

Міністерство освіти і науки України
Національний технічний університет України „КПІ”
Факультет інформатики та обчислювальної техніки

Кафедра автоматизованих систем обробки
інформації та управління

ЗВІТ

з лабораторної роботи № 3
дисципліни
“ТЕХНОЛОГІЇ ПАРАЛЕЛЬНОГО ПРОГРАМУВАННЯ В УМОВАХ
ВЕЛИКИХ ДАНИХ”
на тему:

„Реалізація на основі Apache Hadoop”

Виконали:
студенти групи ІТ-01мн
Панасюк Станіслав
Лесогорський Кирило

Перевірив:
доц. Жереб К. А.

Київ – 2021

Зміст

1. Постановка задачі	3
2. Обрані інструменти.....	3
3. Опис роботи програмного забезпечення.....	3
4. Отримані результати	4
5. Висновки	4

1. Постановка задачі

Обрану задачу необхідно реалізувати, використовуючи Apache Hadoop. У якості задачі було обрано побудову системи пошуку схожих зображень. У ядрі системи лежатиме використання D-hash для знаходження хешу зображення. D-hash дозволяє точно та швидко шукати схожі зображення. Він стійкий до скейлінгу зображення, але погано справляються з обрізаними та повернутими під кутом зображеннями. Тому цю техніку аугментовано за допомогою наступного прийому: при завантаженні зображення воно буде аугментовано за допомогою декількох фільтрів, при цьому для кожного фільтру буде згенеровано хеш і збережено у базу даних. При пошуку зображення буде використовуватись оператор XOR для знаходження зображень зі схожими хешами.

2. Обрані інструменти

Для виконання третьої лабораторної роботи буде використано стандартні інструменти Java із доданою до них бібліотекою Apache Hadoop. У нашому випадку третя та четверта лабораторні роботи взаємопов'язані, адже четверта робота, використовуючи Apache Spark, проводить підготовку даних, а вже третя робота працює із ними за рахунок написаної MapReduce програми.

3. Опис роботи програмного забезпечення

Код для підрахунків D-hash використовується ідентичний до попередніх робіт, тому у даній роботі розглянемо лише нові зміни.

Було написано цілком новий застосунок, який Apache Hadoop добре розуміє. Таким застосунком є MapReduce job, що читає попередньо підготовані через Apache Spark дані з Sequence File, обраховує хеші кожного зображення та зберігає їх у файл.

Застосунок складається з п'яти основних частин:

1. Мапери (це Map частина нашої MapReduce job). Вони реалізують інтерфейс Mapper та описують логіку, за якою будуть оброблятися окремі часточки даних (у нашому випадку це окремі зображення). Тут ми кожне зображення перетворюємо на його хеші;
2. Ред'юсери (це Reduce частина нашої MapReduce job). Вони реалізують інтерфейс Reducer та описують логіку, за якою усі підготовані часточки даних будуть агрегуватись згідно до ключа, який надає мапер із кожним обробленим значенням. У нашому випадку ключ – ім'я файла. Тут ми зберігаємо хеші згідно з ключами до файлу;

3. Утиліти – тут ми скопіювали ту саму логіку підрахунків D-hash;
4. Модель – у ній ми створили кастомний тип даних, який передаємо від мапера до ред'юсера, що дозволяє передати масив LongWritable даних (наших хешів);
5. Драйвер – основний клас, який відповідає за конфігурацію MapReduce job.

4. Отримані результати

В результаті виконання роботи було отримано наступні результати з процесінгу:

```

1  images-1-test-jpg [16318923691360512, -2192927142707191776, -1017266019380444992, -1017319835064860672, -146075649512046080, -4481...
2  images-10-test-jpg [-8407285498379239940, 1295199966510739488, -3341776924914615812, -3622801300206321536, 72110915634648248, 119393...
3  images-11-test-jpg [3565709413589274960, -1010546079462047712, -65970806595080, -149197261786480512, -3038287551325502004, -10218760...
4  images-12-test-jpg [8133995509015765440, -436408047080808416, -1027482974336979576, -3318640615735836544, -1017259010510699580, -101...
5  images-13-test-jpg [-4062284278396265360, -2167927288720498656, -578794225215342104, -1035284918998826880, -4101176226267260816, -91...
6  images-14-test-jpg [5788374979791298852, 3490404509388136480, -7007888186518239104, -435167452735733632, 7921706296733554630, -47688...
7  images-15-test-jpg [-434662205589942776, 2416497918101450784, -1600526357730960400, -2179195616857358208, -4738215194517290624, 8094...
8  images-16-test-jpg [-35009053671262112, -1328596826932117504, -55212191727140, -1123525177270845440, -1763159285185316864, -14816087...
9  images-17-test-jpg [144680518553599870, 3490417761182044224, 1098666377134180350, -9134007612513066880, 289979246216937214, 38436235...
10 images-18-test-jpg [667171596415486, -1107870708334051232, -5787710547593803794, -5193391883301044096, 130320722476597240, -80344080...
11 images-19-test-jpg [-515167866835070704, -1041020383983026176, -2184375639226202240, -882263999463178112, 6762605095775474560, -4503...
12 images-2-test-jpg [-751481895034388168, -3316529078564012000, -3173942676027227264, -2025667889297555392, -8243502729620537270, -921...
13 images-20-test-jpg [-217931116812648250, 8204887308287111168, -440858234270646868, -440928893675650944, 8355271534166770178, 3546637...
14 images-21-test-jpg [-287672897986617284, 1157725760412733472, -750412571412662796, -432915137270103872, 1148417898478436096, 1261440...
15 images-22-test-jpg [-38720126329981600, 2369035935787532288, -1018418549950931572, -465566338440380288, -1925913500882418688, -22422...
16 images-23-test-jpg [-583703278350109354, -1012850957016866720, -6516807695208693364, -4196922641285611392, 6143270536974065728, 8097...
17 images-24-test-jpg [-574758489277325184, 4945438901297152, -3760569710479021592, -1617468220503949312, -293013239665623168, 12298188...
18 images-25-test-jpg [-9223225810433111284, -2177489262203609056, -3493133591528895040, -1745834180133928832, 4339395266663415800, -50...
19 images-26-test-jpg [-2251797047723670524, -1137016029580333024, -137439477762, -144684795465449344, 9185003476860531648, 47160781145...
20 images-27-test-jpg [-5475061383230692790, -8633251239311161296, -5196714543801722890, -2479820272998203200, -5267841992717795510, -5...
21 images-28-test-jpg [-2190701822210, 4629701592428339232, -8536140394624611960, -3631817609906651008, 145995333632, 11279878590460316...
22 images-29-test-jpg [-1524468366491151592, 1283940658268094496, -3635091098431815296, -3635081028336238464, -7144538580789129056, -21...
23 images-3-test-jpg [-141630302194280, 1234515305276919808, -4496358579039503888, -2771550222587723776, -7403513271069442048, -61987103...
24 images-30-test-jpg [7457929925140666262, -440859051815833568, -1028076590429265476, -1161470342779010944, 6363850929097147492, -1686...
25 images-31-test-jpg [-9331000711437300, -4537348064907599872, -9092246487773361784, -3030375053083443072, 4570816750208290816, -21974...
26 images-32-test-jpg [108188107927995324, -1026392465581510624, -9126050969703508, -2477549677329661824, 645001227605113342, -68026499...

```

Рисунок 4.1 - Вихідний файл із ключами у вигляді імені зображення та значеннями у вигляді масиву хешів

5. Висновки

У ході лабораторної роботи було створено MapReduce додаток із використанням Apache Hadoop для досягнення ідентичної до першої та другої лабораторних робіт цілі.