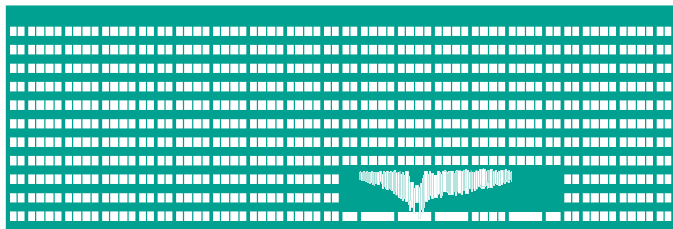


VŠB TECHNICKÁ
UNIVERZITA
OSTRAVA

VSB TECHNICAL
UNIVERSITY
OF OSTRAVA



www.vsb.cz

Algoritmy pro Bioinformatiku

Podobnost sekvencí

Michal Vašínek

VŠB – Technická univerzita Ostrava

FEI/EA443

michal.vasinek@vsb.cz

26. září 2023

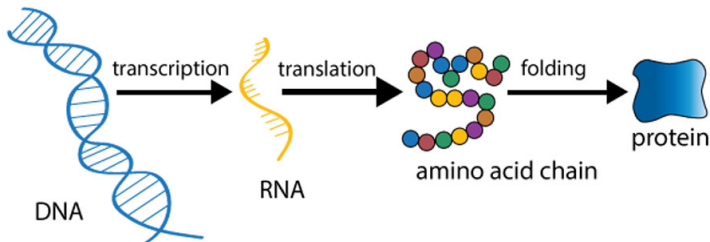
100111100101011110101
ACCGTTACGGA CTTA
CATGGGTAGGGAGGGG
01101100110100111011





- Biotechnologické nástroje:
 - Restrikční enzymy
 - Metoda shotgun
 - PCR reakce
- Měření podobnosti sekvencí:
 - Hammingova vzdálenost
 - Editační vzdálenost

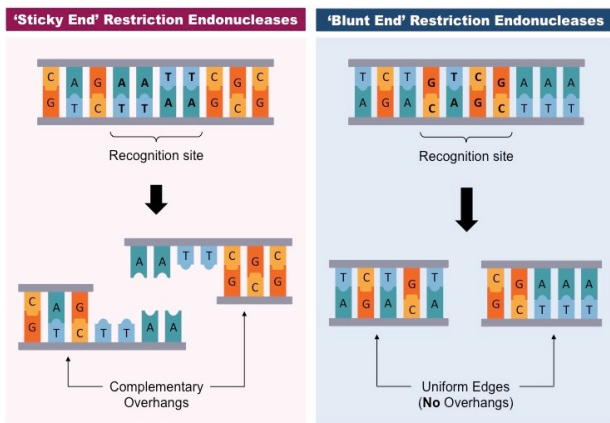
- Centrální Dogma nám říká, jakým způsobem obdržíme protein z konkrétního genu. Tomuto procesu se říká exprese genu.
- Exprese genu se skládá ze dvou kroků:
 - Transkripce: DNA \rightarrow mRNA
 - Translace: mRNA \rightarrow Protein
 - Následně dochází k po-translačním úpravám.
- Každá aminokyselina v proteinu je zakódována trojicí nukleotidů, tzv. codonem.





- Fragmentace DNA
 - Restrikční enzymy
 - Shotgun metoda
- Kopírování DNA
 - Klonování
 - PCR - Polymerase Chain Reaction
- Měření délky DNA
 - Gelová elektroforéza

- Restrikční enzym rozpoznává určitá místa v DNA, kde je schopen se navázat a DNA v tomto místě rozštěpit ¹.



¹https://www.youtube.com/watch?v=Ik_Pxht1LM0

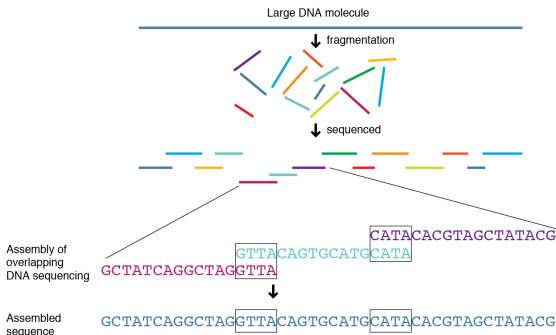


- Objeveny v roce 1968.
- Výskyt u bakterií a archeobakterií.
- Přirozeně se restrikční enzymy používají k přestřížení cizí DNA => předcházení infekcí.
- EcoRI přestřihává DNA v místech, kde se vyskytuje sekvence GAATTC. Je svým vlastním **reverzně komplementárním palindromem**.
- V současnosti známe přes 3000 restrikčních enzymů.
- Využití: optické mapování.²

²<https://www.youtube.com/watch?v=S2ng6glu04I>

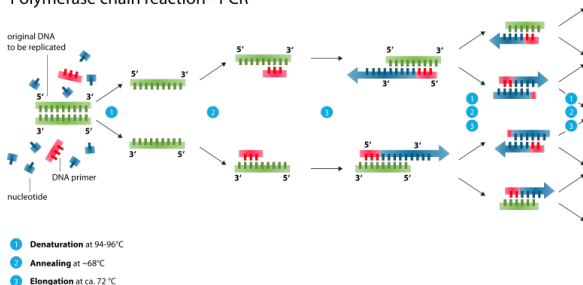


- Náhodná fragmentace DNA na krátké úseky.
- Celý proces proveden několikrát pro zajištění překrývajících se úseků.
- Metoda se používá pro sekvenování celého genomu.
- Krátké úseky pospojovány pomocí algoritmů tzv. de novo sestavení genomu.



- PCR - polymerase chain reaction
- Vynalezeno v roce 1984.
- Umožňuje extrémně rychlou replikaci vybraných úseků DNA.
- Zjišťování mutací ve vybraných oblastech DNA.

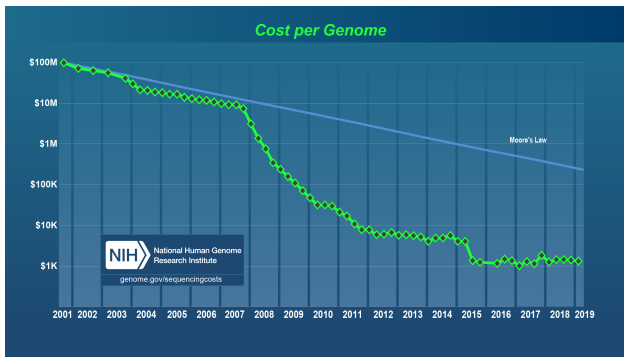
Polymerase chain reaction - PCR



³<https://www.youtube.com/watch?v=wBrNbbAIAFo>



- Lidský genom přečten a publikován 2000/2001.
- Stále zpřesňován. Poslední verze 2022, T2T-CHM13.



Měření podobnosti sekvencí



- Určování biologických funkcí různých genů. Vycházíme z předpokladu, že dva geny mající podobnou sekvenci mají i podobnou funkci.
- Nalezení evoluční vzdálenosti mezi dvěma organismy.^{4 5}
- Sestavení genomu.
- PCR - amplikonové sekvenování - přiřazení sekvence ke správnému místu v referenční sekvenci.

⁴Dva lidé se liší přibližně v 0,6% DNA

⁵Člověk a šimpanz sdílí 99% DNA



- Hammingova vzdálenost:
 - Dva řetězce X a Y stejné délky.
 - Hammingova vzdálenost je počet rozdílných symbolů.

$A = \text{acgtacgt}$

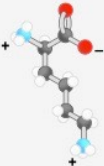
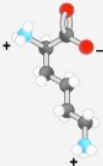
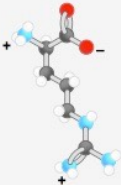
$B = \text{aggtcgat}$

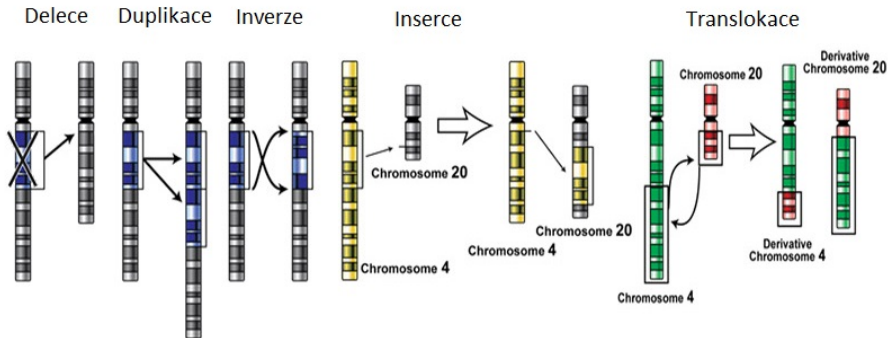
acgtacgt

aggtcgat

- Hammingova vzdálenost $d_H(A, B) = 4$.
- Editační vzdálenost - viz. další slidy



	No mutation	Point mutations		
		Silent	Missense	Nonsense
DNA	TTC	TTT	TCC	ATC
mRNA	AAG	AAA	AGG	UAG
Protein	Lys	Lys	Arg	STOP
				





- Mějme dva řetězce A a B , pomocí kolika následujících operací dokážeme z A obdržet B :
 - Nahrazení písmene za jiné písmeno - R .
 - Vynechání písmene, delece - D .
 - Vložení písmene, inserce - I .
 - Shoda písmem - M .
- Editací vzdálenost \leq Hammingova vzdálenost.

$A = \text{acgtacgt}$

$B = \text{aggtcgat}$

acgtacg_t

aggt_cgat

editační vzdálenost $d_E(A, B) = 3$



- Každá operace s řetězcem M , R , I , D má přiřazenu svou cenu.
- Cíl : nalézt nejmenší počet operací resp. minimální editační vzdálenost.

	_	a	c	g	t
_		1	1	1	1
a	1	0	1	1	1
c	1	1	0	1	1
g	1	1	1	0	1
t	1	1	1	1	0

acgtacg_t

aggt_cgat

$$0 + 1 + 0 + 0 + 1 + 0 + 0 + 1 + 0 = 3$$



- Existence nerozhodnutelných posloupností operací.

$$A = acgtac$$

$$B = actgac$$

acgtac

acgt_ac

ac_gtac

actgac

ac_tgac

actg_ac

- Všechny alternativy mají stejnou editační vzdálenost.

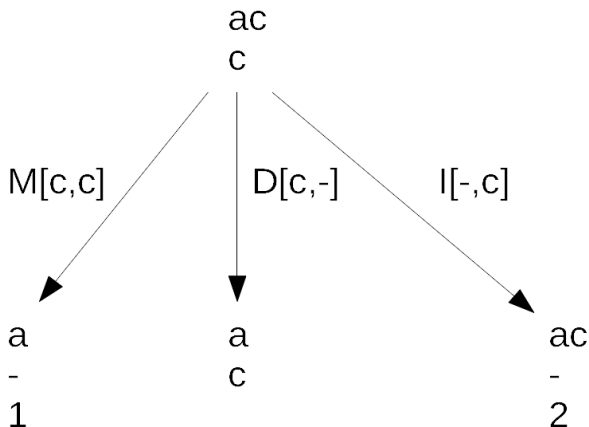


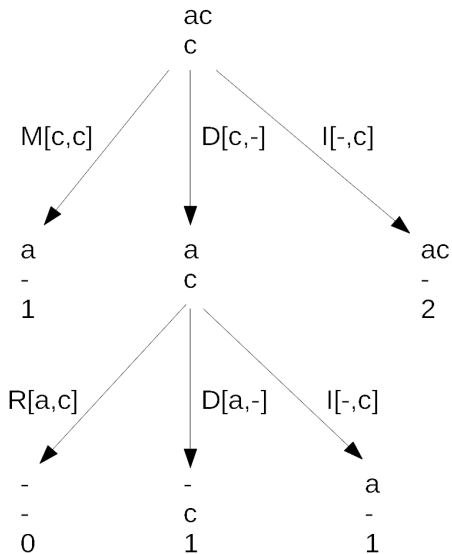
- Uvažujme editační vzdálenost mezi dvěma řetězci s_1 a s_2 , necht' $|s_1| > 0$ a $|s_2| = 0$. Jaká je jejich editační vzdálenost? $d_E = |s_1|$
- Necht' i, j jsou dva prefixy řetězců s_1 a s_2 .
- $d_E[i, j]$ je potom editační vzdálenost mezi dvěma prefixy. $d_E[i, 0] = i$ a $d_E[0, j] = j$.

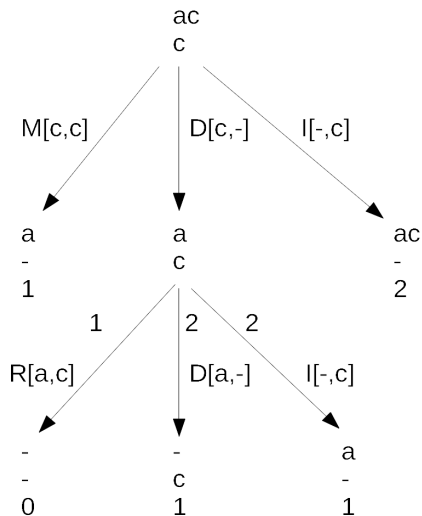


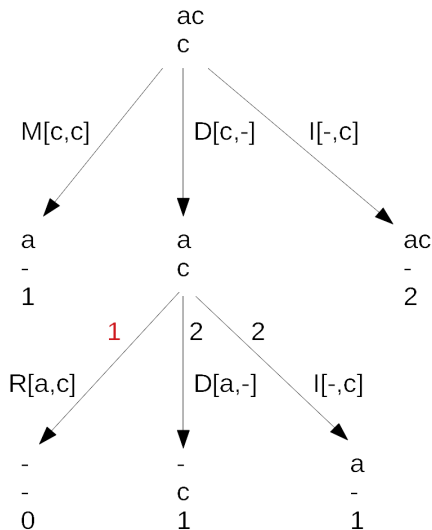


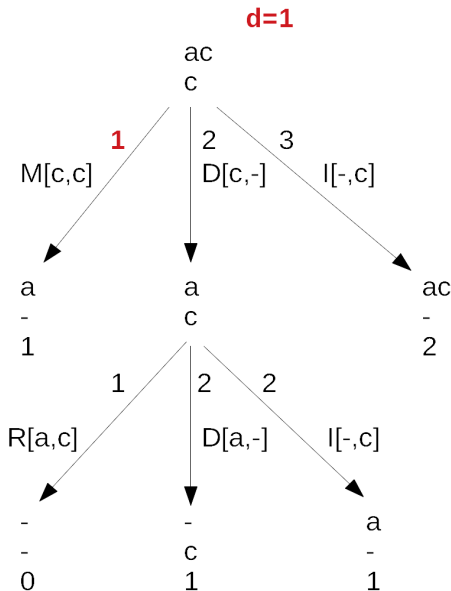
- V příkladu je rekurzivní postup aplikován od konce sekvencí. Na směru však nezáleží.
- $op[s_1[i], s_2[j]]$ - op - jedna z operací R,M,I,D, s_1, s_2 jsou porovnávané sekvence.













- Z definice okrajových podmínek:

$$d_E[0, i] = i$$

$$d_E[j, 0] = j$$

- Předpoklad, že poslední znak v řetězcích je buďto shodou, neshodou nebo insercí:

$$d_E[i, j] = \min \begin{cases} d_E[i-1, j] + 1 \\ d_E[i, j-1] + 1 \\ d_E[i-1, j-1] + \delta(s_1[i-1], s_2[j-1]) \end{cases}$$

- Funkce, která nám pomáhá rozlišit mezi shodou a neshodou:

$$\delta(a, b) = 0 \text{ pokud } a = b, \text{ jinak } \delta(a, b) = 1$$

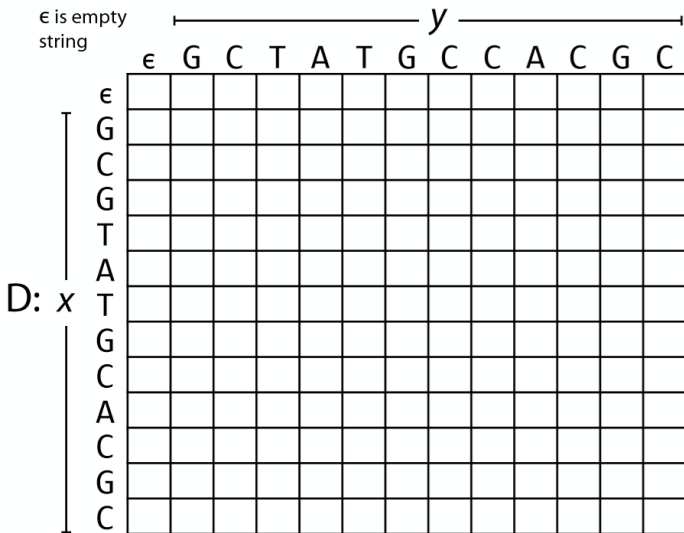
- Rekurentní předpis \Rightarrow rekurzivní funkce.



Mějme dva řetězce:

$$y = \text{GCTATGCCACGC}$$
$$x = \text{GCGTATGCACGC}$$

- Vypĺňujeme tabulku D , podle definic z předchozího slidu.
- Každá buňka tabulky obsahuje vzdálenost mezi dvěma prefixy řetězců x a y .





Doplníme hodnoty $D[0, i] = i$ a $D[j, 0] = j$

	ε	G	C	T	A	T	G	C	C	A	C	G	C
ε	0	1	2	3	4	5	6	7	8	9	10	11	12
G	1												
C	2												
G	3												
T	4												
A	5												
T	6												
G	7												
C	8												
A	9												
C	10												
G	11												
C	12												



Dle vyznačeného směru dopočítáme do tabulky hodnoty:

	ε	G	C	T	A	T	G	C	C	A	C	G	C
ε	0	1	2	3	4	5	6	7	8	9	10	11	12
G	1	→											
C	2	← →											
G	3	← →											
T	4	← →											
A	5					etc							
T	6												
G	7												
C	8												
A	9												
C	10												
G	11												
C	12												



	ε	G	C	T	A	T	G	C	C	A	C	G	C
ε	0	1	2	3	4	5	6	7	8	9	10	11	12
G	1	?											
C	2												
G	3												
T	4												
A	5												
T	6												
G	7												
C	8												
A	9												
C	10												
G	11												
C	12												

$$D[i, j] = \min \begin{cases} D[i-1, j] + 1 & \leftarrow \text{upper} \\ D[i, j-1] + 1 & \leftarrow \text{left} \\ D[i-1, j-1] + \delta(x[i-1], y[j-1]) & \leftarrow \text{upper-left} \end{cases}$$



	ε	G	C	T	A	T	G	C	C	A	C	G	C
ε	0	1	2	3	4	5	6	7	8	9	10	11	12
G	1	0	1	2	3	4	5	6	7	8	9	10	11
C	2	1	0	1	2	3	4	5	6	7	8	9	10
G	3	2	1	1	2	3	3	4	5	6	7	8	9
T	4	3	2	1	2	2	3	4	5	6	7	8	9
A	5	4	3	2	1	2	3	4	5	5	6	7	8
T	6	5	4	3	2	1	2	3	4	5	6	7	8
G	7	6	5	4	3	2	1	2	3	4	5	6	7
C	8	7	6	5	4	3	2	1	2	3	4	5	6
A	9	8	7	6	5	4	3	2	2	2	3	4	5
C	10	9	8	7	6	5	4	3	2	3	2	3	4
G	11	10	9	8	7	6	5	4	3	3	3	2	3
C	12	11	10	9	8	7	6	5	4	4	3	3	2

← Editační vzdálenost

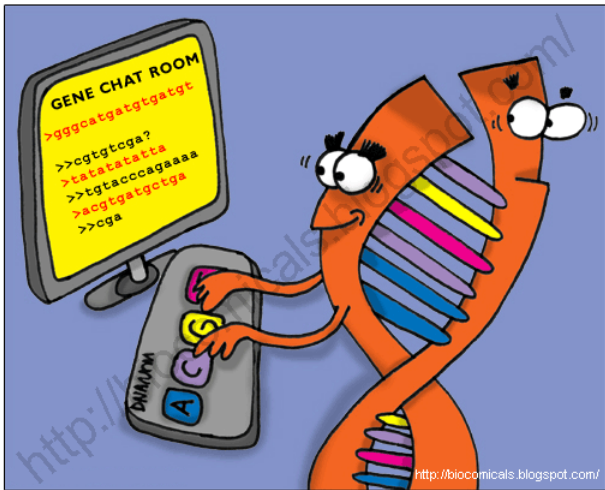


Tabulku lze vyplňovat více způsoby:

	ε	G	C	T	A	T	G	C	C	A	C	G	C
ε	0	1	2	3	4	5	6	7	8	9	10	11	12
G	1												
C	2												
G	3												
T	4												
A	5												
T	6												
G	7												
C	8												
A	9												
C	10												
G	11												
C	12												



- Vyplňujeme tabulku o m řádcích a n sloupcích.
- Editací vzdálenost nalezneme v pravém dolním políčku tabulky D .
- Prostorová a časová náročnost je $O(mn)$.
- Využití k určení podobnosti dvou sekvencí.





- Jak se zarovnávají sekvence vůči referenční sekvenci.
- Příklady algoritmů.
- Jak se modifikuje tabulka pro ohodnocení změn a proč.

DĀŽkuji za pozornost

Michal Vařinek

VřB – Technická univerzita Ostrava

FEI/EA443

michal.vasinek@vsb.cz

26. zřř 2023