

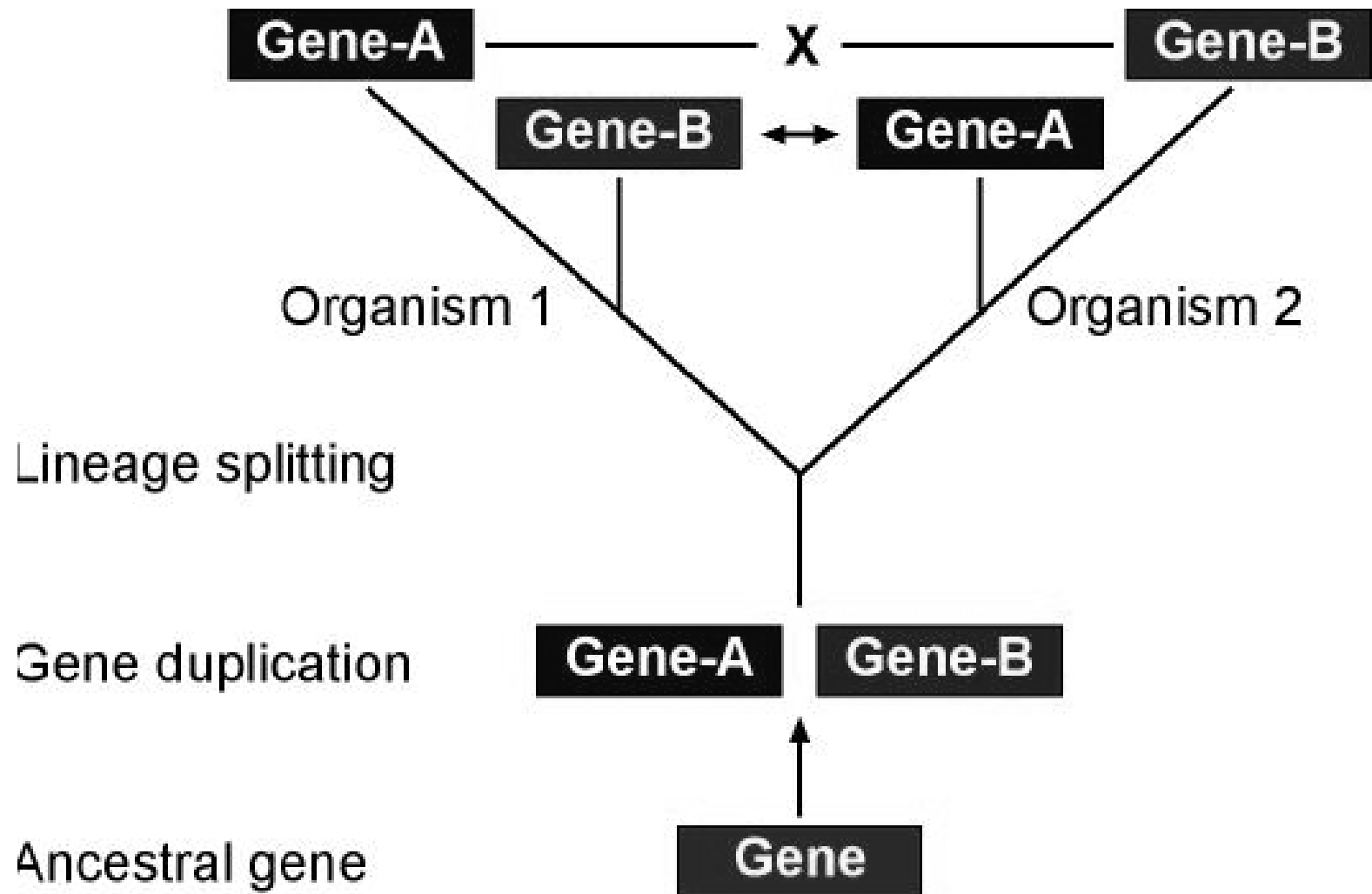
Podstawy bioinformatyki dla biotechnologów

Wykład 3 alignment

Porównywanie sekwencji

Homologia, podobieństwo i analogia

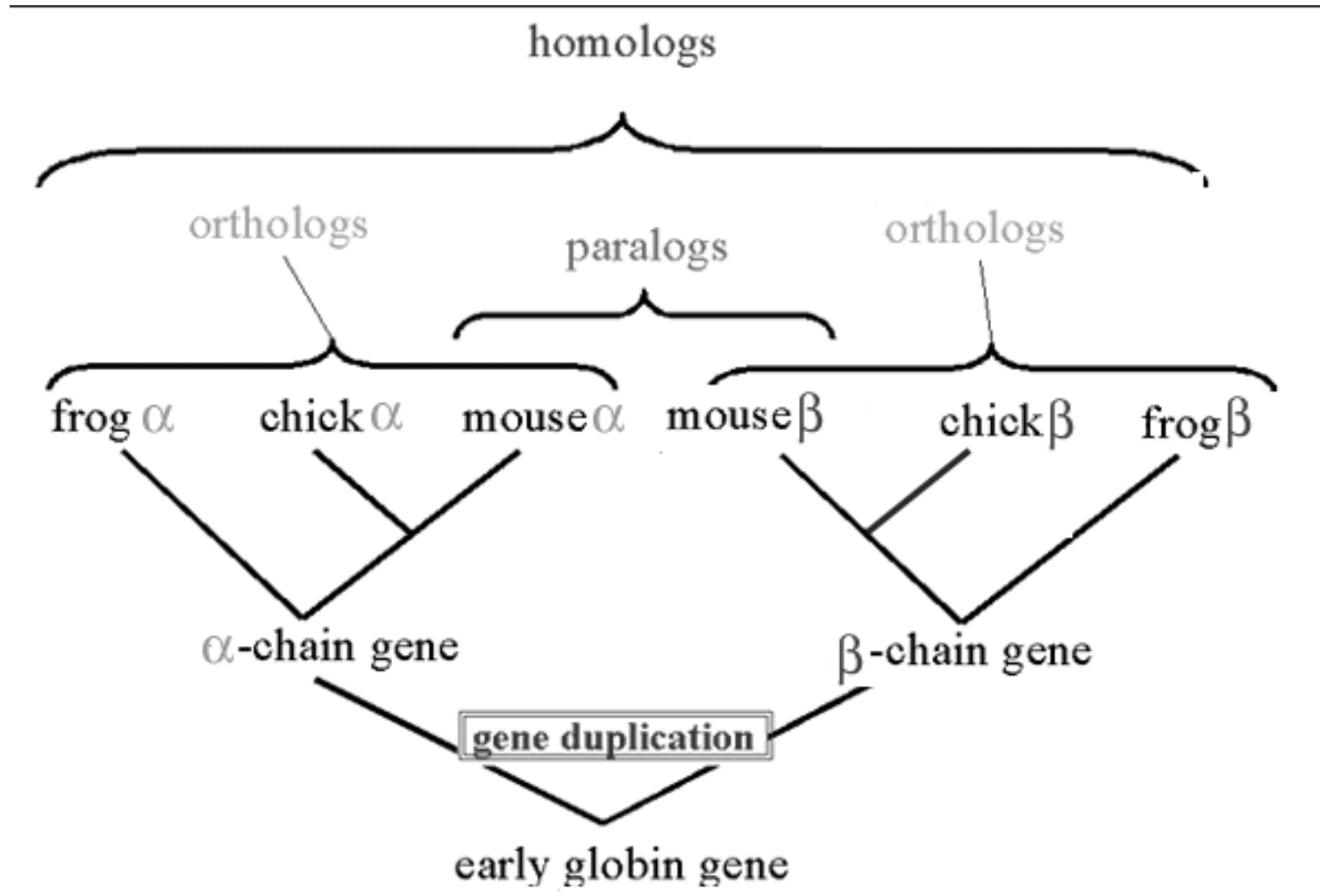
Duplikacja, specjacja



Homologi

- **Ortologi** – homologiczne geny, których rozdzielenie nastąpiło na skutek specjacji, czyli rozdzielenia gatunków, lub rzadziej horyzontalnego transferu genu. Geny ortologiczne mają zwykle taką samą, albo zbliżoną funkcję.
- **Paralogi** – geny pochodzące od wspólnego przodka, rozdzielone w wyniku duplikacji genu. Paralogi mają często różne funkcje w organizmie. Przykładem mogą być mioglobina i hemoglobina u człowieka.

Homo-, para-, orto-, analogi



dopasowanie sekwencji

- Dopasowanie/porównywanie
- Uliniowanie
- Alignment

W bioinformatyce, dopasowanie sekwencji jest sposobem dopasowania struktur pierwszorzędowych DNA, RNA, lub białek do zidentyfikowania regionów wykazujących podobieństwo, mogące być konsekwencją funkcjonalnych, strukturalnych, lub ewolucyjnych powiązań pomiędzy sekwencjami. Zestawione sekwencje nukleotydów lub aminokwasów są zazwyczaj przedstawione jako wiersze macierzy. Pomiędzy reszty wprowadzane są przerwy, tak że reszty zbliżonych do siebie sekwencji tworzą kolejne kolumny.

Jeśli dwie dopasowywane sekwencje mają wspólne pochodzenie, niedopasowania mogą być interpretowane jako mutacje punktowe, a przerwy jako indele (mutacje polegające na delecji lub insercji), które zaszły w jednej lub obu liniach od czasu, kiedy obie sekwencje uległy rozdzieleniu. W przypadku dopasowywania sekwencji białek, stopień podobieństwa pomiędzy aminokwasami zajmującymi konkretną pozycję, może stanowić zgrubną miarę tego, jak konserwatywny jest dany region lub motyw. Brak substytucji lub obecność jedynie konserwatywnych substytucji (tj. zamiany reszty na inną, ale o podobnych właściwościach chemicznych) w określonym regionie sekwencji sugeruje, że jest on ważny strukturalnie lub funkcjonalnie. Dopasowywanie sekwencji może być także stosowane dla sekwencji pochodzenia poza biologicznego, np. danych finansowych lub sekwencji występujących w językach naturalnych.

Masur i inni, Dopasowanie sekwencji, Wikipedia 11.2009

alignment

```
AAB24882 TYHMCQFHCRYVNNHSGEKLIECNERSKAFSCPSHLQCHKRRQIGEKTHEHNQCGKAFPT 60
AAB24881 -----YECNQCCKAFAQHSSSLKCHYRTHIGEKPYECNQCCKAFSK 40
          ****: .***: * *:* ** * :*****.:* *****..
```

```
AAB24882 PSHLQYHERTHTGEKPYECHQCQAFKKCSLLQRHKRTHHTGEKPYE-CNQCCKAFAQ- 116
AAB24881 HSHLQCHKRTHHTGEKPYECNQCCKAFSQHGLLQRHKRTHHTGEKPYMNVINMVKPLHNS 98
          **** * :*****:***:**. : .*****: : *.: :
```

```

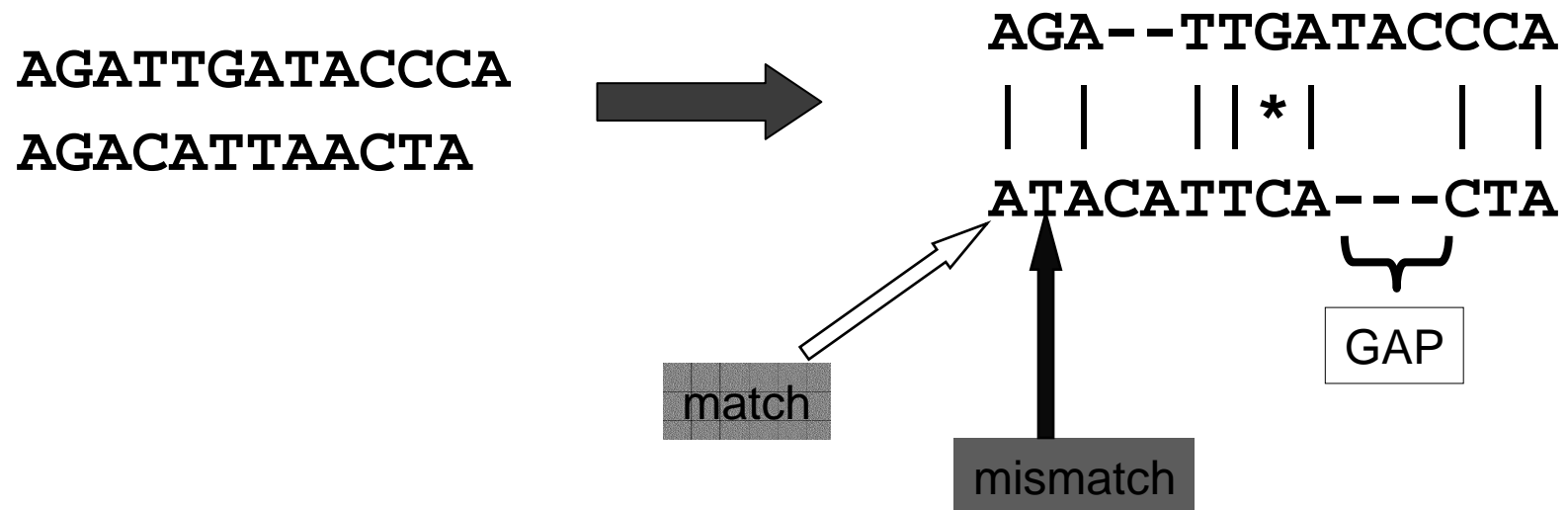
          *           :           *           : : :
Q5E940_BOVIN -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMKNTMMRKAIRGHLENN--PALE 76
RLA0_HUMAN -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMKNTMMRKAIRGHLENN--PALE 76
RLA0_MOUSE -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMKNTMMRKAIRGHLENN--PALE 76
RLA0_RAT -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMKNTMMRKAIRGHLENN--PALE 76
RLA0_CHICK -----MPREDRATWKSNYFMKIIQLLDDYPKCFVVGADNVGSKOMQQIRMSLRGK-AVVLGMKNTMMRKAIRGHLENN--PALE 76
RLA0_RANSY -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMKNTMMRKAIRGHLENN--SALE 76
Q7ZUG3_BRARE -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQTIIRLSLRGK-AVVLGMKNTMMRKAIRGHLENN--PALE 76
RLA0 ICTPU -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQTIIRLSLRGK-AIVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0_DROME -----MVRENKAAWKAQYFIKVVLELDFEPKCFIVGADNVGSKOMQNIIRTSLRGL-AVVLGMKNTMMRKAIRGHLENN--PQLE 76
RLA0_DICDI -----MSGAG-SKRKLFIEKATKLFTTYDKMIVAEADLVGSSQLQKIRKSIIRGI-GAVLMGKNTMIRKVIIRDLADSK--PELD 75
Q54LP0_DICDI -----MSGAG-SKRKNVFIKATKLFTTYDKMIVAEADLVGSSQLQKIRKSIIRGI-GAVLMGKNTMIRKVIIRDLADSK--PELD 75
RLA0_PLAF8 -----MAKLSKQKKQMYIEKLSLIQQYSKILIVHVDNVGSKOMASVRKSLRGK-ATILMGKNTIRIRTAALKKNLQAV--PQIE 76
RLA0_SULAC -----MIGLAVTTTKIAKWKVDEVAELTEKLKTHKTIIIANIEGFPADKLHEIRKKLRGK-ADIKVTKNNLFNIALKNAG----YDTK 79
RLA0_SULTO ----MRIMAVITQERKIAKWKIEEVKELEQKLRHYHTIIIANIEGFPADKLHDIRKKMRGM-AEIKVTKNTLFGIAAKNAG----LDVS 80
RLA0_SULSO ----MKRLALALKQKVASWVLEEVKELTELKNSNTILIGNLEGFPADKLHEIRKKLRGK-ATIKVTKNTLFGIAAKNAG----IDIE 80
RLA0_AERPE MSVVSIVGQMYKREKPIPEWKTMLRELEELFSKHRVVFADLTGTPTFVVQVRVKKLWKK-YPMMAVAKKRIILRAMKAAGLE---LDDN 86
RLA0_PYRAE -MMLAIGKRRYVRTQYPARVKIVSEATELLQKYPYVFLFDLHGLSSRILHEYRYRLRY-GVIKIIKPTLFKIAFTKVYGG---IPAE 85
RLA0_METAC -----MAEERHHTTEHIPQWKKDEIENIKELIQSHKVFVGMVIEGILATKMQKIRRDLDKDV-AVLKVSRLNTEALNQLG----ETIP 78
RLA0_METMA -----MAEERHHTTEHIPQWKKDEIENIKELIQSHKVFVGMVRIEGLATKIQKIRRDLDKDV-AVLKVSRLNTEALNQLG----ESIP 78
RLA0_ARCFU -----MAAVRGS---PEYKVRAVEEIKRMISSKPVVAIVSFRNVPAGQMQKIRREFRGK-AEIKVVKNTLLEALDALG----GDYL 75
RLA0_METKA MAVKAKGQPPSGYEPKVAEWRREVKELKELMDEYENVGLVDLEGIAPAQEQEIRAKLRERDTIIRMSRNTLMRIAEEKLDER--PELE 88

```

alignment

Ułożenie dwóch lub więcej sekwencji biopolimerów (DNA, RNA lub białka) w celu zidentyfikowania regionów podobieństwa istotnego ze względów ewolucyjnych, strukturalnych lub funkcjonalnych (procedura oraz jej efekt).

- dwie sekwencje - pairwise alignment
- wiele sekwencji - multiple sequence alignment



Znaczenie dopasowania

Podobieństwo porównywanych sekwencji (*similarity*) może świadczyć o:

- podobnej strukturze białek
- podobnej funkcji sekwencji
- wspólnej historii ewolucyjnej sekwencji

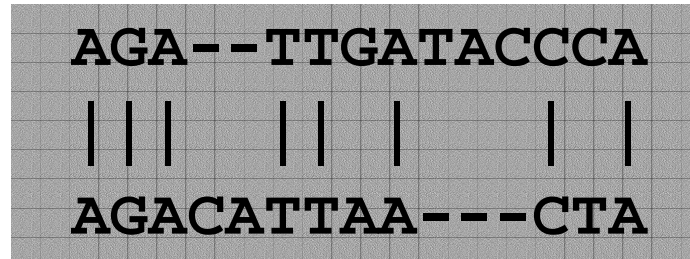
Podobieństwo porównywanych sekwencji (*similarity*) może wynikać z:

- **homologii** - pochodzeniu sekwencji (homologicznych) od wspólnego przodka; sekwencje mogą, ale nie muszą pełnić te same funkcje
- **konwergencji** - podobne motywy, które wyewoluowały w obu sekwencjach (analogicznych) niezależnie; np. chymotrypsyna i subtylizyna - różna struktura 3D, ale podobne centrum aktywne (histrydina, seryna, kwas asparaginowy)

{... Problem rozróżnienia odległej homologii od analogii }

Skąd te różnice

różnice między sekwencjami świadczą o mutacjach, które zaszły po rozdzieleniu się sekwencji od wspólnego przodka



AGA--TTGATACCCA

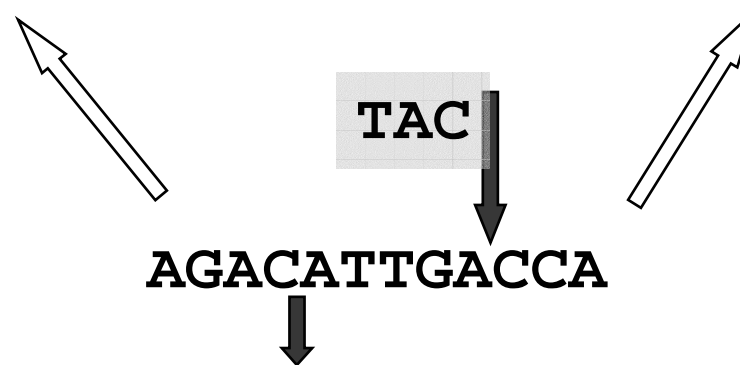
Insercja +TAC

Delecja -CA

AGACATTAA---CTA

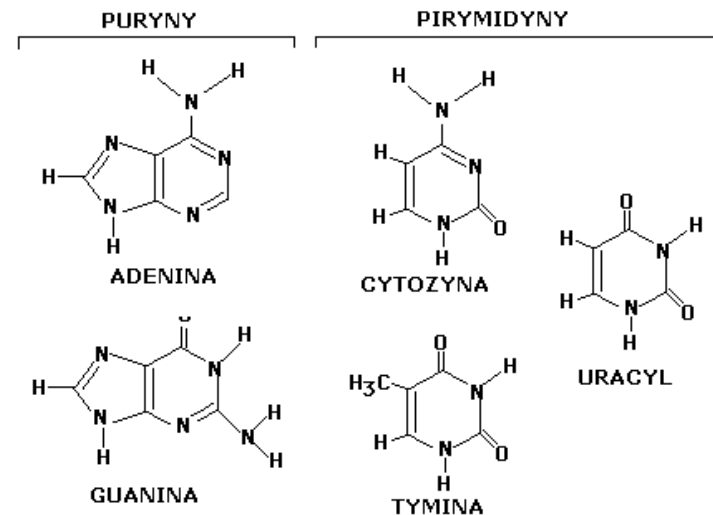
G->A C->T

substytucje



Substytucje nukleotydowe

- **Tranzycja** - okres przejściowy między systemem politycznym, który był, a tym który nastąpi. Proces ten jest krótszy i łatwiejszy od konsolidacji systemu politycznego. Tranzycja kończy się gdy pojawiają się ogólne ramy funkcjonowania nowego systemu. Przykładem są wszystkie państwa byłego bloku wschodniego, w tym Polska. (Czy o to chodzi?)
- **Transwersja** - mutacja genowa, punktowa zmiana chemiczna w obrębie nici DNA, w której zasada purynowa ulega zamianie na pirymidynową lub odwrotnie. Mutacja taka może nie spowodować żadnej zmiany lub zmianę kodu genetycznego (UUU → UUA) albo też skróconą syntezę białka (UCG → UCA).



Zastosowanie alignmentu

- poszukiwaniu oraz określaniu funkcji i struktury (białek) dla „nowych” sekwencji (nieznanych nam do tej pory)
- określaniu powiązań filogenetycznych między sekwencjami - homologii między sekwencjami oraz w analizach ewolucyjnych

Metody dopasowania

dopasowanie par sekwencji (*pairwise alignment*)

- **Macierz punktowe** - dot matrix, dotplot
 - **Programowanie dynamiczne (DP)**
 - **Metody słów (k - tuple methods)** - szybkie metody stosowane przy przeszukiwaniu baz danych sekwencji z wykorzystaniem programów FASTA i BLAST
-
- **dopasowanie wielu sekwencji (*multiple alignment*)**

Etapy dopasowywania sekwencji

1 zestawienie (0 identycznych, 0% podobieństwa)

$x = \text{długość sekwencji (30)}$

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

$y = \text{długość sekwencji (20)}$

2 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

3 zestawienie (1 identyczna, 33% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

4 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

5 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

6 zestawienie (2 identyczne, 33% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

7 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

X-2 zestawienie (3 identyczne, 15% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

x zestawienie (19 identycznych, 95% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

X+1 zestawienie (1 identyczna, 5,26% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

X+2 zestawienie (3 identyczne, 16,67% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

X+Y-4 zestawienie (1 identycznych, 25% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

X+Y-3 zestawienie (1 identycznych, 33,3% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

X+Y-2 zestawienie (0 identyczne, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

X+Y-1 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

Etapy dopasowywania sekwencji

1	RVCPKILMECKKSDCLAECICLEHGYCG	0
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	0.0%
2	RVCPKILMECKKSDCLAECICLEHGYCG	0
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	0.0%
3	RVCPKILMECKKSDCLAECICLEHGYCG	0
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	0.0%
4	RVCPKILMECKKSDCLAECICLEHGYCG	1
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	25.0%
5	RVCPKILMECKKSDCLAECICLEHGYCG	0
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	0.0%
•		
$n-1$	RVCPKILMECKKSDCLAECICLEHGYCG	1
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	3.6%
n	RVCPKILMECKKSDCLAECICLEHGYCG	18
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	62.1%
$n+1$	RVCPKILMECKKSDCLAECICLEHGYCG	5
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	17.2%
$n+2$	RVCPKILMECKKSDCLAECICLEHGYCG	2
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	6.9%
•		
$n+m-3$	RVCPKILMECKKSDCLAECICLEHGYCG	1
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	33.3%
$n+m-2$	RVCPKILMECKKSDCLAECICLEHGYCG	0
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	0.0%
$n+m-1$	RVCPKILMECKKSDCLAECICLEHGYCG	0
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	0.0%
n	RVCPKILMECKKSDCLAECICLEHGYCG	22
	MVCPKILMKCKHSDCLLDCVCLEDIGYCGVS	73%
m		

Za zgodą
dr. Jacka Leluka

Kryteria szacowania podobieństwa sekwencji

1) Zawartość % pozycji identycznych

PKILMECKKD 8
PKILMKCKHD 80%
spokrewnione

PKILMECKKD 2
SDCLLDVCL 20%
niespokrewnione

2) Długość porównywanych sekwencji

LCE 1
WCG 33.3%
nieznaczące

MVEICIEPKIRCIKVCTKDERITCLILDET 8
MYYWCPRRFMHCVHLKAGGCTCWCLRLDYY 26%
znaczące

3) Rozmieszczenie identycznych pozycji wzdłuż porównywanych sekwencji

MVEMICIEPKIRCIKVCTKDERITL 5
HYYWRPERFMHTVKLKAGGCRWL 20%
przypadkowa

MVEMIMAGDARCIKVCTKDERITCL 5
HYYWMAGDAHTVQLKAGGCWCWAG 20%
nieprzypadkowa

4) Typ reszt w pozycjach konserwatywnych

MVCPKILMKCKKHDSDCLLDCVCLLED
EDEGKRRTKREHFKE SNLAAAFKEQ
nieznaczące

MVCPKILMKCKKHDSDTLLDCVCLLED
QNCPPPREWCFTTRMNDSSCACPQT
znaczące

5) Podobieństwo strukturalne/genetyczne aminokwasów w nieidentycznych pozycjach

Kryterium identyczności

MV PKILMK KHDSDLLDV LED
RL RRLVKR RKETE IVE I IDE

Kryterium podobieństwa strukturalnego

MV PKILMK KHDSDLLDV LED
RL RRLVKR RKETE IVE I IDE

Kryterium podobieństwa genetycznego

MV PKILMK KHDSDLLDV LED
RL RRLVKR RKETE IVE I IDE

Za zgodą

dr. Jacka Leluka

Kryteria szacowania podobieństwa sekwencji

- Procent identyczności (względny udział odpowiadających sobie pozycji obsadzonych tymi samymi resztami)
- Długość porównywanych sekwencji (liczba porównywanych pozycji)
- Rozmieszczenie identycznych pozycji wzdłuż porównywanych sekwencji
- Typ reszt okupujących pozycje konserwatywne (sekwencje białkowe)
- Relacje genetyczne/strukturalne między resztami znajdującymi się w odpowiadających sobie nieidentycznych pozycjach (sekwencje białkowe)

Procedura oszacowania stopnia podobieństwa porównywanych sekwencji

Bardzo często oszacowanie stopnia podobieństwa porównywanych sekwencji sprowadzane jest jedynie do określenia względnego udziału pozycji identycznych. Pozostałe kryteria analizy zazwyczaj nie są w ogóle brane pod uwagę (np. bezwzględna długość sekwencji, dystrybucja identycznych pozycji wzdłuż łańcucha). Podejście takie jest niekompletne i stwarza ryzyko błędnej interpretacji otrzymanych wyników.

Przedstawiona niżej metoda oparta jest na prawdopodobieństwie przypadkowego pojawienia się zadeklarowanego stopnia identyczności. Uwzględnia ona podstawowe parametry mające znaczenie dla opisu faktycznego związku między porównywanymi sekwencjami.

Liczbę wszystkich możliwych stopni identyczności dla danych dwóch sekwencji opisuje poniższe r

$$T = x^{2n} = \sum_{a=0}^n \binom{n}{a} x^a (x(x-1))^{n-a}$$

Gdzie:

x – ilość rodzajów jednostek występujących w sekwencjach (20 dla białek; 4 dla kwasów nukleinowych)

n – długość sekwencji (liczba porównywanych par pozycji)

a – ilość pozycji identycznych

Local vs. Global

Global alignment – znajduje najlepsze dopasowanie dla **CAŁYCH** dwóch sekwencji

(Needleman-Wunsch algorithm)

```
ADLGAVFALCDRYFQ
|||||         |||
ADLGRTON-CDRYYQ
```

Global alignment: forces alignment in regions which differ

Local alignment – poszukuje podobnych regionów we **FRAGMENTACH** sekwencji

(Smith-Waterman algorithm)

```
ADLG          CDRYFQ
||||         |||||
ADLG          CDRYYQ
```

Local alignment will return only **regions** of good alignment

Global - local

Globalne dopasowanie dla pary $S_1, S_2 \in \mathcal{G}^*$

to każda taka para $(S_1^Y, S_2^Y) \in (\mathcal{G} \cup \text{'-'})^* \times (\mathcal{G} \cup \text{'-'})^*$,
która spełnia warunki:

- S_1 otrzymuje się z S_1^Y przez usunięcie wszystkich '-',
- $|S_1^Y| = |S_2^Y|$,
- $\forall i \in 1 \dots |S_1^Y| : (S_1^Y(i) = \text{'-'}) \Rightarrow S_2^Y(i) \neq \text{'-'}$

Lokalne dopasowanie sekwencji S_1, S_2 to każde globalne uliniowanie dla pewnych podciągów $s_1 \in S_1^*$, $s_2 \in S_2^*$.

Pairwise alignment

```
AAGCTGAATTCGAA  
AGGCTCATTCTGA
```

Tylko jeden możliwy alignment

```
AAGCTGAATT-C-GAA  
AGGCT-CATTCTGA-
```

This alignment includes:
2 mismatches
4 indels (gap)
10 perfect matches

Kilka możliwych rozwiązań:

AAGCTGAATTCGAA
AGGCTCATTTCTGA

A-AGCTGAATTC--GAA
AG-GCTCA-TTTCTGA-

AAGCTGAATT-C-GAA
AGGCT-CATTTCTGA-

scoring system:

- Perfect match: +1
- Mismatch: -2
- Indel (gap): -1 (*kara za przerwy*)

AAGCTGAATT-C-GAA
AGGCT-CATTCTGA-

A-AGCTGAATTC--GAA
AG-GCTCA-TTTCTGA-

$$\text{Score:} = (+1) \times 10 + (-2) \times 2 + (-1) \times 4 = 2$$

$$\text{Score:} = (+1) \times 9 + (-2) \times 2 + (-1) \times 6 = -1$$



Zadanie 1

- Jaki jest **score** tego alignmentu??

dopasowanie: +1

niedopasowanie: -1

przerwa: -2

```
---bardzo---lubiebioinformatyke
| | | | | | | | | | * | | | | | | | | | | *
niebardzonielubiębioinformatyki
```

Kara za przerwy (gap costs)

Kara za otwarcie przerwy – G

Kara za przedłużenie przerwy – L

$$\text{Kara} = G + Ln$$

gdzie:

n – długość przerwy

Standardowo:

$$G = 10 - 15$$

$$L = 1 - 2$$

Zadanie 2

Kara za otwarcie przerwy – G

Kara za przedłużenie przerwy – L

Kara = G + Ln

gdzie:

n – długość przerwy

Standardowo dla aa:

G = 10 - 15

L = 1 - 2

```
-GAGCTGAA-----GAA
AGAGCTCAATTTCTGA-
```

G = 10

L = 1

Kara = (10 + 5*1),

czy

Kara = (10 + 1*1) + (10 + 5*1) + (10 + 1*1)

Zadanie 3

Wiemy, że w toku ewolucji z danej sekwencji wyskoczyła jedna cała stosunkowo duża domena. Jakie wartości G i L dla kary za przerwy należy ustawić?

```
nie lubie bardzo-----bioinformatyki
      ||||| |||||          ||||| ||||| ||||| ||||| *
-----lubie bardzozobardzobioinformatyke
```

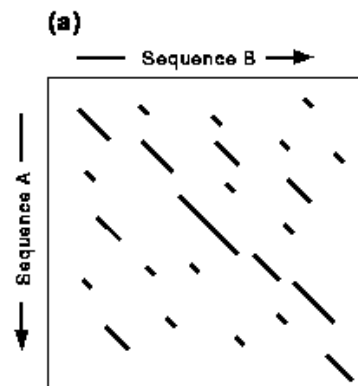
Metody dopasowania

dopasowanie par sekwencji (*pairwise alignment*)

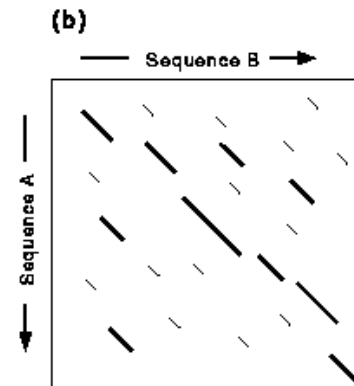
1. **Metody słów (k - tuple methods)** - szybkie metody stosowane przy przeszukiwaniu baz danych sekwencji z wykorzystaniem programów FASTA i BLAST
2. **Macierz punktowe** - dot matrix, dotplot
3. **Programowanie dynamiczne (DP)**

1. „słowa” - FASTA

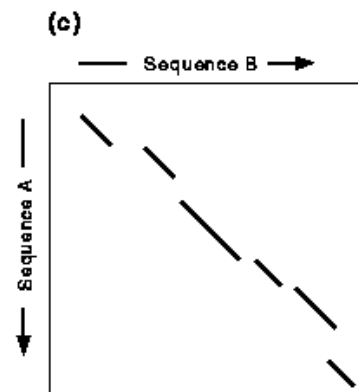
FASTA Algorithm



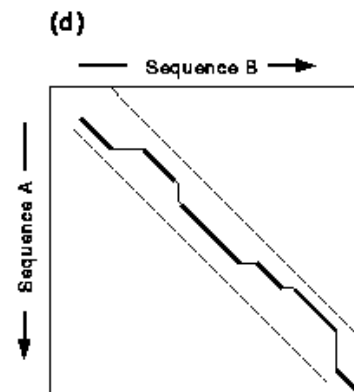
Find runs of identical words



Re-score using PAM matrix
Keep top scoring segments



Join segments using gaps,
eliminate other segments

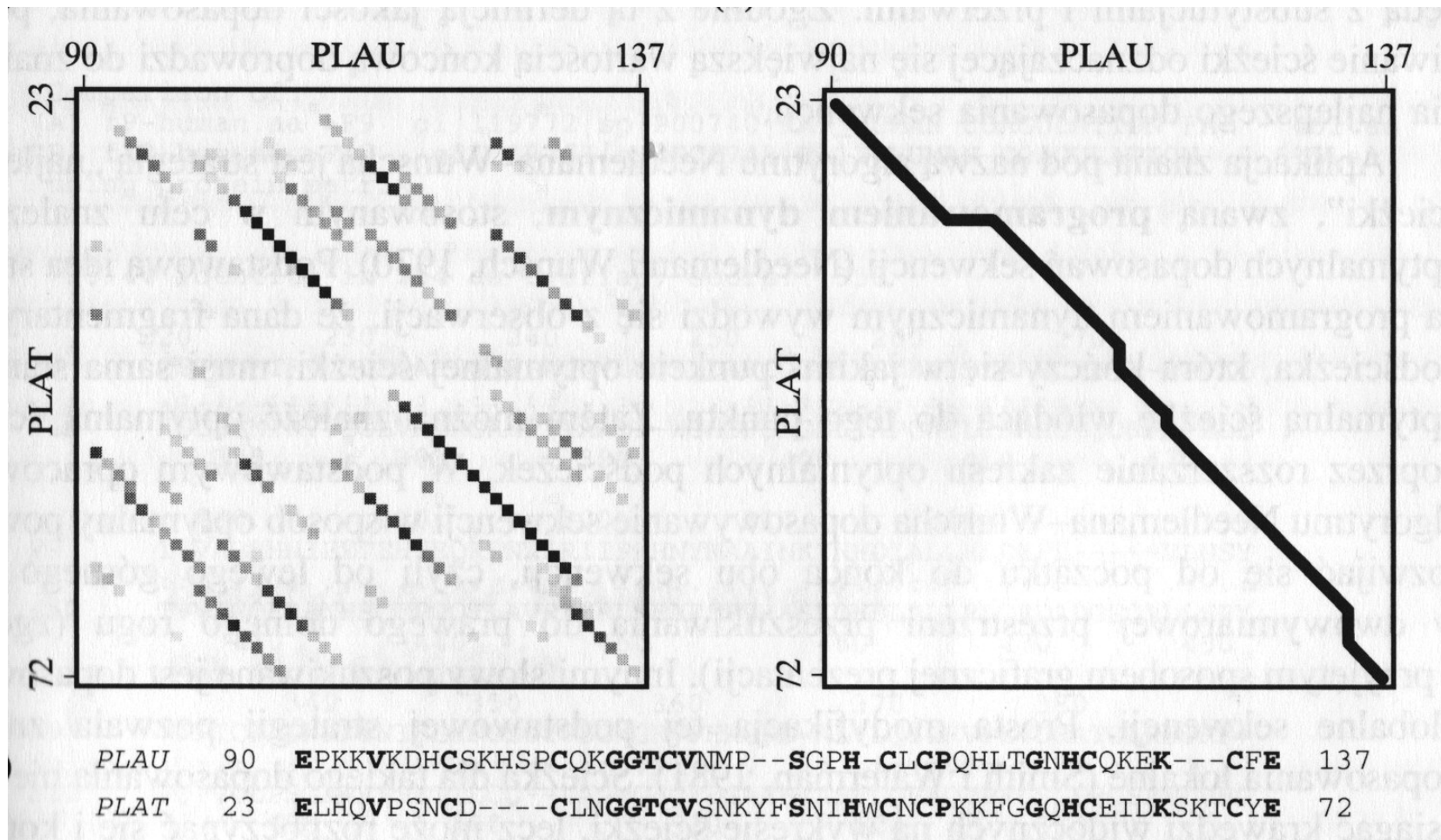


Use dynamic programming to
create an optimal alignment

1. „słowa” - BLAST vs. FASTA

BLAST	FASTA
może podawać więcej niż jeden region o wysokiej punktacji	podaje tylko jedno najlepsze dopasowanie
lepszy dla sekwencji białek niż DNA	lepszy dla sekwencji DNA niż białek
szybszy niż FASTA	wolniejszy niż BLAST
mniej czuły niż FASTA przy użyciu domyślnych ustawień	bardziej czuły niż BLAST
daje gorsze rozróżnienie między prawdziwymi i fałszywymi homologami	daje lepsze rozróżnienie między prawdziwymi i fałszywymi homologami

2. Macierze punktowe



2. Dot-matrix

BioEdit Sequence Alignment Editor - [Dot plot pairwise sequence comparison]

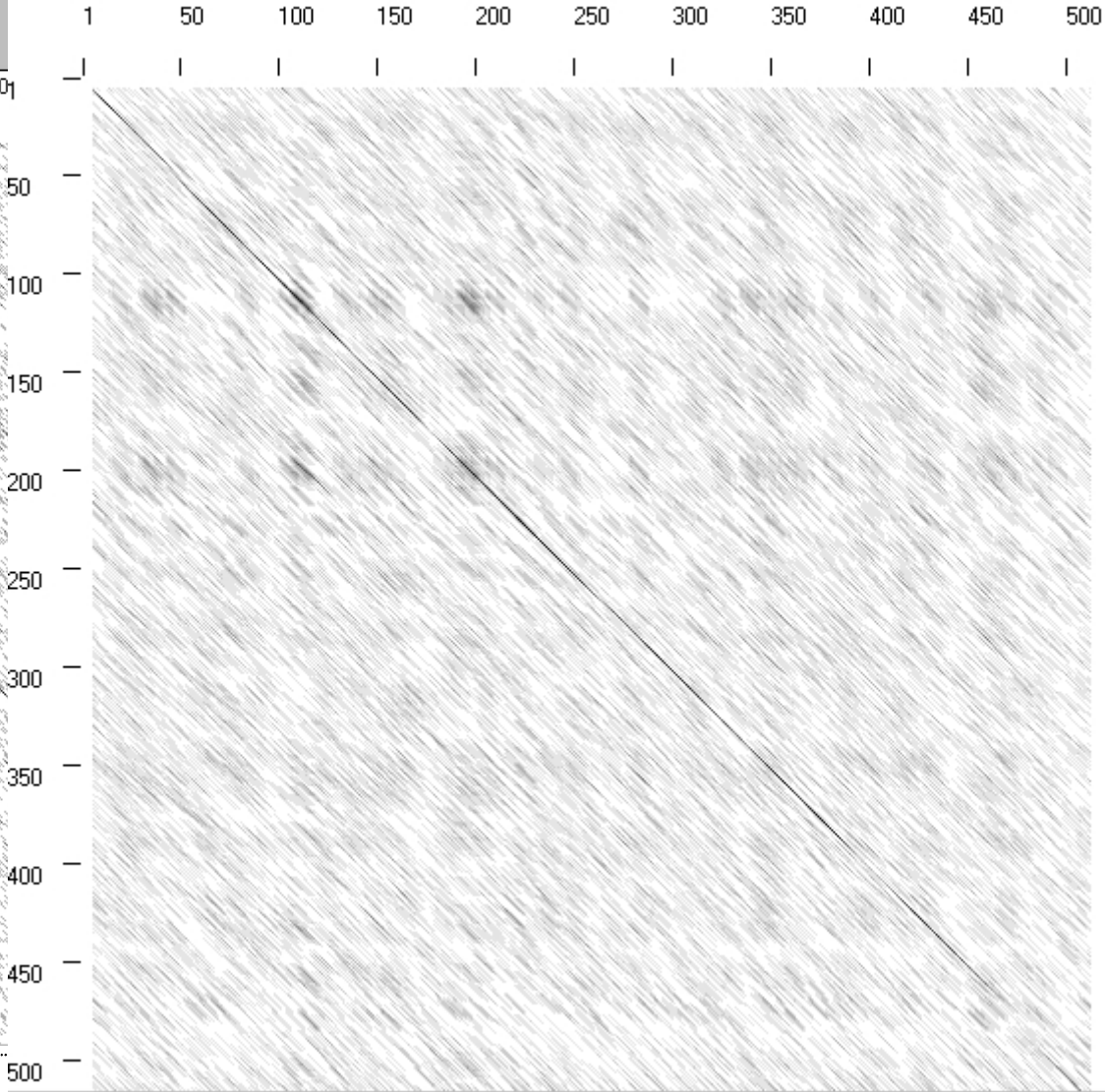
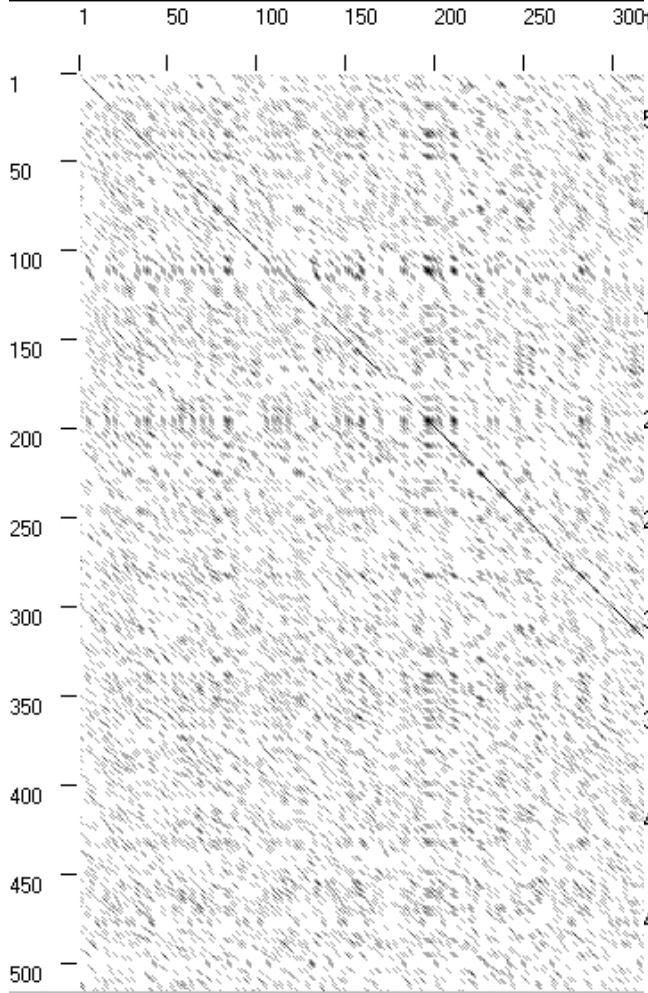
File Plot Edit View File Accessory Application RNA World Wide Web Options Window Help

Zoom: 100 Re-plot

Force point size 4 Pixels

Low threshold: 0.0000

High threshold: 3.0118



3. Programowanie dynamiczne

opiera się na podziale rozwiązywanego problemu na podproblemy względem kilku parametrów.

		<i>B</i>					
		Δ	1	2	3	4	
<i>A</i>	<i>H</i>	Δ	A	L	P	Q	
		Δ	0	-1	-2	-3	-4
	1	A	-1	1	0	-1	-2
	2	D	-2	0	1	0	-1
	3	L	-3	-1	1	1	0
	4	P	-4	-2	0	2	0
5	Q	-5	-3	-1	1	3	

Alignment *A* ... A D L P Q
 B ... A Δ L P Q
 Similarity Score +1 -1 +1 +1 +1 = 3

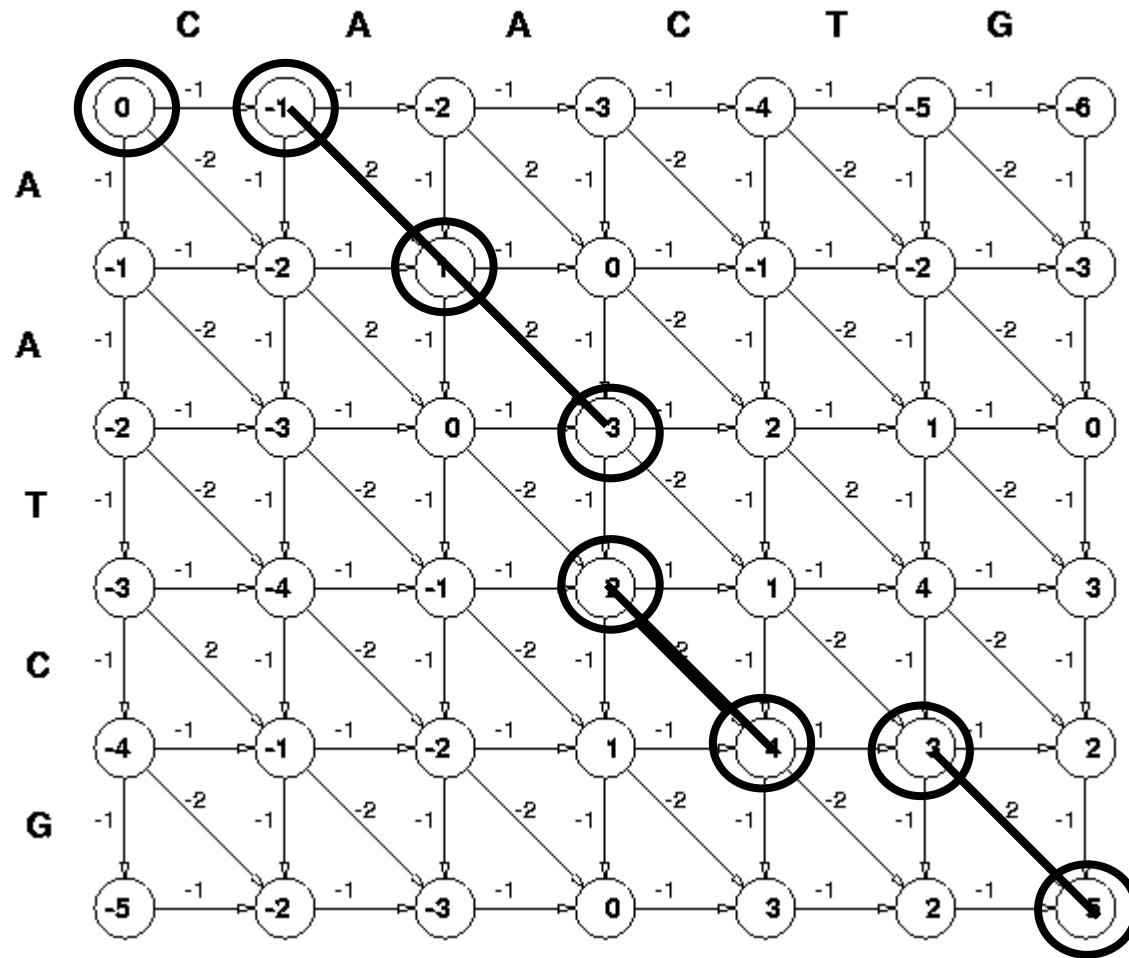
Dynamic programming matrix:

		<i>j</i> → (sequence <i>y</i>)								
		0	1	2	3	4	5	6	7	8 = N
		T	G	C	T	C	G	T	A	
<i>i</i> ↓ (sequence <i>x</i>)	0	0	-6	-12	-18	-24	-30	-36	-42	-48
	1 T	-6	5	-1	-7	-13	-19	-25	-31	-37
	2 T	-12	-1	3	-3	-2	-8	-14	-20	-26
	3 C	-18	-7	-3	8	2	3	-3	-9	-15
	4 A	-24	-13	-9	2	6	0	1	-5	-4
	5 T	-30	-19	-15	-4	7	4	-2	6	0
M = 6 A	-36	-25	-21	-10	1	5	2	0	11	

Optimum alignment scores 11:

T	-	-	T	C	A	T	A
T	G	C	T	C	G	T	A
+5	-6	-6	+5	+5	-2	+5	+5

3. Programowanie dynamiczne



Scoring matrix

- Reprezentuje system punktowania jako tabela lub macierz $n \times n$ (n jest liczbą liter, które zawiera alfabet. $n=4$ dla DNA, $n=20$ dla białek)
- Macierz punktowania jest symetryczna

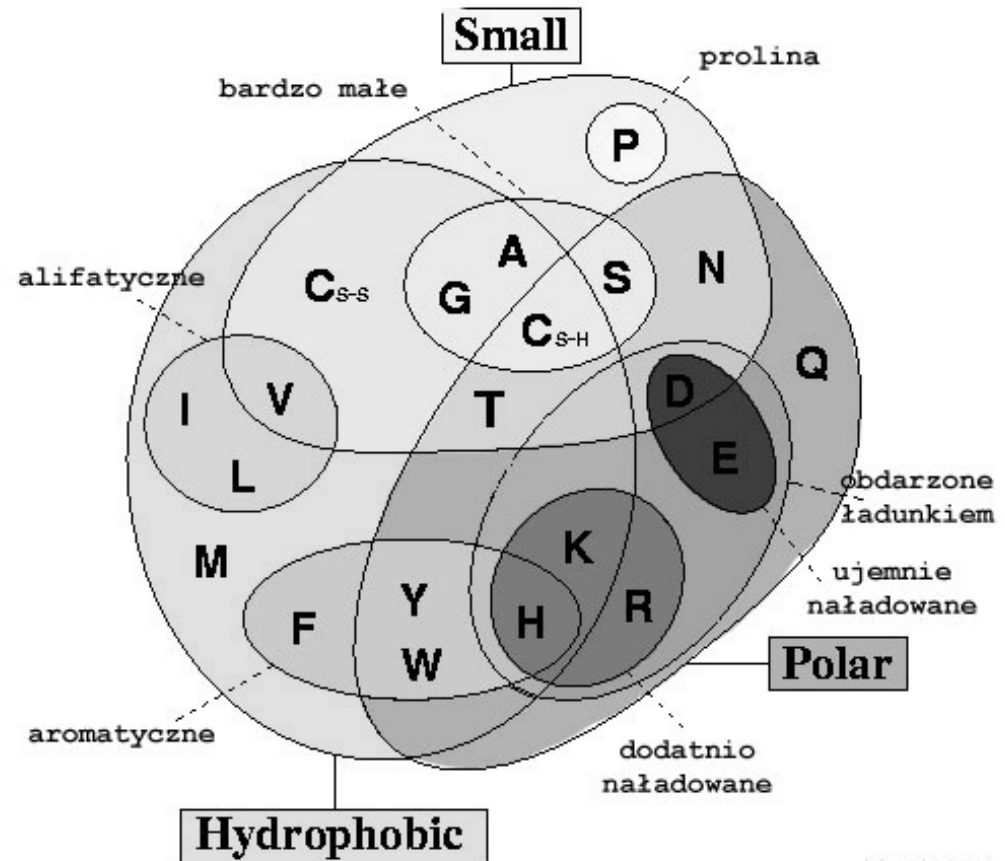
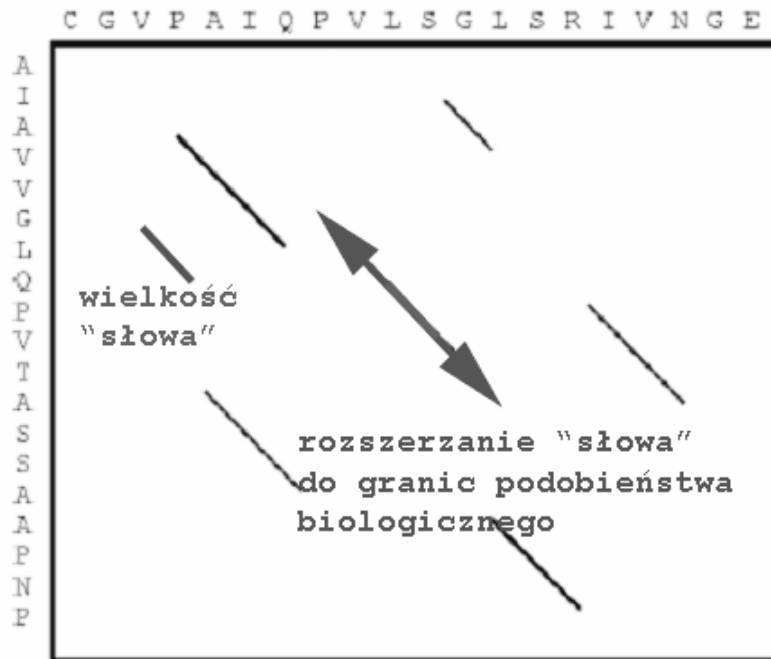
	A	G	C	T
A	2			
G	-6	2		
C	-6	-6	2	
T	-6	-6	-6	2

Mismatch

Match

Podobieństwa biochemiczne i biofizyczne aminokwasów

Diagram Venn-a



(c) ebiolog.pl

Macierze substytucji (podstawień)

- Jak za pomocą liczby określić podobieństwa biochemiczne i biofizyczne poszczególnych aminokwasów tak, aby liczba ta wyrażała jednocześnie realny wpływ na całe białko podstawienia danego aminokwasu w łańcuchu polipeptydowym i była uniwersalna dla wszystkich sekwencji?
- Przede wszystkim należy bazować na danych empirycznych
- Należy stworzyć alignment bardzo wielu blisko spokrewnionych sekwencji – na tyle podobnych, aby bez wątpliwości można było jednoznacznie i precyzyjnie określić częstotliwość substytucji poszczególnych aminokwasów w konkretnych pozycjach.

W kolumnie 4 E i D występują z częstotliwością w 4/8

M	G	Y	D	E
M	G	Y	D	E
M	G	Y	E	E
M	G	Y	D	E
M	G	Y	E	E
M	G	Y	D	E
M	A	Y	E	E
M	A	Y	E	E

PAM Matrix – Point/Percent Accepted Mutations

*n*PAM (*n* Percent Accepted Mutations)

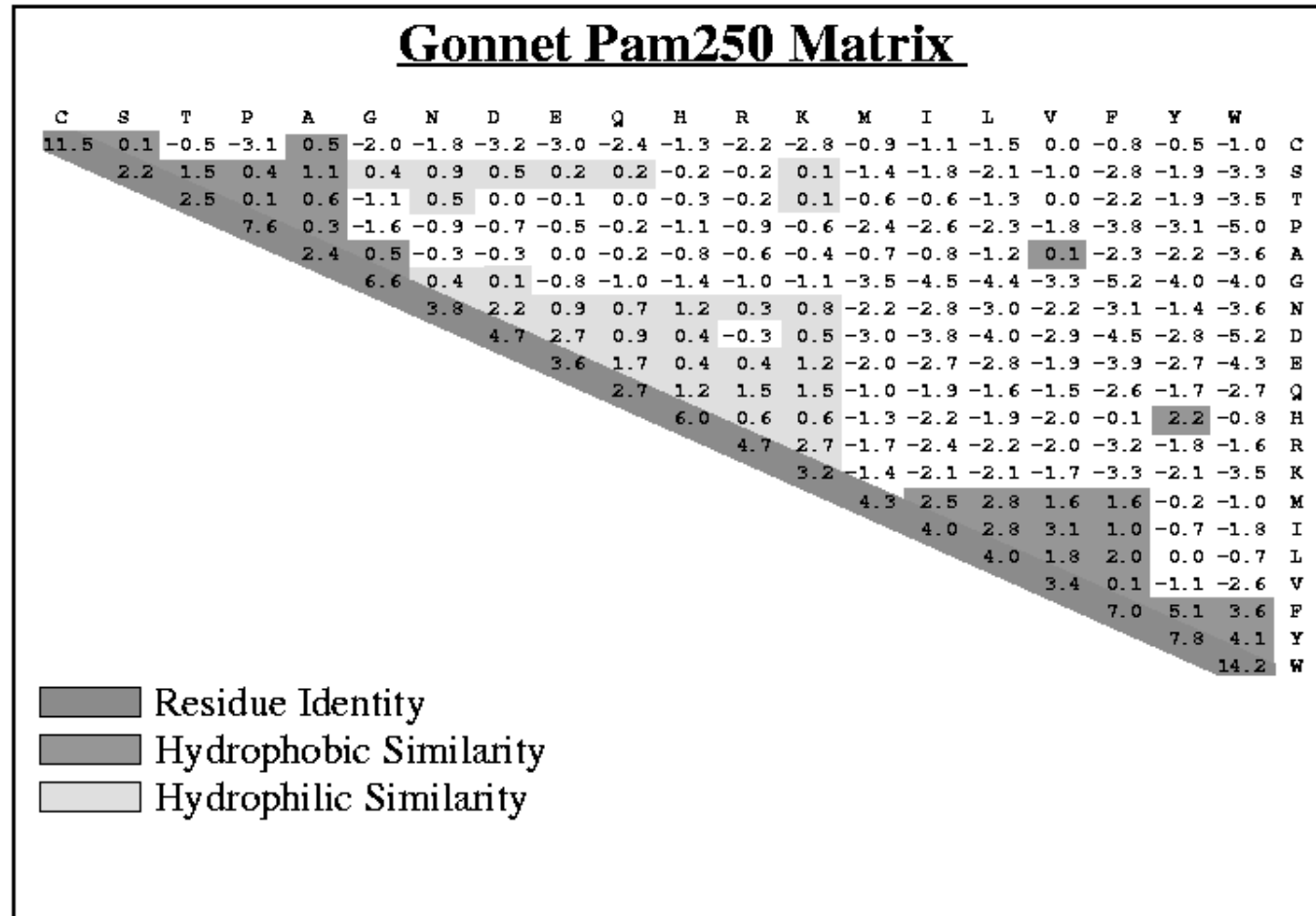
S_1, S_2 różnią się o jednostkę n PAM, jeśli S_2 można otrzymać z S_1 w ciągu akceptowalnych mutacji punktowych takich, że średnia liczba nieletalnych mutacji na 100 wynosi n . Najpopularniejsza jest 250PAM.

- Based on a database of 1,572 changes in 71 groups of closely related proteins (85% identity)
 - Alignment was easy

PAM Matrices

- Family of matrices PAM 80, PAM 120, PAM 250
- The number on the PAM matrix represents evolutionary distance
- Larger numbers are for larger distances

PAM



PAM - limitations

- Only one original dataset - PAM 1
- Examining proteins with few differences (85% identity)
- Bazuje głównie na małych białkach globularnych więc macierz jest nieco stronnicza

BLOSUM

- Henikoff i Henikoff (1992) stworzyli zestaw matryc bazujących na większej ilości danych empirycznych

BLOSSUM_n (*Block Substitution Matrix n*)

Oparta na bazie białek BLOCKS, gdzie są one podzielone na grupy tak, że dwa białka są zaliczane do jednej, jeśli można przejść od jednego do drugiego używając białek pośrednich tak, że dwa każde kolejne białka w tym przejściu mają skład identyczny w co najmniej $n\%$.

Popularne są BLOSUM 50, BLOSUM 62.

- BLOSUM observes significantly more replacements than PAM, even for infrequent pairs

BLOSUM: Blocks Substitution Matrix

- Based on BLOCKS database
 - ~2000 blocks from 500 families of related proteins
 - Families of proteins with identical function
- Blocks are short conserved patterns of 3-60 aa long without gaps

AABCDA	----	BBCDA
DABCDA	----	BBCBB
BBBCDA	AA-	BCCAA
AAACDA	A--	CBCDB
CCBADA	---	DBBDCC
AAACAA	---	BBCCC

BLOSUM

- Each block represent sequences alignment with different identity percentage
- For each block the amino-acid substitution rates were calculated to create BLOSUM matrix

BLOSUM Matrices

- BLOSUM n is based on sequences that shared at least n percent identity
- BLOSUM62 represents closer sequences than BLOSUM45

BLOSUM (62)

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
C	9																				C
S	-1	4																			S
T	-1	1	5																		T
P	-3	-1	-1	7																	P
A	0	1	0	-1	4																A
G	-3	0	-2	-2	0	6															G
N	-3	1	0	-2	-2	0	6														N
D	-3	0	-1	-1	-2	-1	1	6													D
E	-4	0	-1	-1	-1	-2	0	2	5												E
Q	-3	0	-1	-1	-1	-2	0	0	2	5											Q
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8										H
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5									R
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5								K
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5							M
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4						I
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4					L
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4				V
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			F
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7		Y
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11	W
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	

BLOSUM / PAM

Wszystkie macierze na podstawie danych empirycznych	Tylko PAM1 na podstawie danych empirycznych, pozostałe macierze z interpolacji
Opracowywane na podstawie sekwencji o dalszym pokrewieństwie	Opracowane na podstawie bardzo blisko spokrewnionych sekwencji
Podobieństwo sekwencji rośnie wraz ze wzrostem indeksu	Podobieństwo sekwencji maleje wraz ze wzrostem indeksu
Bezpośrednie podobieństwo sekwencji tu i teraz	Poniekąd reprezentuje dystans ewolucyjny (model ewolucyjny akceptowanych mutacji punktowych)
Macierz symetryczna (im wyższa wartość tym łatwiejsza substytucja)	Macierz symetryczna (im wyższa wartość tym łatwiejsza substytucja)
Nie uwzględnia bezpośrednio ani właściwości fizykochemicznych aminokwasów, ani podobieństwa genetycznego (podobieństwa kodonów)	Nie uwzględnia bezpośrednio ani właściwości fizykochemicznych aminokwasów, ani podobieństwa genetycznego (podobieństwa kodonów)

PAM vs. BLOSUM



PAM100	~	BLOSUM90
PAM120	~	BLOSUM80
PAM160	~	BLOSUM60
PAM200	~	BLOSUM52
PAM250	~	BLOSUM45

Sekwencje bardziej odległe

Uwarunkowania genetyczne substytucji aminokwasowych



Podstawy genetyczne algorytmów do zestawień aminokwasów?

Replacement	PAM250	BLOSUM62
Arg/Lys	3	2
Lys/Gln	1	1
Arg/Gln	1	1
Lys/Glu	0	1
Arg/Glu	-1	0

Diagram of amino acid genetic relationships

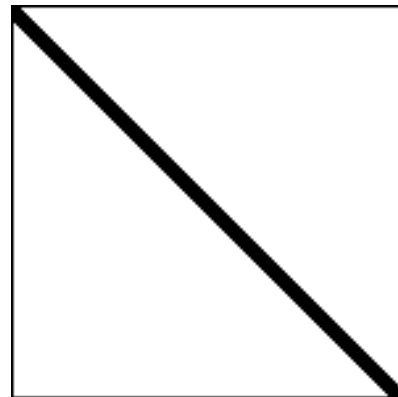
Algorytm semihomologiczny



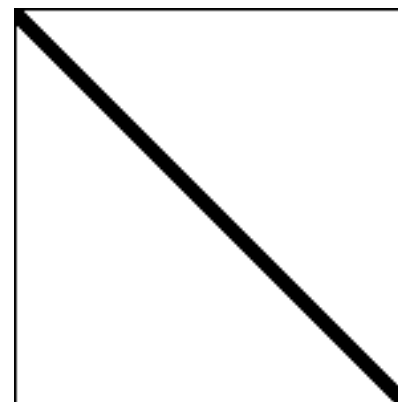
Dot matrix pairwise alignment

Internal homology (gene multiplication)

BLAST 2 SEQUENCES

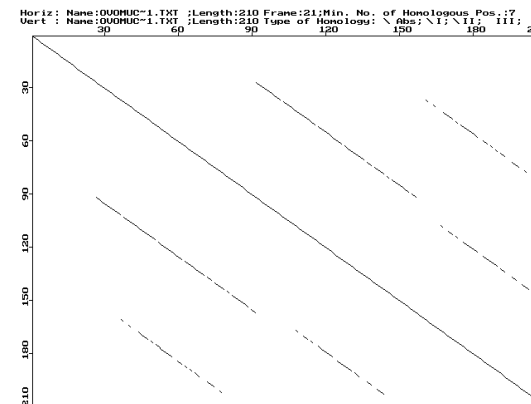
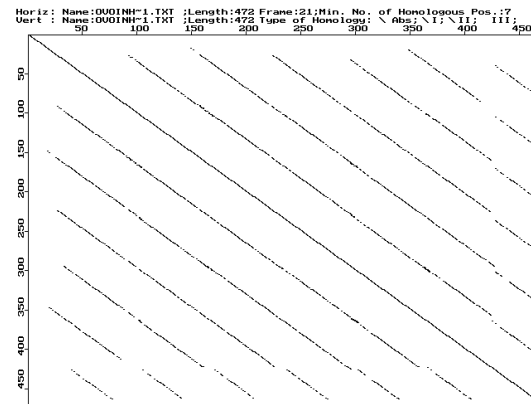


Chicken
ovoinhibitor
precursor
(7 domains)



Chicken
ovomucoid
precursor
(3 domains)

SEMIHOM




```

VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFSS--YAMYWVQRAPG
LSLTCTVSGTSFDD--YYSTWVQRPPG
PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPGA--VTVWKADS--
AALGCLVKDYFPEP--VTVWNSG---
VSLTCLVKGFYPSD--IAVEWWSNG--

```

Tak jak pairwise alignment ALE zestawienie n sekwencji zamiast 2

W rzędach ustawione są poszczególne sekwencje

W kolumnach ustawia się „te same”/”odpowiadające sobie” pozycje (pozycje konserwatywne); grupy pozycji konserwatywnych tworzą bloki konserwatywne (w blokach dozwolone są mutacje – insercje, delecje, substytucje – reprezentowane jako przerwy lub różne pozycje w kolumnach)

MSA & Evolution

MSA może dawać obraz sił kształtujących ewolucję !!!

- Ważne aminokwasy lub nukleotydy (pozycje w sekwencjach) mutują „niechętnie”
- Mniej ważne pozycje dla struktury i funkcji mogą wykazywać większą zmienność w kolumnach porównywanych sekwencji

Pozycje konserwatywne

- Kolumny, gdzie wszystkie sekwencje zawierają takie same aminokwasy lub nukleotydy (lub w większości takie same – pozycje konserwatywne) są bardzo ważne (kluczowe) dla funkcji lub struktury.

```
VTISCTGSSSNIGAG-NHVKWYQQLPG  
VTISCTGSSSNIGS--ITVNWYQQLPG  
LRLSCTGSGFIFSS--YAMYWYQQAPG  
LSLTCTGSGTSFDD-QYYSTWYQQPPG
```

Sekwencja konsensusowa

- W **sekwencji konsensusowej** zachowane są pozycje o największej częstotliwości występowania w każdej z kolumn alignmentu (The consensus sequence holds the most frequent character of the alignment at each column)
- Jest to sposób reprezentowania wyników multiple alignment, gdzie pokrewne sekwencje są porównywane każda do każdej, aby odnaleźć funkcjonalnie podobne motywy sekwencji (domeny białek). Sekwencja konsensusowa obrazuje które pozycje są konserwatywne, a które zmienne.

A	T	C	T	T	G	T
A	A	C	T	T	G	T
A	A	C	T	T	C	T



A	A	C	T	T	G	T
---	---	---	---	---	---	---

*	:	*	*	*	:	*
---	---	---	---	---	---	---

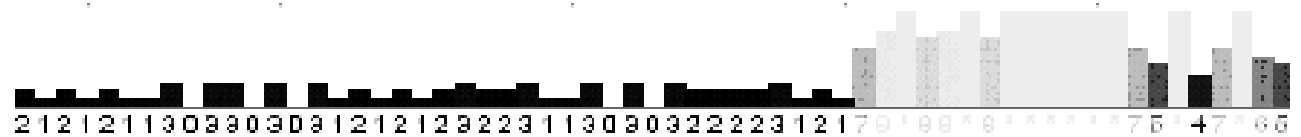
Sekwencja konsensusowa

...*****
 ...

P0_Homo_sapiens|human|P25189|1-209
P0_Pan_troglodytes|chimpanzee|1-258
P0_Musca_domestica|house|NP_011248
P0_Equus caballus|horse|NP_011248
P0_Sus_scrofa|pig|CAU5376.1|1-176
P0_Rattus_norvegicus|rat|NP_1-248
P0_Mus_musculus|mouse|NP_03211-209
P0_Canis_familiaris|dog|XP_571-258
P0_Bos_taurus|bovine|P10522|1-219
P0_Gallus_gallus|chicken|P371-209
P0_Salmo_gairdneri|fish|AA834399.1|1-202
P0_Danio_rerio|zebrafish|NP_1-209

1	I	V	V	T	D	R	E	V	H	G	A	V	G	S	R	V	T	L	H	C	S	F	W	S	S	E	W	S	D	D	I	S	F	T	W	R	Y	Q	P	E	G	G	R	D	A	I	S	I	F	H	T	A	K	G	O	P	
1	I	V	V	T	D	R	E	V	H	G	A	V	G	S	R	V	T	L	H	C	S	F	W	S	S	E	W	S	D	D	I	S	F	T	W	R	Y	Q	P	E	G	G	R	D	A	I	S	I	F	H	T	A	K	G	O	P	
1	I	V	V	T	D	R	E	V	H	G	A	V	G	S	R	V	T	L	H	C	S	F	W	S	S	E	W	S	D	D	I	S	F	T	W	R	Y	Q	P	E	G	G	R	D	A	I	S	I	F	H	T	A	K	G	O	P	
1	I	V	V	T	D	R	E	V	H	G	A	V	G	S	R	V	T	L	H	C	S	F	W	S	S	E	W	S	D	D	I	S	F	T	W	R	Y	Q	P	E	G	G	R	D	A	I	S	I	F	H	T	A	K	G	O	P	
1	I	V	V	T	D	R	E	V	H	G	A	V	G	S	R	V	T	L	H	C	S	F	W	S	S	E	W	S	D	D	I	S	F	T	W	R	Y	Q	P	E	G	G	R	D	A	I	S	I	F	H	T	A	K	G	O	P	
1	I	V	V	T	D	R	E	V	H	G	A	V	G	S	R	V	T	L	H	C	S	F	W	S	S	E	W	S	D	D	I	S	F	T	W	R	Y	Q	P	E	G	G	R	D	A	I	S	I	F	H	T	A	K	G	O	P	
1	I	H	V	T	P	R	E	V	Y	G	T	V	G	S	H	V	T	L	S	C	S	F	W	S	S	E	W	I	S	E	D	I	S	Y	T	W	H	F	Q	A	E	G	G	R	D	A	I	S	I	F	H	T	G	K	G	O	P
1	I	V	I	T	T	G	W	E	R	H	A	L	V	G	S	D	I	R	L	S	C	S	F	F	S	W	R	W	T	S	D	D	V	T	F	S	W	S	Y	R	P	D	G	A	R	D	A	I	S	I	F	H	T	G	G	A	P
1	V	V	N	T	D	R	E	K	H	A	L	V	G	S	D	V	R	L	S	C	S	F	F	S	W	Q	W	T	S	P	E	Y	S	F	T	W	H	Y	R	P	D	G	A	R	D	A	I	S	I	F	H	T	G	G	E	A	

Conservation



Quality



Consensus



Alignment methods

- Progressive alignment (Clustal)
- Iterative alignment (mafft, muscle)

- All methods today are an approximation strategy (**heuristic algorithm**), yield a possible alignment, but not necessarily the best one

Praca domowa

- iteracja (np. pętle w programowaniu)
- heurestyka (głównie w informatyce)
- Alignment progresywny

Jak wyświetlić na ekranie liczby od 1 do 5000 za pomocą 2 linijek kodu?

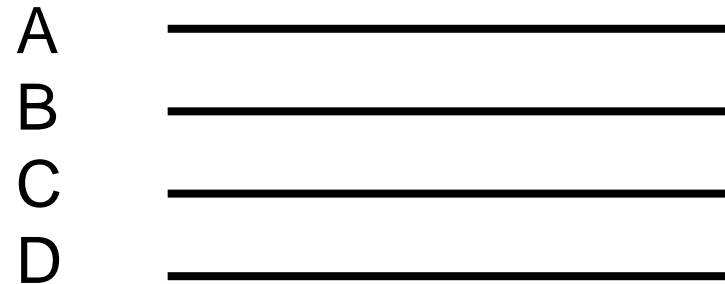
```
<?php
  for( $x = 1; $x = 5000; $x++ )
    echo $x. "<br /> ";
?>
```

dla (zmiennej x początkowo równej 1; aż do momentu kiedy
zmienna x osiągnie wartość równą 5000; z każdym
krokiem powiększając wartość zmiennej x o +1)
wyświetl wartość zmiennej i przejdź do nowej linii;

Iteracja (łac. iteratio ‘powtórzenie’) to czynność powtarzania (najczęściej wielokrotnego) tej samej instrukcji (albo wielu instrukcji) w pętli. Mianem iteracji określa się także operacje wykonywane wewnątrz takiej pętli.

Progressive alignment

First step:



Compute the pairwise alignments for all against all (6 pairwise alignments) the similarities are stored in a table

	A	B	C	D
A				
B	11			
C	3	1		
D	2	2	10	

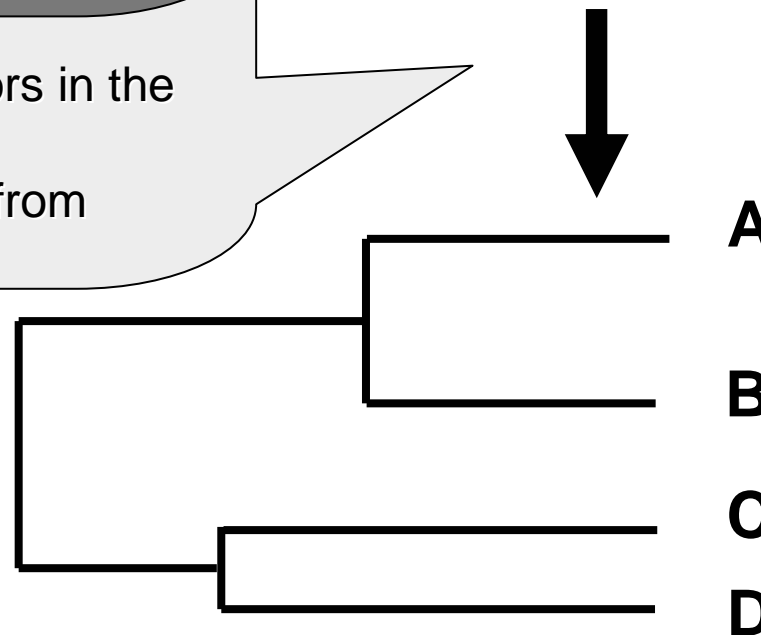
Second step:

The guide tree is imprecise and is NOT the tree which truly describes the relationship between the sequences!

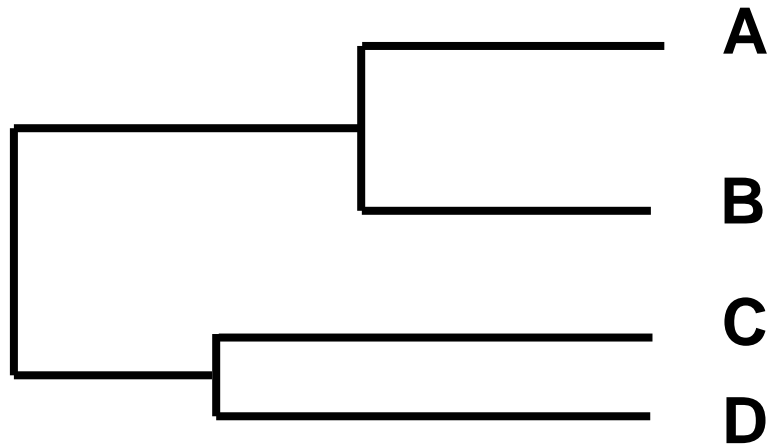
sequences are to be aligned

- similar sequences are neighbors in the tree
- distant sequences are distant from each other in the tree

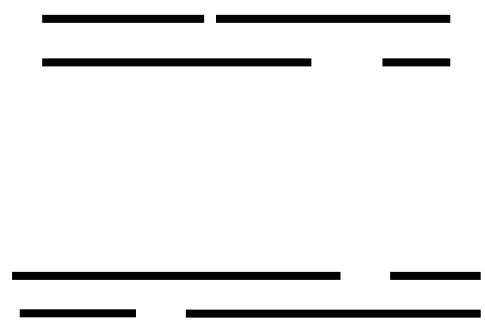
	A	B	C	D
A				
B	11			
C	3	1		
D	2	2	10	



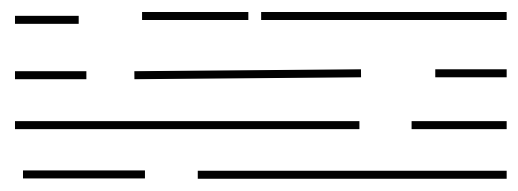
Third step:



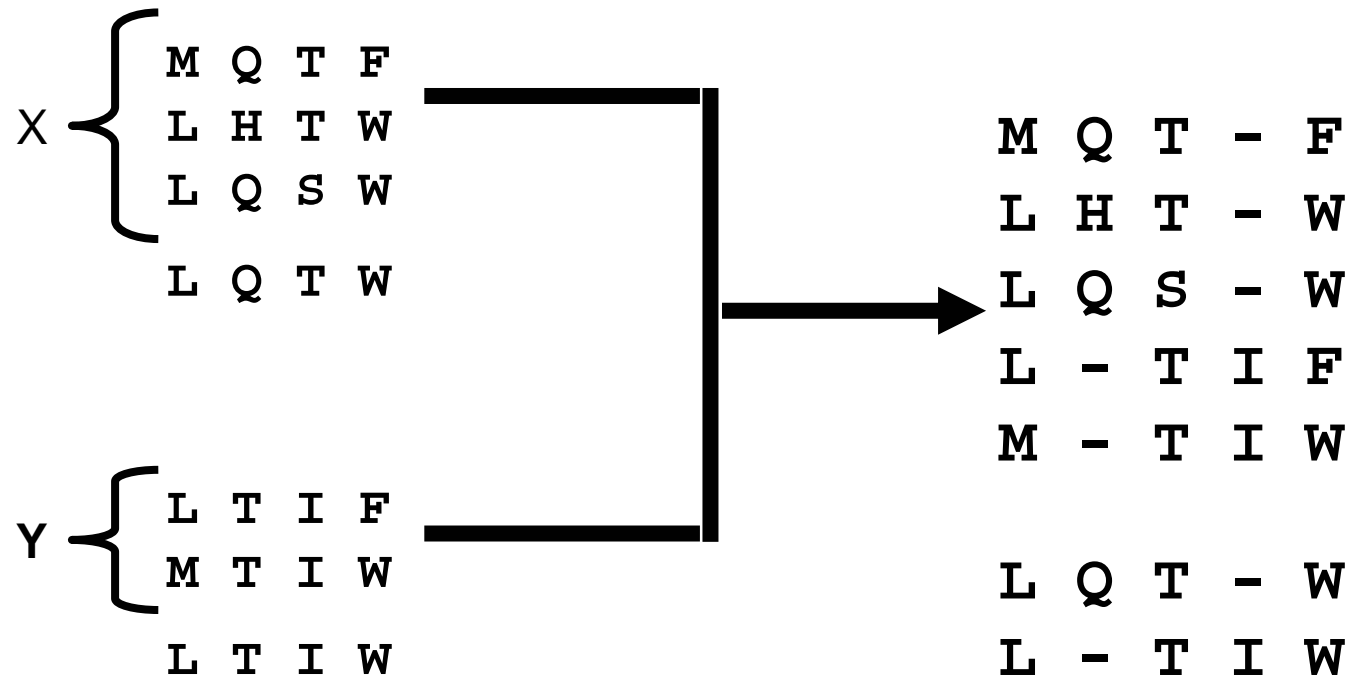
Align most similar pairs



Align the alignments as if each of them was a single sequence (replace with a single consensus sequence or use a profile)

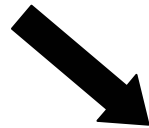
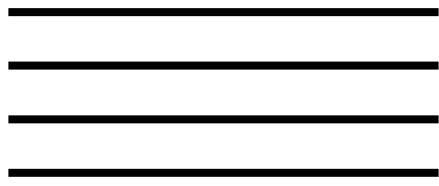


Alignment of alignments



Iterative alignment

A
B
C
D

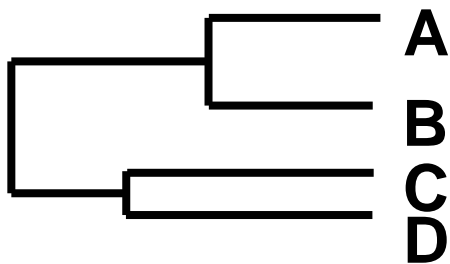


Pairwise distance table

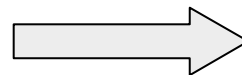
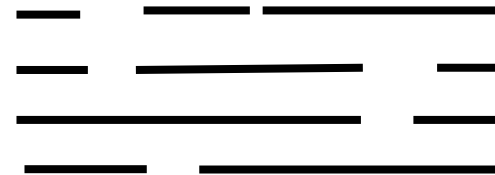
	A	B	C	D
A				
B	11			
C	3	1		
D	2	2	10	

Iterate until the MSA doesn't change

Guide tree



MSA



Searching for remote homologs

- Sometimes BLAST isn't enough.
- Large protein family, and BLAST only gives close members. We want more distant members
- PSI-BLAST
- Profile HMMs

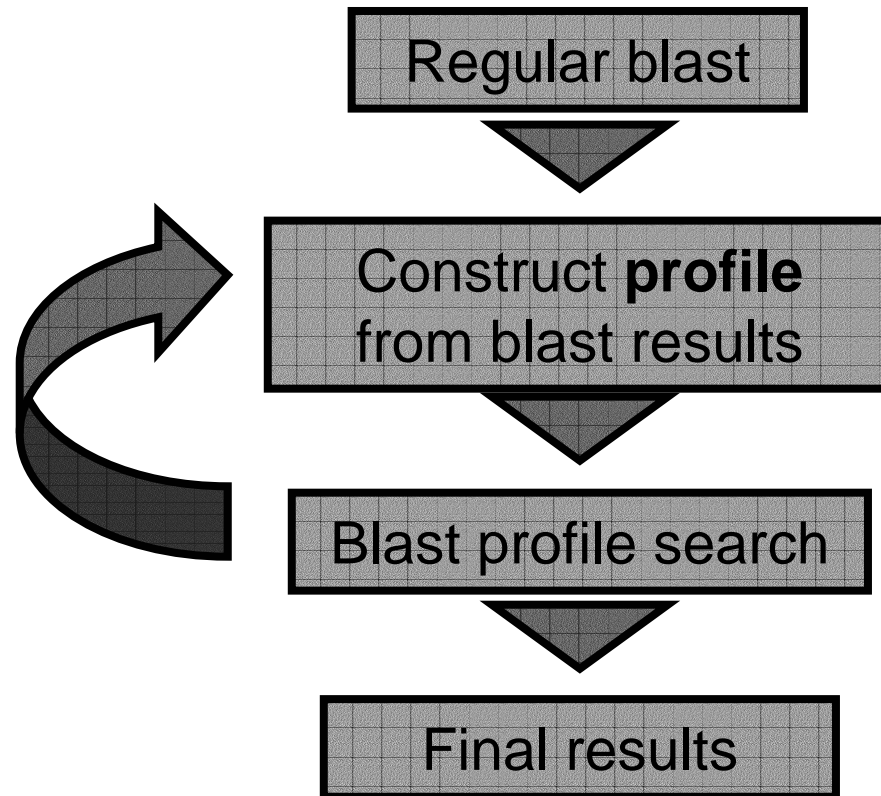
Profile

		1	2	3	4	5	6	
A	T	C	T	T	G	T		
A	A	C	T	T	G	T		
A	A	C	T	T	C	T		
	A	T	C	T	G	T		
	T							
	C							
	G							

Profile =
PSSM – Position Specific Score Matrix

PSI-BLAST

- **Position Specific Iterated BLAST**



PSI-BLAST

- **zalety:** PSI-BLAST looks for seq.s that are close to ours, and learns from them to extend the circle of friends
- **wady:** if we found a **WRONG** sequence, we will get to unrelated sequences (contamination). This gets worse and worse each iteration

Profile HMM

- Similar to PSI-BLAST: also uses a profile
- Takes into account:
 - Dependence among sites (if site n is conserved, it is likely that site $n+1$ is conserved \rightarrow part of a domain)
 - The probability of a certain column in an alignment

PSI BLAST vs profile HMM

PSI BLAST

Profile HMM

Less exact

Faster

More exact

Slower