

Aminokyseliny

glycin	alanin	valin	leucin	ileucin	asparagová kys.	asparagin	glutamová kys.	glutamin	arginin	lysin	histidin	fenylalanin	serin	threonin	tyrosin	tryptofan	metionin	cytein	prolin	seleocytein	pyrolysin
G	A	V	L	I	D	N	E	Q	R	K	H	F	S	T	Y	W	M	C	P	U	O

Třídění aminokyselin

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

CC(C)C(N)C(=O)O
Isoleucine

CC(C)C(C)C(N)C(=O)O
Leucine

Nukleové báze

Adenine
Nc1ncnc2[nH]cnc12

Cytosine
Nc1cc[nH]c(=O)n1

Guanine
Nc1nc2[nH]cnc2c(=O)[nH]1

Thymine
Cc1c[nH]c(=O)[nH]c1=O

Uracil
O=c1cc[nH]c(=O)[nH]1

Nukleová báze
Adenin

Nukleosid
Adenosin

Nukleotid
Adenosinmonofosfát
AMP

Nukleotid
Adenosintrifosfát
ATP

adenin	cytosin	guanin	thymin	uracil
A	C	G	T	U

Polysacharidy

Komplikované sekvence – alignment se neprovádí

Polymer	Protein	Nukleová kyselina	Polysacharid
Počet druhů základních stavebních jednotek	20 (22)	4 (DNA) 4 (RNA)	desítky
Počet typů vzájemných vazeb	1	1	2 x 4 (pro hexosu)

Struktura proteinů (NK)

ADSQTSSNRAGEFSIPPNTDFRAIF
FANAEEQQHKLFIGDSQEPAAAYHK
LTTDRGPREATLNSGNGKIRFEVSV
NGKPSATDARLAPINGKKSDBGSPF
TVNFGIVVSEDDGHSDYNDGIVVL
QWPIG

primární (sekvence)

Cont:

sekundární

terciární

kvarterní

Alignment

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

Význam alignmentu

- Identifikace sekvence v databázi
- Hledání podobných sekvencí v databázi
- Detekce mutací
- Hledání konzervovaných částí sekvence
- Odhalování příbuzenských vztahů
- Předpověď funkce makromolekuly
- Předpověď vyšších struktur

Typy alignmentu

Pairwise alignment – dvě sekvence

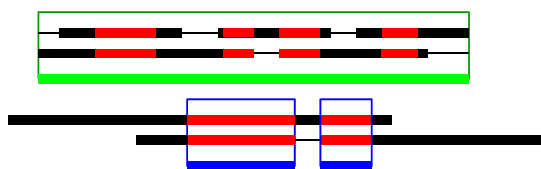
```
WLA K A K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M
W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M
```

Multiple sequence alignment – více sekvencí

```
W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M
W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M
W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M
W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M
W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M W L A K A L K Y L M E T A Q A S S I S T E L A R H H P R A V D A K R K S E M K R K T A M
...
```

Pair-wise alignment

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



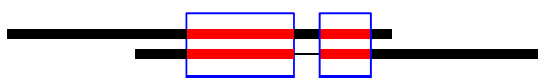
Global alignment

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 příliš neodpovídají.



Local alignment

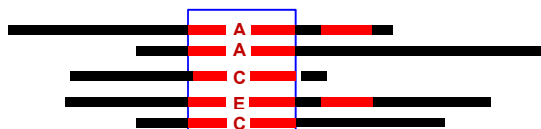
Hledá úseky dvou sekvencí, které si podle zvolených kritérií dobře odpovídají. Nesnaží se zahrnout celé sekvence, pokud si jejich některé části neodpovídají.



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 známých 3D struktur

- BLOSUM62 – znamená, že ke konstrukci matrice byly použity proteiny s průměrnou identitou 62%.



- A - C = 4
- A - E = 2
- C - E = 2
- A - A = 1
- C - C = 1

- výskyt každého páru AA v každém sloupci každého bloku je sečten
- čísla získána ze všech bloků slouží pro výpočet BLOSUM matricí

Číslování BLOSUM jde v obráceném pořadí oproti PAM

– čím menší číslo, tím odlišnější sekvence byly použity

Matrix	Best use	Similarity (%)
Pam40	Short highly similar alignments	70-90
PAM160	Detecting members of a protein family	50-60
PAM250	Longer alignments of more divergent sequences	~30
BLOSUM90	Short highly similar alignments	70-90
BLOSUM80	Detecting members of a protein family	50-60
BLOSUM62	Most effective in finding all potential similarities	30-40
BLOSUM30	Longer alignments of more divergent sequences	<30

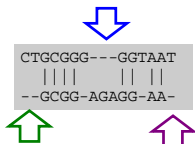
Similarity column gives range of similarities that the matrix is able to best detect.

Mezery (Gaps)

Příčiny vzniku mezer:

- Bodová mutace** (velmi častá příčina)
- Nepřesný crossover při meióze (inzerce nebo delecce 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“, často 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:

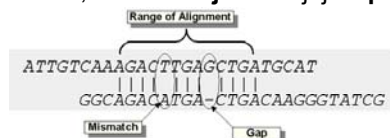
```
ATCTTCAGTGTTCCTCCCTGTTTGGCCC-ATTTAGTTCGCTC
|||||
ATCTTCAGTGTTCCTCCCTGTTTGGCCCATTTAGTTCGCTC
```

Dlouhá mezera:

```
ATCTTCAGTGTTCCTCCCTGTTTGGCCC-----ATTTAGTTCGCTC
|||||
ATCTTCAGTGTTCCTCCCTGTTTGGCCCGCCCCCCCCCCCCCCCCCAATTTAGTTCGCTC
```

Skóre

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



$$S = \sum(\text{identities, mismatches}) - \sum(\text{gap penalties})$$

$$\text{Score} = \text{Max}(S)$$

Čím vyšší je skóre, tím vyšší je podobnost.

Podle použité matice může být skóre i záporné.

Příklad výpočtu

AAEECCDDEF
AADDKKKEFGG

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

Například, celkové pozitivní skóre na úrovni jednotlivých aa

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

Celkové skóre 32 – 24 = 8

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.

Multiple sequence alignment - MSA

(mnohonásobné přiložení)

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 – náročné na paměť a čas, nevhodné pro více než 3-4 sekvence (n=rozměrný prostor)
- Progresivní alignment** (progressive sequence alignment) – nejčastěji používaný k vytvoření alignmentu; využívá **fylogenetické informace** – hierarchický, nejdříve identifikuje nejpodobnější sekvence a následně inkorporuje ostatní
- Iterativní alignment** (iterative sequence alignment) – odstraňuje problémy progresivního alignmentu, který je závislý na prvotním přiložení nejpodobnějších sekvencí pomocí **opakování alignmentu** pro podskupiny sekvencí následující po globálním alignmentu
- Hledání motivů – nalezení částí konzervovaných sekvenčních motivů pomocí globálního přiložení a následně „hodnocení“ těchto úseků nezávisle na celé sekvenci

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

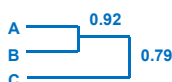


Progresivní multiple alignment

- Používá ho většina programů
- Vznik – 1987
Feng, D.-F. and Doolittle, R.F. (1987) J. Mol. Evol. 25, 351-360.

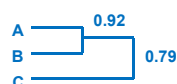
1) sestavení příbuzenského stromu (guide tree) na základě distanční matice (distance matrix) z nepřiložených sekvencí

A	-		
B	0.92	-	
C	0.65	0.79	-
	A	B	C



Počet exaktně stejných shod dělená celkovou délkou sekvence (ignoruje mezery)

Progresivní multiple alignment



Nejdříve provede pairwise alignment A a B
Pak přidá sekvenci C do předešlého alignmentu
(inzerce mezer, pokud je potřeba)

2) tvorba párových alignmentů postupně podle příbuznosti (topologie guide tree)

- Dnes obsahuje často iterativní smyčky

Guide tree vs. phylogenetic tree

- **Guide tree** je vypočítán na základě matice vzdáleností (distance matrix) vytvořené podle skóre pairwise alignmentů. Výstupem je .dnd soubor. [NEMA fylogenetický význam](#)
- **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).

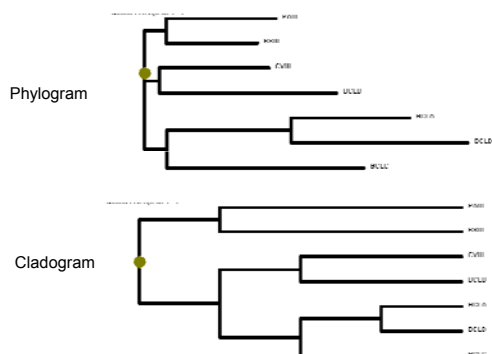
.dnd soubor

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

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



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 kotvících bodů v místě dobré shody
- Následný globální alignment pak tyto odpovídající úseky sekvencí zahrnuje (využito např. v ClustalW2)

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 dynamické programování
- umožňují vložení mezer

Needle – globální pairwise alignment, Needleman-Wunsch algoritmus

Water – lokální pairwise alignment, Smith-Waterman algoritmus

Nelze však spoléhat na zdánlivě dobrá řešení

```

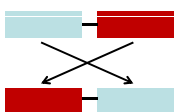
PLLSASIVSAPVVTSETYVDIPLGYLDVAKAGIRDGKLVILNVPTPYATGNNFPGIYFAIATNQGVADGCFYSSKY
PESTRGRMPFTLVATIDVSGVTFVKGQWKSVRGSAMHIDSYASLSAIWGTAAAPSSQSGNQGAETGGTAGNIG
GGGERDGFNLPPHIKFGVLTALHAANDQTIIDDDPKPAATFKGAGAQDQNLGKVLDSGNGRVRVIMMANGR
PSRLGSRQVDIFKKSYPFGIGSEGDADDYNDGIVFLNWPLG

```

```

ERDGTFLNPPHIKFGVLTALHAANDQTIIDDDPKPAATFKGAGAQDQNLGKVLDSGNGRVRVIMMANGRPSR
LGSRQVDIFKKSYPFGIGSEGDADDYNDGIVFLNWPLGPLLSASIVSAPVVTSETYVDIPLGYLDVAKAGIRDGKLV
ILNVPTPYATGNNFPGIYFAIATNQGVADGCFYSSKYPESTRGRMPFTLVATIDVSGVTFVKGQWKSVRGSAM
HIDSYASLSAIWGTAAAPSSQSGNQGAETGGTAGNIGGGKLAALAEIKRASQPELAPEDPEVHHHHHHH

```



EMBOSS_001	1	...	100
EMBOSS_001	1	ERDGTFLNPPHIKFGVLTALHAANDQTIIDDDPKPAATFKGAGAQDQNLGKVLDSGNGRVRVIMMANGRPSR	100
EMBOSS_001	1	PLLSASIVSAPVVTSETYVDIPLGYLDVAKAGIRDGKLVILNVPTPYATGNNFPGIYFAIATNQGVADGCFYSSKY	100

BLAST algoritmus

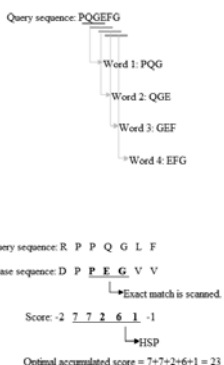
BLAST (Basic Local Alignment Search Tool)

The BLAST Search Algorithm

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 dané 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.

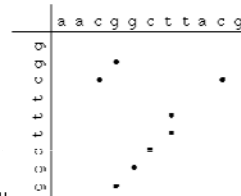


Novější verze BLASTu (BLAST2) má mj. níže nastavenou hladinu pro hledání podobných slov, což rozšiřuje možnost nalezení vzdálenějších homologů.

FASTA algoritmus

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). 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 GGCTTCGG a AACGGCTTACC

MSA „programy“

- Za posledních 15 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

```

I P P N T D F R A I F F A N A A E Q Q H I K L F I G D S Q E P A A Y H K L T T R D G E R E -- A T L N S G N G K I R F E
L P P N T A F K A I F Y A N A A D R Q D L K L F I D D A P E P A A T F V G N S E D G V R L -- F T L N S K G G K I R I E
L P P N I A F G V T A L V N S S A P Q T I E V F V D D N P K P A A T F Q G A G T Q D A N L N T Q I V N S G K G K V R V V
L P P H I K F G V T A L T H A A N D Q T I D I Y I D D D P K P A A T F K G A G A Q D C N L G T K V L D S G N G R V R V I
: * * : * . . . . . : : : : * . * * * . . . . . : : . : : * * : * :
  
```

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**.

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 !!!**

