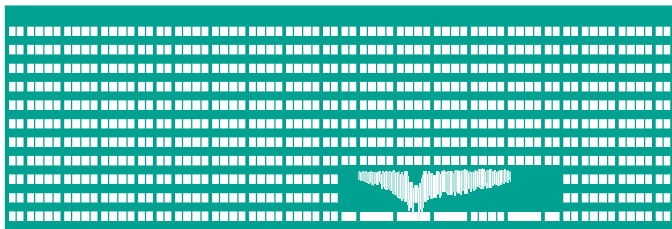


VŠB TECHNICKÁ
UNIVERZITA
OSTRAVA

VSB TECHNICAL
UNIVERSITY
OF OSTRAVA



www.vsb.cz

Algoritmy pro Bioinformatiku

Indexace Genomu

Michal Vašínek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

2. října 2019



- Trie
- Suffixový strom
- Suffixové pole



Vyhledávání řetězců

Mějme text T o délce $n = |T|$ a řetězec P o délce $m = |P|$. Základní úlohy:

- Find - nalezněte pozice všech výskytů P v T .
- Exists - vyskytuje se P v T ?
- Occ - Kolikrát se vyskytuje P v T ?

Řešení standardními algoritmy pro prohledávání textů:

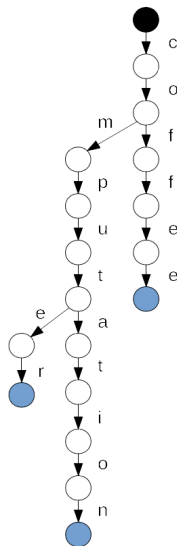
- Hrubou silou $O(nm)$.
- Boyer-Moore algoritmus - předzpracování vzorku $O(n + m)$.

Lidský genom $|G| = 3.2$ mld symbolů. FastQ soubor 1 mil sekvencí.
Použitím Boyer-Moorova algoritmu $10^6|G|$ porovnání symbolů.

Indexace Textu Trie

$T = \{\text{computer, computation, coffee}\}$

- Trie (ang. výslovnost "try") je stromová struktura reprezentující množinu řetězců:
 - Každá hrana je označena symbolem $c \in \Sigma$.
 - Pro libovolný uzel stromu platí, že uzel má nanejvýš jednu hranu označenou symbolem c .
 - Vložený řetězec jsme schopni odečíst průchodem od kořene stromu k listu.



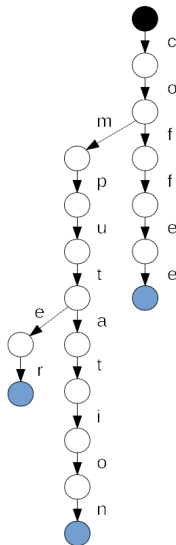

 $T = \{\text{computer, computation, coffee}\}$

Použití mimo bioinformatiku:

- Nahrazení hashovací tabulky => vyhledávání klíčů.
- Repräsentace slovníku => doplňování slov v mobilních klávesnicích.

Bioinformatika:

- Indexace suffixů.
- Indexace k-merů.

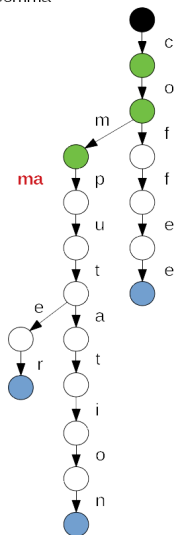


- Při ověření existence vzoru P o délce $|P| = m$ procházíme nejvýše m hran.
- Pokud do trie vložíme w řetězců o celkové délce W , pak počet hran v trie je nejvýše W .

Implementace:

- Hashovací tabulka \Rightarrow předpokládáme složitost $O(1)$ pro nalezení symbolu v uzlu, pak vyhledání řetězce má složitost $O(m)$.
- Setřazené pole \Rightarrow binárním vyhledáváním nalezneme symbol v čase $\log \sigma$, tudíž řetězec nalezneme v čase $O(m \log \sigma)$.

$T = \{\text{computer, computation, coffee}\}$
 $P = \text{comma}$





- Dosud jsme do trie vkládali oddělené řetězce.
- Pro vyhledávání v dlouhém textu T vkládáme do trie všechny suffixy textu T .

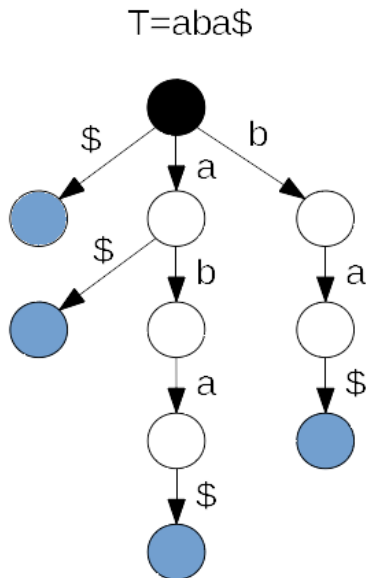
```

T:  c o f f e e $
    c o f f e e $
      o f f e e $
        f f e e $
          f e e $
            e e $
              e $
                $
  
```

Celkem mají všechny suffixy $n(n+1)/2$ symbolů. Prostorová složitost $O(n^2)$.

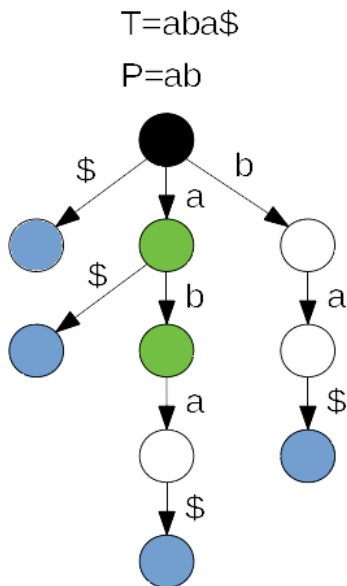
Každý podřetězec T je prefixem nějakého suffixu.

- K textu T přidáváme na konec speciální terminální symbol $\$,$ menší, než všechny ostatní symboly.
- Speciální symbol $\$$ nám zajistí dvě důležité vlastnosti:
 - Kratší řetězce začínající stejným prefixem jsou při setřizení vepředu. např "e" a "ee" ve slově coffee.
 - Žádný suffix T nemůže být zároveň prefixem jiného suffixu.
- Trie sestavujeme od nejdelšího po nejkratší suffix.



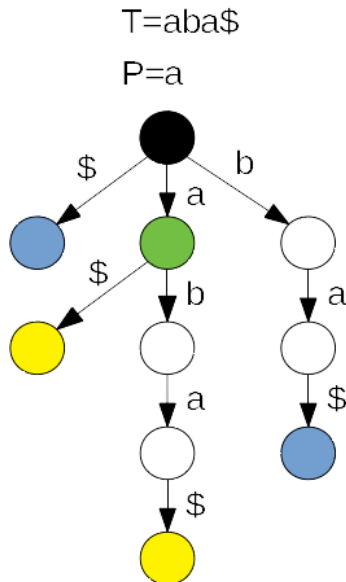


- Každá cesta od kořene k listu reprezentuje jeden suffix.
- Ověření výskytu řetězce P v T provedeme opět průchodem od kořene.



Výpočet počtu výskytů:

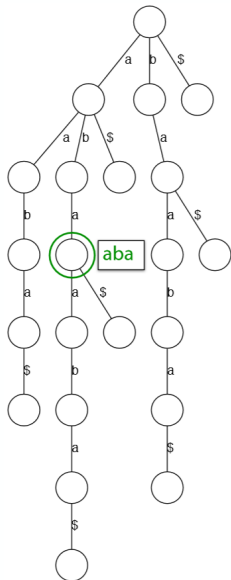
- Postupujeme od kořene k uzlům.
- Pokud se řetězec P v trie nevyskytuje, pak je počet roven 0, jinak je počet výskytů P dán počtem listových uzlů odpovídajících podstromu začínajícímu v uzlu, do kterého jsme dorazili průchodem přes symboly z P .
- Průchod do hloubky.





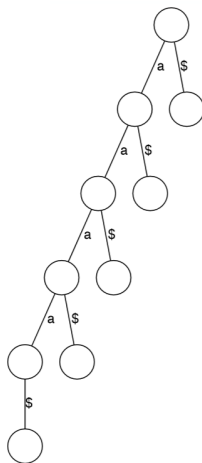
Výpočet nejdelšího opakujícího se podřetězce:

- Najděte uzel s alespoň dvěma potomky, ke kterému je nejdelší cesta od vrcholu.



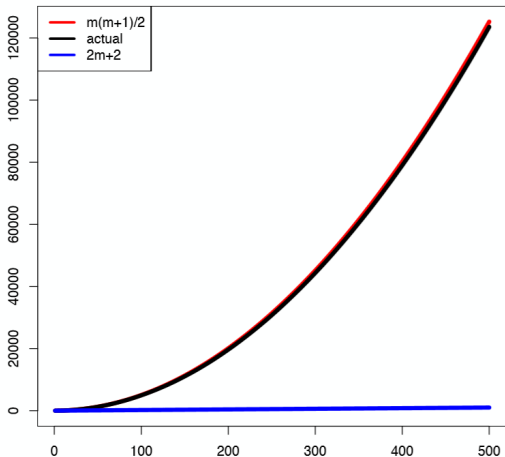


- Uvažujme řetězec $T = a^n$, tedy n -opakování symbolu a .
- Počet uzlů:
 - 1 kořenový
 - n uzlů se vstupující hranou "a".
 - $n + 1$ uzlů se vstupující hranou "\$".
 - Celkem $2n + 2$ uzlů.
- V nejhorším případě $n(n + 1)/2 = O(n^2)$ uzlů.





Lambda fág genom. Počet uzlů roste přibližně rychlostí nejhoršího případu.



Suffixový strom

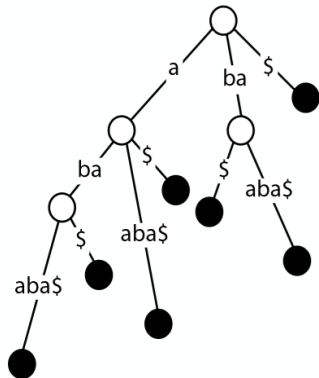


- Víme, že ověření zda řetězec P je obsažen v T lze pomocí trie vyřešit v $O(m)$.
- Prostorová složitost nicméně $O(n^2)$. Lidský genom $|G| = 3.2 \times 10^9$
 $\Rightarrow |G|^2$.
- Cíl: zmenšit paměťové nároky \Rightarrow suffixový strom.

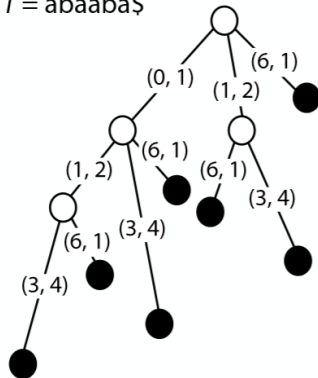


Myšlenka 2

Uložit si společně se stromovou strukturou samotný text T a převést popisky hran na dvojice (offset,délka) v T .

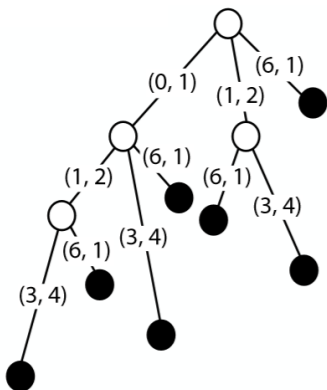


$T = \text{abaaba\$}$





- Offset a délka nabývají maximální hodnoty $|T| = n$ a lze je uložit pomocí $2\lceil \log n \rceil$ bitů.
- Celková prostorová složitost: $2n - 1$ uzlů vede k uchování $2n - 2$ hran. Počet dvojic (offset, délka) $\Rightarrow (2n - 2)2 \log n = O(n \log n)$.

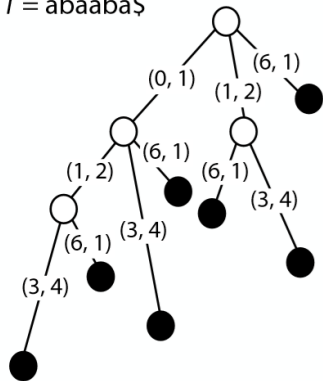




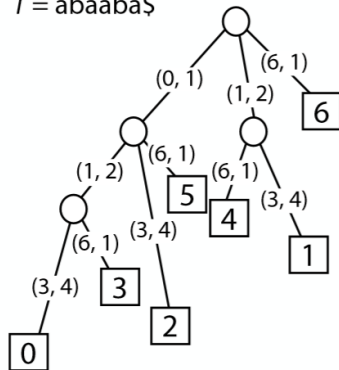
Myšlenka 3

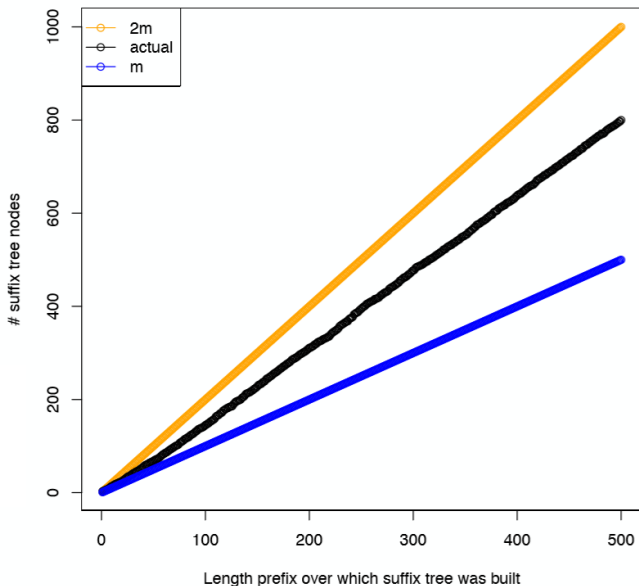
Do listových uzlů uložíme začátek suffixu v textu T .

$T = \text{abaaba}\$$



$T = \text{abaaba}\$$





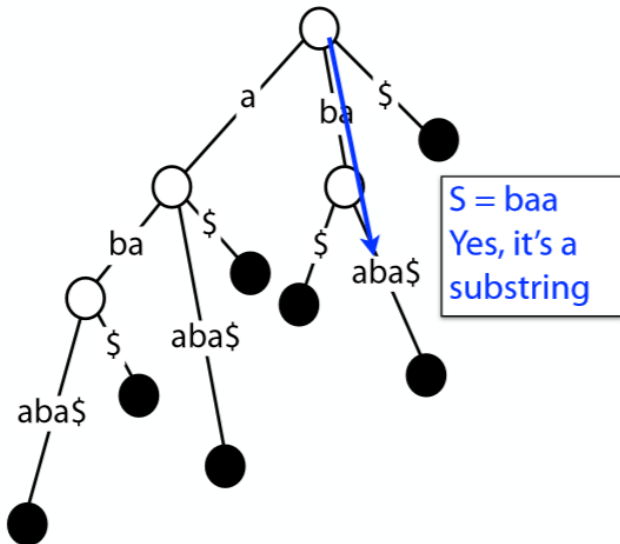


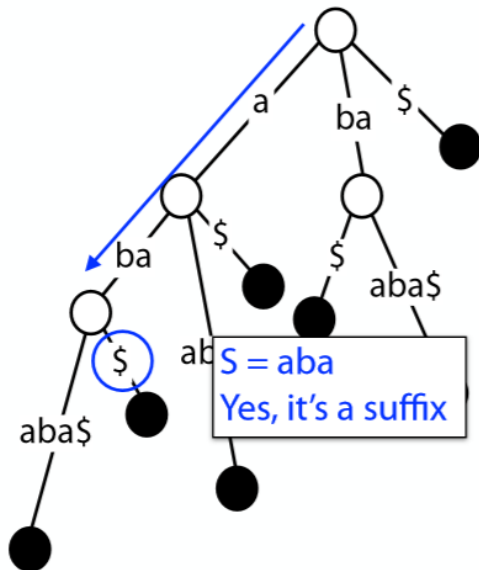
Dva algoritmy k sestavení:

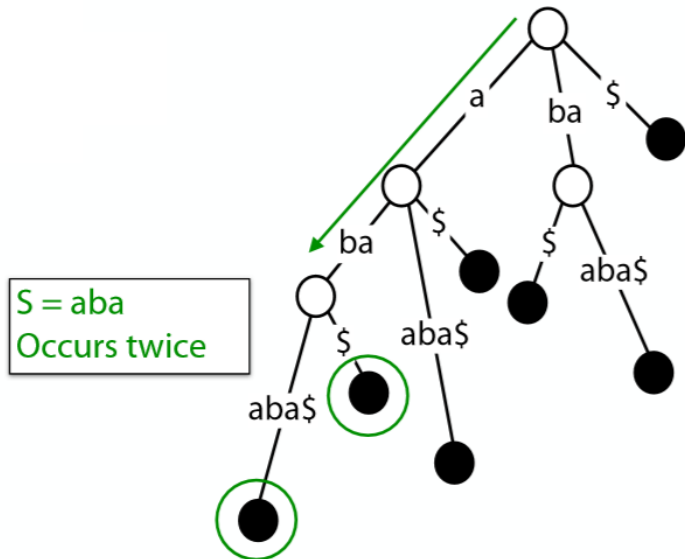
- Trie \Rightarrow Suffixový strom - $O(n^2)$ čas a prostor.
 - 1 Sestavit trie
 - 2 Spojit nevětvící se uzly a přeznačit označení hran.
- Přímé sestavení suffixového stromu - $O(n^2)$ čas, $O(n\sigma \log n)$
 - 1 Sestavte strom pro nejdelší suffix.
 - 2 Upravit strom, aby obsahoval druhý nejdelší suffix.
 - 3 Upravit strom, aby obsahoval třetí a další suffixy...



- Weiner's algorithm [1973]
- McCreight [JACM 1976]
- Ukkonen [Algorithmica 1995] - $O(m)$ čas a prostor.
- Farach [FOCS 1997]
- Hon, Sadakane, Sung [FOCS 2003]







Suffixové pole



- Suffixový strom pro lidský genom je sice teoreticky $O(m)$, nicméně v realitě to znamená $> 45\text{GB}$.

$T = \text{abaaba\$}$
 0123456

Suffixové pole

Suffixové pole je pole kladných celých čísel, které reprezentují počáteční pozice suffixů v lexikografickém pořadí.

$SA(T) =$

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$



- K řetězci o délce n přidáváme opět terminální symbol.
- Pole tak obsahuje $n + 1$ indexů.
- Každý index vyžaduje $\lceil \log n \rceil$ bitů.
- Prostorová složitost $O(n \log n)$.
- Lidský genom cca 12 GB plus původní text.

$T = \text{abaaba\$}$
 0123456

SA(T) =

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$



- Hledáme vzor P .
- Jestliže se P vyskytuje v T , pak P musí být prefixem nějakého suffixu.
- Navíc, suffixy sdílející stejný prefix jsou setřizeny v suffixovém poli za sebou.
- Nalezení P binárním vyhledáváním.

$T = \text{abaaba\$}$
0123456

SA(T) =

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$



- $P = ba$.
- Začneme na pozici $SA[n/2] = 3$.
 $T[SA[3]..] = aba\$$.
- $aba\$ < ba$ prohledáváme spodní polovinu suffixového pole.
- $SA[(6 + 4)/2] = 5$.
 $T[SA[5]..] = ba\$$. Shoda P a suffixu od pozice 5.
- Prohledáme pozice nad a pod indexem číslo pět, abychom zjistili rozsah pro prefix $ba \Rightarrow [4,1]$ indexy v textu.

$T = abaaba\$$
0123456

$SA(T) =$

6	$\$$
5	$a \$$
2	$a a b a \$$
3	$a b a \$$
0	$a b a a b a \$$
4	$b a \$$
1	$b a a b a \$$



- Provádíme $\log n$ půlení intervalu.
- V každém kroku navíc provedeme až m porovnání symbolů $\Rightarrow O(m)$ při každém půlení
- Nalezení prvního výskytu je tedy operace s $O(m \log n)$ časovou složitostí.
- Nalezení všech k výskytů vyžaduje dalších až $k + 2$ porovnání, tudíž dohromady $m \log n + mk \Rightarrow O(m(k + \log n))$ složitost.

$T = \text{abaaba\$}$
 0123456

SA(T) =

6	\$
5	a \$
2	a a b a \$
3	a b a \$
0	a b a a b a \$
4	b a \$
1	b a a b a \$



- S použitím běžných třídících algoritmů $O(n^2 \log n)$.
- V praxi 80% řetězců délky 20 jsou unikátní \Rightarrow
 $O(20n \log n) = O(n \log n)$
- Například Quicksort, Heapsort.
- Karkainen [2003] - $O(n)$ algoritmus.
- Ko, Aluru [2003] - SAIS, $O(n)$ algoritmus.



- Spočítáme si počty stejných prefixů o konstantní délce p u všech suffixů. $O(pn) = O(n)$
- Rozdělíme si suffixové pole do přihrádek podle počtu suffixů se stejným prefixem - RadixSort.
- Každou přihrádku, pak seřídíme pomocí Quicksortu.
- Dostatečně rychlé řešení pro indexaci genomu.
- Teoreticky nejlepší výkon pro $p = 3$, cca 10x rychlejší, než samotný Quicksort.



When you need to make a lexicographically ordered suffix array

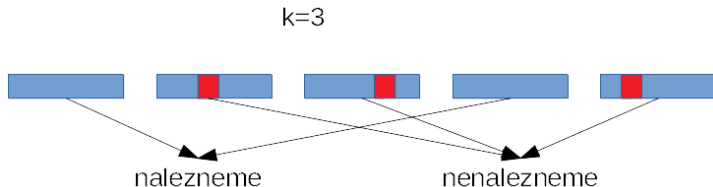




- Struktury umožňují nalézt pouze exaktní shody vzoru a textu => mutace?.
- V angl. literatuře Pigeonhole principle.
- Používá se např. v aligneru Bowtie.

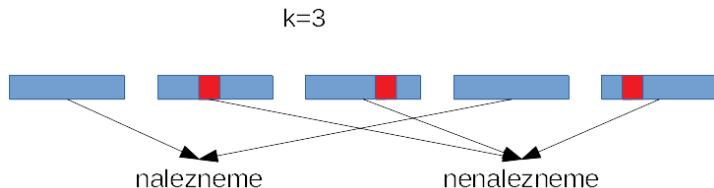
Dirichletův princip

Uvažujme, že v našem systému tolerujeme k mutací. Pak pokud vstupní sekvenci rozdělíme na $k + l$ nepřekrývajících se řetězců, máme jistotu, že alespoň l řetězců nebude obsahovat mutaci.





- 80% sekvencí délky 20nt jsou unikátní.
- Každou hledanou sekvenci rozdělíme na $k + l$ řetězců.
- Každý řetězec zkusíme nalézt v referenčním genomu a pozice jeho výskytu si uložíme do seznamu.
- Seznamy spojíme do jednoho seznamu/pole a setřídíme.
- Pokud se nějaké pozice vyskytnou alespoň l -krát, pak na danou sekvenci od dané pozice aplikujeme dynamické programování a vybereme pozici s nejvyšším skóre.





- Ukkonenův algoritmus pro sestavení suffixového stromu v $O(n)$.
- SAIS algoritmus pro sestavení suffixového pole v $O(n)$.



- Backtracking v suffixovém poli
- FM-Index
- Dokončení alignmentu - kombinace FM-Indexu a dynamického programování.

DĚKUJI za pozornost

Michal Vašinek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

2. října 2019