083);
)
```

.dst soubor

7							
PAIL	0.000	0.649	0.684	0.930	0.904	0.939	0.982
RSIL	0.649	0.000	0.377	0.912	0.895	0.912	0.939
CVIL	0.684	0.377	0.000	0.912	0.939	0.930	0.947
BCLA	0.930	0.912	0.912	0.000	0.921	0.904	0.965
BCLB	0.904	0.895	0.939	0.921	0.000	0.956	0.921
BCLC	0.939	0.912	0.930	0.904	0.956	0.000	0.965
BCLD	0.982	0.939	0.947	0.965	0.921	0.965	0.000

Phylogram a cladogram

- **Phylogram** (phylogeny tree) – je rozvětvený diagram (strom), který naznačuje fylogenezi (postupný vývoj). Délka jednotlivých větví je úměrná **velikosti změny** v průběhu evoluce.
- **Cladogram** – rovněž strom, v němž však všechny větve mají **stejnou délku**. Ukazuje tak sice „společné předky“ pro jednotlivé sekvence, ale ne množství změn, jež od té doby prodělaly (evoluční dobu).

Phy

Cladogram

Iterativní přístup

(Gotoh, 1996; Notredame & Higgins, 1996)

Vzniklý strom i alignment jsou následně **optimalizováni** do konvergence. Jinak jsou chyby vzniklé při prvním alignmentu (tvorba stromu) zachovány i ve výsledku.

Nezaručuje nalezení nejlepšího výsledku, ale – na rozdíl od deterministických alternativ – je dostatečně **robustní** a dobře použitelný i pro velký počet sekvencí.

Kombinace local a global alignment

- S výhodou lze kombinovat lokální a globální alignment.
- Lokální alignment může být reprezentován sadou kotvicích bodů v místě dobré shody
- Následný globální alignment pak tyto odpovídající úseky sekvencí zahrnuje (využito např. v ClustalW2)

Výstup

Výstupem je sada sekvencí (případně s vloženými mezerami)
Různé formáty, nejčastěji používán **.aln soubor**, ale též **.fasta**, aj.

Mnoho programů sloužících pro zobrazení a/nebo editaci

- Bioedit
- JalView
- CINEMA 2.1...
- JavaShade
- ...

Výstup - .aln soubor

```

CLUSTAL 2.0.10 multiple sequence alignment

PA1IL -----ATQGVFF
RS1IL -----AQQGVFF
CV1IL -----AQQGVFF
BCLB --LREKLPQGVVWDLKATIPSPDQNSQWVKTDAASRVACTVFMGASQVYLGQAA
BCLC A1ATNQQVVAQQCFYVSKVPESTGRMPFLWATLWVIGVGVFFVQKNSVYRGSAMHIDS
BCLA -----
BCLD LBETALAAEAVSVLFIKFAKLDAGIVAPIELVVDAAVFPADQLLHPGCRPLKDHVY

PA1IL -----ATQGVFF
RS1IL -----AQQGVFF
CV1IL -----AQQGVFF
BCLB KFGVGVVFN-----YFSKATPQVQVPAFV-----TQDGERDGIPT
BCLC YASLSA1WG-----TAAPSQQSSQNSQASRTQGVAGNIQQGGERDRTFM
BCLA -----ASQGF-----SRRRAGEFF
BCLD RSDVLAAGATTCTADFAVCDRDQTVSGYFRWETSLEIAGSQPQTKQPFKSSDRNRFIS
* *

PA1IL LPANTFQVTFANSSQTQVNVLVNNEA--ATFSDQSTNSNAVIGTVLNSQSSKQVQV
RS1IL LPANTFQVTFANSSQTQVNVLVNNEA--ATFSDQSTNSNAVIGTVLNSQSSKQVQV
CV1IL LPARINFQVTVLNSAATQVVEIPEDEER--AAFSDVDTQDNLNTQVINSQ--QVYV
BCLB LPNRIAPQVTVLNSAATQVVEIPEDEER--AAFSDVDTQDNLNTQVINSQ--QVYV
BCLC LPNRIAPQVTVLNSAATQVVEIPEDEER--AAFSDVDTQDNLNTQVINSQ--QVYV
BCLA IPPNTPFRAIFPANAARQQIKLFIQDSQKPAVHKLTRDQRE--ATLNSQ--QKIRP
BCLD LPNTPAFKAIFFYANAARQQIKLFIQDSQKPAVHKLTRDQRE--ATLNSQ--QKIRP
* * * * *

```

Programové balíčky



- Existují programy pro pairwise alignment i pro MSA
- Využívají lokální nebo globální alignment nebo příp. kombinaci obou
- Neexistuje univerzální „nejlepší“ program – záleží na konkrétním použití

Pairwise alignment „programy“

Oblasti použití:

- Přímé porovnání dvou sekvencí
- Vyhledávání podobných sekvencí v databázích

emboss Needle & Water

- vytvořeny 1970
Needleman S.B. and Wunsch C.D. (1970) A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology* 48: 443-453.
- využívají dynamické programování,
- umožňují vložení mezer

Needle – globální pairwise alignment, Needleman-Wunsch algoritmus

Water – lokální pairwise alignment, Smith-Waterman algoritmus

T-Coffee

http://www.tcoffee.org/Projects_home_page/t_coffee_home_page
(Tree-based Consistency Objective Function for alignment Evaluation)

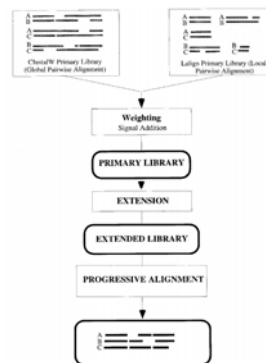


- Pomalejší ale výrazně přesnější než ClustalW
- Je schopen kombinovat data z více předchozích alignmentů, které mohly být vytvořeny různými postupy (lokální, globální, strukturní podobnost,...)

Hlavním rozdílem oproti tradičním metodám progresivního alignmentu je použití pozičně specifického skórovacího schématu (**extended library**) namísto substituční matice.

T-Coffee

- 1) Provedení pairwise alignmentů pro všechny dvojice sekvencí pomocí globálního a pomocí lokálního alignmentu (dve primární knihovny).
- 2) Jednotlivým pairwise alignmentům je přiřazena váha podle poměru počtu identických residuí k celkovému počtu residuí.
- 3) Kombinace obou knihoven. Pokud je rozdíl v globálním a lokálním alignmentu, jsou zachovány oba s příslušnou vahou. Vzniká pozičně specifická matice (extended library), která je dále použita pro vlastní progresivní alignment.



Zlepšení přesnosti – strukturní informace

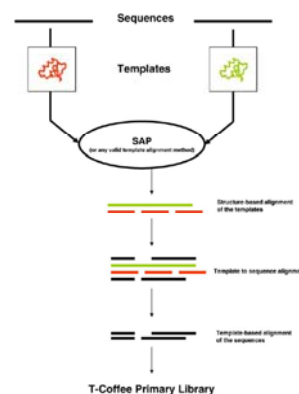
- Sekvence s vyšší homologií (>40%) – vysoká přesnost alignmentu
- Bez homologie – nepoužitelné
- Tzv. twilight zone – málo podobné sekvence (nižší než 20% homologie) = špatná (méně než 30%) přesnost alignmentu

Řešení: nejčastěji využití znalosti strukturní podobnosti (2D nebo 3D), která se během evoluce zachovává více než sekvence AK.

Rozšíření konzistentního modelu

Template-based alignment metody – vytváří alignment na základě strukturní informace známých homologů jednotlivých sekvencí v PDB databázi nebo získava "profil" sekvence na základě homologních sekvencí (pomocí BLASTu)

Výhoda: vyšší přesnost



Expresso

- Je založeno na 3DCoffee
- Expresso je MSA server, který srovnává sekvence za užití strukturní informace. Po zadání sekvencí vyhledá v databázi struktur (PDB) pomocí BLASTu homology a použije je jako templáty pro následný alignment zadaných sekvencí pomocí metod MSA založených na struktuře (např. SAP, Fugue).

Benchmark (srovnávací testy)

BALIBASE - První vytvořená sada benchmarkových testů pro multiple alignment programy (Thompson et al., 1999) – byla vytvořena pomocí manuálně provedeného alignmentu

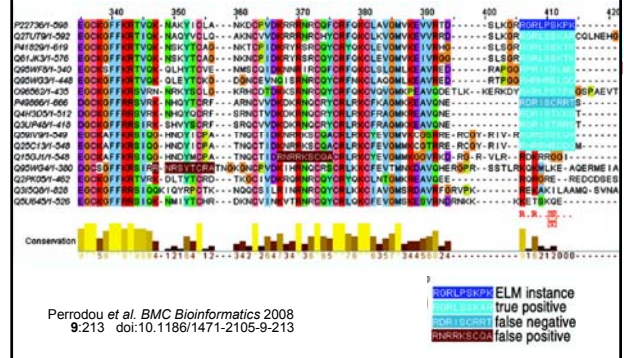
Na základě srovnání 3D struktur byly vytvořeny další sety:

- HOMSTRAD [Mizuguchi *et al.*, 1998].
- OxBench [Raghava *et al.*, 2003]
- PREFAB [Edgar, 2004]

Existují i specificky zaměřené benchmarkové sady, např.

IRMBASE [Subramanian *et al*, 2005] – náhodné (nepřiložitelné) sekvence s vloženými motivy. Slouží k testování metod pro lokální alignment

BaliBASE – ukázka alignmentu



Zopakování / shrnutí

- ▼ **Alignment** – přiložení sekvencí (2 nebo více) na základě podobnosti
- ▼ **Využití** pro hledání příbuznosti sekvencí, tvorba profilů proteinových rodin, aj.
- ▼ Řada **programů** využívajících rozdílné přístupy – použití závisí na vstupních datech a účelu
- ▼ Nejčastěji používaný (ClustalW) neznamená nepřesnější – každý program je **kompromisem mezi přesností a rychlostí**
- ▼ Každý alignment potřebuje **lidskou kontrolu !!!**



Local alignment

- For two-sequence comparisons, there is the well-known Smith and Waterman (1981) algorithm. Here we use Lalign
- For multiple sequences, the Gibbs sampler (Lawrence *et al.*, 1993) and Dialign2 (Morgenstern, 1999) are the main automatic methods. These programs often perform well when there is a clear block of ungapped alignment shared by all of the sequences. They perform poorly, however, on general sets of test cases when compared with global methods

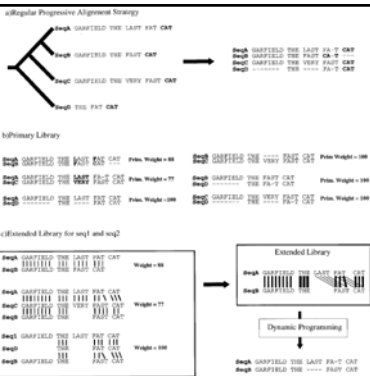


Figure 2. The library extension. (a) Progressive alignment. Four sequences have been designed. The tree indicates the order in which the sequences are aligned when using a progressive method such as ClustalW. The resulting alignment is shown, with the novel C/A mismatched. (b) Primary Library. Each pair of sequences is aligned using ClustalW. In these alignments, each pair of aligned residues is associated with a weight equal to the average identity among matched residues within the complete alignment (numbers are indicated in bold type). (c) Library extension for a pair of sequences. The three possible alignments of sequence A and B are shown (A and B, A and B through C, A and B through D). These alignments are combined, as explained in the text, to produce the position-specific library. This library is resolved by dynamic programming to give the correct alignment. The thickness of the lines indicates the strength of the match.

BALIBASE [Thompson *et al.*, 1999] contains eight reference sets, each dealing with a different type of alignment problem. Ref1 deals with test cases containing small numbers of equivalent sequences, and is further subdivided by percent identity. Ref2 alignments contain "topical" or unrelated sequences. Ref3 test cases contain a pair of divergent substrates, with less than 25% identity between the two groups. Ref4 is concerned with long terminal extensions, while Ref5 test cases contain large internal insertions and deletions. Test sets from references 6-8 deal with problems like transmembrane regions, inverted domains, and repeat sequences. In previous versions of BALIBASE, test cases were confined to homologous regions. In practice, the boundaries of such regions may be unknown. The current version [Thompson *et al.*, 2005] now also provides 6 disparate test cases containing full-length sequences. Only the first five reference sets are used here, as they have been corrected and verified in the latest release.

OnBench [Moguch *et al.*, 2003] comprises 3 related datasets. Test cases in the MASTER set deal with isolated domains derived exclusively from sequences of known structure. The FULL set was generated from suitable MASTER test cases, using full-length sequence data. High scoring homologous sequences were added to each MASTER test case to generate the EXTENDED set. The results from this third set, however, are not used here. It was found that some of the test cases in the EXTENDED set proved too large for some programs, and aborted due to excessive memory requirements. Of the 276 test cases selected from EXTENDED, T-COFFEE returned 295 alignments, and Align in was only able to align 167, using a single processor with 4GB of RAM.

PREFAB [Edge, 2004] test cases are generated by taking a pairwise alignment of sequences of known 3D structure, and adding up to 24 high scoring homologues for each sequence. Accuracy is assessed on the structural alignment of the original pair alone.

KALIBRAK [van Isterik *et al.*, 2005] is divided into two subsets. Each test group in the SLIPSTREAM V set represents a SCOP superfamily, whose sequences are 25-100% identical. Each test group in the TRILIGHT set represents a common SCOP fold and sequences are 0-20% identical. In addition, these two subsets are also provided with non-homologous (these positive) sequences included within each group. Instead of a single alignment acting as a reference, S40marks provides multiple pairwise references for each test, and it is the average score from each of these references that is taken here as a score for each test case.

IRMBASE [Subramanian *et al.*, 2005] test cases contain a number of simulated motifs [Bhaya *et al.*, 1998] inserted into otherwise random (unalignable) sequences, and so are entirely different to the other benchmarks used in this study. Test cases are designed to examine whether a method can detect isolated motifs within sequences, and so are tailored to a local alignment approach.

HOMSTRAD [Moguch *et al.*, 1998] is a database exclusively based on protein structures derived from the PDB, arranged into homologous protein families. It was not specifically designed as a benchmark database, although it is regularly employed as such.

Method	Score	Templates	Validation Values		Server
			Prefab	HOMSTRAD	
ClustalW [14]	Mean	—	61.86 [12]	—	http://www.ebi.ac.uk/blast/
Julgln	Mean	—	63.00 [16]	—	http://ma.zgib.li/ser/
MUSCLE [8]	Mean	—	68.00 [16]	63.0 [8]	http://www.drive5.com/muscle/
T-Coffee [10]	Consistency	—	69.97 [12]	44.0 [9]	http://www.tcoffee.org/
ProbCons [7]	Consistency	—	70.54 [12]	—	http://probcons.stanford.edu/
MAFFT [6]	Consistency	—	72.30 [12]	—	http://align.genome.jp/mafft/
MCoffee [12]	Consistency	—	72.91 [12]	—	http://www.tcoffee.org/
HOMSTRAD [16]	Consistency	—	73.10 [16]	—	http://prodna.som.mcgill.ca/homstrad/
DIClustal [24]	Profiles	—	—	—	http://bioparsis.org/align/align/
PRANK [6]	Mean	Profiles	—	62.2 [6]	http://sourceforge.net/projects/prankserver/
PROBALS [14]	Consistency	Profiles	70.00 [16]	—	http://prodna.som.mcgill.ca/probals/
SPM [28]	Mean	Profiles	77.00 [28]	—	http://geneticsinformatics.cup.edu.edu/Software/Services_files/spm.htm
Expres [13]	Consistency	Structures	—	71.18 [13]*	http://www.tcoffee.org/
T-Less [26]	Consistency	Structures	—	—	http://www.molku.bioinf.de/tless/

Validation values were compiled from several sources, and selected for comparability. Prefab validations were made using Prefab version 3. HOMSTRAD validations were made on datasets having less than 30% identity. The source of each value is indicated by the accompanying reference citation. The Expres value comes from a slightly more demanding subset of HOMSTRAD (HOMD) made of sequences less than 25% identical. doi:10.1371/journal.pcbi.1002022.t001

Table 1: Programs used in this investigation

Method	OVERVIEW
Align-Q (2)	http://bioinformatics.vub.ac.be/bioinformatics/align-q.html Local, specialised for highly divergent sequences. [Van Walle et al., 2004]
ClustalW (1.8)	http://www.ebi.ac.uk/blast/ Global, progressive alignment package. [Thompson et al., 1994]
Dialign2 (2.2)	http://bioinformatics.toronto.utoronto.ca/dialign/ Local, aligns segments of sequences rather than individual residues. [Morjanter, 1999]
Dialign-Q (1.3)	http://align-f.genetics.de/ Local, progressive alignment. Recent re-implementation of Dialign2. [Hiltenen et al., 2005]
MAFFT (6.521)	http://www.genome.jp/mafft/align/mafft/align/mafft/ Suite of alignment programs. [Katoh et al., 2002]
FFTS	Global, uses Fast Fourier Transforms to generate tree.
FFTNS	As FFTS, but with iteration step to refine alignment.
NWAS	Global, uses traditional Needleman-Wunsch algorithm.
NWAS	As NWAS, but with iteration step to refi
FNIS	Local, iterative, uses local pairwise alg
GNIS	Global, iterative, uses global pairwise alignment information.
MUOCLC (3.0)	http://www.drive5.com/muscle/ [Edgar, 2004] Global, iterative, progressive alignment program that uses Log Expectation as scoring function.
ProbCons (1.00)	http://probcons.stanford.edu/ [Do et al., 2005] Global, uses posterior probabilities from HMMs and pairwise alignment consistency.
PCMA (2.0)	http://hpc.som.mcgill.ca/PCMA/ [Wei et al., 2003] Global, switches alignment strategies dependent on sequence data. ClustalW is used to align highly similar sequences and to form pre-aligned groups. T-COFFEE is used to align the more divergent groups.
POA (v2)	http://www.bioinformatics.ucla.edu/poa/ [He et al., 2002] Local, uses Partial Order graphs.
T-COFFEE (1.37)	http://align-server.cbi.msi.ru/~chirov/projects_home_page/t_coffee_home_page.html [Notredame et al., 2000] Combines both global and local methods; uses consistency.

Blackshield 2006 oznacil ProbCons jako nejlepsi na zaklade 6 benchmarkovych testu