

Definice problému

Vyhledání podřetězce v textu (Substring matching). Cílem je nalezen všechny ch v skytu hledaného řetězce v zadaném textu.

Naivní algoritmus

Datově nezávislý, se zaručenou složitostí $\Theta(m \cdot n)$. Spočívá v průchodu všechny ch možností vzájemných pozic řetězců a jejich následném porovnání.

Předčasné ukončení vnitřního cyklu může v ideálním případě urychlit v počtu až na $\Omega(n)$, horní složitost zůstává $\Theta(m \cdot n)$.

Knuth - Morris - Pratt algoritmus

Algoritmus používá předzpracování vzorku pro urychlení vyhledání. Dobu vyhledání zkracuje pomocí posunu o více než jednu pozici v závislosti na struktuře dat. Algoritmus nejprve vyplní tabulku posunů a v druhé fázi pak spustí samotný vyhledání.

V případě neshody znaku se podle tabulky a provede předpokládaný posun v závislosti na struktuře vzorku.

Pro v počtu tabulky potřebujeme $\Theta(m)$ operací a pro samotný vyhledání $\Theta(n)$. Paměťová složitost je $\Theta(m) = \Omega(m)$ pro tabulku posunů.

Měření

Měření času bylo provedeno na serveru star.fit.cvut.cz, pomocí knihovny funkce `omp_get_wtime`. Každé měření proběhlo jednotným způsobem, po načtení vstupních hodnot programem byl změřen početný čas, spuštěn algoritmus v počtu a po jeho dokončení spočten rozdíl čas v sekundách. Přibližná hodnota MIPS byla měřena během zvláštního běhu programu, aby nebyl ovlivněn čas měření.

Optimalizace

Protože oba algoritmy jsou silně datově závislé a oba mají různou složitost, měřil jsem jejich časy odděleně a na jiných instancích problému. Pro algoritmus KMP jsem použil přibližně 16x větší vstupní řetězec než pro naivní implementaci.

V sledky

Po zapnutí optimalizací se průměrný čas běhu zkrátil přibližně na polovinu v důsledku optimalizací kompilátoru. Do srovnávacích cyklů jsem musel vložit nadbytečnou operaci, aby nedošlo k jejich odstranění v důsledku agresivní optimalizace.

Vylepšením, aplikovaným na oba algoritmy bylo rozdělení vstupního textu na části, které se vejdou do cache procesoru. Po změření jsem nezaznamenal větší zrychlení, protože už původní verze měla hodnotu cache-miss velice nízkou. A to i přes značně veliký vstupní řetězec, zatím pouze odhaduji, že v nové verzi bude na vřady trpět více.

Dalším vylepšením použitým na oba přístupy byla změna datového typu řetězce z `uint8_t` na `uint64_t`. Po tom průměrný čas klesl v závislosti na vstupu přibližně 6 – 8 krát pro oba algoritmy.

Na zrychlenou verzi naivního algoritmu jsem aplikoval ruční unrolling nejvnitřnějšího cyklu. Pro některé vstupy se čas téměř nezměnil, u jiných nastalo cca 30% zrychlení.

Mimo tyto opravy jsem změnil vnitřní cyklus naivního postupu kvůli vektorizaci.

Nevhodou je nalezen všech rozdílů vůči textu mnoho prvních, kvůli nemožnosti větvení ve vektorizovaném cyklu. Tím se průměrná složitost dostala až na $\Omega(m^2)$, urychlení je dosaženo pouze díky vektorizaci, která násobně snížila počet průchodů cyklem.

Protože ale došlo k logické změně algoritmu, nedá se v sledek porovnat s ostatními.

V sledky jsou vyneseny v příložených grafech.

Navíc jsem vyhodnotil přibližně zkrácení doby v počtu během v počtu MIPS. Měření došlo k průměrnému nárůstu času o 10 – 15%.