

Міністерство освіти і науки, молоді та спорту України

Національний Технічний Університет України

«Київський Політехнічний Інститут ім. Ігоря Сікорського»

Навчально-науковий комплекс

Кафедра системного проектування

Курсова робота

з дисципліни

“Паралельні обчислення”

«Створення інвертованого індексу»

Виконав:

студент ДА-71

Кузнєцов Олексій

Київ – 2020

Завдання

1. Вивчити теорію з глави 4 «Побудова індексу».
2. Розробити модель системи для побудови індексу в паралельному середовищі, заклавши можливість варіювати кількість потоків (від 1 до N).
3. Завантажити вхідні текстові дані згідно до варіанту.
4. Створити проект Github (або у будь-якому іншому репозиторії).
5. Створити програмну реалізацію моделі з п.2, правильно використовуючи систему контролю версій (бажано, щоб один commit був однією атомарною зміною в проекті); описувати кожен commit в одному стилі.
6. Виконати тестування розробленої системи.
7. Зробити графічне порівняння часу роботи системи, задаючи різну кількість потоків (від 1 до N).
8. Перевірити правильність результату роботи послідовного та паралельних рішень.
9. Зробити форматування коду в одному стилі.
10. Зробити файл з інструкцією як зібрати та запустити проект, додати його в репозиторій.
11. Зробити протокол з результатами роботи, що має містити інформацію з попередніх пунктів.

Теоретична основа та алгоритм роботи

Інвертований індекс (англ. Inverted index) - структура даних, в якій для кожного слова колекції документів у відповідному списку перераховані всі документи в колекції, в яких воно зустрілося. Інвертований індекс використовується для пошуку за текстами.

Приклад :

Нехай у нас буде колекція з 3 текстів $T_0 = \text{"it is what it is"}$, $T_1 = \text{"what is it"}$ і $T_2 = \text{"it is a banana"}$, тоді інвертований індекс буде виглядати наступним чином :

"a": {2}

"banana": {2}

"is": {0, 1, 2}

```
"what": {0, 1}
```

```

graph LR
    A[Файли з директорії] --> B[Потік 1]
    A --> C[Потік n]
    B -- "Індекси" --> D[Об'єднання індексів]
    C -- "Індекси" --> D
    D --> E[Вивід інвертованих індексів]
  
```

Реалізація та тестування програми

Алгоритм побудови індексу реалізований у цій роботі мовою програмування JAVA. Були використані вбудовані бібліотеки цієї мови для роботи з паралелізмом та структурами даних. Приклад роботи програми продемонстрований на малюнку нижче.

```
C:\Program Files\JetBrains\IntelliJ IDEA Community Edition 2019.3.3\bin\java.exe "C:\Program Files\JetBrains\IntelliJ IDEA Community Edition 2019.3.3\lib\idea_rt.jar=S4380;C:\Program Files\JetBrains\IntelliJ IDEA Community Edition 2019.3.3\bin\java.exe" -javaagent:C:\Program Files\JetBrains\IntelliJ IDEA Community Edition 2019.3.3\lib\idea_rt.jar=S4380;C:\Program Files\JetBrains\IntelliJ IDEA Community Edition 2019.3.3\bin\java.exe
```

```
output final result  
word = "", occurrences = [test\neg\2000_1.txt, test\neg\2001_1.txt, test\neg\2002_1.txt, test\neg\2003_1.txt, test\neg\2004_1.txt, test\neg\2005_1.txt, test\neg\2006_1.txt, test\neg\2007_3.txt, test\neg\2008_1.txt, test\neg\2009_1.txt]  
word = "a", occurrences = [test\neg\2000_1.txt, test\neg\2002_1.txt, test\neg\2003_1.txt, test\neg\2004_1.txt, test\neg\2005_1.txt, test\neg\2006_1.txt, test\neg\2007_3.txt, test\neg\2008_1.txt, test\neg\2009_1.txt]  
word = "abandoned", occurrences = [test\neg\2011_1.txt]  
word = "aberrations", occurrences = [test\neg\2196_1.txt]  
word = "abilities", occurrences = [test\neg\2211_2.txt]  
word = "ability", occurrences = [test\neg\2088_1.txt, test\neg\2092_2.txt, test\neg\2134_4.txt, test\neg\2211_2.txt, test\neg\2245_1.txt, test\neg\2250_1.txt]  
word = "able", occurrences = [test\neg\2014_1.txt, test\neg\2021_1.txt, test\neg\2093_1.txt, test\neg\2100_1.txt, test\neg\2125_2.txt, test\neg\2135_1.txt, test\neg\2136_3.txt, test\neg\2170_1.txt, test\neg\2178_1.txt]  
word = "ably", occurrences = [test\neg\2228_4.txt]  
word = "abner", occurrences = [test\neg\2224_3.txt]  
word = "aboard", occurrences = [test\neg\2016_2.txt, test\neg\2091_2.txt, test\neg\2094_3.txt]  
word = "abolished", occurrences = [test\neg\2013_1.txt]  
word = "abound", occurrences = [test\neg\2136_3.txt]  
word = "about", occurrences = [test\neg\2002_1.txt, test\neg\2004_1.txt, test\neg\2008_1.txt, test\neg\2009_1.txt, test\neg\2010_1.txt, test\neg\2013_1.txt, test\neg\2014_1.txt, test\neg\2016_2.txt, test\neg\2017_2.txt]  
word = "abouts", occurrences = [test\neg\2009_1.txt]  
word = "above", occurrences = [test\neg\2011_1.txt, test\neg\2024_1.txt, test\neg\2036_1.txt, test\neg\2042_3.txt, test\neg\2078_4.txt, test\neg\2103_1.txt, test\neg\2125_2.txt, test\neg\2167_2.txt]  
word = "abraham", occurrences = [test\neg\2014_1.txt]  
word = "abrupt", occurrences = [test\neg\2011_1.txt, test\neg\2013_1.txt]  
word = "absent", occurrences = [test\neg\2079_1.txt, test\neg\2129_4.txt, test\neg\2224_3.txt, test\neg\2246_1.txt]  
word = "absolutely", occurrences = [test\neg\2005_1.txt, test\neg\2011_1.txt, test\neg\2014_1.txt, test\neg\2016_2.txt, test\neg\2025_2.txt, test\neg\2028_1.txt, test\neg\2051_1.txt, test\neg\2085_3.txt]  
word = "absorbed", occurrences = [test\neg\2079_1.txt]  
word = "abstract", occurrences = [test\neg\2049_1.txt]  
word = "absurd", occurrences = [test\neg\2077_4.txt, test\neg\2087_3.txt, test\neg\2163_1.txt]  
word = "absurdities", occurrences = [test\neg\2190_3.txt]
```

Скріншот роботи програми

В процесі розробки програми була використана система контролю версій Git за допомогою веб-сервісу Github.

Сторінка проекту:

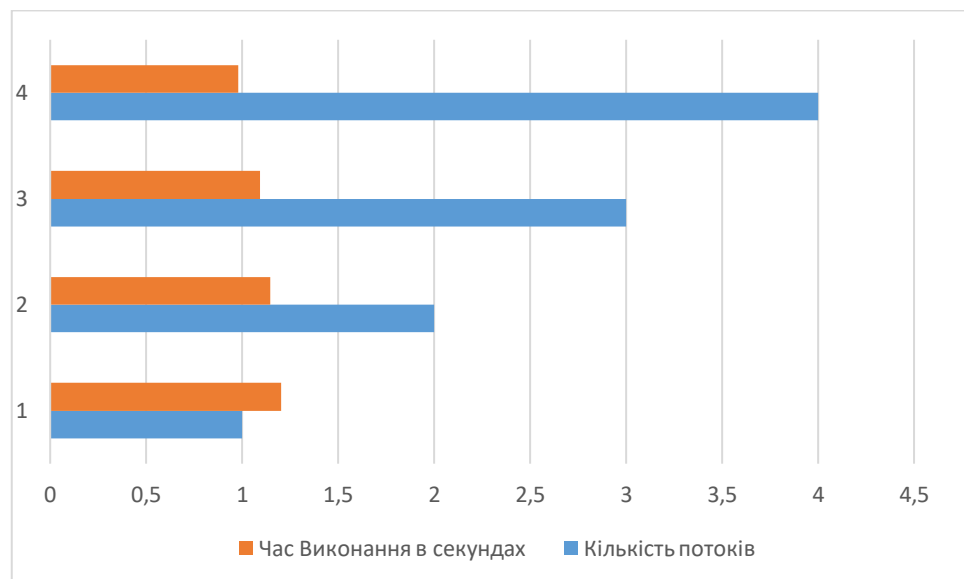
https://github.com/Kinelan/course_work. На ній знаходиться інструкція по компіляції та використанню проекту.

Перевіримо залежність роботи програми від кількості потоків, які ми виділимо на роботу програми. В якості тестових прикладів використаємо базу новин, в якій присутні наступні папки:

1. aclImdb\test\neg – N=250 файлів
2. aclImdb\test\pos – N=250 файлів
3. aclImdb\train\neg – N=250 файлів
4. aclImdb\train\pos – N=250 файлів
5. aclImdb\train\unsup – N=1000 файлів

Результати тестування ось такі :

Кількість потоків	Час виконання в нс	Час Виконання в секундах
1	1203752600	1,2037526
2	1145782000	1,145782
3	1092709500	1,0927095
4	981291400	0,9812914



Швидкодія програми

Як можемо побачити, завдяки використанню паралелізму, вдалося скоротити час виконання програми. Цей ефект буде більше помітний при використанні більшого об'єму інформації, тому що

алгоритм має найвищу ефективність при роботі з великими та надвеликими масивами даних.

Висновки

У ході виконання курсової роботи було досліджено алгоритм побудови інвертованого індексу. На етапі теоретичного ознайомлення була розроблена модель програми.

На етапі проектування я реалізував алгоритм та створив консольний додаток. Результатом розробки програми є закріплення навичок проектування паралельних алгоритмів та використання засобів мови програмування JAVA у паралельних додатках. Також була показана ефективність використання паралельних алгоритмів для збільшення швидкодії програми.