# CS205 Project 4 Report

11812804 董正

# 1 Introduction

The documentation of class `matrix` is hosted on [my GitHub](#).



Alternatively, view this report online [here](#).

## 1.1 Project Description

This project is to design a class for matrices.

1. The class should contain the data of a matrix and related information.
2. The class support different data types.
3. Do not use memory hard copy if a matrix object is assigned to another.
4. Implement some frequently used operators.
5. Implement region of interest (ROI) to avoid memory hard copy.
6. Test the program on X86 and Arm platforms, and describe the differences.

## 1.2 Development Environment

- x86_64

    - Windows 10 Home China x86_64
    - Kernel version `10.0.19042`
    - `Intel i5-9300H (8) @ 2.400GHz`
    - `g++.exe (tdm64-1) 10.3.0`
    - C++ standard: `c++11`

- Arm64

    - `macOS 12.0.1 21A559 arm64`
    - Darwin Kernel Version `21.1.0`
    - Apple M1 Pro (10-cores)
    - `Apple clang version 13.0.0 (clang-1300.0.29.3)`

- C++ standard: `c++11`

# 2 Design and Implementation

Header files and macros used in this section:

```
1  #include <fstream>
2  #include <iostream>
3
4  #define STRASSEN_LOWER_BOUND 128
```

## 2.1 Fields Of Matrix Class

To support different data types, I used template class.

`matrix` class has three basic fields:

```
1  template <typename T>
2  class matrix {
3      int nrows;
4      int ncols;
5      T* data;
6
7      ...
8  }
```

`data` is a 1D dynamic array containing all the elements of the matrix.

To avoid memory copy, I added a field named `ref_count`:

```
1  int* ref_count;
```

See details in section 2.2 below.

To support ROI, I used another pointer:

```
1  const matrix<T>* parent_matrix;
```

## 2.2 Constructor: Avoid Mem Copy

As the project required, it is forbidden to use hard copy if a matrix object is assigned to another.

If use soft copy directly, there will be an error when freeing `data` array:
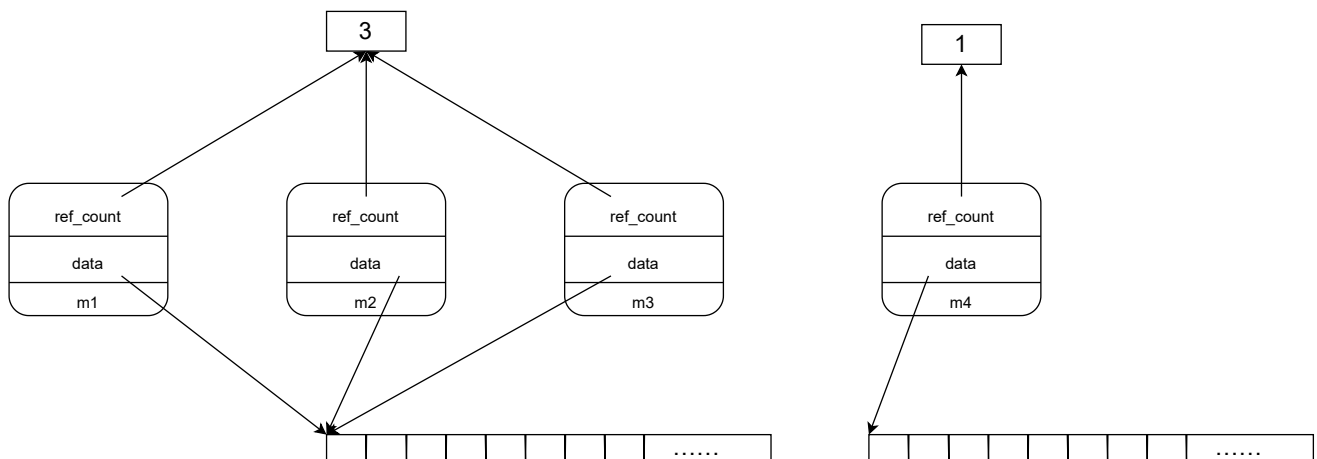
```
1  matrix<int> m2 = m1; // will set m2.data = m1.data
```

Therefore, when destructing `m1` and `m2`, since their `data` field is pointing to the same address, the array will be freed twice.

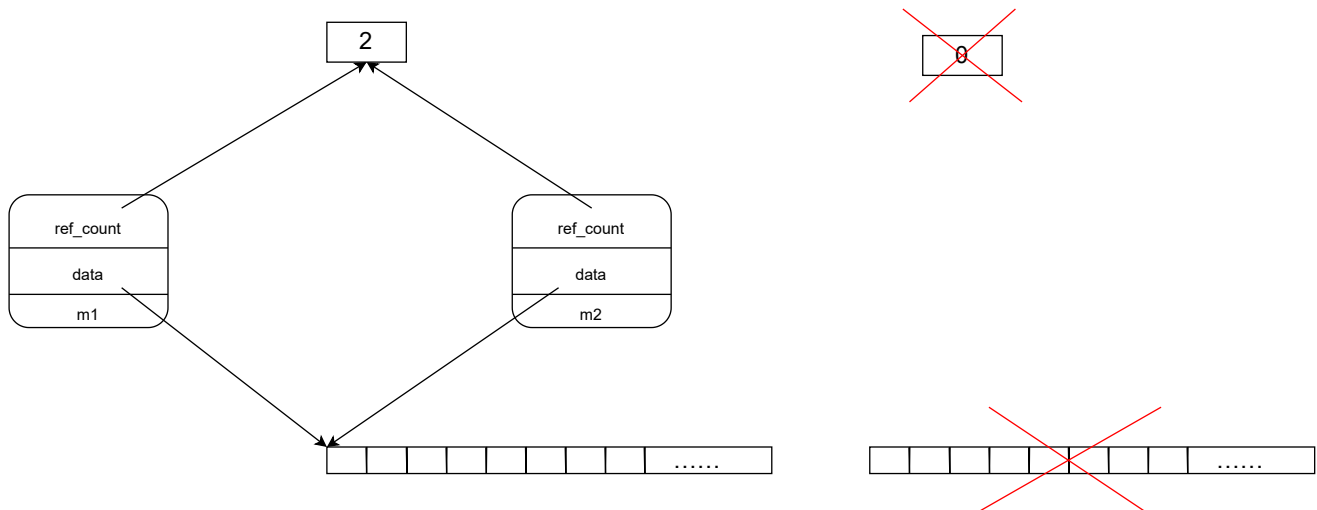A solution is to override destructor and do not free `data`. But it will cause memory leak.

To solve this, I referenced `cv::Mat` and added a field `int* ref_count`. It serves like a global integer for each allocated array.

The idea is described in the following picture:



When destructing a matrix, just decrease the integer that `ref_count` points to.

Example: delete `m3, m4`.



When `ref_count` decreases to 0, the array will be finally freed.

Therefore, the destructor is:

```
1  template <typename T>
2  inline matrix<T>::~matrix() { // version 1
3      *(this→ref_count) -= 1;
4      if (*(this→ref_count) == 0 && this→data ≠ nullptr) {
5          delete[] this→data;
6          delete this→ref_count;
7      }
8  }
```

And for copy constructor, just do soft copy, and increase `ref_count` by 1.

```
1  template <typename T>
2  inline matrix<T>::matrix(const matrix<T>& other) {
3      this→nrows = other.nrows;
4      this→ncols = other.ncols;
5      this→parent_matrix = other.parent_matrix;
6
7      this→data = other.data;
8      this→ref_count = other.ref_count;
9      *(this→ref_count) += 1;
10 }
```

For assignment operator, be sure to free data if it is the last reference of the array.

```
1  template <typename T>
2  inline matrix<T>& matrix<T>::operator=(const matrix<T>& other) {
3      this→nrows = other.nrows;
4      this→ncols = other.ncols;
5      this→parent_matrix = other.parent_matrix;
6
7      *(this→ref_count) -= 1;
8      if (*(this→ref_count) == 0 && this→data ≠ nullptr) {
9          delete this→ref_count;
10         delete[] this→data;
11     }
12
13     this→data = other.data;
14     this→ref_count = other.ref_count;
15     *(this→ref_count) += 1;
16
17     return *this;
18 }
```
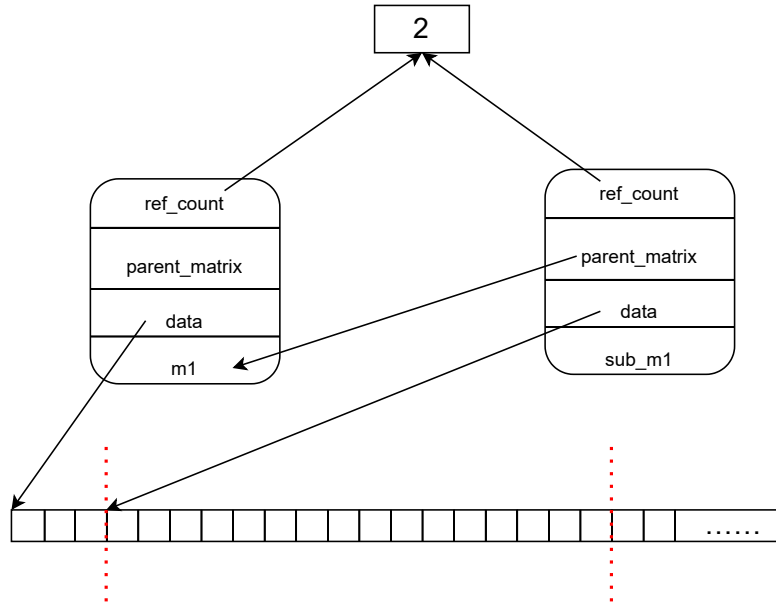
# 2.3 ROI For Submatrix

In my understanding, ROI is a concept of submatrix that shares data with its parent matrix.

To use ROI correctly, we need to know the shape of its parent matrix's size. Because when we get the `[i][j]` position element, we need to know the columns of its parent matrix in order to translate `i` and `j` to the correct position in the whole array.

i.e. `sub[i][j]=sub.data[i * parent.ncols + j]`.

So the design is very clear. Use a `parent_matrix` pointer to store its parent:
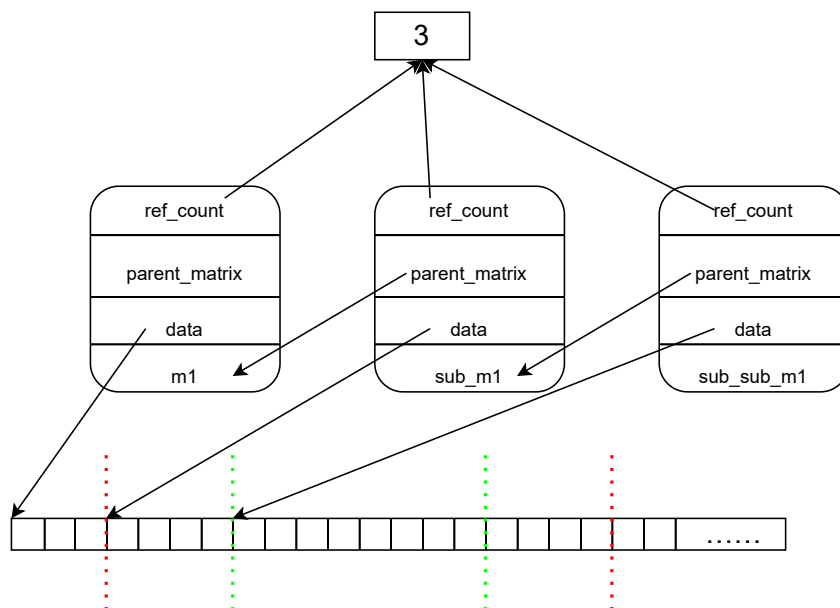


However, here comes two problems with this design.
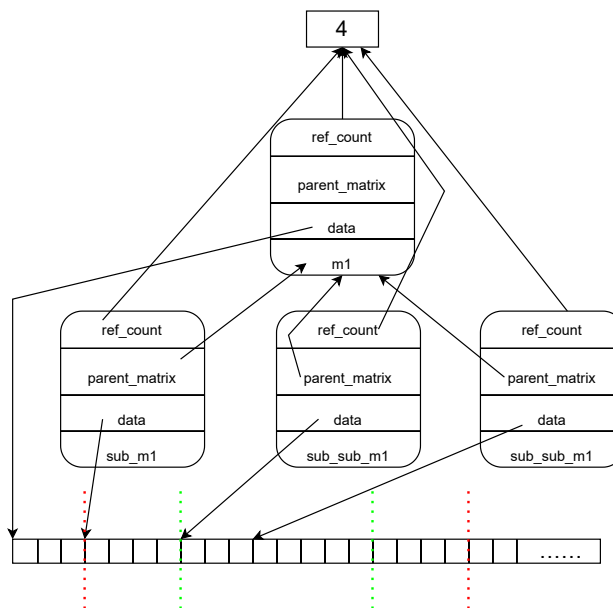
## 2.3.1 Multi-level Submatrix

The first is, multi-level submatrix.

One choice is to concatenate them like a linked-list:



It is unacceptable because `sub_sub_m1[i][j]=sub_sub_m1.data[i * sub_sub_m1->parent_matrix->ncols + j]=sub_sub_m1.data[i * sub_m1.ncols + j]` and it is wrong.

Therefore, we must let all the submatrices point to the root matrix.



Thus, the submatrix construction function is:
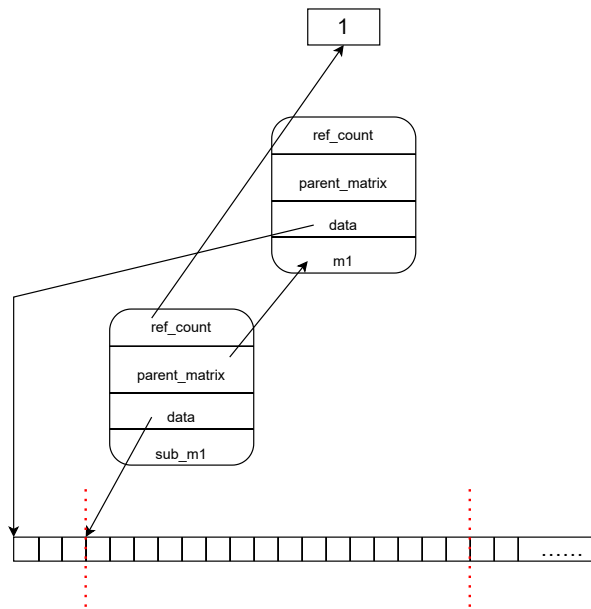
```
1   template <typename T>
2   matrix<T> matrix<T>::submatrix_ROI(int row_start, int row_end, int
    col_start, int col_end) {
3       matrix<T> res;
4       res.nrows = row_end - row_start;
5       res.ncols = col_end - col_start;
6
7       if (this→parent_matrix == nullptr) {
8           res.data = this→data + row_start * this→ncols + col_start;
9           res.parent_matrix = this;
10      } else {
11          res.data = this→data + row_start * this→parent_matrix-
    >ncols + col_start;
12          res.parent_matrix = this→parent_matrix;
13      }
14
15      res.ref_count = this→ref_count;
16      *(res.ref_count) += 1;
17
18      return res;
19  }
```

## 2.3.2 Submatrix Deletion

The second problem is the case that delete parent matrix before delete all submatrices.

When delete `sub_m1`, we should delete `m1.data`, not `sub_m1.data` because it is not the head address.

Therefore, `~matrix()` should be modified:

```cpp
template <typename T>
inline matrix<T>::~matrix() {
    *(this→ref_count) -= 1;
    if (*(this→ref_count) == 0 && this→data ≠ nullptr) {
        if (this→parent_matrix == nullptr) {
            delete[] this→data;
        } else {  // it is a submatrix
            delete[] this→parent_matrix→data;
        }
        delete this→ref_count;
    }
}
```

### 2.3.3 ROI Adjustment

Change the address of `data` pointer.

```cpp
template <typename T>
void matrix<T>::adjust_ROI(int row_start, int row_end, int col_start,
int col_end) {
    if (this→parent_matrix == nullptr) {
        cout << "ROI adjustment error: not a submatrix." << endl;
        exit(EXIT_FAILURE);
    }

    if (row_start < 0 || row_end > this→parent_matrix→nrows ||
col_start < 0 || col_end > this→parent_matrix→ncols) {
        cout << "ROI adjustment error: array index out of bound." <<
endl;
        exit(EXIT_FAILURE);
```

```
11       }
12
13       this→nrows = row_end - row_start;
14       this→ncols = col_end - col_start;
15       this→data = this→parent_matrix→data + row_start * this-
   >parent_matrix→ncols + col_start;
16     }
```

## 2.4 Matrix Operations

### 2.4.1 Index Operator

Since the data is stored in a 1d array, we need an API to access elements like a 2d array. Be careful of submatrices.

```
1  template <typename T>
2  inline T* matrix<T>::operator[](int i) {
3      if (i > this→nrows - 1) {
4          cout << "Index error: array index out of bound." << endl;
5          exit(EXIT_FAILURE);
6      }
7
8      if (this→parent_matrix ≠ nullptr) {
9          return this→data + i * this→parent_matrix→ncols;
10      }
11      return this→data + i * this→ncols;
12  }
13
14  template <typename T>
15  inline T& matrix<T>::operator()(int i, int j) {
16      if (i > this→nrows - 1 || j > this→ncols - 1) {
17          cout << "Index error: array index out of bound." << endl;
18          exit(EXIT_FAILURE);
19      }
20      if (this→parent_matrix ≠ nullptr) {
21          return this→data[i * this→parent_matrix→ncols + j];
22      }
23
24      return this→data[i * this→ncols + j];
25  }
```

### 2.4.2 Other Functions

```
1  bool operator==(matrix<T>& other);
2  bool operator≠(matrix<T>& other);
3  matrix<T> operator+(matrix<T>& other);
4  matrix<T> operator-(matrix<T>& other);
5  matrix<T> operator*(matrix<T>& other);
6  matrix<T> operator*(T coef);
```

```cpp
 7   template <typename U> friend matrix<U> operator*(int coef, matrix<U>&
     m);
 8   matrix<T> operator^(int expo); // fast power algo
 9   matrix<T>& operator*=(matrix<T>& other);
10   matrix<T>& operator*=(T coef);
11   matrix<T>& operator+=(matrix<T>& other);
12   matrix<T>& operator-=(matrix<T>& other);
13   matrix<T> multiply_elements(matrix<T>& other);

15   static matrix<T> merge_matrix(matrix<T>& C11, matrix<T>& C12,
     matrix<T>& C21, matrix<T>& C22);

17   int get_nrows();
18   int get_ncols();
19   int* shape();
20   void print_shape();
21   bool shape_equals(matrix<T>& other);

23   matrix();
24   matrix(int nrows, int ncols);
25   matrix(int nrows, int ncols, T fill);
26   matrix<T> copy(); // hard copy

28   static matrix<T> create_row_vec(int ncols, T fill);
29   static matrix<T> create_col_vec(int nrows, T fill);
30   static matrix<T> create_diagonal(int nrows, T fill);
31   static matrix<T> read_matrix(const char* file_name);
32   void print(const char* file_name);
33   void print();
```

Check their description and usage in the [document](document).

# 3 x86 and Arm

## 3.1 Metrics: CPU Cycles

Speaking to differences between x86 and Arm, the first thought in my head is instructions. x86 is a CISC instruction set and Arm is a RISC instruction set. So to compare them, the most obvious metrics is the number of instructions executed. We expect that Arm will execute much more instructions than x86 on the same program.

However, counting the number of instructions need CPU and operating system's API support. Inside most CPU there is a PMU (Performance Monitoring Unit) that counts information on machine-level. Luckily Linux has a tool `perf` to access PMU.

The command is `perf stat <program>`. However, I tried it on three different Linux computers (my PC, my lab's server, Huawei's server) and it just did not work.

```
 1  [dongzheng@ecs001-0021-0032 ~]$ perf stat -v ls
 2  Using CPUID 0×00000000480fd010
 3  Warning:
 4  cycles event is not supported by the kernel.
 5  Warning:
 6  instructions event is not supported by the kernel.
 7  Warning:
 8  branches event is not supported by the kernel.
 9  Warning:
10  branch-misses event is not supported by the kernel.
11  htop-2.2.0-8.fc32.aarch64.rpm  neofetch-7.1.0-1.5.noarch.rpm
    openEuler_aarch64.repo  test.cpp
12  task-clock: 740580 740580 740580
13  context-switches: 0 740580 740580
14  cpu-migrations: 0 740580 740580
15  page-faults: 61 740580 740580
16  failed to read counter cycles
17  failed to read counter instructions
18  failed to read counter branches
19  failed to read counter branch-misses
20
21   Performance counter stats for 'ls':
22
23              0.74 msec task-clock                #    0.797 CPUs
    utilized
24                 0      context-switches          #    0.000 K/sec
25                 0      cpu-migrations            #    0.000 K/sec
26                61      page-faults               #    0.082 M/sec
27    <not supported>      cycles
28    <not supported>      instructions
29    <not supported>      branches
30    <not supported>      branch-misses
```

```
31
32          0.000929716 seconds time elapsed
33
34          0.000965000 seconds user
35          0.000000000 seconds sys
```

It all showed `<not supported>` on these three different Linux systems, which is very weird. Then I spent all day finding solutions on foreign websites, and I failed. All the configs are right and god knows why it doesn't work.

```
1  [dongzheng@ecs001-0021-0032 ~]$ cat /usr/src/kernels/4.19.90-
   2003.4.0.0036.oe1.aarch64/.config | grep CONFIG_HW_PERF_EVENTS
2  CONFIG_HW_PERF_EVENTS=y
3  [dongzheng@ecs001-0021-0032 ~]$ cat /usr/src/kernels/4.19.90-
   2003.4.0.0036.oe1.aarch64/.config | grep CONFIG_PERF_EVENTS
4  CONFIG_PERF_EVENTS=y
```

Then I started to read Intel and Arm's PMU manual to get information about perf events.

In Intel(R) 64 and IA-32 Architectures Software Developer's Manual Volume 3B, section 19.1, there is a table:

Table 19-1. Architectural Performance Events

| Event Num. | Event Mask Name | Umask Value | Description |
|---|---|---|---|
| 3CH | UnHalted Core Cycles | 00H | Counts core clock cycles whenever the logical processor is in C0 state (not halted). The frequency of this event varies with state transitions in the core. |
| 3CH | UnHalted Reference Cycles[1] | 01H | Counts at a fixed frequency whenever the logical processor is in C0 state (not halted). |
| C0H | Instructions Retired | 00H | Counts when the last uop of an instruction retires. |
| 2EH | LLC Reference | 4FH | Accesses to the LLC, in which the data is present (hit) or not present (miss). |
| 2EH | LLC Misses | 41H | Accesses to the LLC in which the data is not present (miss). |
| C4H | Branch Instruction Retired | 00H | Counts when the last uop of a branch instruction retires. |
| C5H | Branch Misses Retired | 00H | Counts when the last uop of a branch instruction retires which corrected misprediction of the branch prediction hardware at execution time. |

To get these counters, use `perf stat -e r<umask><event#> <program>`. And I found the cycles counter (0x003c) work.

On my lab's server (Intel Xeon Gold 6240)

```
1  undergrad1@s001:~/dz$ perf stat -e r003c ls
2  LightGBM.tar.gz  log  supersegment
3
4   Performance counter stats for 'ls':
5
6        1,742,296        r003c
7
8      0.000919820 seconds time elapsed
9
```

So I found a way to count CPU cycles.

And for Arm's CPU, look up its ISA first.

On Arm server:

```
1  [dongzheng@ecs001-0021-0032 ~]$ cat /proc/cpuinfo
2  processor        : 0
3  BogoMIPS         : 200.00
4  Features         : fp asimd evtstrm aes pmull sha1 sha2 crc32 atomics
   fphp asimdhp cpuid asimdrdm jscvt fcma dcpop asimddp asimdfhm
5  CPU implementer : 0×48
6  CPU architecture: 8
7  CPU variant      : 0×1
8  CPU part         : 0×d01
9  CPU revision     : 0
10
11 ...
```

Referring to `MIDR` register. The CPU of this server is HiSilicon's Kunpeng-920 r1p0, ARMv8.

Then find an ARMv8 architecture CPU on Arm's document, eg. Cortex-A53. In PMU events section, there is another [table](#):

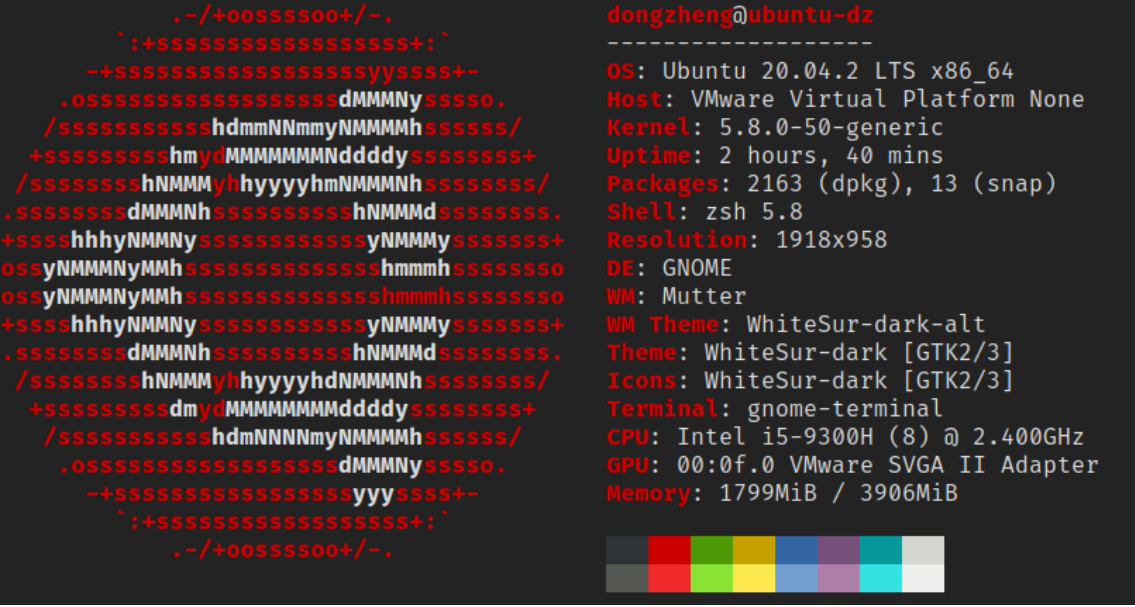The event number of CPU cycles is `0×11`.

Therefore, to count CPU cycles, use `perf stat -e r11 <program>`, and it did work.

```
1  [dongzheng@ecs001-0021-0032 ~]$ perf stat -v -e r11 ls
2  Using CPUID 0×00000000480fd010
3  htop-2.2.0-8.fc32.aarch64.rpm   neofetch-7.1.0-1.5.noarch.rpm
   openEuler_aarch64.repo   test.cpp
4  r11: 2147483647 728390 728390
5
6   Performance counter stats for 'ls':
7
8       2,147,483,647         r11
9
10        0.000918686 seconds time elapsed
11
12        0.000957000 seconds user
13        0.000000000 seconds sys
14
```
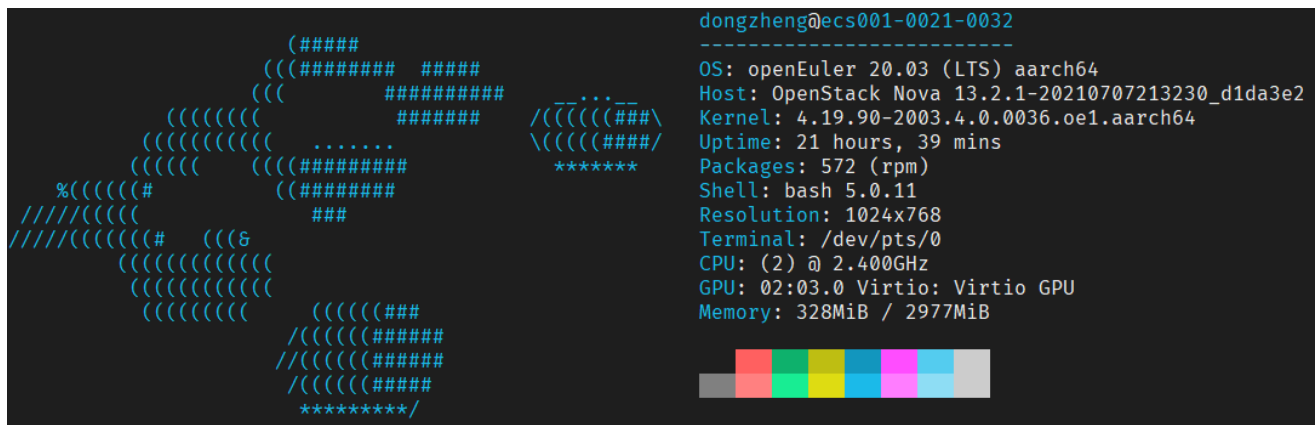
Finally, I found a way to compare CPU cycles.

# 3.2 Test Platform

x86_64:



Arm64:

## 3.3 Dataset & Test Cases

Use the same dataset as **project3**, the dimension of matrices are 32, 64, 128, 256, 512, 1024 and 2048.

The test program is still matrix multiplication, compiled with `set(CMAKE_BUILD_TYPE "Release")` to get maximum compiler optimization.
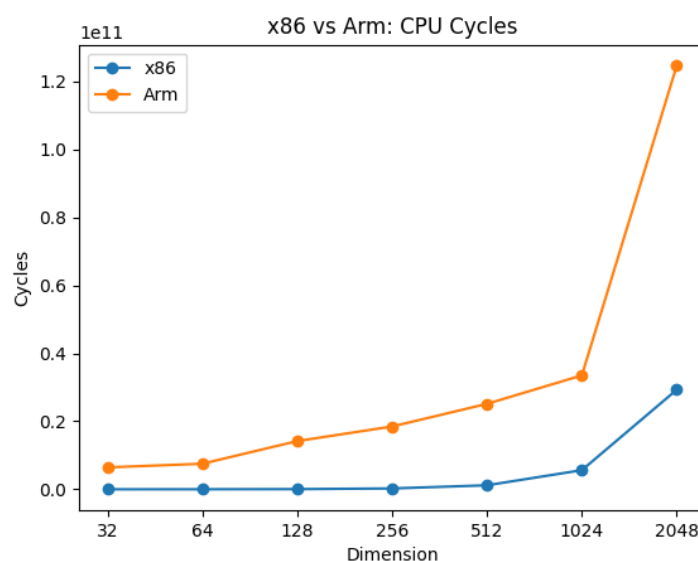
For each dimension, run it for 10 times and calculate average to improve accuracy.

The command (on x86) I use is:

```
1   perf stat -e r003c -x, -r 10 ../matmul.out ../data/mat-A-{dim}.txt
    ../data/mat-B-{dim}.txt ./out/out-{dim}.txt 2>>res_x86.csv
```

## 3.4 Comparison On Matrix Multiplication

Here is the result:



Raw data:

| Dim | x86 | Arm |
|-----|-----|-----|

| Dim | x86 | Arm |
|---|---|---|
| 32 | 5,437,450 | 6,442,450,941 |
| 64 | 15,888,681 | 7,516,192,764 |
| 128 | 54,962,123 | 14,173,392,070 |
| 256 | 248,184,298 | 18,468,359,364 |
| 512 | 1,177,740,654 | 25,125,558,670 |
| 1024 | 5,620,732,185 | 33,500,744,893 |
| 2048 | 29,277,503,823 | 124,554,051,526 |

From the result, we can see that Arm used much more CPU cycles, about $10^1$ larger than x86's, as we expected.

From my test results, on same programs, on average, Arm uses **6.3126** times of x86's CPU cycles (sum of arm / sum of x86). Which means Arm is short for complex programs.

# 3.5 References

[How can I get the number of instructions executed by a program?](#)

[why does perf stat show "stalled-cycles-backend" as <not supported>?](#)

[Why can't I find hardware cache event in my perf list?](#)

[Perf Wiki Tutorial](#)

[Intel(R) 64 and IA-32 Architectures Software Developer's Manual Volume 3B: System Programming Guide](#)

[Using the perf utility on Arm](#)

[PERF tutorial: Counting hardware performance events](#)

[Arm CPU Vendor 及 Part ID 映射关系](#)

[Arm Cortex-A53 MPCore Processor Technical Reference Manual](#)

# 4 Conclusion

In this project, I learnt to design a complete class in C++, especially overriding the operators.

I learned how to manage memory when using soft copy. And a new concept ROI and its manipulations.

In addition, I learned how to use `doxygen` to generate docs for C++ source codes.

And I compared the performance of x86 and Arm ISA. Thus, I had a better understanding on different architectures. And I need an OS engineer to fix these perf events and a doctor to fix me.