

Přibližné vyhledávání nad vlastními indexy

Vedoucí: Jan Holub

Lukáš Hrbek

Problém:

- Přibližné vyhledávání s nejvýše k chybami
- Chyby definované pomocí Hammingovy, či Levenshteinovy vzdálenosti
- Zobrazení všech výskytů a jejich zarovnání
- Vyhledávání v komprimovaných datech

Řešení:

- Zobecněný filtrační algoritmus založený na pigeonhole principu
- Matematický model pro vhodné nastavení algoritmu
- Sublineární složitost vyhledávání při vhodném nastavení
- Použití index FM-Index verze 2

Motivace:

Přibližné vyhledávání ve velkých textech je využíváno při výzkumu DNA, např. pro hledání nových léků.

Přibližné vyhledávání:

Dovoluje vyhledat také výskyty, které se od hledaného vzorku mírně liší. Rozdíly mohou být definovány takto:

- Replace CAG → CTG
- Insert CAG → CATG
- Delete CAG → CG

Vlastní index:

Indexování původního textu umožňuje rychleji provádět vyhledávání (např. složitost nebude záviset na délce prohledávaného textu). Vlastní index navíc umožňuje rekonstrukci původního textu - ten proto nemusí být uložen. Protože indexový soubor je menší, dochází ke kompresi původního textu.

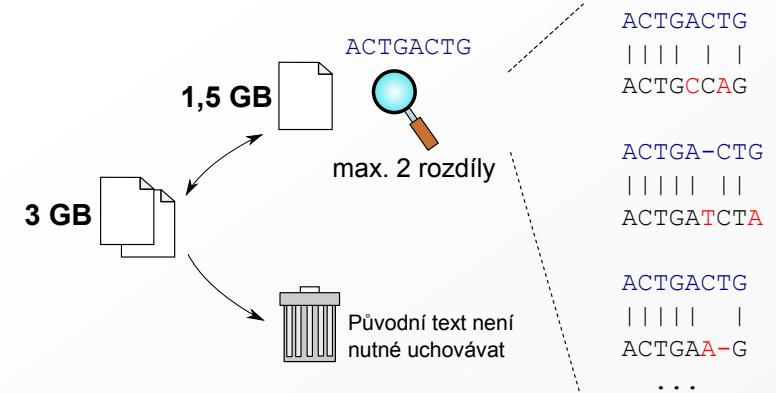
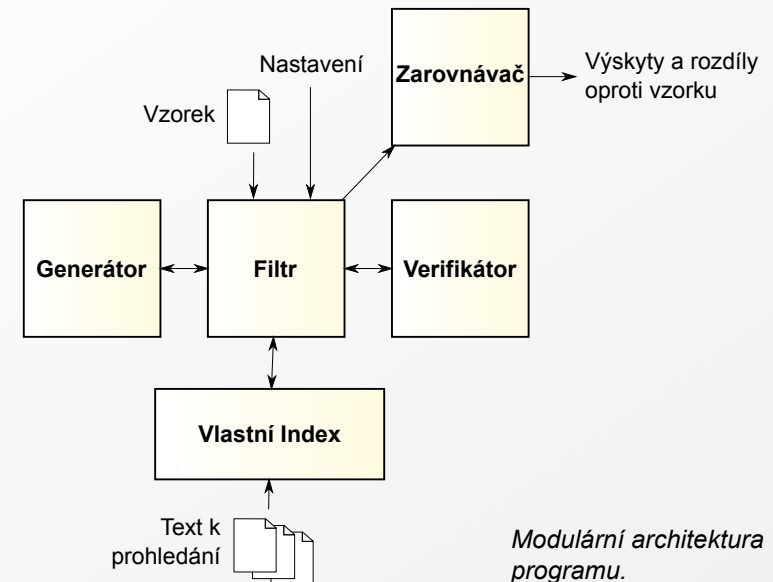
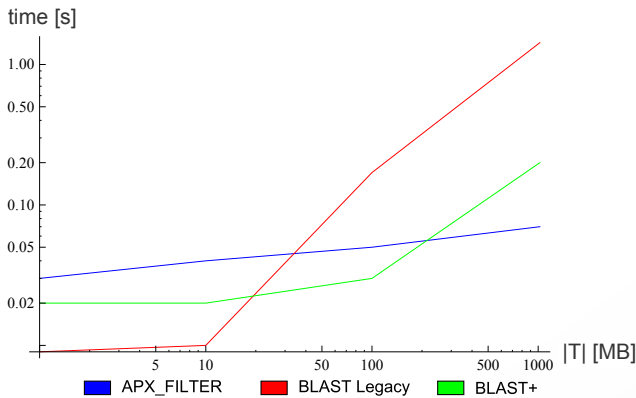


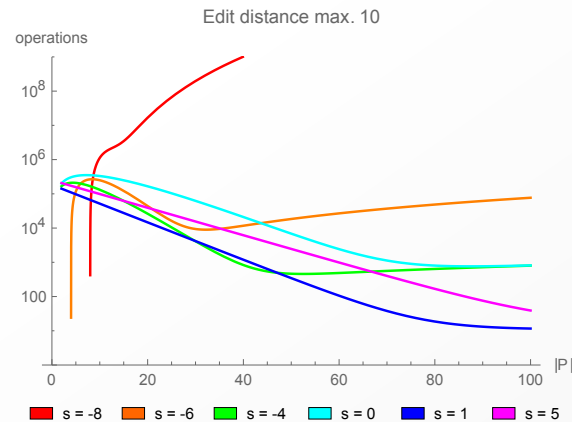
Schéma typického použití programu.



Modulární architektura programu.



Logaritmický graf vlivu velikosti původního textu (|T|) na čas vyhledávání.



Analytický model odhadující dobu výpočtu v závislosti na délce vzorku (|P|) pro různé nastavení filtračního algoritmu.

Srovnání:

Tyto problémy řeší BLAST. Jedná se však jen o heuristickou metodu specializovanou pro DNA. Nové řešení je pro velké soubory rychlejší než BLAST a pro menší soubory využívá významně méně paměti.