# Project 3: Improved Matrix Multiplication in C

| SID | NAME |
| --- | --- |
| 12212108 | Bai Junyan |

# Introduction

In Project 2, we have compared the performance of matrix multiplication written in C and Java and reach a conclusion that, generally, C outperforms Java due to better memory management, higher load of CPU and enhanced compiler optimization. In this project, we will focus on improving the overall performance of matrix multiplication in C as much as possible and compare our best participant against OpenBLAS (Open Basic Linear Algebra Subprograms).

So as to "win" the battle(though it's unlikely), we shall know detailed message about OpenBLAS --- why is it so fast?
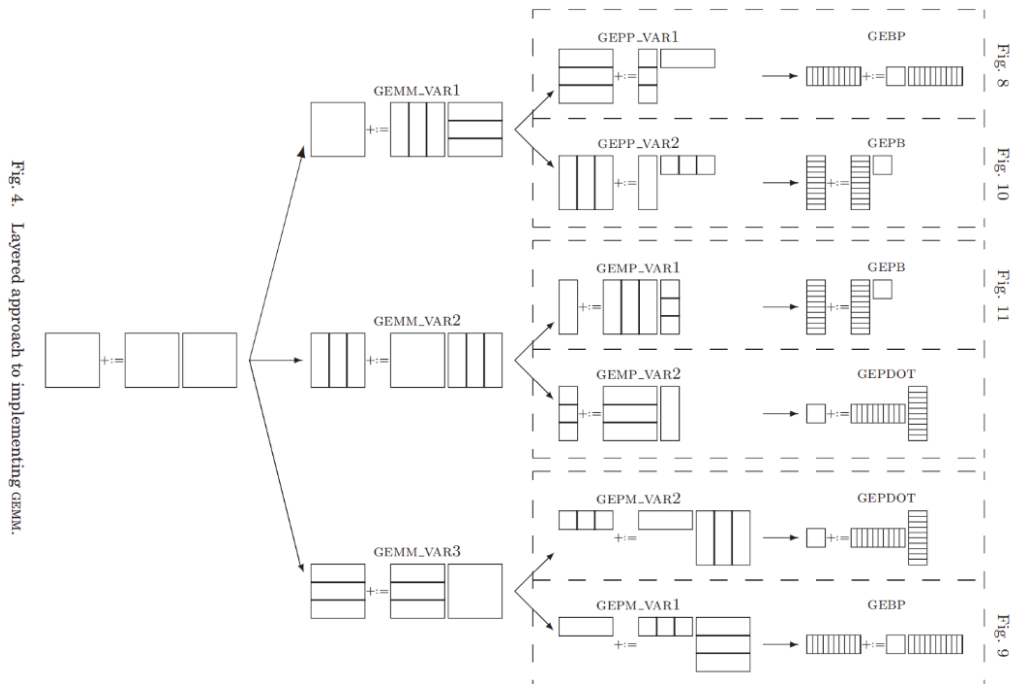
*Why OpenBlAS is so fast in Matrix Multiplication?*

OpenBLAS, an open-source implementation of the BLAS interface, is highly optimized specially for nearly every common computer architecture by hundreds of devoted engineers, it makes the most out of the nowadays "Computer Memory Hierarchy" and combine skills including Multi-threading, Cache-aware Algorithm, SIMD Instructions, Compiler Optimizations etc.

The core idea of GEMM in OpenBLAS is based on GotoBLAS, of which the concept is clarified in the essay written by KAZUSHIGE GOTO - "Anatomy of High-Performance Matrix": **If all matrix operands are $n * n$ in size, $O(n^3)$ floating point operations are performed with $O(n^2)$ data so that the cost of moving data between memory layers (main memory, the caches, and the registers) can be amortized over many computations.** Detailed explanation is as follow:

## Key Ideas

Divide GEMM into smaller sub problem GEBP,GEPB or GEPDOT, by optimizing the lowest level(the smaller block-size matrix) kernel, then high performance can be acquired in other circumstances.

Fig. 4. Layered approach to implementing GEMM.

The pseudocode is as follow:

```
1   for p = 1:K
2       for i = 1:M
3           for j = 1:N
4               Cij += Aip * Bpj
5           endfor
6       endfor
7   endfor
```

To optimize GEPB problems, core idea is to consider the multi-level of Computer memory hierarchy.
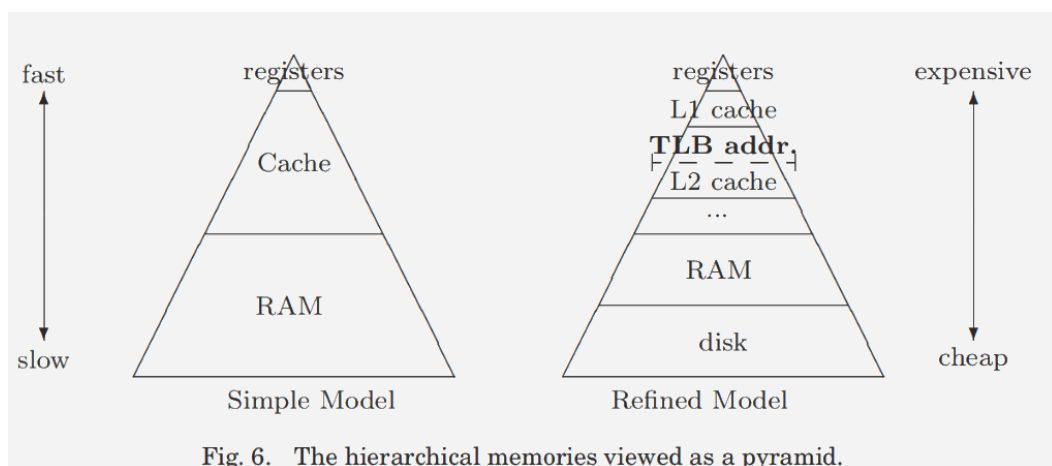


Fig. 6. The hierarchical memories viewed as a pyramid.

Using cache, the pseudocode can be modified to:

```
1  Load A into cache    (m_c*k_c memops)
2  for j = 0,.... N-1
3      Load B_j into cache (k_c * n_r memops)
4      Load C_j into cache (m_c * n_r memops)
5      calculate GEPB  (C_j := AB_j + C_j)
6      Store C_j into memory (m_c * n_r memops)
7  endfor
```

By analyzing the concept of **amortizing**, the problem is clear:

- Which level of cache to put matrix A in: as far as possible, to make $m_c * k_c$ as large as possible, as it turns out, L2 cache is our choice. And $B_j, C_j$ are stored into cache. We can measure each $R_{comp}, R_{load}$ as the rate of floating point arithmetic by CPU and loading data from L2 stream to register groups. We need to ensure that

$$n_r \geq \frac{R_{comp}}{2R_{load}}$$

- TLB miss is way more unaffordable than cache miss(cache miss can be compensated by prefetch), need to make sure that:

    - A,B,C can be found by TLB, so that there is no TLB miss when calculating $C_j := AB_j + C_j$

    - A is repeatedly accessed, A should be stay in TLB as long as possible.

Then our solution is:

- Packing: pack A into continuous accessed memory, then choose argument $m_c, k_c$ to fulfill the above requirement. Also, packing B or C is better.

- Store all of the matrix by column major order, which would make the multiplication $C_{aux} := A^- B_j$ consecutive in memory.

- Choosing $m_r * n_r$

- Choosing $k_c$: To amortize, $k_c$ should be as large as possible.

    - $B_j$' $k_c n_r$ floating point number should only take up below half the space of L1 cache.

    - $m_c k_c$(the space of $A^-$) floating point number should take up much of $L_2$.

## My Implementation

Based on the approach proposed by GOTO and following the critical steps of another essay - "BLISlab: A Sandbox for Optimizing GEMM" (BLIS is based on OpenBLAS and performs nearly as good as it), I conduct my experiment in the following steps:

- Implementing plain matrix multiplication and plain SIMD optimization first, and compare their performance.

- Implementing blocked plain matrix multiplication and compare the performance difference.

- Implementing the scheme for GEPB and test its performance.

- Optimizing for the $8 * 4$ `micro-bench` and test its performance shift.

- Go parallelizing utilizing `OpenMP`.

- Testing performance on the server.

Note that some famous algorithm such as Strassen are not implemented since they are "theoretically" fast, which is in practice not ideally fast especially for GEMM in rather smaller dimension. I mainly focus using Computer itself for the most.

# Experiment

## Local Experiment Environment and Project Structure

Hardware Environment:

```
 1  Architecture:            x86_64
 2  CPU op-mode(s):          32-bit, 64-bit
 3  Byte Order:              Little Endian
 4  Address sizes:           39 bits physical, 48 bits virtual
 5  CPU(s):                  20
 6  On-line CPU(s) list:     0-19
 7  Thread(s) per core:      2
 8  Core(s) per socket:      10
 9  Socket(s):               1
10  Vendor ID:               GenuineIntel
11  CPU family:              6
12  Model:                   154
13  Model name:              12th Gen Intel(R) Core(TM) i7-12700H
14  Stepping:                3
15  CPU MHz:                 2688.011
16  BogoMIPS:                5376.02
17  Virtualization:          VT-x
18  Hypervisor vendor:       Microsoft
19  Virtualization type:     full
20  L1d cache:               480 KiB
21  L1i cache:               320 KiB
22  L2 cache:                12.5 MiB
23  L3 cache:                24 MiB
```

Software Environment:

```
1  GNU Make 4.2.1
2  gcc version 9.4.0 (Ubuntu 9.4.0-1ubuntu1~20.04.2)
```

The project structure is as follow, `dclock.c` is a time measuring function adapted from the original `bl2_clock()` routine in the BLIS library.

```
1  ├── dclock.c
2  ├── m.c
3  ├── main.c
4  ├── makefile
5  ├── output.t
6  └── parameter.h
```

## Step0: Plain Matrix Multiplication

In step0, plain matrix multiplication is implemented as follow, note that **macro is defined to make sure matrices are stored in column-major**. The linear array at address $C$ is used to store elements $C_{i,j}$ so that $i, j$ elements is mapped to location `j * ldc + i`.

```
1  #define A(i,j) a[ (j)*lda + (i) ]
2  #define B(i,j) b[ (j)*ldb + (i) ]
3  #define C(i,j) c[ (j)*ldc + (i) ]
```

`a,b,c` are matrices `m * k, k * n, m * n` each, while all of them are embedded in a larger array that has each `lda, ldb, ldc` rows, which are called "leading dimension".

```
1  void matmul_plain(int m,int n,int k,float *a,int lda,float *b,int
   ldb,float *c,int ldc){
2      if(a==NULL||b==NULL||c==NULL) return;
3      int i,j,p;
4      for(i=0;i<n;i++){
5          for(j=0;j<k;j++){
6              for(p=0;p<m;p++){
7                  C(p,i) = C(p,i) + A(p,j)*B(j,i);
8              }
9          }
10     }
11 }
```

Compile this with `-O3 -march=native` and compare the GFLOPS with OpenBLAS library with largest matrix dimension as 1024, the result is as follow:



OpenBLAS literally kills the plain implementation. Let's take a closer look at the plain GEMM. With the increase of the matrix size, GLOPS drop significantly with plain implementation. Since the order of loop has already been changed, cache miss is significantly reduced due to continuous memory access.

## Step0.5: Plain SIMD

Now's let's use plain SIMD to optimize the code, since my CPU is Intel's supporting `AVX2`, then 256-bit register `__m256` can be used in GEMM which can calculate 8 elements of `c` within one iteration. The key point here is the use of `_mm256_set1_ps` which store `B(j,i)` to all elements of `y`, avoiding storing data into register in the inner `p` loop repeatedly (in the inner loop, `B(j,i),B(j+1,i),...,B(j+7,i)` all needs to be loaded into register).
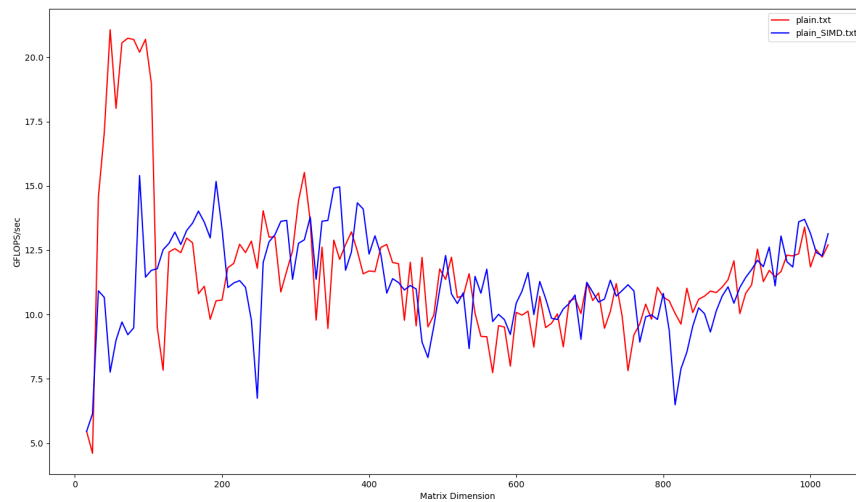
```
1   void matmul_plain_SIMD(int m,int n,int k,float *a,int lda,float *b,int
    ldb,float *c,int ldc){
2       if(a==NULL||b==NULL||c==NULL) return;
3       if(m%8!=0||n%8!=0||k%8!=0) fprintf(stderr,"Input size must be
    multiplier of 8!");
4       int i,j,p;
5       __m256 x,y;
6       __m256 gamma=_mm256_setzero_ps();
7       for(i=0;i<n;i++){
8           for(j=0;j<k;j++){
```

```
9              y=_mm256_set1_ps(B(j,i)); //store B[i,j] to all element of
   y(the jth column)
10             for(p=0;p<m;p+=8){
11                 x=_mm256_loadu_ps(&A(p,j));
12                 gamma=_mm256_add_ps(gamma,_mm256_mul_ps(x,y));
13                 _mm256_storeu_ps(&C(p,i),gamma);
14             }
15         }
16     }
17 }
```

The performance comparison between `matmul_plain` and `matmul_plain_SIMD` is as follow:



Surprisingly, the overall performance between two of them are nearly the same with `matmul_plain` outperforms `matmul_plain_SIMD` at smaller dimensions. The problem is located at the `-O3` compiler optimization, in the assembly code, instruction with prefix `-v` indicates that SIMD is used:
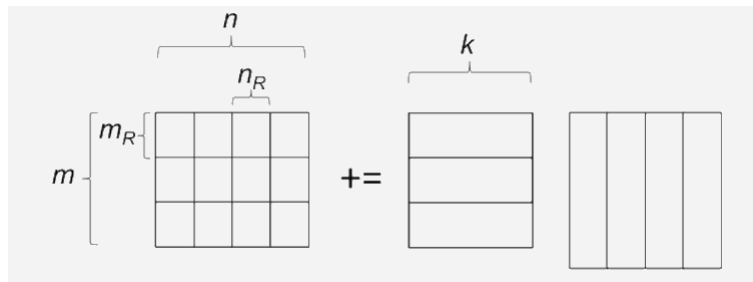
```
1  .L75:
2      vmovups (%rax,%rdi), %xmm3
3      vmovups (%r12,%rdi), %xmm4
4      vinsertf128 $0x1, 16(%rax,%rdi), %ymm3, %ymm0
5      vinsertf128 $0x1, 16(%r12,%rdi), %ymm4, %ymm1
6      vfmadd132ps %ymm2, %ymm1, %ymm0
7      vmovups %xmm0, (%r12,%rdi)
8      vextractf128    $0x1, %ymm0, 16(%r12,%rdi)
9      addq    $32, %rdi
10     cmpq    %r13, %rdi
11     jne .L75
12     cmpq    %r14, %r15
```

After manually implementing SIMD instructions, the compiler cannot grasp what the program is essentially doing, leading to performance decreasing.

## Step1: GEPB Prepared

From Goto's essay, it's essential to optimize GEPB on `micro kernal`. In step1, our "modest first goal" is to utilize this idea and conduct GEMM under GEPB. Essentially, what we want is to implement GEPB so that $m_R * n_R$ blocks of $C$ are kept in registers.



Implementation is as follow, two functions `AddDot` and `AddMM` are defined to calculate matrix multiplication for each $m_R * n_R$ blocks. And fundamental optimizations are implemented, such as loop unrolling and pointer points to the top of the column. `DGEMM_NRr` and `DGEMM_MRr` are predefined as macro which are equal to $m_R$ and $n_R$, they are set to 4 for convenience at first.

```
1   void matmul_tech(int m,int n,int k,float *a,int lda,float *b,int
    ldb,float *c,int ldc){
2       if(a==NULL||b==NULL||c==NULL) return;
3       int i,j;
4       for(j=0;j<n;j+=DGEMM_NRr){
5           for(i=0;i<m;i+=DGEMM_MRr){
6               AddMM(k,&A(i,0),lda,&B(0,j),ldb,&C(i,j),ldc);
7           }
8       }
9   }
10  void AddDot(int k,float *a,int lda,float *b,int ldb,float *gamma){
11      if(a==NULL||b==NULL||gamma==NULL) return;
12      int p;
13      for(p=0;p<k;p++){
14          *gamma += A(0,p)*B(p,0);//A's row multiply by B's column
15      }
16  }
17  void AddMM(int k,float *a,int lda,float *b,int ldb,float *c,int ldc){
18      if(a==NULL||b==NULL||c==NULL) return;
19      int h,l;
20      for(h=0;h<DGEMM_NRr;h++){
21          for(l=0;l<DGEMM_MRr;l++){
22              AddDot(k,&A(l,0),lda,&B(0,h),ldb,&C(l,h));
```
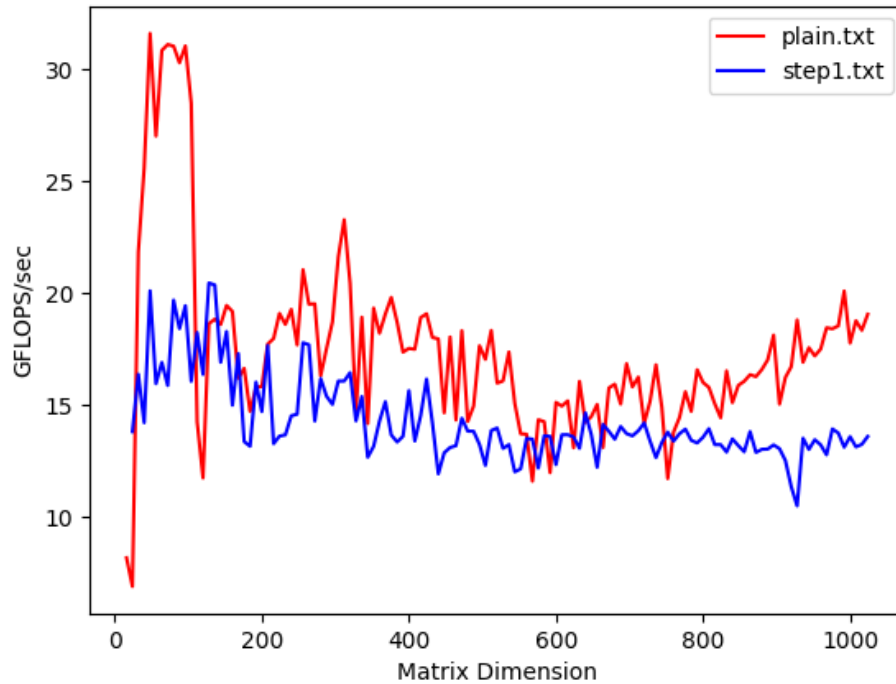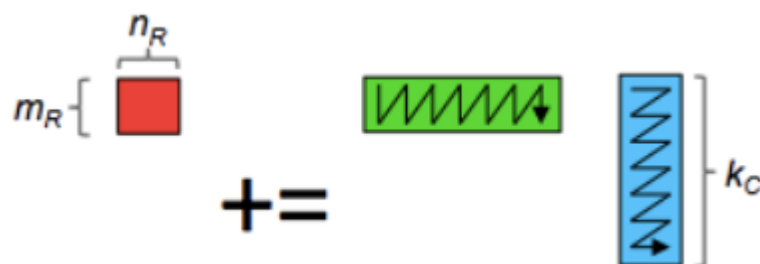
```
23            }
24        }
25  }
```

However, with `-O3 march=native` option turned on, there is essentially no difference between `matmul_plain` and `matmul_tech` (former one outperforms the latter one) which entails the fact again that how much optimization the compiler has already done and what we do is only messing up with it. But we are moving towards the right direction for further optimization.
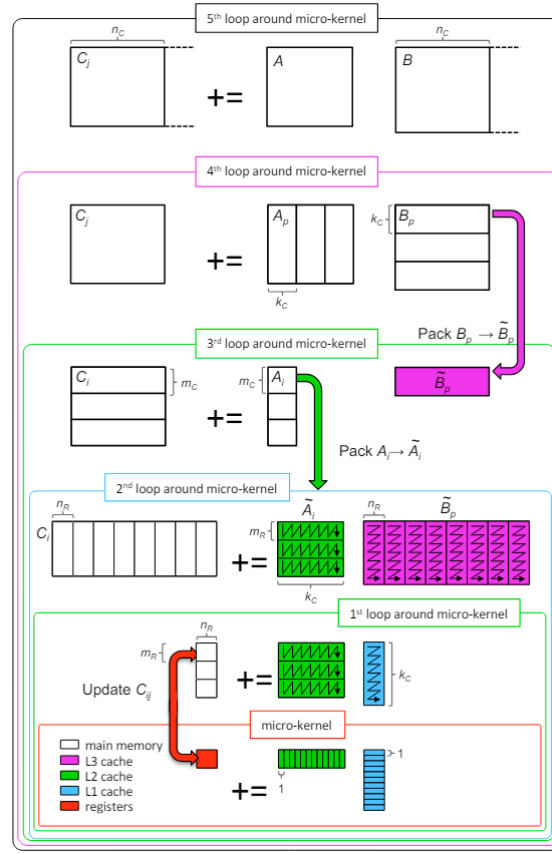


## Step2: GEPB Blocking

Now it's time to implement the approach of GEPB, which is blocking for multiple levels of cache. To amortize the cost of memory loading on computing of floating point, the bottom `micro-kernal` should be highly optimized. And in order to get the best performance, it helps if all data is accessed contiguously, which is why at some point prior to reaching the micro-kernel, data is packed in the order indicated by the arrows:

In the Intro section, it has been discussed that since each column of the block of $A$ in the above picture is multiplied by each element in the corresponding row of the block of $B$ (which are so-called `micro-panels`), the latency to the L2 cache (the time required to bring in an element of the micro-panel of $A$ from that cache) can be amortized over $2n_R$ flops. And by `prefetching` techniques (slightly different from `prefetching` in computer memory system, but the idea is the same), we can put continuously accessed columns on each of $A's$ `micro-panel` onto registers to amortize even further. And the `micro-kernal` of B are kept in L1 cache, which we have to choose parameter $k_c$ carefully. Ideally, this approach may reach the so-called `locally optimzal`, in the sense that assuming data is in a certain memory layer in the hierarchy, the proposed blocking at that level optimally amortizes the cost of data movement with the next memory layer.



Currently, the whole scheme for this approach is established, so the only remaining task to do is to optimize for the calculation of in `micro-kernal`, which would tell exactly the computer to store $A$ into L2 cache, $B$ into L3 cache, continuous $A$ into registers and $B_j$ into L1 cache, just as we want. There are two ways to complete this task: use `avx2` in C or directly write x86 assembly code. Though in Computer Organization class we learn about RISC-V, I'm not so famailiar with X86. I will use `avx2` to write optimize for the `micro-kernal`. Note that the code is adapted from `double` implementation by "BLIS Framwork".

## Step2.1: GEPB Scheme

In this step, we will leave the optimization for `micro-kernal` aside temporarily and complete the Scheme for GEPB proposed in the paper. Critically, there are several factors we need to take into consideration.

- The parameters `DGEMM_NC`, `DGEMM_KC` ,`DGEMM_MC`, `DGEMM_MRr` and `DGEMM_NRr` which each stands for $n_c$,$k_c$,$m_c$, $m_r$ and $n_r$ in the above graph. As what we have discussed in **Step2.1**, these parameters represent the blocksize splited on the way to `micro-kernal`, and all of them have some restrictions, considering a single core on my CPU:

  - $k_c * n_c$ packed $B$ should be kept in L3 cache, which is `24MiB`.

  - $m_c * k_c$ packed A should be kept in L2 cache for most of the time, which is `1.25MiB`.

  - $n_R * k_c$ `micro-panels` $B_j$ should be kept in L1 cache, which is `48KiB`.

We first take $m_R * n_R$ as $8 * 4$ which means 32 single floating points number, each SIMD register is 256-bit, then we need four registers to store them.(This is the size for `micro-kernal`). Since in Goto's paper, experiment shows that $B_j$ should only take up less than half of L1 cache to ensure it would stay there when calculating, thus $k_c$ should be set to 256. Then we set $m_c$ to be 640 to make sure $A$ take more than half of the L2 cache. Finally, $n_c$ can be set up to 24567 since the shared L3 cache is very large, we will set it to 8192.

To pack $A$, $B$ into continuous space, we can simply set a temporary pointer pointing to the discrete space in $A$ and $B$, like `ptr = A(i+di,j+dj)`, however, the discontinuous space will cause a lot of trouble, we could inform the compiler to optimize it by use pointers. Also, we declare the two packing functions to be `inline` function. Notice the slight difference between the order data is scanned in $A$ and $B$ in `micro-panel`.

```
1   inline void packA(int m,int k,float *x,int ldx,int offset,float *pack)
    {//data should be accessed by column first
2       if(x==NULL||pack==NULL) return;
3       int i,p;
4       float *a_ptr[DGEMM_MRr];
5       for(i=0;i<m;i++){a_ptr[i]=x+(offset+i);}
6       for(i=m;i<DGEMM_MRr;i++){a_ptr[i]=x+(offset+0);}
7       for(p=0;p<k;p++){
8           for(i=0;i<DGEMM_MRr;i++){
9               *pack = *a_ptr[i];
10              pack ++;
11              a_ptr[i] = a_ptr[i]+ldx;
12          }
```

```
13        }
14   }
15   inline void packB(int n,int k,float *y,int ldy,int offset,float *pack)
     {//data should be accessed by row first
16       if(y==NULL||pack==NULL) return;
17       int j,p;
18       float *b_ptr[DGEMM_NRr];
19       for(j=0;j<n;j++){b_ptr[j]=y+ldy*(offset+j);}
20       for(j=n;j<DGEMM_NRr;j++){b_ptr[j]=y+ldy*(offset+0);}
21       for(p=0;p<k;p++){
22           for(j=0;j<DGEMM_NRr;j++){
23               *pack++ = *b_ptr[j]++;
24           }
25       }
26   }
```

When allocating space for the packed datas, aligned memory is  allocated to optimize for `micro-kernal`.

```
1    float *malloc_aligned(int m,int n,int size){
2        float *ptr;
3        int e;
4        e = posix_memalign((void**)&ptr,(int)32,size*m*n);//needs a pointer
     to the pointer to reallocate the memory
5        if(e){
6            printf("aligned memory aligned fails");
7            exit(1);
8        }
9        return ptr;
10   }
```

Then it's time to write the 6 loops to approach GEPB, we add the previous `packA` and `packB` each to Loop3 and Loop4 outside of the `micro-kernal`, one thing to emphasize is while entering the next loop, we need to check whether our preset `DEMM_Mx` is out of boarder of the matrix size. We always choose the smaller one by calling `pb=min(k-pc,DGEMM_KC)`.

```
1    for(jc=0;jc<n;jc+=DGEMM_NC){//Loop5
2            jb=min(n-jc,DGEMM_NC);//ensure that the remaining doesn't exceed
     the columns of B when packing
3            for(pc=0;pc<k;pc+=DGEMM_KC){//Loop4
4                pb=min(k-pc,DGEMM_KC);
5                for(j=0;j<jb;j+=DGEMM_NRr){//PackB,iterate through column of
     B
6                    packB(min(jb-
     j,DGEMM_NRr),pb,&b[pc],ldb,jc+j,&_B_[j*pb]);
7                }
```

```
8                    for(ic=0;ic<m;ic+=DGEMM_MC){//Loop3
9                        ib=min(m-ic,DGEMM_MC);
10                       for(i=0;i<ib;i+=DGEMM_MRr){//PackA,iterate through row
   of A
11                           packA(min(ib-
   i,DGEMM_MRr),pb,&a[pc*lda],m,ic+i,&_A_[i*pb]);
12                       }
13                       for(jr=0;jr<jb;jr+=DGEMM_NRr){//Loop2
14                           jrb=min(jb-jr,DGEMM_NRr);
15                           for(ir=0;ir<ib;ir+=DGEMM_MRr){//Loop1
16                               irb=min(ib-ir,DGEMM_MRr);
17                               for(kr=0;kr<pb;kr++){//Loop0--micro-kernal
18                                   C(ir,jr) = A(ir,kr)*B(kr,jr) + C(ir,jr);
19                               }
20                           }
21                       }
22                   }
23               }
24           }
```
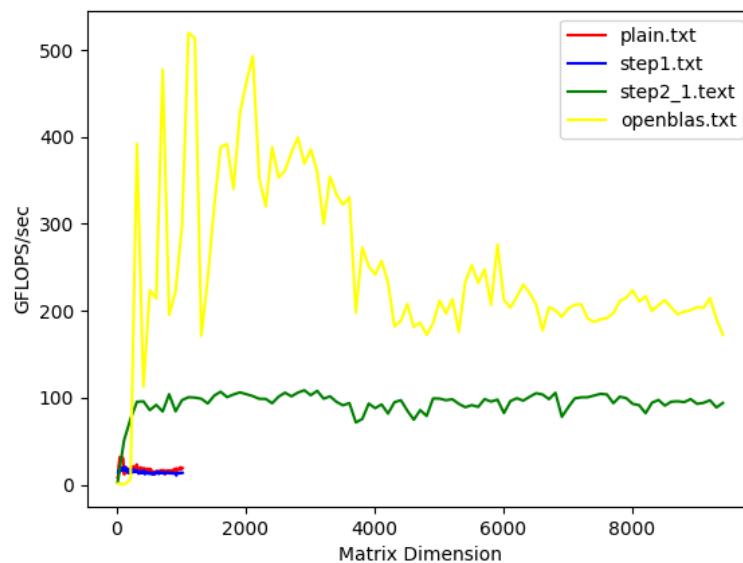
The result is as follow:



It is unsurprising to see that when the matrix size is small, let's say, smaller than $3000 * 3000$, OpenBLAS still way outperforms the GEPB. That's because, well, remember when setting the essential parameter, they are set as follow:
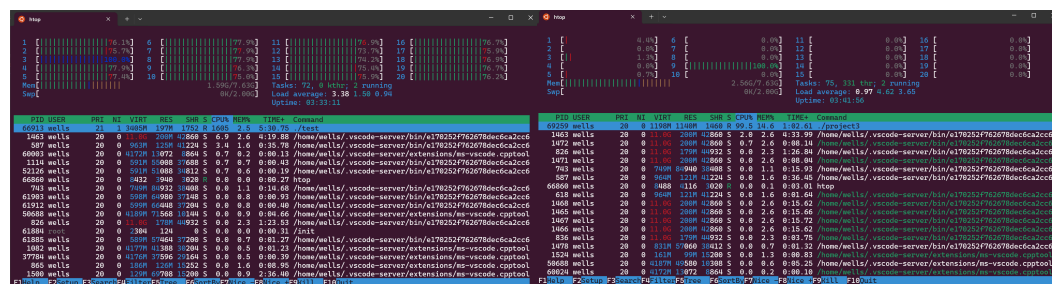
```
1   #define LDA 10000
2   #define LDB 10000
3   #define LDC 10000
4   #define DGEMM_MRr 8
5   #define DGEMM_NRr 4
6   #define DGEMM_NC 8192
7   #define DGEMM_MC 640
8   #define DGEMM_KC 256
```

$m_c, n_c, k_c$ are set to a rather larger value to optimize for caches operation(modern computer has larger and larger cache), and when the matrix size is small, GEPB method is not actually used since we choose the smaller parameter when entering the inner loop.

Still, large performance shift can be seen (more than 10 times faster). And when the matrix size is large, around $8000 * 8000$, we are getting closer to the result of OpenBLAS! That's approximately half of performance, and what's better is that the GFLOPS remain rather stable even the size is getting large.
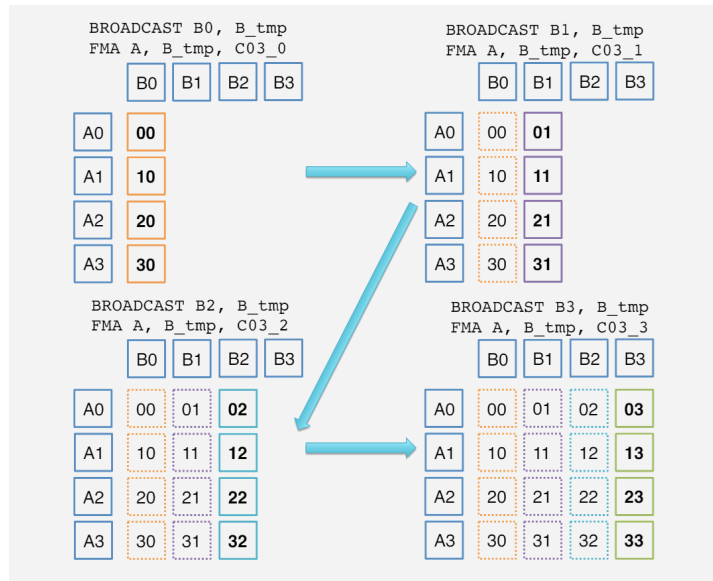
Before optimizing for the `micro-kernal` manually, let's dive deeper why OpenBLAS is faster than the current implementation. Using `htop` on Shell, we can get (OpenBLAS on the left):



OpenBLAS is draining every logical cores of CPU, thus multithreading technology is used, while my plain GEPB implementation only use one single core, that explains the difference in performance. Let's move on to optimize the floating point operation on `micro-kernal`.

## Step2.2: `Micro-Kernal` Optimization

Let's implement a simpler `micro-kernal` using `avx2`  first based on the simple `matmul_plain_SIMD` we implemented in step0. Simpler means we plainly use `avx2` for quicker floating point arithmetic but not optimizing for cache use. The implementation can be abstracted to this graph:
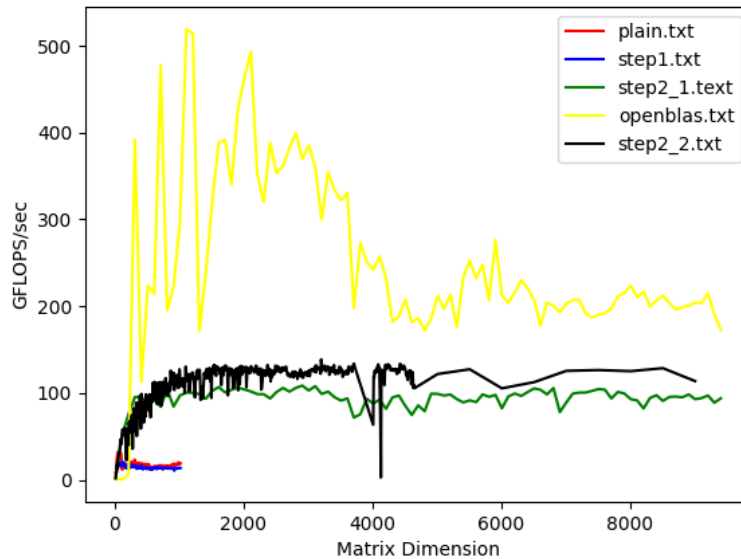
Because my CPU doesn't support `__m256 _mm256_fmadd_ps (__m256 a, __m256 b, __m256 c)`, I can only use `_mm256_add_ps` and `_mm256_mul_ps` combined, which will increase the cost of computation.

```
for (int p = 0; p < k; p++) {
        __m256 a_0p_a_1p_vreg = _mm256_load_ps(a);
        __m256 a_2p_a_3p_vreg = _mm256_load_ps(a + 2);
        a += 4;

        __m256 b_p0_vreg = _mm256_broadcast_ss(b);
        __m256 b_p1_vreg = _mm256_broadcast_ss(b + 1);
        __m256 b_p2_vreg = _mm256_broadcast_ss(b + 2);
        __m256 b_p3_vreg = _mm256_broadcast_ss(b + 3);
        b += 4;
        c_00_c_10_vreg = _mm256_add_ps(c_00_c_10_vreg,
_mm256_mul_ps(a_0p_a_1p_vreg, b_p0_vreg));
        c_01_c_11_vreg = _mm256_add_ps(c_01_c_11_vreg,
_mm256_mul_ps(a_0p_a_1p_vreg, b_p1_vreg));
        c_02_c_12_vreg = _mm256_add_ps(c_02_c_12_vreg,
_mm256_mul_ps(a_0p_a_1p_vreg, b_p2_vreg));
        c_03_c_13_vreg = _mm256_add_ps(c_03_c_13_vreg,
_mm256_mul_ps(a_0p_a_1p_vreg, b_p3_vreg));

        c_20_c_30_vreg = _mm256_add_ps(c_20_c_30_vreg,
_mm256_mul_ps(a_2p_a_3p_vreg, b_p0_vreg));
        c_21_c_31_vreg = _mm256_add_ps(c_21_c_31_vreg,
_mm256_mul_ps(a_2p_a_3p_vreg, b_p1_vreg));
        c_22_c_32_vreg = _mm256_add_ps(c_22_c_32_vreg,
_mm256_mul_ps(a_2p_a_3p_vreg, b_p2_vreg));
        c_23_c_33_vreg = _mm256_add_ps(c_23_c_33_vreg,
_mm256_mul_ps(a_2p_a_3p_vreg, b_p3_vreg));

    }
```

And the result is as follow:



We can see that after slight optimization on `micro-kernal`, the performance is getting closer to OpenBLAS, reaching nearly 64% of its performance. One thing to notice is that my there are actually two modes of my computer -- "Quiet Mode" and "Game Mode",which differs in CPU Clock Frequency. In this case, I only use the "Quiet Mode" for reference. But when doing parallelizing, CPU Clock Frequency matters a lot,(So OpenBLAS's performance right now it's demined), we will see it in step3.

## Step3: Parallelizing with `OpenMP`

Since there are 20 logical cores in my CPU, parallelizing is definitely approachable taken that there are multiple embedded loops. In this project, for convenience and thread safety concern, I choose to use `OpenMP` library to go parallelizing. The question is, which loop to parallelize. By reading another paper "Anatomy of High-Performance Many-Threaded Matrix Multiplication", we can get some insights.

- Parallelizing within Loop0: This is parallelizing when calculating the `micro-kernal`, which is not suggested. There are two factors:
    - We need multiple thread to calculate a single block of C, increasing the cost of the `reduce` operations of threads.
    - The calculation is just too little for every thread since the matrix size is rather small within Loop0, `amortation` is not implemented.
- Parallelizing within Loop1: This is parallelizing in loop `ir`, which is calculating every `micro-panel` of $A$ multiply by the same column $B_j$ which is in L1 cache. To amortize the cost of loading $B_j$ from L3, factor $\frac{mc}{mr}$ must be large for us to parallelize.

- Parallelizing within Loop2: This is parallelizing within loop `jr`, which is calculating every `micro-panel` of $B$ multiply by the same row of $A_i$ that is stored in L2 cache. To amortize the cost of packing the shared $A_i$ on every core, factor $\frac{n_c}{n_r}$ must be large for us to parallelize.

- Parallelizing within Loop3: This is parallelizing within loop `ic`, which is splitting the $m$ rows of $A$ by $m_c$. So we need $\frac{m}{m_c}$ large, which is satisfied since we set $m$ to $640$.

- Parallelizing within Loop4: This is parallelizing within loop `pc`, which is calculating the exact $C(ir, jr)$. Since each thread need to update for it, thread competence exits, will likely lead to poor performance.

- Parallelizing within Loop5: This is parallelizing within loop `jc`, which is parallelizing for the whole matrix multiplication. This involves shared L3 cache fetch, which is better for multi socket machine with `NUMA`.

We will experiment parallelizing within Loop3, Loop4 and 5. To achieve assigning nearly the same workload for each thread, we need to separate the task to fit in different threads, we define an inline function called `thread_shared`:

```
1   inline void thread_shared(int n,int bf,int *start,int *end){
2       int way =omp_get_num_threads();th
3       int id = omp_get_thread_num();
4       //thread partitioning
5       //we partition the space between all_start and
6       // all_end into n_way partitions, each a multiple of block_factor
7       // with the exception of the one partition that recieves the
8       // "edge" case (if applicable).
9       int all_start=0;
10      int all_end=n;
11      int size=all_end-all_start;
12      int whole=size/bf;
13      int left=size%bf;
14      int lo=whole/way;int hi=whole/way;
15      //differentiate between edge cases
16      int b_th_lo=whole%way;
17      if(lo!=0) lo++;
18      int size_lo=lo*bf,size_hi=hi*bf;
19      // Precompute the starting indices of the low and high groups.
20      int lo_start=all_start,hi_start=all_start+b_th_lo*size_lo;
21      // Compute the start and end of individual threads' ranges
22      // as a function of their work_ids and also the group to which
23      // they belong (low or high).
24      if(id<b_th_lo){
25          *start=lo_start+id*size_lo;
26          *end=lo_start+(id+1)*size_lo;
27      }
```
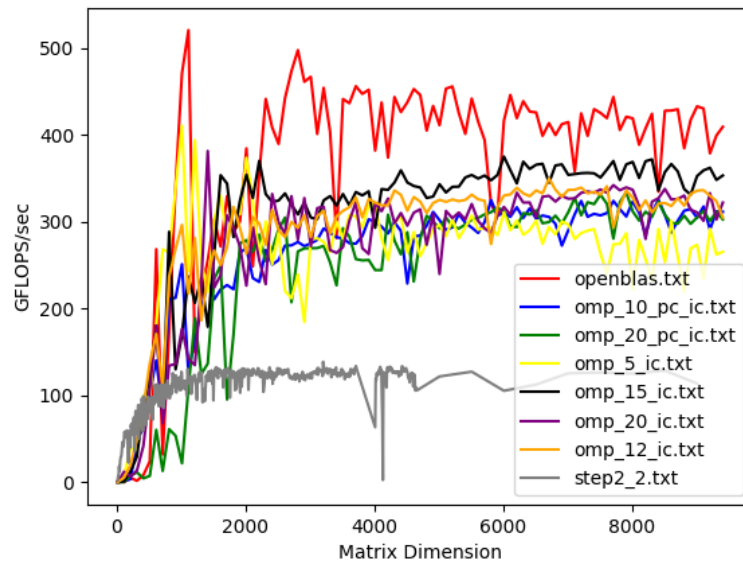
```
28        else{*start=hi_start+(id-b_th_lo)*size_hi;
29        *end=hi_start+(id-b_th_lo+1)*size_hi;
30        if(id==way-1) *end+=left;}
31   }
```
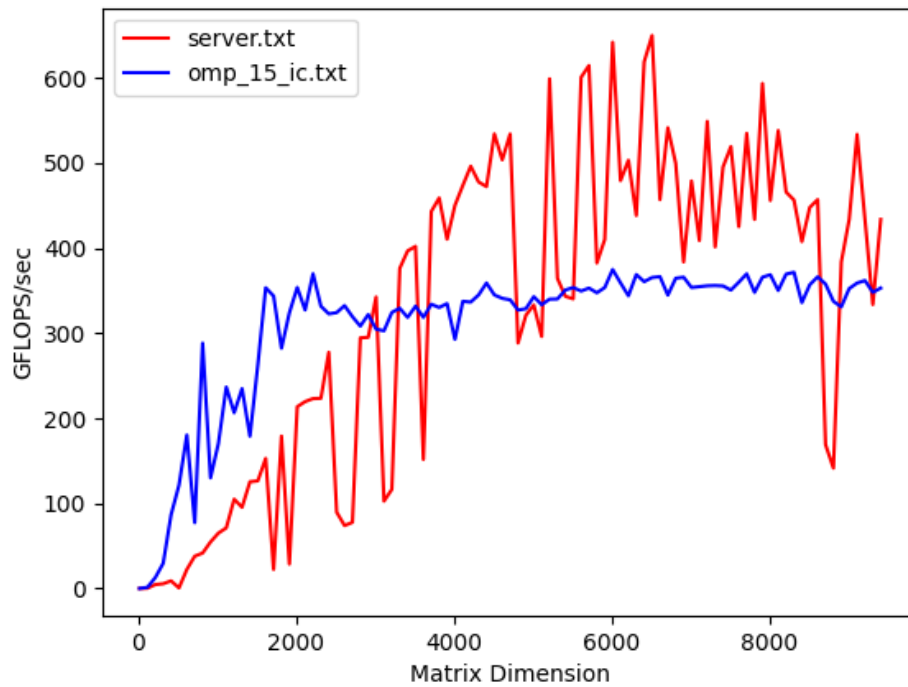
Then experiment parallelizing within different loops, the result compared to OpenBLAS is as follow. One thing to note is that `Game Mode` of my computer is turned on to enable higher CPU frequency. It can be seen that both OpenBLAS and my implementation increases):



Compared with single thread, GFLOPS is around 2.5 times bigger. And right now we are reaching more than 75% the performance of OpenBLAS! And the result shows that only parallelizing within Loop3 has the best performance among them. Adding parallelizing within Loop2(`pc`) is indeed worsen the performance, in correspondence to the previous theory.

## Step4: Testing on the Server

Testing result on the server, we can get

As expected, it outperfoms local performance since the better CPU and more cores.

Running the $64000 * 64000$ matrix multiplication, the result is as follow:

| My | Server |
|---|---|
| 3.451305e+02 GFLOPS/sec | 4.8999e+02 GFLOPS/sec |

# Conclusion

By applying the method proposed by Goto's paper, I reach nearly 75% of the performance of OpenBLAS at last. This project really makes me realize the importance of diving deep into the computer memory and optimize our code according to it. I find it amazing when this project reflects close connection to what I've learned in another course Computer Organization, which encourages me to learn more about the bottom level knowledge. For example, I read the famous article "What every programmer should know about computer memory", which inspires me a lot.

There are still things I want to do but fail to accomplish, and I wish to complete them in the future:

- optimize the `micro-bench` by writing assembly code, taking advantage of the cache memory.

- Using GPU to run the code and see its performance.

- Fully understand how MKL and OpenBLAS works so well on different platforms by learning the source code.

# Reference

1. https://github.com/OpenMathLib/OpenBLAS/wiki

2. https://www.cs.utexas.edu/~flame/pubs/GotoTOMS_revision.pdf

3. https://www.cs.utexas.edu/~flame/pubs/blis3_ipdps14.pdf

4. https://github.com/flame/blislab

5. https://github.com/flame/how-to-optimize-gemm