

CS301-HPC Lab 2 (2021)

Assignment-1 Deadline: 3rd Feb

CPU architecture, Vector Triad, Measuring Performance and Code Optimization Techniques

To be included in the report: Part- A,B,C

A) Vector triad benchmarking

As seen in the previous exercise, the CPU memory is divided into 4 parts:

1. L1 cache
2. L2 cache
3. L3 cache
4. RAM Memory

These are listed in the order of the bandwidth/ latency at which they operate.

The goal of this exercise is to explore the impact of the memory hierarchy on the computation speed of a CPU. For this, we will use the Vector Triad benchmarking.

This benchmark consists of performing a vector sum between 3 vectors and storing it in another vector. Consider 4 vectors **A, B, C, D** of size **N** each.

Pseudocode for Vector triad is as follows:

```
for i=1 to N
    A[i] = B[i] + C[i]*D[i]
```

For benchmarking purposes we will measure the time taken by the vector triad for different problem sizes **N**.

```
minSize = 2^8
maxSize = 2^29
total = maxSize
for size =minSize; size<maxSize; size*=2
```

```

RUNS = total/size //initialise arrays
double A[size], B[size], C[size], D[size]
start_time = clock()
for j=1 to RUNS
    for i=1 to size
        A[i] = B[i] + C[i]*D[i]
    end_time = clock() - start_time
    throughput = (sizeof(double)*2*total)/end_time
print size, throughput

```

Using python/matlab/gnuplot make the necessary plots (e.g [problem size vs performance](#)) and investigate the following:

- Make the necessary modifications in the code to measure the run-time accurately. What are the couple of clock functions that can be used to measure the elapsed time? Understand how the clock works and how it is used to measure runtime.
- Make modifications in the code to approximately measure the time taken by the computation part.
- Make modifications in the code to approximately measure the time taken by the data access part.
- How does your measured performance compare with the peak-theoretical performance of the computing system?
- Can you show which operation is more expensive, addition or division by making some modifications in the code? Provide explanations for your result.
- Is it possible to fix the clock speed to 2.0 GHz? If yes, can you interpret the observable changes to the results in throughput?

B) Which optimization strategies would you suggest for the piece of pseudo-code below?

```

double mat[N][N], s[N][N], val;
int i, j, v[N];
//
// ... v[] and s[][] --Ⓜ assumed to contain valid data
//
for(i=0; i<N ; ++i) {
    for(j=0; j<N; ++j) {
        val = (double)(v[i] % 256);
        mat[j][i] = s[j][i]*(sin(val)*sin(val)-cos(val)*cos(val));
    }
}

```

Do not make any major assumptions about the size of N.

Compare the performance of the basic implementation and your proposed optimized implementation

- Does the optimization strategy depend on the problem size?
- Can you predict an upper performance limit on the given lab system?

Assignment-2 (to be combined with assignment 1)

C) High Performance Matrix Multiplication on a CPU

Goal is to optimize matrix multiplication to run on a single core and understand how fast we can perform important linear algebra kernels/ routines.

We will start with matrix multiplication which is a level-3 BLAS subroutine (https://en.wikipedia.org/wiki/Basic_Linear_Algebra_Subprograms)

Matrix-matrix multiplication computes $C = C + A * B$, where C is m -by- n , A is m -by- k and B is k -by- n .

Assume we have two levels of the memory hierarchy, fast (cache) and slow (RAM), and that all data initially resides in slow memory.

we can count

m = number of memory references to slow memory needed just to read

the input data from slow memory, and write the output data back

c = number of floating point operations

$CMA = c/m$ = average number of flops per slow memory reference

The significance of CMA : for each word/ byte read from slow memory (the expensive operation), one can hope to do at most CMA ratio operations on it (on average) while it resides in fast memory. The higher the CMA is, the more the algorithm will operate at the top of the memory hierarchy, where it is most efficient.

Some important points/ assumptions to be taken into account as discussed during the lecture, particularly during Block Matrix Multiplication lecture:

1. There are just two levels in the hierarchy, fast and slow.
2. The small, fast memory has size " c " bytes, where $c \ll n^2$, so we can only fit a small part of an entire n -by- n matrix, but $c \geq 4*n$, so we may fit several whole rows or columns depending on the size of the matrix.
3. In the code - Each word/ element is read from slow memory individually (as discussed, in practice, larger groups of words/ elements are read, such as cache lines or memory pages).

In this assignment, we will consider 2 implementations of matrix multiply and compute CMA for each and also do memory/ computational complexity analysis as performed during the last lecture.

Algorithm 1: Simple matrix multiply (Unblocked) $C=A*B$

Look into the pattern - in this case, the innermost loop is doing a dot product of row i of A and column j of B .

What is the CMA ratio in this case?

for a standard code and for matrix size of $n \times n$

$$CMA = (2*n^3)/(n^3 + 3*n^2) \sim \text{around two (2)}$$

$$\begin{aligned} m = \# \text{ slow memory access} &= n^3 \quad \text{read each column of } B \text{ } n \text{ times} \\ &+ n^2 \quad \text{read each row of } A \text{ once for each } i, \\ &\quad \text{and keep it in fast memory during} \\ &\quad \text{the execution of the two inner loops} \\ &+ 2*n^2 \quad \text{read/write each entry of } C \text{ once} \\ &= n^3 + 3*n^2 \end{aligned}$$

Algorithm 2: Square blocked matrix multiply (also called 2D blocked MM)

Now consider C to be an N -by- N matrix of n/N -by- n/N subblocks C_{ij} , with A and B similarly partitioned. As discussed during the lecture.

Show that : CMA in this case is n/N .

Compare the computational throughput of the above two matrix multiplication algorithms (algorithm1 and algorithm2) using several computational experiments with different problem sizes

(changing the size of the matrix, taking into account the Cache sizes and RAM size) and support your observations quantitatively.

To multiply two matrices, we use 3 nested loops: Perform the following numerical experiments, report your observations and support it with **quantitative justifications** in terms of cache performance, data access pattern, reuse of cached data, cache blocking, temporal/ spatial locality of memory accesses etc.

1. Try all 6 different loop orderings. Which ordering perform best for (a) 512 by 512 and (b) 2048-by-2048 matrices? Which ordering is the worst? Support your observations in terms of striding through the matrices with respect to the innermost loop?
2. Why does performance drop for large values of matrices (start from 256 by 256 and try till 2048 by 2048)? How much is the drop in terms of compute throughput?
3. What is the effect on performance if elements of the matrices are changed from integer to double precision? Perform a quantitative analysis for the case 1024 by 1024 and 2048 by 2048.