

Zwielowatkowanie problemu Longest Common Sequence (LCS)

Maciej Sanocki

November 2023

1 Wprowadzenie

Znajdowanie najdłuższych podciągów dwóch tekstów okazuje się kluczowe w rozważaniu wielu problemów z różnych dziedzin. Często jest, że interesuje nas porównanie bardzo długich ciągów, na przykład kodów genetycznych, w jak najkrótszym czasie co tłumaczy chęć optymalizacji czasu wykonania.

Rozważę dwa algorytmy dynamiczne rozwiązujące problem LCS:

- Najczęściej stosowany algorytm dynamiczny, obliczający LCS w czasie $O(n^2)$, i pamięci $O(n^2)$
- Algorytm Hirschberga, który również działa w czasie $O(n^2)$, ale za to potrzebuje liniowej pamięci.

2 Programy

Wszystkie programy na wejściu przyjmują liczbę testów z i długość każdego ciągu znaków w teście n , a następnie z par ciągów, każdy o długości n .

Sekwencyjne: LCS_01.cpp Algorytm wyliczający LCS, przechodzący po kolejnych wierszach tablicy wynikowej.

LCS_02.cpp Algorytm wyliczający LCS, przechodzący po kolejnych przekątnych tablicy wynikowej, jednak tablica wynikowa nadal jest ustawiona wierszami.

LCS_03.cpp Algorytm wyliczający LCS, przechodzący po kolejnych przekątnych tablicy wynikowej, specjalnie ułożonej.

LCS_04.cpp Algorytm wyliczający LCS, przechodzący po kolejnych wierszach tablicy wynikowej, ale używający trick-u (numer paragrafu) pozwalającego na nie odwoływanie się do tego samego wiersza w tablicy wynikowej.

Hirschberg_0.cpp Algorytm Hirschberg-a.

Zrównoległone: LCS_threads1.cpp Bardzo zła implementacja programu wielowatkowego, głównie przez tworzenie wątków niepotrzebnie wiele razy.

LCS_threads2.cpp zwielowatkowana wersja programu programu LCS_02.cpp używająca `std::barrier` do synchronizacji przed przejściem do kolejnej przekątnej.

LCS_omp2.cpp zwielowatkowana wersja programu programu LCS_02.cpp używająca biblioteki `omp`.

LCS_threads3_1.cpp i LCS_0threads3_2.cpp zwielowatkowana wersja programu programu LCS_03.cpp używająca `std::barrier` do synchronizacji przed przejściem do kolejnej przekątnej, Obydwa programy różnią się sposobem chodzenia przez wątki po przekątnej, wersja 3.1 skacze o liczbę wątków, a w wersji 3.2, każdy wątek oblicza spójny kawałek przekątnej. (wersja 3.2 z bardzo dziwnych powodów kompletnie nie działa dla więcej niż 4 wątków, czyli ilości rdzeni na mojej maszynie).

LCS_omp3.cpp zwielowatkowana wersja programu programu LCS_03.cpp używająca biblioteki `omp`.

LCS_omp4.cpp zwielowatkowana wersja programu programu LCS_04.cpp używająca biblioteki `omp`.

Hirschberg_threads.cpp Zwielowatkowany algorytm Hirschberg-a z użyciem techniki dziel i zwyciężaj dla początkowych głębokości.

Hirschberg_omp.cpp Zwielowatkowany algorytm Hirschberg-a z użyciem techniki dziel i zwyciężaj z użyciem kolejki zadań z biblioteki `omp`.

3 Problem

Znalezienie najdłuższego wspólnego podciagu najczęściej jest rozwiązywane przez użycie algorytmów dynamicznych. Problem dla ciągów $X = (x_1 x_2 \dots x_n)$, $Y = (y_1 y_2 \dots y_m)$, oraz X_0, X_1, \dots, X_n i Y_0, Y_1, \dots, Y_n jako odpowiednio ich prefixy, możemy zapisać rekurencyjnie w postaci:

$$LCS(X_i, Y_j) = \begin{cases} \epsilon, & \text{if } i = 0 \text{ or } j = 0 \\ LCS(X_{i-1}, Y_{j-1})x_i, & \text{if } x_i = y_j \\ \max(LCS(X_i, Y_{j-1}), LCS(X_{i-1}, Y_j)), & \text{else} \end{cases}$$

Analizując te funkcje widzimy, że pola w tablicy wynikowej są zależne od wartości w polach:

- w tym samym wierszu, ale kolumnie wcześniej.
- w tej samej kolumnie, ale wiersz wyżej.
- kolumnie wcześniej i wiersz wyżej.

4 Optymalizacja algorytmów sekwencyjnych

W przypadku optymalizacji dla programów sekwencyjnych możemy zauważyć, jak wielką rolę odgrywa dostęp do pamięci. W szczególności jak porównamy programy LCS_03.cpp i LCS_02.cpp. W obydwu programach przechodzimy po kolejnych przekatnych, ale tablica wynikowa w programach jest odpowiednio ułożona wzdłuż przekatnych i wzdłuż wierszy. W rezultacie możemy zauważyć, że program LCS_03.cpp działa najszybciej wśród wszystkich programów sekwencyjnych i aż trzy razy szybciej od LCS_02.cpp na sekwencjach długości 10 tysięcy znaków. Jest też znacząco szybszy od algorytmu chodzącego po wierszach (LCS_01.cpp), który też *cache-uje* pamięć dosyć optymalnie, ale wyniki w tablicy są zależne od wartości obliczanych bezpośrednio przed nimi, co ogranicza możliwość stosowania operacji takich jak chociażby SIMD.

Nie jest też niczym dziwnym, że algorytm Hirschberg-a radzi sobie gorzej, to znaczy jest 2 razy wolniejszy od podstawowego LCS_01.cpp.

5 Zwielowatkowanie klasycznego algorytmu dynamicznego

Metoda narzucająca się podczas analizy funkcji z sekcji 3. jest iterowanie się po tablicy wynikowej wzdłuż jej przekatnych ponieważ komórki macierzy należące do tej samej przekatnej są od siebie niezależne. Jeśli chcielibyśmy się iterować wzdłuż kolejnych wierszy, to natrafilibyśmy na problem ponieważ każda komórka zależy od wartości w tym samym wierszu.

Program LCS_threads1.cpp ma za zadanie pokazać jak drogie potrafi być tworzenie wątków będąc 8-10 razy wolniejszym od implementacji sekwencyjnej na średniej wielkości testach.

W przypadku wielowatkowych implementacji, zarządzanie pamięcią staje się jeszcze ważniejsze, o czym świadczą większe różnice w czasach wykonania, to znaczy program LCS_threads2.cpp jest 4 razy wolniejszy od programu LCS_threads3_1.cpp na średnich testach.

Możemy też zauważyć, że *parallel for* z biblioteki openMP radzi sobie wyraźnie gorzej z dużą ilością stosunkowo małych i prostych zadań, będąc parę razy gorszym w porównaniu z implementacjami na zwykłych threadach z barierą.

Dość sporym minusem korzystania z przekatnych jest fakt, że przekątne są różnej wielkości, co jest problemem zwłaszcza na początku i na końcu wywoływania programu gdzie wątki wykonują bardzo mało pracy, a i tak muszą się synchronizować zabierając czas. Dlatego spróbowałem w LCS_omp4.cpp i LCS_04.cpp implementacji triku który obchodzi ten problem używając dodatkowej tablicy pamiętającej ostatnie wystąpienia kolejnych elementów z alfabetu, co pozwala na odwoływanie się tylko do wiersza wyżej. Wprowadza to jednak dodatkową złożoność w obliczeniach radząc sobie istotnie gorzej z obliczaniem problemu.

6 Zwielowatkowanie algorytmu Hirschberg-a

Algorithm 1 bardzo ogólny schemat algorytmu Hirschberga

if $m \leq 2$ **or** $n \leq 2$ **then**
| **zwróć** rozwiązanie obliczone algorytmem pamięciochłonnym
end

oblicz ostatni wiersz tablicy wynikowej dla pierwszej połowy pierwszej sekwencji i drugiej sekwencji.(1)

oblicz ostatni wiersz tablicy wynikowej dla odwróconej drugiej połowy pierwszej sekwencji i odwróconej drugiej sekwencji.(2)

Znajdź taki punkt podziału k.(3)

zwróć sumę dla wywołań na mniejszych podzadaniach według podziału k.(4)

Analizując algorytm Hirschberg-a możemy, go podzielić na części które mogą być wykonane asynchronicznie i części wykonywane synchronicznie.

Tak więc części (1) i (2) są od siebie niezależne i mogą być wywołane równolegle, ale muszą się zsynchronizować przed wywołaniem (3), po czym (4) część również może być wykonana równolegle.

Zdecydowałem się na spróbowanie dwóch rozwiązań zrównoleglania dla tego algorytmu:

threads obliczamy operacje (1) i (2), oraz (4) równolegle tylko dla pierwszych paru głębokości rekursji.

OpenMP dodając jako podzadania do puli wątków operacje które mogą być wykonane równolegle.

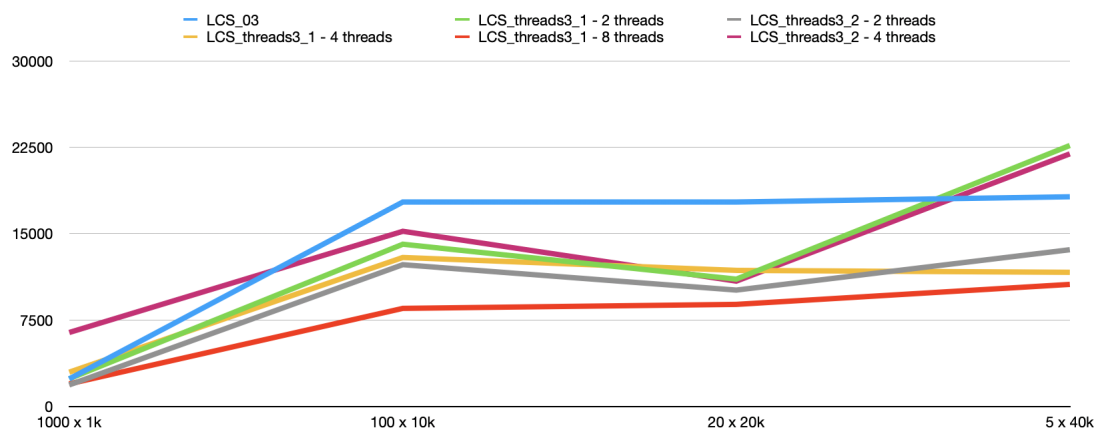
Co było dla mnie dość sporym zaskoczeniem, wersja threads okazała się bardzo dobrze zrównoleglająca się, osiągając wyniki konkurujące z najlepszymi programami urzywającymi pamięci kwadratowej na średnich testach i miażdżąc je na bardzo dużych testach. Jest to po części zasługą tego, że wątki rzadko odwołują się do tych samych miejsc w pamięci, oraz z faktu, że środowisko wykonania posiada bardzo mało pamięci RAM(16 GB).

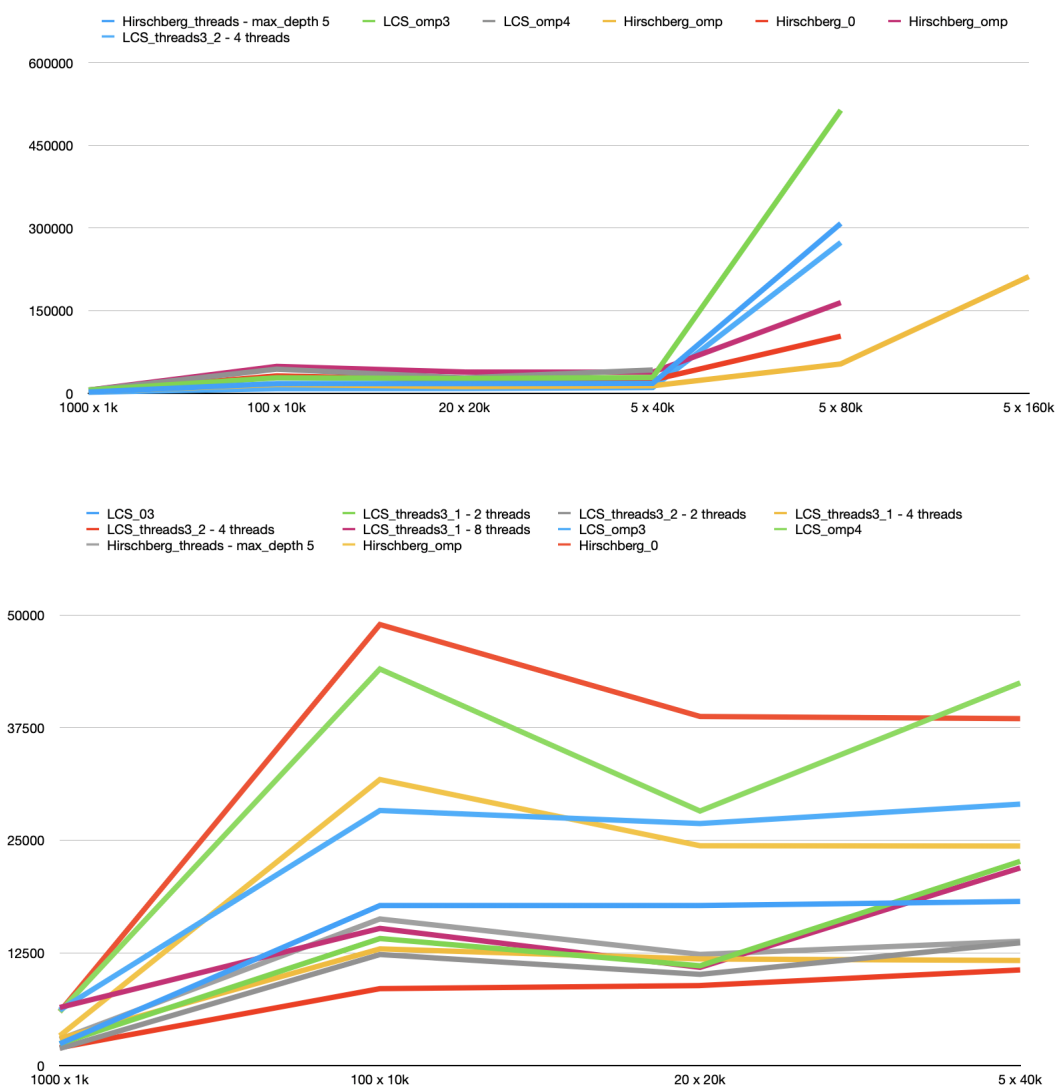
Z drugiej jednak strony wersja używająca openMP, która na pierwszy rzut oka mogła lepiej zarządzać obciążeniem wszystkich wątków okazała się znacząco gorsza, chociaż i tak poprawiła czas o 36% względem wersji sekwencyjnej.

7 Porównanie wyników

Wyniki testów prezentują się następująco. Możemy zauważyć, że bardzo duża rola w rozwiązywaniu problemu odgrywa dostępna pamięć, ponieważ programy oparte na algorytmie Hirschberga osiągają znacząco lepsze rezultaty.

	1000 x 1k	100 x 10k	20 x 20k	5 x 40k	5 x 80k	5 x 160k
LCS_03	2388	17744	17744	18201	308269	
LCS_threads3_1 - 2 threads	2402	14081	11049	22652	330317	
LCS_threads3_2 - 2 threads	1868	12304	10099	13604	320884	
LCS_threads3_1 - 4 threads	2973	12928	11810	11642	302726	
LCS_threads3_2 - 4 threads	1958	8511	8858	10595	273794	
LCS_threads3_1 - 8 threads	6419	15205	10875	21925	381821	
LCS_omp3	6130	28302	26852	29010	513836	
LCS_omp4	5950	44020	28241	42477	b.d	
Hirschberg_threads - max_depth 5	2883	16249	12312	13775	53459	211923
Hirschberg_omp	3274	31742	24384	24353	103771	
Hirschberg_0	5959	48958	38746	38504	164715	





8 Dalsza Możliwość poprawy

Wszystkie programy, używają instrukcji warunkowych *if* co powoduje, że kompilator nie stosuje instrukcji SIMD, których zastosowanie mogłoby skutkować dalszą redukcją czasu wykonania.