

ĐẠI HỌC BÁCH KHOA HÀ NỘI
TRƯỜNG CÔNG NGHỆ THÔNG TIN & TRUYỀN THÔNG
CÔNG NGHỆ THÔNG TIN VIỆT PHÁP



PROJECT 1

Truy xuất dữ liệu lớn với bộ nhớ nhỏ

GVHD: Nguyễn Đức Anh
SV thực hiện: Nguyễn Hữu Mạnh – 20205213

Hà Nội, Tháng 1/2024

Contents

1	Problem Statement	2
2	Design	2
3	Implementation	2
3.1	Generate Arrays	2
3.2	Access Arrays	4
3.3	Code	5
4	Testing	5
4.1	Compile Code	5
4.2	Generate Arrays	5
4.3	Access Arrays	5
4.4	Test	6
5	Conclude	7

1 Problem Statement

Tạo một chương trình truy xuất dữ liệu lớn trên bộ nhớ kích thước nhỏ:

1. Tạo dữ liệu nhân tạo bằng cách: Tạo 3,000,000 (3 triệu) mảng interger, kích thước mỗi mảng là một giá trị nguyên ngẫu nhiên từ 100 tới 10000. Lưu trữ tùy ý (Gợi ý, lưu trữ qua file).
2. Tạo một hàm truy xuất ngẫu nhiên tới mảng thứ i trong 3 triệu mảng trên (Gợi ý, cần lưu thêm giá trị hỗ trợ để truy xuất tốt hơn).

Giới hạn bộ nhớ RAM: Dưới 500 MB.

2 Design

Lưu trữ

- Lưu dữ liệu vào file ở dạng nhị phân để bảo toàn dữ liệu.
- Lưu vị trí của mỗi mảng vào 1 file riêng để truy xuất.

Truy xuất

- Sử dụng file lưu vị trí để truy cập đến mảng nhanh.
- In ra mảng đó.

Để có thể giới hạn bộ nhớ, ta sử dụng `java -Xmx500m`

3 Implementation

3.1 Generate Arrays

1. Tạo biến lưu tổng số array và lưu giới hạn elements trong array

```
1 int totalArr = 3_000_000;  
2 int[] rangeArr = new int[] {100, 10_000};
```

2. Thực hiện ghi vào 2 file **data.dat** và **pos.dat**

File **data.dat** lưu dữ liệu ở dạng nhị phân

File **pos.dat** vị trí của từng array ở dạng nhị phân

```
1 try (BufferedOutputStream binout =  
2     // save data  
3     new BufferedOutputStream(new FileOutputStream("../bin/data.dat"));  
4     BufferedOutputStream posout =  
5     // save position  
6     new BufferedOutputStream(new FileOutputStream("../bin/pos.dat"))) {
```

Khai báo 1 biến pos để lưu vị trí của từng array

Ta cần biến pos là long vì mỗi array có 100 đến 10000 elements nên 3 triệu mảng sẽ có tổng số elements > max value int có thể lưu là 2^{31}

```
1 long pos = 0; // long because total element in 3M > int =  $2^{31}$ .
```

Lặp qua từng array

```
1 for (int i = 1; i <= totalArr; i++) {
```

Khai báo biến size là số elements của array i. Biến size random trong khoảng rangeArr

```
1 int size = random.nextInt(rangeArr[1] - rangeArr[0] + 1) + rangeArr[0];
```

Chuyển value pos sang byte và ghi vào file **pos.dat**

```
1 byte[] posBytes = toBytes(pos, 8);  
2 posout.write(posBytes);  
3 pos += size;
```

Lặp qua từng element trong array, lấy giá trị của mỗi element là random 1 số int
Chuyển value element sang byte và ghi vào file **data.dat**

```
1     for (int j = 0; j < size; j++) {
2         int element = random.nextInt();
3         byte[] elementBytes = toBytes(element, 4);
4         binout.write(elementBytes);
5     }
```

Sau khi ghi xong ta sẽ lưu thêm 1 value pos cuối cùng nữa để tính được số element của array cuối cùng

```
1     if (i == totalArr) {
2         byte[] lastPosBytes = toBytes(pos, 8);
3         posout.write(lastPosBytes);
4     }
```

3. Hàm chuyển int hoặc long sang byte

```
1     private static byte[] toBytes(long value, int a) {
2         byte[] result = new byte[a];
3         for (int i = 0; i <= (a - 1); i++) {
4             result[i] = (byte) (value >> (8 * (a - 1 - i)));
5         }
6         return result;
7     }
```

Code hoàn chỉnh

```
1 import java.io.BufferedOutputStream;
2 import java.io.FileOutputStream;
3 import java.io.IOException;
4 import java.util.Random;
5
6 public class GenerateArrays {
7
8     public static void main(String[] args) {
9
10        int totalArr = 3_000_000;
11        int[] rangeArr = new int[] {100, 10_000};
12
13        try (BufferedOutputStream binout =
14            // save data
15            new BufferedOutputStream(new FileOutputStream("../bin/data.dat"));
16            BufferedOutputStream posout =
17            // save position
18            new BufferedOutputStream(new FileOutputStream("../bin/pos.dat"))) {
19            long pos = 0; // long because total element in 3M > int = 2^31.
20
21            for (int i = 1; i <= totalArr; i++) {
22                int size = random.nextInt(rangeArr[1] - rangeArr[0] + 1) + rangeArr[0];
23
24                byte[] posBytes = toBytes(pos, 8);
25                posout.write(posBytes);
26                pos += size;
27
28                for (int j = 0; j < size; j++) {
29                    int element = random.nextInt();
30                    byte[] elementBytes = toBytes(element, 4);
31                    binout.write(elementBytes);
32                }
33
34                if (i == totalArr) {
35                    byte[] lastPosBytes = toBytes(pos, 8);
36                    posout.write(lastPosBytes);
37                }
38
39                System.out.println("Array: " + i);
40            }
41        } catch (IOException e) {
42            System.err.println("Error write file: " + e.getMessage());
43        }
```

```

43     }
44 }
45
46 private static final Random random = new Random();
47
48 // (int || long) -> byte (a = total byte)
49 private static byte[] toBytes(long value, int a) {
50     byte[] result = new byte[a];
51     for (int i = 0; i <= (a - 1); i++) {
52         result[i] = (byte) (value >> (8 * (a - 1 - i)));
53     }
54     return result;
55 }
56 }

```

3.2 Access Arrays

1. Nhập vào thứ tự của array muốn truy xuất

```

1 Scanner sc = new Scanner(System.in);
2 System.out.print("Input index: ");
3 int index = sc.nextInt();
4 sc.close();

```

2. Thực hiện truy xuất array

```

1     getArray(index);

```

3. Hàm truy xuất getArray

Đọc dữ liệu từ 2 file **data.dat** và **pos.dat**

Ta sử dụng **RandomAccessFile** để có thể truy nhập file bằng byte.

```

1     public static void getArray(int index) throws IOException {
2         try (RandomAccessFile pos = new RandomAccessFile("../bin/pos.dat", "r");
3             RandomAccessFile dat = new RandomAccessFile("../bin/data.dat", "r")) {

```

Thực hiện lấy vị trí của array cần truy xuất

Dùng **seek** để có thể truy nhập file vào đúng vị trí cần truy nhập

```

1         index--;
2         pos.seek(index * 8); // pos.dat lưu du lieu long

```

Đọc vị trí và kích thước của array

```

1         long posData = pos.readLong();
2         long size = pos.readLong() - posData;

```

Truy nhập vào vị trí của array cần truy xuất

```

1         dat.seek(posData * 4); // data.dat lưu du lieu int

```

Dùng vòng lặp **for** để đọc từng **element** của array cần truy xuất và in ra

```

1         for (long i = 1; i <= size; i++) {
2             int element = dat.readInt();
3             System.out.print(element + " ");
4         }

```

Code hoàn chỉnh

```

1 import java.io.IOException;
2 import java.io.RandomAccessFile;
3 import java.util.Scanner;
4
5 public class AccessArrays {
6
7     public static void main(String[] args) throws IOException {
8         Scanner sc = new Scanner(System.in);
9         System.out.print("Input index: ");

```

```

10     int index = sc.nextInt();
11     sc.close();
12
13     getArray(index);
14 }
15
16 public static void getArray(int index) throws IOException {
17     try (RandomAccessFile pos = new RandomAccessFile("../bin/pos.dat", "r");
18         RandomAccessFile dat = new RandomAccessFile("../bin/data.dat", "r")) {
19
20         index--;
21         pos.seek(index * 8); // pos.dat lưu du lieu long
22         long posData = pos.readLong();
23         long size = pos.readLong() - posData;
24
25         // System.out.print(posData + " " + size + "\n");
26         dat.seek(posData * 4); // data.dat lưu du lieu int
27         for (long i = 1; i <= size; i++) {
28             int element = dat.readInt();
29             System.out.print(element + " ");
30         }
31     }
32 }
33 }

```

3.3 Code

[Source code toàn bộ của project₁](#)

4 Testing

4.1 Compile Code

Thư mục **src/** chứa source code, thư mục **bin/** chứa các file **.class** và các file dữ liệu.

Giả sử ta đang ở thư mục cha của thư mục **src/** và thư mục **bin/**

```
javac -d ./bin src/*.java
```

Vào thư mục bin

```
cd ./bin
```

4.2 Generate Arrays

```
java GenerateArrays
```

Sau khi chạy ta được 2 file **data.dat** và **pos.dat** được lưu ở cùng thư mục bin

4.3 Access Arrays

Chạy AccessArrays với flag **-Xmx500m** để giới hạn bộ nhớ còn **500Mb**

```
java -Xmx500m AccessArrays
Input index: // nhập 1 số i bất kì
```

> In ra console **array i** cần truy xuất

Thời gian truy xuất với 3 triệu mảng, mỗi mảng từ 100 đến 10000 nghìn được thực hiện chỉ tính bằng mili giây, ta có thể kiểm tra bằng

```
time java -Xmx500m AccessArrays
```

lệnh này sẽ in ra thời gian thực hiện câu lệnh.

4.4 Test

Vì số lượng array + element quá lớn nên không thể kiểm tra được truy xuất có chính xác không nên ta sẽ tạo 1 file **test.java** với số array 100000, mỗi array từ 10-20 element để có thể kiểm tra việc truy xuất có đúng không

```
1 int totalArr = 100000;
2 int[] rangeArr = new int[] {10, 20};
```

Ghi dữ liệu dạng text vào file **data.txt**

```
1 try (BufferedWriter writer = new BufferedWriter(new FileWriter("../bin/data.txt"));
```

Trong vòng lặp for khi thực hiện random element và ghi vào **data.dat** ta cũng sẽ lưu dữ liệu vào **data.txt**, mỗi array sẽ được lưu vào 1 dòng

```
1     for (int j = 0; j < size; j++) {
2         int element = random.nextInt();
3         byte[] elementBytes = toBytes(element, 4);
4         binout.write(elementBytes);
5
6         writer.write(String.valueOf(element) + " ");
7     }
8     writer.newLine();
```

Code hoàn chỉnh

```
1 import java.io.BufferedOutputStream;
2 import java.io.BufferedWriter;
3 import java.io.FileOutputStream;
4 import java.io.FileWriter;
5 import java.io.IOException;
6 import java.util.Random;
7
8 public class test {
9
10     public static void main(String[] args) {
11
12         int totalArr = 100000;
13         int[] rangeArr = new int[] {10, 20};
14
15         try (BufferedWriter writer = new BufferedWriter(new FileWriter("../bin/data.txt"));
16             BufferedOutputStream binout =
17                 new BufferedOutputStream(new FileOutputStream("../bin/data.dat"));
18             BufferedOutputStream posout =
19                 new BufferedOutputStream(new FileOutputStream("../bin/pos.dat"))) {
20             long pos = 0; // long because total element in 3M > int = 2^31.
21
22             for (int i = 1; i <= totalArr; i++) {
23                 int size = random.nextInt(rangeArr[1] - rangeArr[0] + 1) + rangeArr[0];
24
25                 // writer.write(String.valueOf(i) + " ");
26
27                 byte[] posBytes = toBytes(pos, 8);
28                 posout.write(posBytes);
29                 pos += size;
30
31                 for (int j = 0; j < size; j++) {
32                     int element = random.nextInt();
33                     byte[] elementBytes = toBytes(element, 4);
34                     binout.write(elementBytes);
35
36                     writer.write(String.valueOf(element) + " ");
37                 }
38                 writer.newLine();
39
40                 if (i == totalArr) {
```

```

41     byte[] lastPosBytes = toBytes(pos, 8);
42     posout.write(lastPosBytes);
43 }
44 }
45 System.out.println("Done");
46 } catch (IOException e) {
47     System.err.println("Error write file: " + e.getMessage());
48 }
49 }
50
51 private static final Random random = new Random();
52
53 // (int || long ) -> byte (a = total byte)
54 private static byte[] toBytes(long value, int a) {
55     byte[] result = new byte[a];
56     for (int i = 0; i <= (a - 1); i++) {
57         result[i] = (byte) (value >> (8 * (a - 1 - i)));
58     }
59     return result;
60 }
61 }

```

Biên dịch tương tự giống như trên

Tạo dữ liệu nhỏ hơn với 100000 arrays, mỗi array từ 10->20 element

```
java test
```

Sau khi chạy ta được 3 file **data.dat**, **pos.dat** và **data.txt** lưu dữ liệu ở dạng văn bản có thể đọc, mỗi mảng được lưu 1 dòng cách nhau bởi " ". Được lưu ở cùng thư mục bin

Truy xuất thử 1 array bất kì, ví dụ với array thứ 100000

```
java -Xmx500m AccessArrays
Input index: 100000
```

> In ra console array 100000 cần truy xuất

Để kiểm tra truy xuất có chính xác không ta kiểm tra file **data.txt**

Ta sử dụng command ``sed'` để in ra màn hình array 100000

```
sed -n '100000p' data.txt
```

Lệnh này sẽ in ra dòng thứ 100000 trong file **data.txt**

> Đối chiếu kết quả của 2 lệnh thấy giống nhau => truy xuất chính xác.

5 Conclude

Qua bài trên em đã học được cách tổ chức lưu trữ 1 dữ liệu lớn để có thể tối ưu trong việc truy xuất.

Cách để có thể truy xuất vào 1 vị trí bất kì trong 1 file dữ liệu lớn trong thời gian ngắn gần như bằng 0 với bộ nhớ cực kì nhỏ.