

Chemoinformatika a bioinformatika

Sequence alignment



Osnova

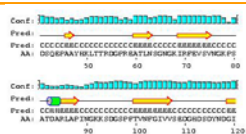
1. Struktura biomakromolekul – sekvence
2. Alignment a jeho typy
3. Užívané algoritmy
4. Multiple sequence alignment
5. Programové balíky
6. Benchmark – porovnávání alignmentů

Struktura proteinů (NK)

```

ADSQTSSNRAGEFSIPPNTDFRAIF
FANAAEQQHIKLFIGDSQEPAAAYHK
LITRDGPREATLNSGNGKIRFEVSV
NGKPSATDARLAPINGKKSDDGSPF
TVNFGIVVSEDDGHSDYNDGIVVL
QWPIG
    
```

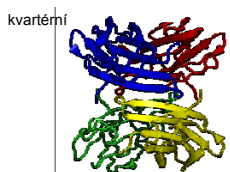
primární
(sekvence)



sekundární

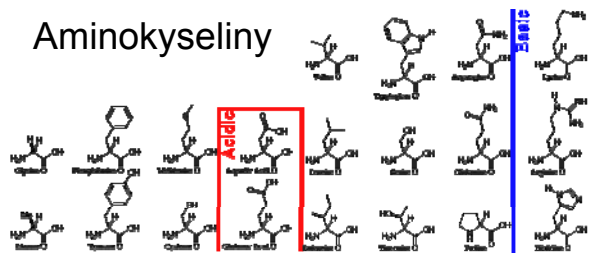


terciární



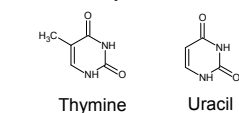
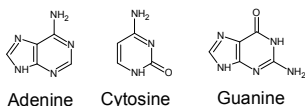
kvartérní

Aminokyseliny



glyc	alan	val	leu	ile	asp	asn	glu	gln	arg	lys	his	phe	ser	thr	tyr	trp	met	cys	pro	sec
G	A	V	L	I	D	N	E	Q	R	K	H	F	S	T	Y	W	M	C	P	U

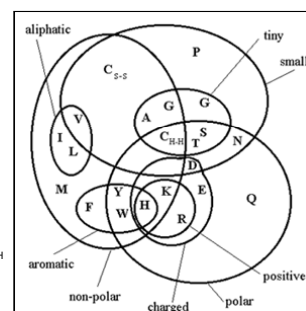
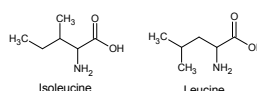
Nukleové báze



adenin	cytosin	guanin	thymin	uracil
A	C	G	T	U

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

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

- Téměř 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 struktury

Vstupní data

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

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

- **FASTA formát**

```
>název_popis dle vlastní volby_1
SEKVENCESEKVENCESEKVENCESEKVENCE
ESEKVENCESEKVENCE
```


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 substitute.

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

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

Krátká mezera:

```
ATCTTCAGTGTTCCTCCCTGTTTGGCC...ATTAGTTCGCTC
|||||
ATCTTCAGTGTTCCTCCCTGTTTGGCCGATTAGTTCGCTC
```

Dlouhá mezera:

```
ATCTTCAGTGTTCCTCCCTGTTTGGCC.....ATTAGTTCGCTC
|||||
ATCTTCAGTGTTCCTCCCTGTTTGGCCGCCCCCCCCCCCCCCCCATTAGTTCGCTC
```

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.

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

Metody MSA

- Dynamické programování (dynamic programming) – rozšíření pairwise alignmentu
- **Progresivní alignment** (progressive sequence alignment) – nejčastěji používaný k vytvoření alignmentu; využívá fylogenetické informace
- Iterativní alignment (iterative sequence alignment) – odstraňuje problémy progresivního alignmentu pomocí opakování alignmentu pro podskupiny sekvencí

Dynamické programování

- **Simultánní alignment všech sekvencí** - analogické pairwise alignmentu
- Programové balíky: MSA (Lipman et al., 1989) a DCA (Stoye et al., 1997), založené na Carrilově a Lipmanově algoritmu (1988)
- Využívá skórovací matice, ale vytváří n-rozměrný prostor (n = počet sekvencí)
- Extrémně **náročný na výpočetní kapacity**
- I při zjednodušení nepoužitelné pro více než cca 20 sekvencí



Progresivní multiple alignment

- Používá ho většina programů
- Vznik – 1987
Feng, D.-F. and Doolittle, R.F. (1987) J. Mol. Evol. 25, 351-360.
- 1) sestavení příbuzenského stromu (guide tree) z nepřiložených sekvencí

Guide tree vs. phylogenetic tree

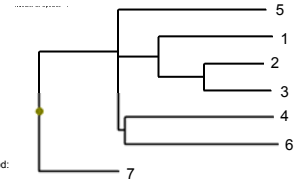
- Guide tree** je vypočítán na základě matice vzdáleností (distance matrix) vytvořené podle skóre pairwise alignmentů. Výstupem je .dnd soubor.
- Phylogenetic tree** je vypočten na základě vytvořeného MSA. Vzdálenosti mezi sekvencemi jsou vypočteny a uloženy jako .ph soubor. Následně je možno je využít pro konstrukci fylogenetického stromu (soubory .nj, .ph, .dst) pomocí zvolené metody (nj, phylip, dist).

.dnd soubor

```
(
(
(PAII:0.16435,
RSII:0.13654)
:0.03384,
(
(CVII:0.16563,
BCLB:0.26800)
:0.02264,
(
(BCLA:0.17899,
BCLD:0.26633)
:0.18717,
(BCLC:0.29707)
:0.03484);
);
);
);
```

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 = 0.9386; 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

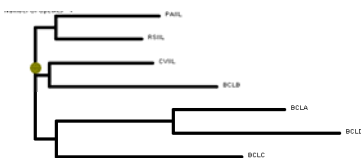
.nj soubor



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 (0.44258) 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)

.ph soubor

```
(
(
(PAII:0.34101,
RSII:0.17807,
CVII:0.19912)
:0.13706,
(
(
(BCLA:0.44518,
BCLD:0.45833)
:0.01864,
(
(BCLB:0.44298,
BCLD:0.47807)
:0.02083);
);
);
);
```



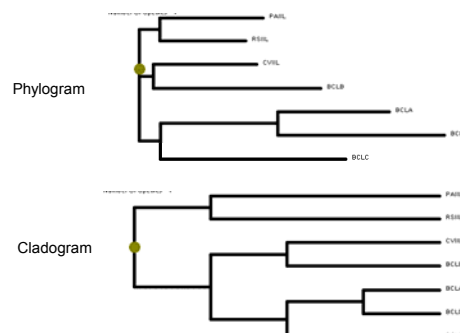
.dst soubor

```
7
PAII 0.000 0.649 0.684 0.930 0.904 0.939 0.982
RSII 0.649 0.000 0.377 0.912 0.895 0.912 0.939
CVII 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).

Phylogram a cladogram



Progresivní multiple alignment

- Používá ho většina programů
 - Vznik – 1987
Feng, D.-F. and Doolittle, R.F. (1987) J. Mol. Evol. 25, 351-360.
 - 1) sestavení příbuzenského stromu (guide tree) z nepříložených sekvencí
 - 2) tvorba párových alignmentů postupně podle příbuznosti (topologie guide tree)
- Dnes obsahuje často iterativní smyčku

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 -----
RS1IL -----
CV1IL -----
BCLB --LVEKLPQYVFDIATIPYDFWGSQKQVKTDAAGEVACTVFA.GAPVLP.GAAA
BCLC AIATNQDVADOCFFYSKVVPESTGRMPPFLVATIDVIGSIVTPYVQKMSVRSGLMIIDB
BCLA -----
BCLD LREITALRAEIVLPIEFALKDAGIVPIELVEVDAATVADADLLHPGCRPLEKHVM

PA1IL -----ATQGVFT
RS1IL -----AQQGVFT
CV1IL -----AQQGVFT
BCLB KFGVGVVW-----VFSEKATQPVPQAPVDP-----TQDGERDGIPT
BCLC YASLSALWQ-----TAAPSSQSSQWQAEVYDTQKNIQQGGERDGIPT
BCLA -----ASRQV-----SSRGAQEF
BCLD RSDVLAAGATTCADFAVCIDRSDVSGYFRWETSLEIAGSQVDTKQPFKPSDRGNHFS

PA1IL LPANTRFQVIAFANSSQTVNVLVNNETA--ATFSGQSTNNAVIGTVLNSSSGKVVV
RS1IL LPANTRFQVIAFANSSQTVNVLVNNETA--ATFSGQSTNNAVIGTVLNSSSGKVVV
CV1IL LPARINFDVTLVNSAATQVEIPIVDSNEP--AAFSGVDTQDNLGTVINSGS--DNVIV
BCLB LPVNIAPQVTLVNSAATQVEIPIVDSNEP--AAFSGVDTQDNLGTVINSGS--DNVIV
BCLC LPPELEFQVTLVNSAATQVEIPIVDSNEP--AAFSGVDTQDNLGTVINSGS--DNVIV
BCLA LPVNTFRALIFANAAKQKIKLPIQDSQEPAAVHKLTRDQPR--ATLANSQ--GKIFP
BCLD LPVNTFRALIFANAAKQKIKLPIQDSQEPAAVHKLTRDQPR--ATLANSQ--GKIFP

```

Programové balíky



- 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 dynamic programming,
- umožňují vložení mezer

Needle – globální pairwise alignment, Needleman-Wunsch algoritmus

Water – lokální pairwise alignment, Smith-Waterman algoritmus

Globálně podobné sekvence

```
Needle
PA-IIL 1 ATQGVFTLPANTRFGVTFANSSGTQTVNVLVNNETAATFSQGSTNNAVI 50
      |.|||||.....|.....|.....|.....|.....|.....|.....|.....|
RS-IIL 1 AQQGVFTLPANTSGVTFANANAANTQTIQVLVDNVVKATFPGSGTSDKLL 50
PA-IIL 51 GTQVLNSGSSGKVVQVSVNGRPSDLVSAQVILTNELNLFALVGSSEDTND 100
      |.....|.....|.....|.....|.....|.....|.....|.....|.....|
RS-IIL 51 GSQVLNSG-SGAIKIQVSVNGKPSDLVSNQTLANKLNFMVGSSEDTND 99
PA-IIL 101 DYNDAVVVINWPLG 114
      |.....|.....|.....|
RS-IIL 100 DYNDGIAVLNWLPG 113

Water
PA-IIL 1 ATQGVFTLPANTRFGVTFANSSGTQTVNVLVNNETAATFSQGSTNNAVI 50
      |.|||||.....|.....|.....|.....|.....|.....|.....|.....|
RS-IIL 1 AQQGVFTLPANTSGVTFANANAANTQTIQVLVDNVVKATFPGSGTSDKLL 50
PA-IIL 51 GTQVLNSGSSGKVVQVSVNGRPSDLVSAQVILTNELNLFALVGSSEDTND 100
      |.....|.....|.....|.....|.....|.....|.....|.....|.....|
RS-IIL 51 GSQVLNSG-SGAIKIQVSVNGKPSDLVSNQTLANKLNFMVGSSEDTND 99
PA-IIL 101 DYNDAVVVINWPLG 114
      |.....|.....|.....|
RS-IIL 100 DYNDGIAVLNWLPG 113
```

Lokálně podobné sekvence

```
Needle
PA-IIL 1 ----- 0
BCLB 1 SQPFTHDDLVALQLAGNDATAVQANGDQAVLDRMRQFMTAQLVEKLPQ 50
PA-IIL 1 ----- 0
BCLB 51 DVFVDIATIPVSPDVGSGNKNVKTDAAGEVVACTVWAGAPGVLPGAANK 100
PA-IIL 1 -----ATQGVFTLPANTRFGVTFANSS 22
      |.....|.....|.....|
BCLB 101 FGVGAVNMFYSKATPQPVPAPVPTGGGERDGIPTLPPNIAFGVLTALVNS 150
PA-IIL 23 SGTQTVNVLV--NNETAATFSQGSTNNAVIQTQVLNSGSSGKVVQVSVN 70
      |.....|.....|.....|.....|.....|.....|.....|.....|.....|
BCLB 151 SAPQTIIEVFVDDNPKPAATFQAGTQDANLNTQIVNSG-KGKVRVVVTAN 199
PA-IIL 71 GRPSDLVSAQVILTNELNLFALVGSSEDTNDYNDYNDYNDYNDYNDYND 114
      |.....|.....|.....|.....|.....|.....|.....|.....|.....|
BCLB 200 GRPSKIGSRQVDIFKTKYFGLVGSSEDTNDYNDYNDYNDYNDYNDYND 243

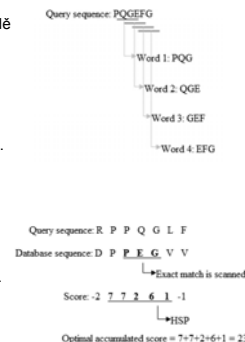
Water
PA-IIL 4 GVFTLPANTRFGVTFANSSGTQTVNVLV--NNETAATFSQGSTNNAVI 51
      |.....|.....|.....|.....|.....|.....|.....|.....|.....|
BCLB 132 GFTLPPNIAFGVLTALVNSAPQTIIEVFVDDNPKPAATFQAGTQDANLN 181
PA-IIL 52 TQVLNSGSSGKVVQVSVNGRPSDLVSAQVILTNELNLFALVGSSEDTND 101
      |.....|.....|.....|.....|.....|.....|.....|.....|.....|
BCLB 182 TQVLNSG-KGKVRVVVTANVNSAPQTIIEVFVDDNPKPAATFQAGTQDANLN 230
PA-IIL 102 DYNDAVVVINWPLG 114
      |.....|.....|.....|
BCLB 231 YNDGIAVLNWLPG 243
```

BLAST algoritmus

Heuristický algoritmus jehož základem je hledání slov (několikapísmenných sekvencí), s dostatečnou podobností (poskytují dostatečně vysoké skóre v substituční matici)



- **Tvorba k-písmenných slov ze vstupní sekvence**
pro proteiny typicky 3-písmenných (v případě DNA 11-písmenných)
- **Porovnání slov na základě substituční matice**
algoritmus BLAST hledá na základě vloženého skóre slova, která jsou podobná každému slovu v zadané sekvenci. Vyhovující slova jsou následně uspořádána.
- **Prohledání databázových sekvencí**
Je hledána shoda s nalezenými vysoce podobnými slovy.
- **Rozšíření slov na segmenty**
Přesné shody slov s databázovými sekvencemi jsou rozšiřovány oběma směry. To pokračuje dokud skóre pro tuto dvojici sekvencí je dostatečně vysoké.



Novější verze BLASTu (BLAST2) má mj. níže nastavenou hranici pro hledání podobných slov, což rozšiřuje možnost nalezení vzdálenějších homologů.

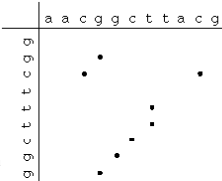
FASTA algoritmus

FastA algoritmus nejprve provádí rychlé prohledání pro nalezení odpovídajících sekvencí, následuje přesnější porovnání zadané sekvence s databázovou sekvencí. Na rozdíl od algoritmu BLAST jsou zde tolerovány mezery.

Proces:

Obě porovnávané sekvence tvoří horizontální a vertikální osu grafu.
Následně jsou jednotlivá slova z jedné sekvence porovnávána se slovy sekvence druhé.
Odpovídající páry pak vytvoří sadu bodů. Body na úhlopříčce signalizují významnou shodu (či podobnost) v daném úseku. Cílem je nalezení nejdelšího shodného úseku (úseku s nejvyšším skóre).

V dalších krocích jsou zahrnuty konzervativní změny pro nejlepší úseky z prvního prohledání. Program pak vyhledává možnost spojení více takových úseků (může mezi nimi být mezera, či jsou na různých diagonálách) a tyto spojené úseky jsou posuzovány z hlediska zadaných kritérií



Příklad porovnání sekvencí
GGCTTTCGG a
AACGGCTTACC

MSA „programy“

- Za posledních 10 let vzniklo přes 50 MSA programových balíčků (Wallace, I. M., O'Sullivan, O., Higgins, D. G. and Notredame, C. (2006) M-Coffee: combining multiple sequence alignment methods with T-Coffee. Nucleic Acids Res. 34: 1692-1699.)
- Clustal W (Thompson et al., 1994)
- Clustal X (Thompson et al., 1997)
- Dialign2 (Morgenstern, 1999)
- T-Coffee (Notredame et al., 2000)
- MAFFT (Kato et al., 2002)
- MUSCLE (Edgar, 2004)
- Kalign (Lassmann, 2005)
- ...

Clustal <http://www.ebi.ac.uk/clustalw/>

- V současné době nejužívanější program
- První verze 1988
Higgins, D.G. and Sharp, P.M. (1988) CLUSTAL: a package for performing multiple sequence alignment on a microcomputer. Gene, 73: 237-244.
- Dnes používané verze:
Clustal W (Thompson et al., 1994)
Clustal X (Jeanmougin et al., 1998)
- Využívá progresivní alignment

ClustalW: Jednotlivým sekvencím přiřazuje váhy (weight – W) podle četnosti zastoupení (čím více jsou si sekvence podobné, tím nižší mají váhu a naopak) a penalizuje přítomnost mezer v závislosti na jejich pozici (position-specific gap penalties)

ClustalW2 – postup

1. Provedení **pairwise alignmentů** pro každou dvojici sekvencí a určení jejich podobnosti – v závislosti na množství neodpovídajících residuí a mezer
2. Sestavení **příbuzenského stromu** (similarity tree)
3. **Kombinace** alignmentů (viz. 1.) v pořadí dle příbuznosti – od nejvíce podobných k nejméně příbuzným (viz. 2.). Jednou vložené mezery jsou zachovány.

Clustal W/Clustal X

Pod alignmentem je uváděn tzv. **consensus** – dohodnuté symboly vyjadřující „konzervovanost“ každého sloupce:

- * - identické residuum ve všech sekvencích
- : - silně konzervovaný sloupec
- . - slabě konzervovaný sloupec

```

IPNTDFRALFFANAAEQOHIKLFIGDSQEPAAHYKLTTRDGER--ATLNSGNGKIRFE
LPNTAFKALFYANAADROQLKLFIDDAPEPAATFVGNSEDGVRLL--FTLNSKGGKIRIE
LPPNIAFGVTALVNSSAPQIIEVVDNPKPAATFGAGTQDANLNTQIVNSGKGVVIVV
LPPHIKFGVTALTHAANDOTIIDIYIDDDPKPAATFKGAGAQDNLGTVLDSGNGRVRVI
::*: . . . . . : . . . . . * . . . . . * . . . . . * . . . . .

```

MUSCLE



(Multiple Sequence Comparison by Log-Expectation)
<http://www.drive5.com/muscle>

Rychlejší určení „vzdálenosti“ dvou sekvencí
Tzv. log-expectation skórovací funkce
Refinement metodou restricted partitioning

Vhodný i pro velký počet sekvencí (5000 seq po 350 bp za 7 min na PC – rok 2004)

Postup:

- Sestavení matice pro každou dvojici sekvencí, určení jejich „vzdálenosti“ a sestavení matice vzdáleností (distance matrix)
- Na základě distance matrix je sestaven první příbuzenský strom (tree1)
- Skládání sekvencí v pořadí dle tree1 od větvi ke kmenu – v každém rozvětvení je vytvořen profil, který při dalším porovnávání nahrazuje původní sekvence – výsledkem je první MSA

Algoritmus MUSCLE (podobne PRRP a MAFFT)

- Přepočítání vzdáleností sekvencí na základě vzniklého MSA1 – tvorba druhé distance matrix (D2)
 - Na základě D2 sestaven vylepšený příbuzenský strom (tree2)
 - Progresivní alignment (viz bod 3) na základě tree2 – vytvoření druhého MSA
-
- Refinement** – rozdělení vzniklého stromu na dvě části a vytvoření MSA pro každou z nich. Pokud je výsledný alignment lepší, je zachován. Toto se opakuje do konvergence (žádná další změna nevede k lepšímu výsledku) nebo do určeného počtu kroků

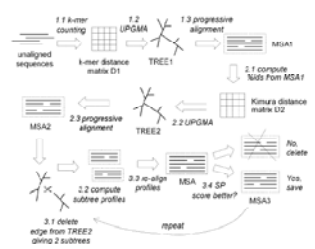


Figure 2. This diagram summarizes the flow of the MUSCLE algorithm. There are three main stages: Stage 1 (draft progressive), Stage 2 (improved progressive) and Stage 3 (refinement). A multiple alignment is available at the completion of each stage, at which point the algorithm may terminate.

Další skórovací schémata (scoring schemes) pro pairwise alignment

Algoritmy založené na matici (matrix-based algorithms) – např. ClustalW, MUSCLE; pomocí substituční matice je příslušné dvojici (AK) přiřazena hodnota. Rozhoduje pouze **identita** těchto dvou AK, případně jejich **nejbližší okolí** (viz. např. BLAST)

Schémat založená na konzistenci (consistency-based schemes) – poprvé v T-Coffee, dále v PCMA, ProbCons, MUMMALS, MAFFT, aj. Vychází z nejlepších možných alignmentů každé dvojice sekvencí. Využívá často i **data z různých zdrojů** (např. strukturní informace). Cílem je dosáhnout maximální konzistence (vnitřní shody). Výsledek je přesnější, ale výpočet je časově náročnější.

T-Coffee

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

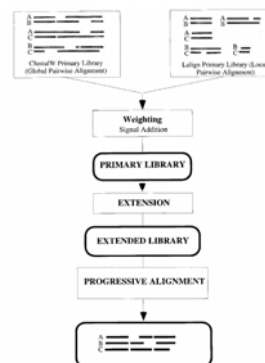


- Pomalejší ale výrazně přesnější než ClustalW
- Je schopen kombinovat data z více předchozích alignmentů, které mohly být vytvořeny různými postupy (lokální, globální, strukturní podobnost, ...)

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

T-Coffee

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



Zlepšení přesnosti – strukturní informace

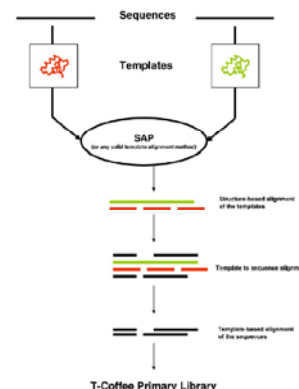
- Sekvence s vyšší homologií (>40%) – vysoká přesnost alignmentu
- Bez homologie – nepoužitelné
- Tzv. twilight zone – málo podobné sekvence (nižší než 20% homologie) = špatná (méně než 30%) přesnost alignmentu

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 homologu a použije je jako šablony pro následný alignment zadaných sekvencí pomocí metod MSA založených na struktuře (např. SAP, Fugue).

Benchmark (srovnávací testy)

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

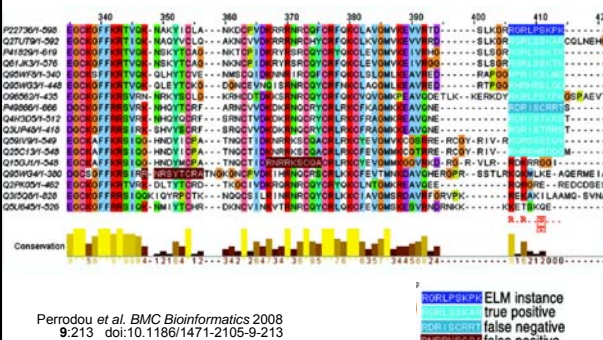
Na základě srovnání 3D struktur byly vytvořeny další sady:

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



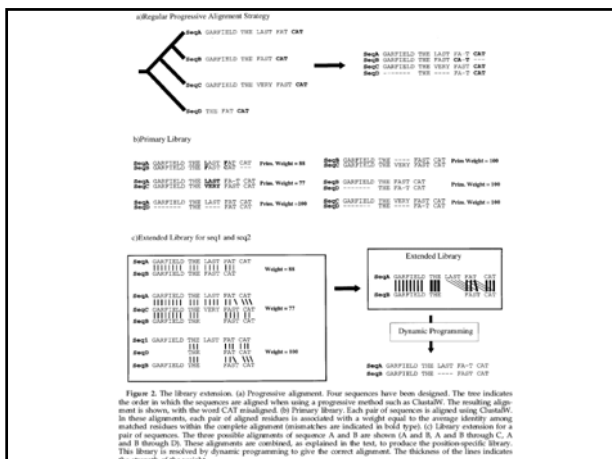
Zopakování / shrnutí

- ▼ **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 nejpřesnější – 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



BAIBASE [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 sequences, with less than 20% identity between the two groups. Ref4 is concerned with long terminal exons, while Ref5 test cases contain large internal insertions and deletions. Last sets from references 6-8 deal with problems like transmembrane regions, inversed domains, and repeat sequences. In previous versions of BAIBASE, 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 duplicate 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.

OxBench [Raghava 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 236 alignments, and Align-IT was only able to align 107, using a single processor with 4GB of RAM.

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

SABmark [Van Walle et al., 2000] is divided into two subsets. Each test group in the SUPERFAMILY set represents a SCOP superfamily, whose sequences are 25%-100% identical. Each test group in the TWILIGHT set represents a common SCOP fold and sequences are 0.25% identical. In addition, these two subsets are also provided with non-homologous (near positive) sequences included within each group. Instead of a single alignment acting as a reference, SABmark 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 [Edwards et al., 2005] test cases contain a number of simulated motifs [Eddy et al., 1998] inserted into otherwise random (unsignal) sequences, and as such is 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 [McGuinly 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	---	61.88 [12]	---	http://www.ebi.ac.uk/blast/
Kalign	Mean	---	63.00 [10]	---	http://ma.rzh.bsu.ru/
MUSCLE [6]	Mean	---	68.00 [16]	45.0 [10]	http://www.drive5.com/muscle/
T-Coffee [17]	Consistency	---	69.97 [12]	44.0 [10]	http://www.tcoffee.org/
ProbCons [9]	Consistency	---	70.54 [12]	---	http://procons.stanford.edu/
MAFFT [8]	Consistency	---	72.25 [12]	---	http://align.genome.jp/mafft/
RCoffee [12]	Consistency	---	72.91 [12]	---	http://www.tcoffee.org/
MEMMA [15]	Consistency	---	73.15 [16]	---	http://pendra.sommed.edu/maemem/ma/
DiscJuret [24]	Profiles	---	---	---	http://bioparis.inria.fr/ProfilAlign/
PRALINE [5]	Mean	Profiles	---	50.2 [10]	http://lsc.usc.edu/~jgarnier/programs/psd/psdnewer/
PROFALS [16]	Consistency	Profiles	79.00 [16]	---	http://pedra.sommed.edu/profals/
SPM [20]	Mean	Profiles	77.00 [20]	---	http://igork.informatics.lupulandia/Softwares/Service/SPM/SPM.htm
Expasiv [13]	Consistency	Structures	77.18 [11]*	---	http://www.tcoffee.org/
ELAND [23]	Consistency	Structures	---	---	https://www.mil.biorh.de/te/LEA/

Validation values were computed from random datasets, not selected for comparability. Prefab validation values were made using Prefab version 3. HOMSTRAD validation values were made on datasets having less than 30% identity. The score of each value is included by the accompanying reference citation. *The Expasiv value comes from a slightly more demanding subset of HOMSTRAD (HOM2) made of sequences less than 22% identical. <http://www.tcoffee.org/2005/03/05/>

Table 1: Programs used in this investigation

Method	OVERVIEW
Align_m[2,3]	http://icr.informatics.vub.ac.be/software/software.html
[Van Walle et al., 2004]	Local, specialised for highly divergent sequences.
ClustalW [1,8]	http://www.ebi.ac.uk/blast/
[Thompson et al., 1994]	Global, progressive alignment package
Dialign2[2,2]	http://ribosome.tachikawa-lab.nu/bioinf/dialign/
[Rognes et al., 1999]	Local, aligns segments of sequences rather than individual residues.
Dialign-[5,1,3]	http://align-l.gubios.de/
[Edwards et al., 2005]	Local, progressive alignment. Recent re-implementation of Dialign2.
MAFFT [6,5,9]	http://www.bioinformatics.jp/~asap/~kashi/program/align/mafft/
[Kato et al., 2002]	Suite of alignment programs:
FFTNS	Global, uses Fast Fourier Transform to generate tree.
FFTNS	As FFTNS, but with iteration step to refine alignment.
NWNS	Global, uses traditional Needleman-Wunsch algorithm.
NWNS	As NWNS, but with iteration step to ref.
FINS	Local, iterative, uses local pairwise align
	GNNS
	Global, iterative, uses global pairwise alignment information.
	MUSCLE [3,6]
	http://www.drive5.com/muscle/
	Global, iterative, progressive alignment program that uses Log Expectation as scoring function.
	ProbCons [1,9]
	http://procons.stanford.edu/
	[Do et al., 2002]
	Global, uses posterior probabilities from HMMs and pairwise alignment consistency.
	PCMA [2,10]
	http://hpc.sommed.edu/pub/PCMA/
	Global, switches alignment strategies dependent on sequence data
	[Wu et al., 2002]
	Local, uses Partial Order graphs
	POA [4,2]
	http://www.informatics.usda.edu/psw/
	Local, uses Partial Order graphs
	T-COFFEE [1,3,7]
	http://lsc.usc.edu/~jgarnier/Projects/home_page/T_Coffee_home_page.html
	[Wolfdame et al., 2000]
	Combines both global and local methods; uses consistency

Blackshield 2006 oznacil ProbCons jako najbolji na zaklade 6 benchmarkovych testu