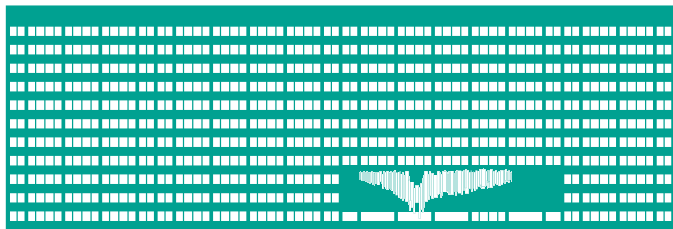


VŠB TECHNICKÁ
UNIVERZITA
OSTRAVA

VSB TECHNICAL
UNIVERSITY
OF OSTRAVA



www.vsb.cz

Algoritmy pro Bioinformatiku

Fylogenetická analýza

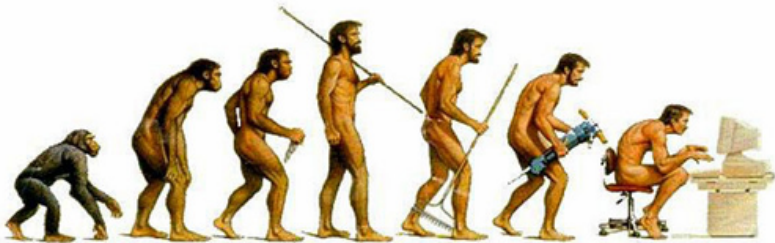
Michal Vašínek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

23. listopadu 2022





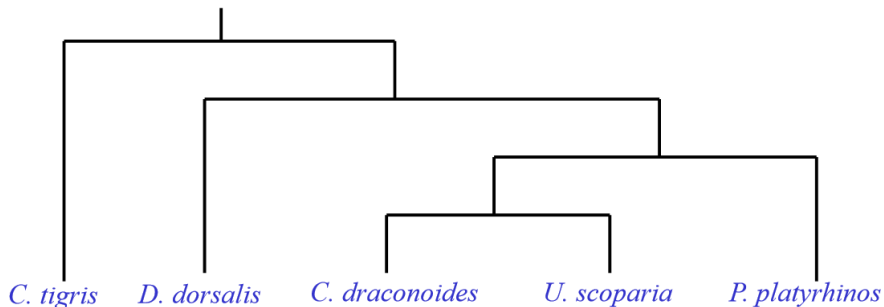
- Úvod do problematiky
- Adam a Eva
- Parsimonie
- Metody založené na vzdálenosti

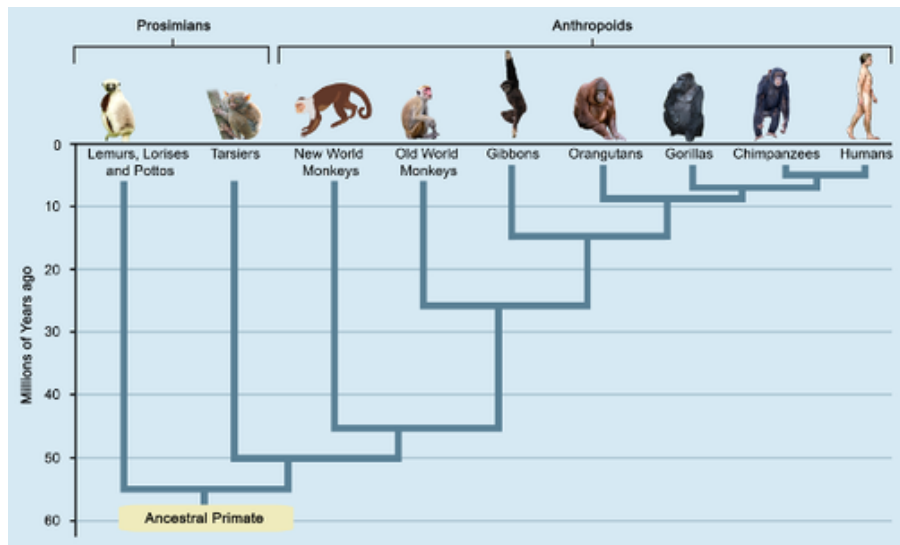


- DNA kóduje informace o živém organismu.
- Organismy předávají informaci v DNA svým potomkům.
- V důsledku mutací se DNA v průběhu času mění.
- V rámci milionů let se vyvinuly různé druhy organismů.
- Fylogenetika se zabývá studiem genetických vztahů mezi různými druhy.



- Porozumění evoluční historii organismů.
- Porozumění tomu, jak se vyvinuly u organismů různé funkce.
- K mapování mutací ve virech.
- K porozumění, co jsou klíčové stavební prvky organismů. Identifikovat je napříč organismy.



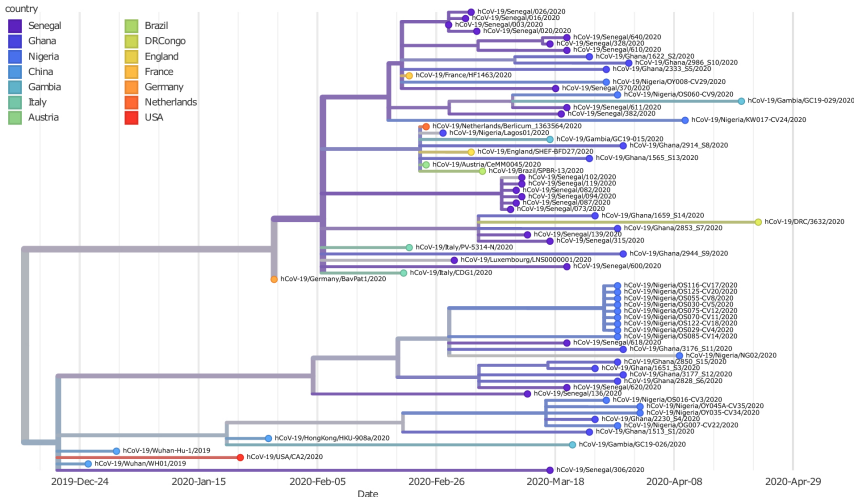




- Chřipková vakcína, založena na studiu vývoje genomu chřipkového viru.
- Studium rapidního vývoje viru HIV.
- Mutace SARS-COV2.
- Fylogenetika umožňuje zjistit zdroj infekce - konkrétní organismus.

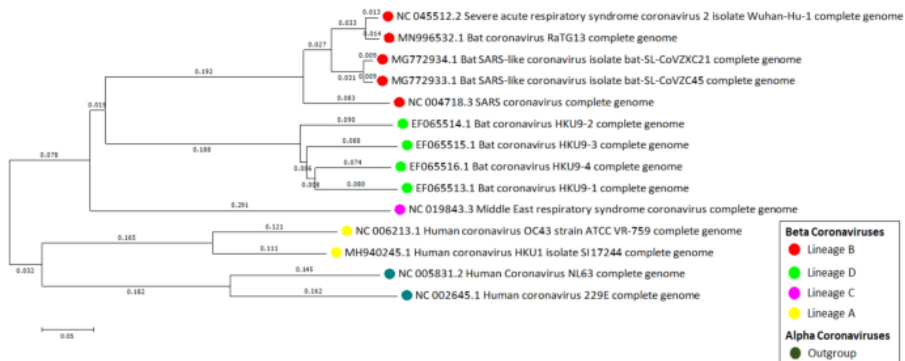
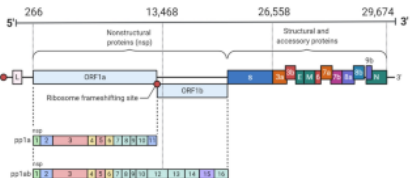
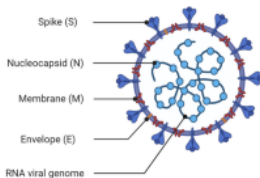


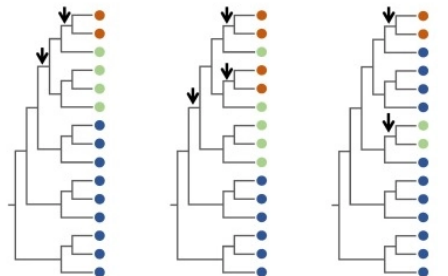
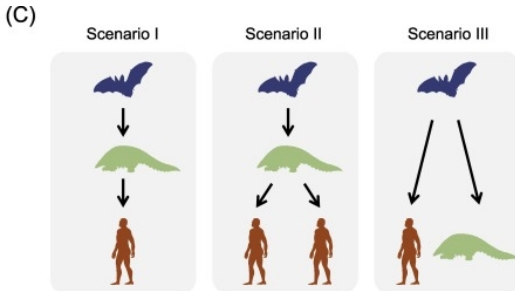
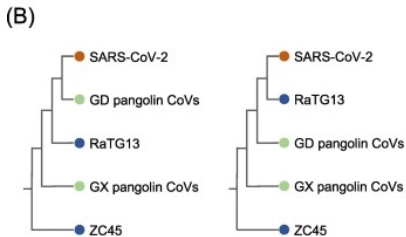
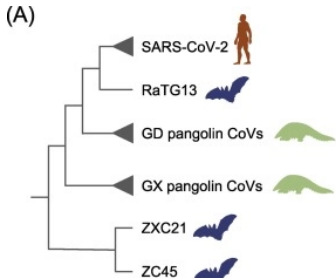
a



b









- Do poloviny 50. let se fylogenetický vývoj hodnotil na základě subjektivního hodnocení expertů.
- Počátkem 60. let snaha o nalezení objektivních kritérií pro konstrukci fylogenetických vazeb.
- Fylogenetika se stala nedílnou součástí biologie:
 - Klasifikace organismů.
 - Porozumění biologickým mechanismům.



- Klasická fylogenetická analýza používala morfologické rysy organismů:
 - Počet končetin, velikost, hmotnost, receptory, ...
- Moderní metody jsou založeny na porovnávání molekul:
 - Sekvence genů.
 - Sekvence proteinů.



- DNA - Velmi citlivá na změny, mutace probíhá různě rychle v různých oblastech genomu.
- cDNA/RNA - Užitečná pro porovnání vzdálených organismů.
- Proteinové sekvence - Porovnání evolučně vzdálených organismů, mutace probíhají rovnoměrněji.
- Příklad: Ribosomální RNA 16S - sekvence existuje ve všech prokaryotních organismech a v určité formě také v eukaryotických organismech, je jedním ze základních stavebních kamenů ribosomů => zachovává se v evoluci.



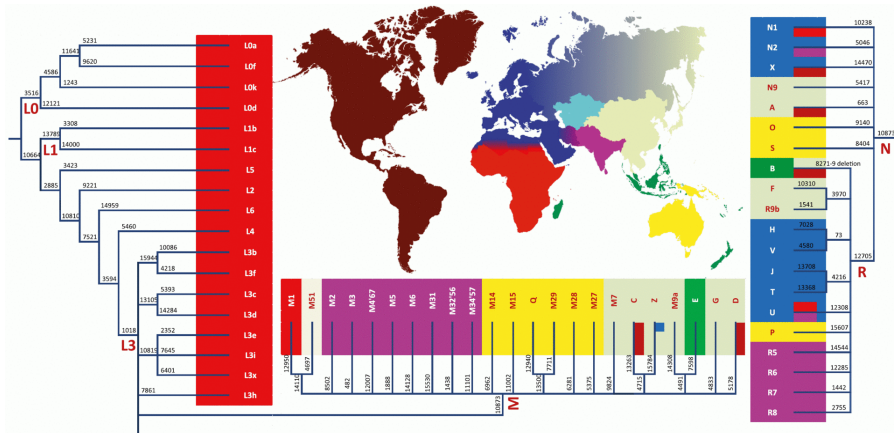
- Fylogenetickou analýzou můžeme zjistit jak probíhala evoluce člověka.
- Zkoumané sekvence:
 - Mitochondriální DNA.
 - Y chromozom.



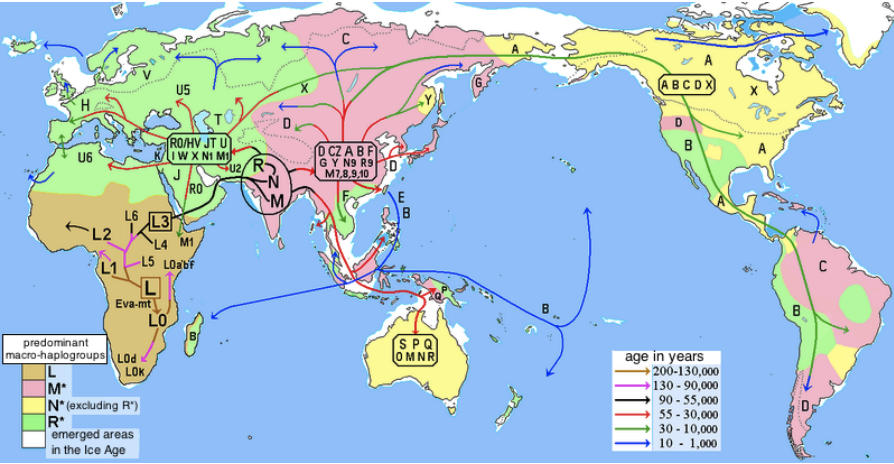
- Lidská mitochondriální DNA (mtDNA)
 - Kruhová dvoušroubovicová DNA sekvence skládající se z cca 16 500 párů bazí.
 - Každý dědíme mtDNA od matky.
 - Rychlost bodových mutací je přibližně 10x větší, než li u jaderné DNA.
 - Nedochází k rekombinaci.
- Všichni jsme podělili mtDNA od jedné původní matky člověka (Evy)!



- Vigilant, L., Stoneking, M., Harpending, H., Hawkes, K. Wilson, A. C. African populations and the evolution of human mitochondrial DNA. Science 253, 1503-1507 (1991).
 - Experiment: vybrali placentální tkáň u 147 žen různých ras z různých zemí.
 - Předpokládali konstantní molekulární hodiny - míra mutací je za určitý čas konstantní.
 - Jejich výzkum vedl k objevu, že Eva se objevila poprvé přibližně před 143 tisíci lety.



Mitochondriální DNA - Eva



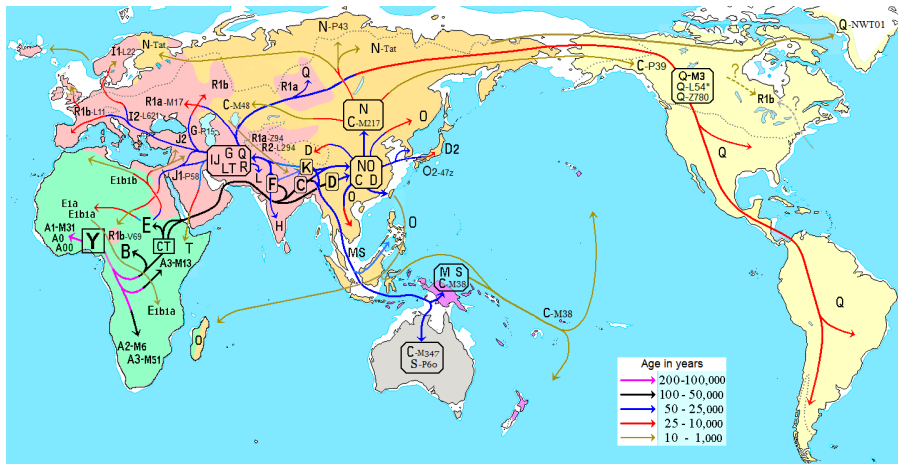


- Chromozóm Y se vyskytuje pouze u mužů => předává se od otce.
- Rychlost mutací je pomalejší, než u mtDNA => potřeba většího počtu vzorků.
- V roce 1997 studie mezi 900 muži a jejich 93 polymorfismů (bodových mutací).
 - 15% obyvatel jihoafrických buší mají konkrétní polymorfismus A
 - 5-10% Etiopanů a Súdánců mají také A
 - Zbylá část afričanů a mimoafrických obyvatel mají T.
 - Bushmen, etiopané a sudánci jsou zřejmě geneticky nejbližšími potomky Adama s chromozomem Y.



- Underhill et al. Y chromosome sequence variation and the history of human populations. Nature Genetic, 26:358-361, 2000.
- Studie 1062 mužů z 22 různých geografických oblastí.
- Identifikováno 167 haplotypů (polymorfismů, rozdílů).
- Stáří společného předka všech 167 haplotypů se odhaduje na 59 000 let.
- https://en.wikipedia.org/wiki/Human_Y-chromosome_DNA_haplogroup

Y chromozom - Adam





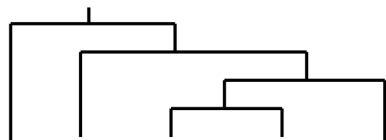
- Sekvenace DNA, RNA či proteinových sekvencí.
- Vytvoření mnohačetných alignmentů sekvencí.
- Spočtení párové vzdálenosti.
- Sestavení fylogenetického stromu.
- Odhad spolehlivosti
- Vizualizace



- Listy stromu reprezentují porovnávané objekty. (geny, jedince, druhy). Používá se pro ně výraz Taxon.
- Vnitřní uzly jsou hypotetiční předci.
- V kořenovém stromu cesta od kořene k listu reprezentuje evoluční cestu, kde kořenový uzel reprezentuje společného prapředka.
- V nekořenovém stromu jsou zobrazeny vztahy mezi organismy, ale ne evoluční cesta.

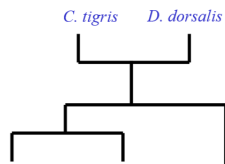


- Fylogenetické stromy jsou kořenovými stromy.
- Protože však určení společného prapředka je na vědecké úrovni velmi obtížná úloha, uvažujeme také řešení s nekořenovými stromy.



C. tigris *D. dorsalis* *C. draconoides* *U. scoparia* *P. platyrhinos*

Kořenový strom



C. draconoides *U. scoparia* *P. platyrhinos*

Nekořenový strom



Taxa	Selected Sequence Positions (sites) and character							
	1	2	3	4	5	6	7	8
1	A	A	G	A	G	T	G	C
2	A	G	C	C	G	T	G	C
3	A	G	A	T	A	T	C	C
4	A	G	A	G	A	T	C	C

- Pozice 1,6,8 nenesou informaci.
- Pozice 2,3,4 nepreferují žádný strom.
- Pozice 5,7 informativní.

Výběr sekvencí

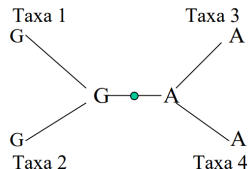
Pravidlo: Vybíráme takové pozice, kde se znaky vyskytují alespoň ve dvou taxonech a jsou zde alespoň dva různé znaky.



Taxa	Selected Sequence Positions (sites) and character							
	1	2	3	4	5	6	7	8
1	A	A	G	A	G	T	G	C
2	A	G	C	C	G	T	G	C
3	A	G	A	T	A	T	C	C
4	A	G	A	G	A	T	C	C

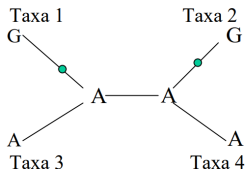
Adapted from Li and Graur 1991

Tree 1



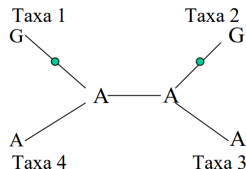
Length = 1

Tree 2



Length = 2

Tree 3



Length = 2



- Velmi podobné sekvence - Metody maximální parsimonie (maximum parsimony).
- Středně podobné sekvence - Metody založené na výpočtu vzdáleností.
- Velmi rozdílné sekvence - Metody maximální věrohodnosti (maximum likelihood)



- Nalezení stromu, který vysvětluje data s minimálním množstvím změn.
- Časově náročná metoda.
- Například programy PHYLIP a PAUP.
- Formální definice:
 - Necht' S je množina sekvencí.
 - Hammingova vzdálenost dvou sekvencí x a y je $H(x, y)$.
 - Nejúspornější (most parsimony) strom je strom T , kde listy jsou označeny sekvencemi z S a každý vnitřní uzel má přiřazenu takovou sekvenci, že platí $H(T) = \sum_{(x,y) \in E(T)} H(x, y)$ je minimální. $H(T)$ značí úspornost T (parsimony length).



- Small Parsimony problem - nalezení délky parsimonie (úspory) a ohodnocení hran daného stromu.
- Large Parsimony problem - nalezení nejvíce parsimonního stromu.



Problém nalezení nejmenší parsimonie

- Vstup: Máme množinu sekvencí S a topologii kořenového stromu T s listy označenými prvky z S .
- Výstup: Délka parsimonie stromu T a označení vnitřních uzlů T .

Problém lze řešit pomocí Fitchova algoritmu.



Fitchův algoritmus

- Vstup: strom T , kde každý list v je označen jedním symbolem v_c .
- Výstup: nejvíce parsimonní strom T s označenými vnitřními uzly.

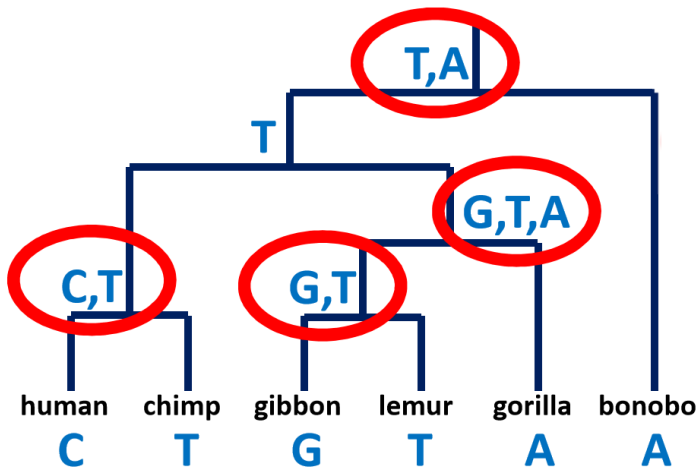
- Řešení ve dvou fázích:
 - Určení možné množiny pojmenování vnitřních uzlů.
 - Výběr symbolu pro vnitřní uzly.



- 1 Pro každý uzel v necht' $S_v = \{v_c\}$.
- 2 Pro každý vnitřní uzel v s potomky u, w necht':

$$S_v = \begin{cases} S_u \cap S_w & \text{if } S_u \cap S_w \neq \emptyset \\ S_u \cup S_w & \text{otherwise} \end{cases}$$

Délka parsimonie je rovna počtu operací sjednocení.

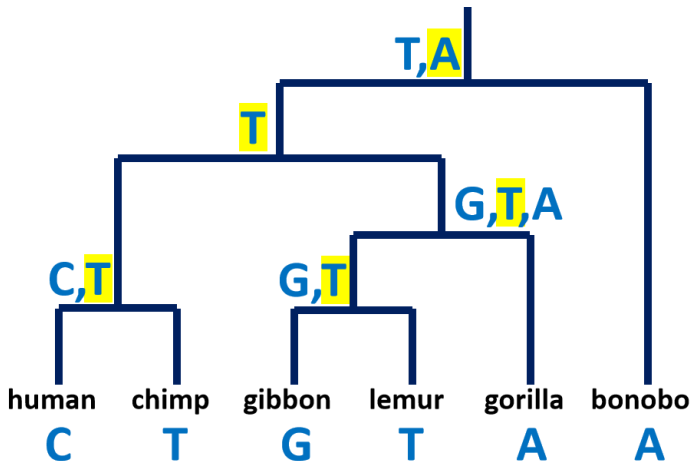




Varianta lineárního programování, volba symbolu pomocí backtrackingu.

- 1 Začneme v kořenovém uzlu, kde volíme libovolný symbol.
- 2 Procházíme stromem od kořene k listům.
- 3 Určíme v_c vnitřního uzlu v s rodičovským uzlem u :

$$v_c = \begin{cases} \text{if } u_c \in S_v \rightarrow v_c = u_c \\ \text{otherwise} \rightarrow \text{vyberme libovolný symbol z } S_v \end{cases}$$





- Symboly na různých pozicích jsou na sobě nezávislé.
- Problém se sekvencemi libovolné délky m lze řešit pomocí m řešení základní úlohy.
- Časová složitost celého postupu je: $O(mnk)$, kde k je velikost abecedy, m je délka sekvence a n je počet sekvencí.



Large parsimony problem

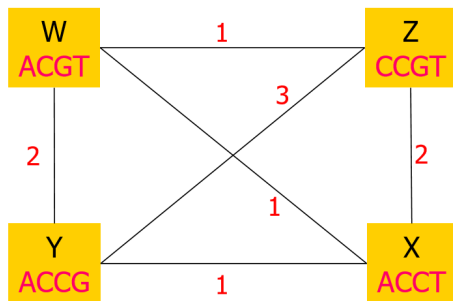
- Vstup: množina sekvencí S
 - Výstup: nejvíce parsimonní strom.
-
- Large Parsimony je NP-těžký problém.
 - Existuje však 2-aproximační polynomiální algoritmus. (Takový algoritmus, jehož řešení je garantováno nejhůře 2x horší, než je optimální řešení)



Aproximace

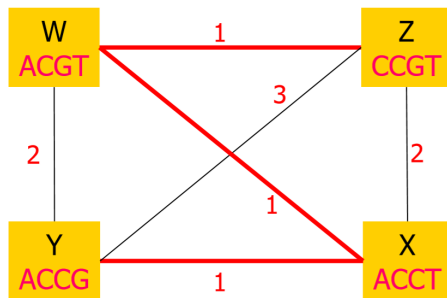
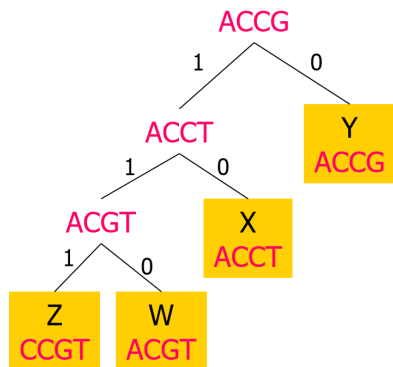
Mějme množinu sekvencí S , definujme $G(S)$ jako vážený úplný graf, jehož uzly jsou označeny sekvencemi z S a každá hrana (i, j) má váhu $H(i, j)$ (Hammingovu vzdálenost).

	1	2	3	4
W	A	C	G	T
X	A	C	C	T
Y	A	C	C	G
Z	C	C	G	T





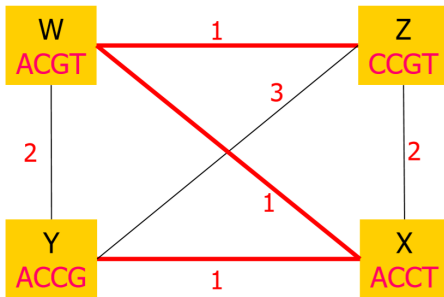
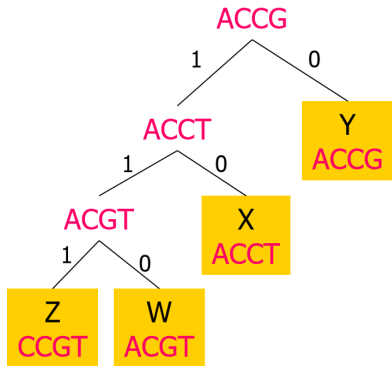
- Necht' T je minimální kostra grafu (spanning tree) $G(S)$.
- Minimální kostra grafu např. Kruskalův nebo Jarníkův algoritmus.





Věta: Aproximace minimální kostrou grafu

Něcht T je minimální kostra grafu $G(S)$. Potom délka parsimonie stromu T je nejvýše dvakrát tak velká, jako nejvíce parsimonní strom.

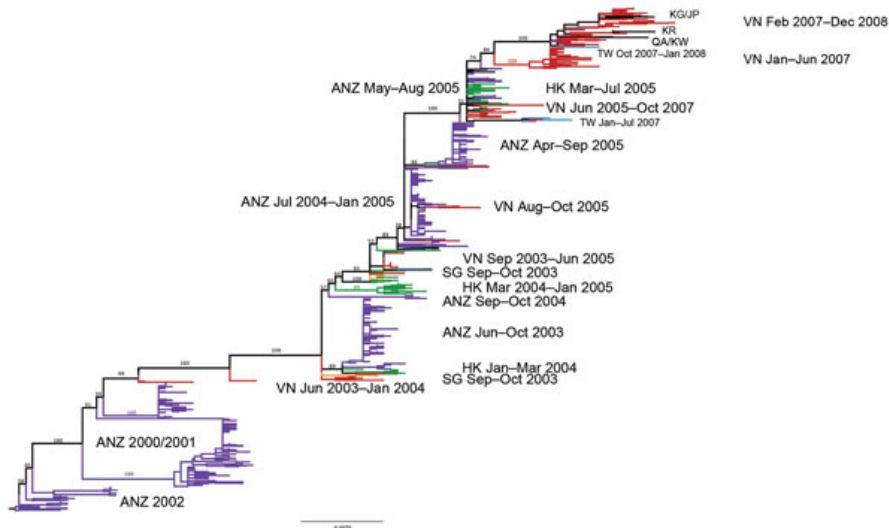




- Bush, R. M., C. A. Bender, et al. (1999) "Predicting the evolution of human influenza A." *Science* 286: 1921-1925.
- Chřipka je rapidně se vyvíjející virus.
- Bush a kolektiv ukázali, že virus lidské chřipky typu A (podtyp H3) může být použit pro predikci evolučního vývoje budoucích chřipkových kmenů.
- Tyto predikované kmeny jsou součástí chřipkové vakcíny při očkování proti chřipce.



- CDC - center for disease control.
- Každý rok osequenují cca 7000 genomů chřipky.
- Určení podobnosti chřipkových genomů.
- Monitorování evoluce chřipkových virů.
- Vyhodnocení účinnosti vakcín.





- Cílem je vygenerovat strom, ve kterém podobné sekvence s krátkou vzdáleností jsou blízko a součet délek hran dvou uzlů je roven jejich vzdálenosti.
- ClustalW
- PAUP
- PHYLIP, DNADIST, PROTDIST pro vygenerování matice vzdáleností.



- UPGMA - unweighted pair group method with arithmetic mean.
- Předpoklad molekulárních hodin (evoluce konstantí rychlostí).
- Vytváří kořenové stromy.
- Podmínka ultrametricky: pro libovolné tři taxony a, b, c platí:

$$d_{ac} \leq \max(d_{ab}, d_{bc})$$



- Inicializace: Definujme množinu všech listových uzlů T , jeden pro každou sekvenci. Nechť výška každého uzlu je 0. Nechť $L=T$.
- Opakujte dokud je v L víc než jeden uzel.
 - 1 Vyberte dva nejbližší uzly (A,B) a vytvořte z nich rodičovský uzel K . Spojte A,B a K . Nastavte výšku K na $d_{ab}/2$. Nastavte délku větve mezi K a $A =$ výška $K -$ výška A , nastavte délku větve mezi K a $B =$ výška $K -$ výška B .
 - 2 Odeberte A,B z L a přidejte K do L . Přepočítejte vzdálenosti mezi K a dalšími uzly v L . Vzdálenost mezi K a dalšími uzly je dána průměrnou vzdáleností potomků K k ostatním uzlům.

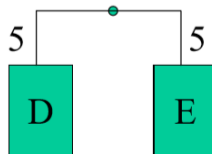


	A	B	C	D	E
A	-	20	26	26	26
B		-	26	26	26
C			-	16	16
D				-	10
E					-



	A	B	C	DE
A	-	20	26	26
B		-	26	26
C			-	16
DE				-

Step 1: Select D and E



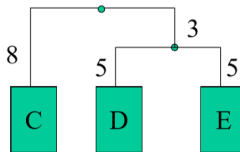


	A	B	C	DE
A	-	20	26	26
B		-	26	26
C			-	16
DE				-

Step 2: Select (DE) and C

↓

	A	B	DEC
A	-	20	26
B		-	26
DEC			-



$$\text{dist}(\text{DEC}, \text{A}) = (d_{\text{DA}} + d_{\text{EA}} + d_{\text{CA}}) / 3 = 26$$

$$\text{dist}(\text{DEC}, \text{B}) = (d_{\text{DB}} + d_{\text{EB}} + d_{\text{CB}}) / 3 = 26$$

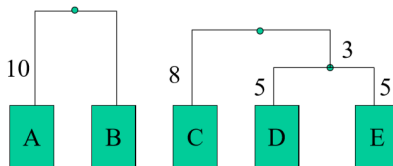


	A	B	DEC
A	-	20	26
B		-	26
DEC			-



	AB	DEC
AB	-	26
DEC		-

Step 3: select A, B

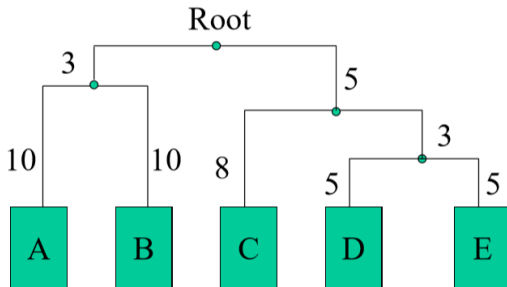


$$\text{dist}(\text{DEC}, \text{AB}) = (d_{\text{DA}} + d_{\text{DB}} + d_{\text{EA}} + d_{\text{EB}} + d_{\text{CA}} + d_{\text{CB}}) / 6 = 26$$



	AB	DEC
AB	-	26
DEC		-

Step 4: select (A,B), (D,E),C

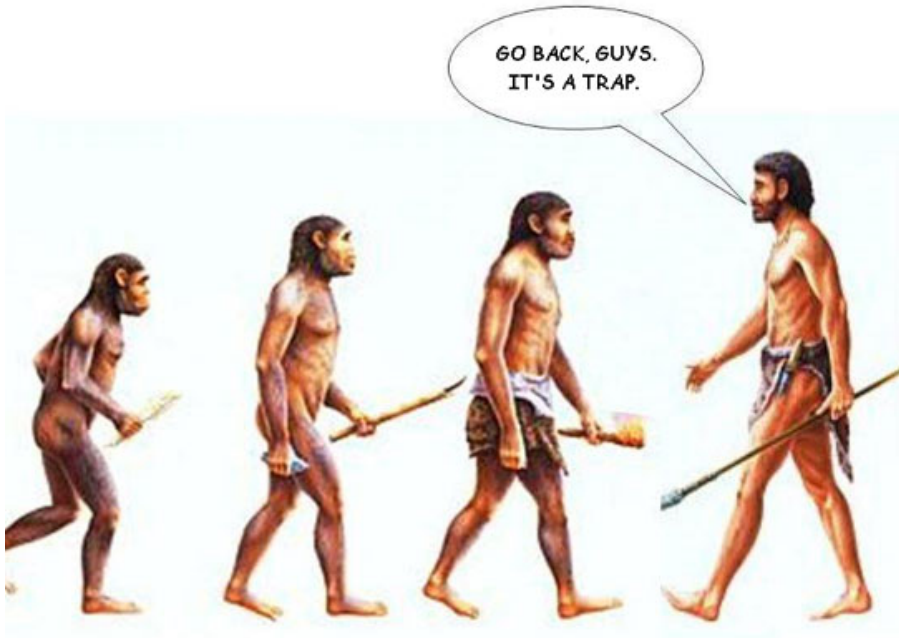




- Celkový počet bodových neshod.
- Celkový počet bodových neshod a indelů.



- Bootstrapping
- Mějme sekvence o délce m , vybíráme podmnožinu symbolů sekvencí.
- Pro každou náhodně vybranou podmnožinu vytvoříme fylogenetický strom.
- Opakujeme požadovaný počet krát.
- Ověříme vztahy mezi dvojicemi sekvencí. Pokud je jejich vazba stabilní \Rightarrow vyskytuje se ve většině stromů, pak jsme nejspíše obdrželi fylogenetický strom blízký skutečnému stromu.



DĚKUJI za pozornost

Michal Vašinek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

23. listopadu 2022