

LINMA2710 - Scientific Computing

Shared-Memory Multiprocessing

P.-A. Absil and B. Legat

☐ Full Width Mode ☐ Present Mode

Table of Contents

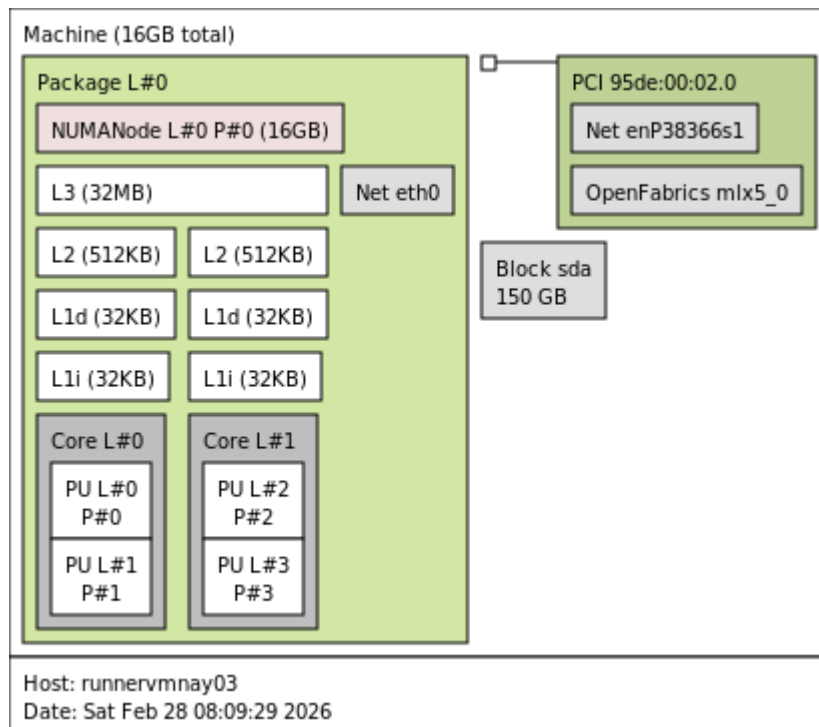
Memory layout

Parallel sum

Amdahl's law

[Eij10] V. Eijkhout. *Introduction to High Performance Scientific Computing*. 3 Edition, Vol. 1 (Lulu.com, 2010).

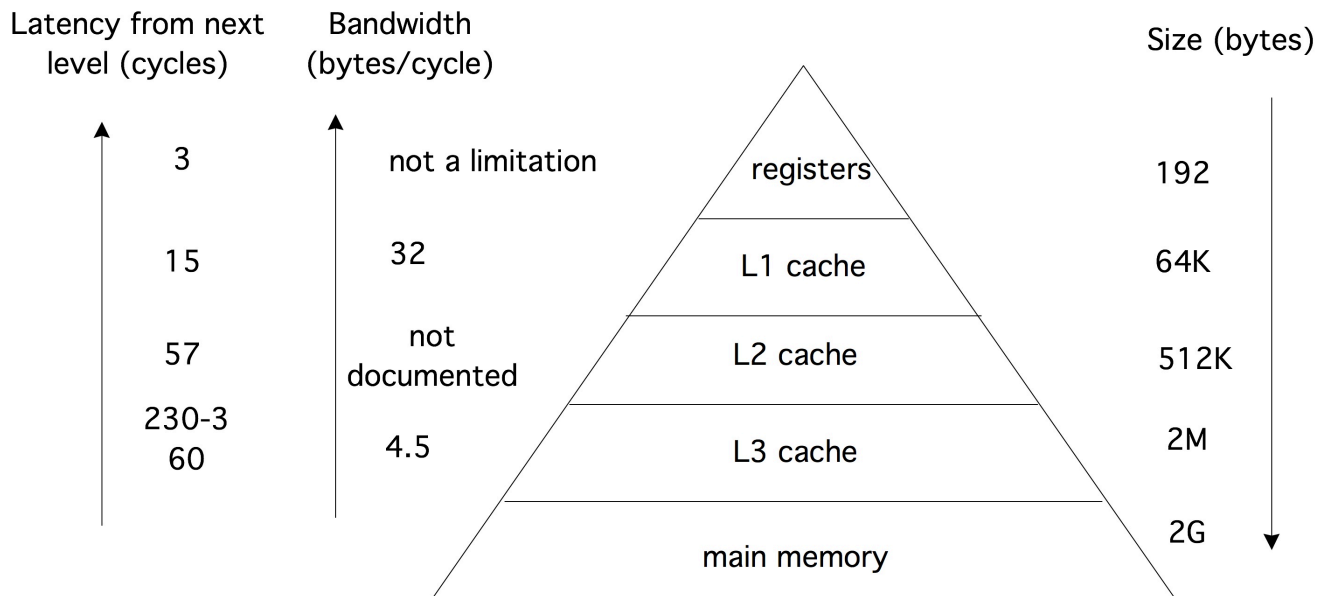
Memory layout ↗



Try it on your laptop!

```
$ lstopo
```

Hierarchy ↷



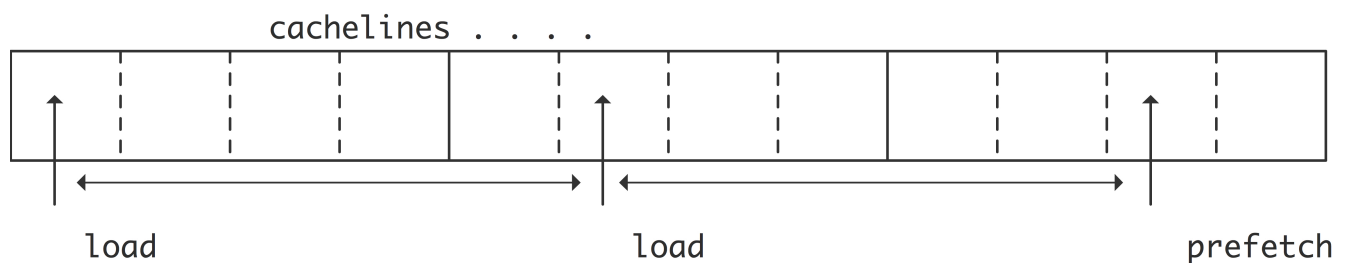
Latency of n bytes of data is given by

$$\alpha + \beta n$$

where α is the start up time and β is the inverse of the bandwidth.

[Eij10; Figure 1.5]

Cache lines and prefetch ↷



- Accessing value not in the cache → *cache miss*
- This value is then loaded along with a whole cache line (e.g., 64 or 128 contiguous bytes)

- Following cache lines may also be anticipated and prefetched

This shows the importance of *data locality*. An algorithm performs better if it accesses data close in memory and in a predictable pattern.

[Eij10; Figure 1.11]

Illustration with matrices

32919.04f0

```
1 @btime c_sum($mat)
```



```
79.337 μs (0 allocations: 0 bytes)
```



```
mat = 256×256 Matrix{Float32}:
```

```
 0.559856  0.059603  0.337242  0.116544  ...  0.127319  0.283428  0.673206
 0.509872  0.50455   0.622709  0.989656  ...  0.97062   0.715898  0.825475
 0.512169  0.148668  0.408      0.0178322 ...  0.129164  0.532762  0.7569
 0.236851  0.293102  0.711055  0.973665  ...  0.393093  0.528149  0.557636
 0.397415  0.85693   0.422944  0.374782  ...  0.38415   0.137626  0.390188
 0.0855426 0.962534  0.575806  0.0184651 ...  0.788814  0.829789  0.297246
 0.482592  0.612828  0.776442  0.492575  ...  0.230199  0.558317  0.24554
  ⋮
 0.0655522 0.418429  0.262089  0.213704  ...  0.83024   0.246336  0.694124
 0.00976539 0.820938  0.687293  0.430617  ...  0.816382  0.57458   0.494722
 0.327859  0.490751  0.689928  0.656891  ...  0.0484409 0.094496  0.669354
 0.36795   0.816392  0.879788  0.0249674 ...  0.0197487 0.139243  0.62569
 0.432604  0.641924  0.76403   0.440289  ...  0.2005    0.590192  0.062048
 0.0963827 0.194705  0.190262  0.848333  ...  0.481382  0.308776  0.203466
```

```
1 mat = rand(Cfloat, 2^8, 2^8)
```

```
1 c_sum(x::Matrix{Cfloat}) = ccall(("sum", sum\_matrix\_lib), Cfloat, (Ptr{Cfloat},
  Cint, Cint), x, size(x, 1), size(x, 2));
```

```
#include <stdio.h>
```

```
float sum(float *mat, int n, int m) {
    float total = 0;
    for (int i = 0; i < n; i++) {
        for (int j = 0; j < m; j++) {
            total += mat[i + j * n];
        }
    }
    return total;
}
```

► What is the performance issue of this code ?

Arithmetic intensity ⇔

Consider a program requiring m load / store operations with memory for o arithmetic operations.

- The *arithmetic intensity* is the ratio $a = o/m$.
- The arithmetic time is $t_{\text{arith}} = o/\text{frequency}$
- The data transfer time is $t_{\text{mem}} = m/\text{bandwidth} = o/(a \cdot \text{bandwidth})$

As arithmetic operations and data transfer are done in parallel, the time per iteration is

$$\max(t_{\text{arith}}, t_{\text{mem}})/o = 1/\min(\text{frequency}, a \cdot \text{bandwidth})$$

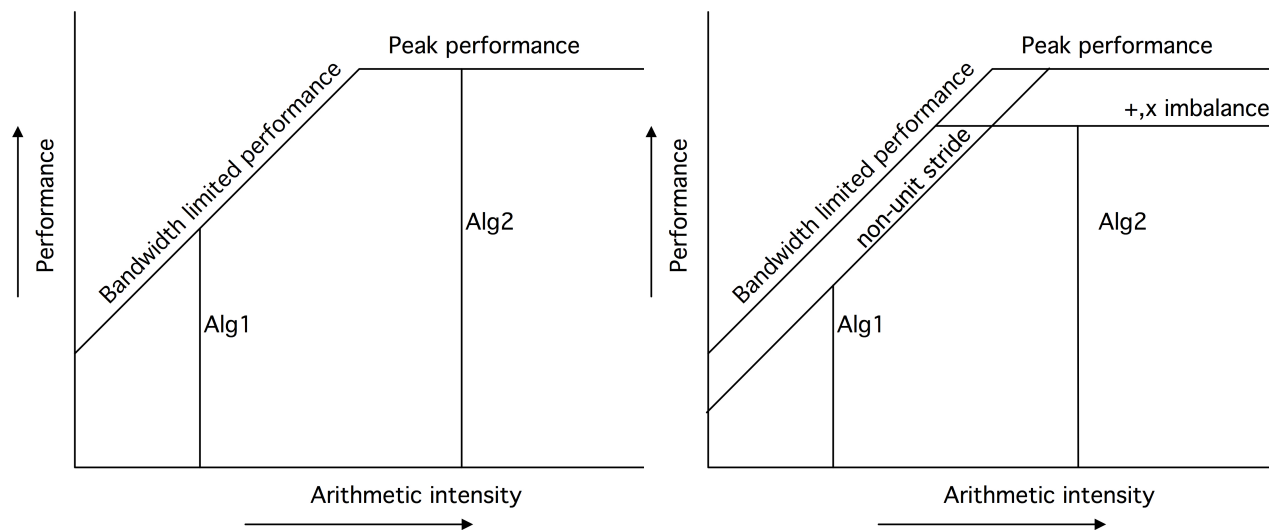
So the number of operations per second is $\min(\text{frequency}, a \cdot \text{bandwidth})$.

This piecewise linear function in a gives the *roofline model*.

Tip

See examples in [Eij10; Section 1.6.1].

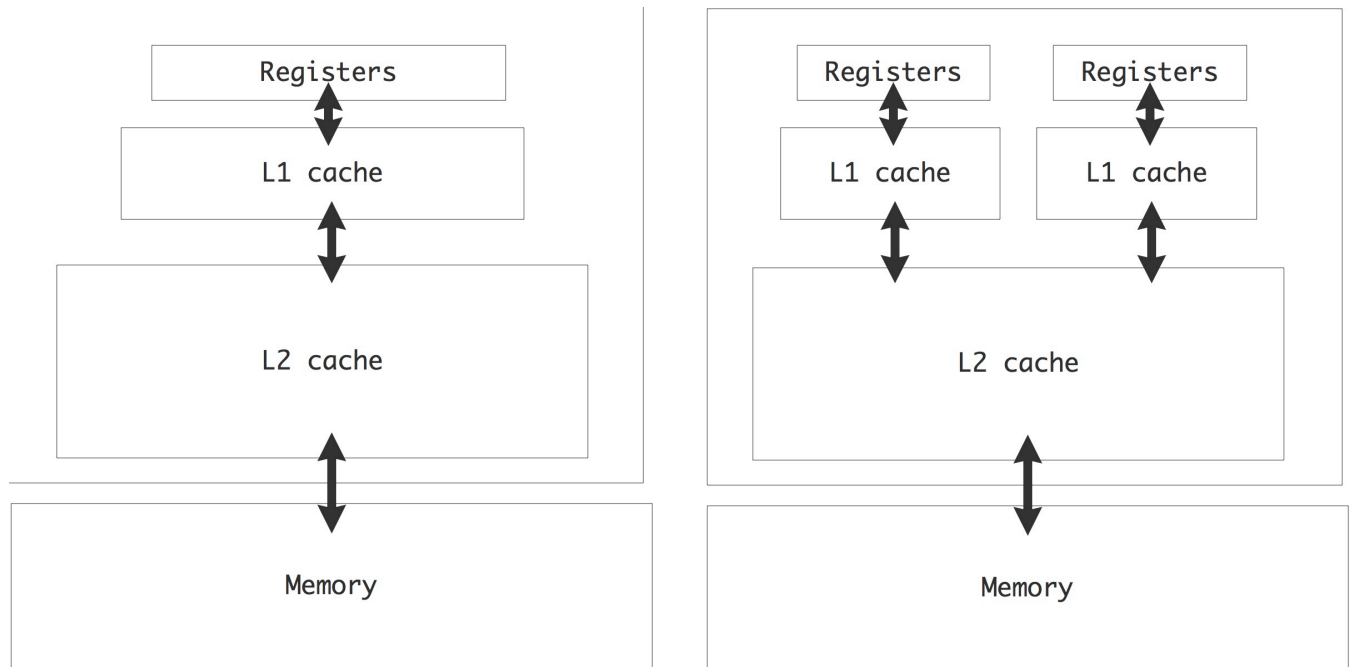
The roofline model



- *compute-bound* : For large arithmetic intensity (Alg2 in above picture), performance determined by processor characteristics
- *bandwidth-bound* : For low arithmetic intensity (Alg1 in above picture), performance determined by memory characteristics
- Bandwidth line may be lowered by inefficient memory access (e.g., no locality)
- Peak performance line may be lowered by inefficient use of CPU (e.g., not using SIMD)

[Eij10; Figure 1.16]

Cache hierarchy for a multi-core CPU ↩



Cache coherence : Update L1 cache when the corresponding memory is modified by another core.

[Eij10; Figure 1.13]

Parallel sum ⇄

```
#include <vector>
#include <stdint.h>
#include <omp.h>
#include <stdio.h>

extern "C" {
float sum(float *vec, int length, int num_threads, int verbose) {
    float total = 0;
    omp_set_dynamic(0); // Force the value 'num_threads'
    omp_set_num_threads(num_threads);
    #pragma omp parallel
    {
        int thread_num = omp_get_thread_num();
        int stride = length / num_threads;
        int last = stride * (thread_num + 1);
        if (thread_num + 1 == num_threads)
            last = length;
        if (verbose >= 1)
            fprintf(stderr, "thread id : %d / %d %d:%d\n", thread_num, omp_get_num_threads(),
                stride * thread_num, last - 1);
        #pragma omp simd
        for (int i = stride * thread_num; i < last; i++)
            total += vec[i];
    }
    return total;
}
```

32803.38f0

```
1 @btime c_sum($vec, num_threads = 1, verbose = 0)
```



2.743 μs (0 allocations: 0 bytes)



32803.38f0

```
1 @btime c_sum($vec; num_threads, verbose = 0)
```



3.969 μs (3 allocations: 80 bytes)



32803.38f0

```
1 @time c_sum(vec; num_threads, verbose = 1)
```

```
thread id : 0 / 2 0:32767
thread id : 1 / 2 32768:65535
    0.000051 seconds (47 allocations: 2.422 KiB)
```

vec =

► [0.628269, 0.854317, 0.517498, 0.63904, 0.38714, 0.965666, 0.638772, 0.706559, 0.316497, 0.638772]

```
1 vec = rand(Cfloat, 2^log_size)
```

```
1 c_sum(x::Vector{Cfloat}; num_threads = 1, verbose = 0) = ccall(("sum", sum_lib),  
Cfloat, (Ptr{Cfloat}, Cint, Cint, Cint), x, length(x), num_threads, verbose);
```

Low level implementation using POSIX Threads (pthreads) covered in "LEPL1503 : Projet 3". We use the high level [OpenMP](#) library in this course.

log_size = 16

num_threads =  2

► Can you spot the issue in the code ?

Many processors ⇄

```
#include <omp.h>
#include <stdio.h>

extern "C" {
void sum_to(float *vec, int length, float *local_results, int num_threads, int verbose) {
    omp_set_dynamic(0); // Force the value 'num_threads'
    omp_set_num_threads(num_threads);
    #pragma omp parallel
    {
        int thread_num = omp_get_thread_num();
        int stride = length / num_threads;
        int last = stride * (thread_num + 1);
        if (thread_num + 1 == num_threads)
            last = length;
        if (verbose >= 1)
            fprintf(stderr, "thread id : %d / %d %d:%d\n", thread_num, omp_get_num_threads(),
stride * thread_num, last - 1);
        float no_false_sharing = 0;
        #pragma omp simd
        for (int i = stride * thread_num; i < last; i++)
            no_false_sharing += vec[i];
        local_results[thread_num] = no_false_sharing;
    }
}

float sum(float *vec, int length, int num_threads, int factor, int verbose) {
    float* buffers[2] = {new float[num_threads], new float[num_threads / factor]};
    sum_to(vec, length, buffers[0], num_threads, verbose);
    int prev = num_threads, cur;
    int buffer_idx = 0;
    for (cur = num_threads / factor; cur > 0; cur /= factor) {
        sum_to(buffers[buffer_idx % 2], prev, buffers[(buffer_idx + 1) % 2], cur, verbose);
        prev = cur;
        buffer_idx += 1;
    }
    if (prev == 1)
        return buffers[buffer_idx % 2][0];
    sum_to(buffers[buffer_idx % 2], prev, buffers[(buffer_idx + 1) % 2], 1, verbose);
    return buffers[(buffer_idx + 1) % 2][0];
}
```

Benchmark ⇄

If we have many processors, we may want to speed up the last part as well:

32803.38f0

```
1 @time many_sum(vec; base_num_threads, factor, verbose = 1)
```



```
thread id : 0 / 2 0:32767
thread id : 1 / 2 32768:65535
thread id : 0 / 1 0:1
0.000102 seconds (17 allocations: 2.062 KiB, 58845.64% compilation time)
```



32737.598f0

```
1 @btime c_sum($many_vec)
```



```
2.866 μs (0 allocations: 0 bytes)
```



32737.59f0

```
1 @btime many_sum($many_vec; base_num_threads, factor)
```



```
8.716 μs (3 allocations: 80 bytes)
```




many_vec =

► [0.630695, 0.525817, 0.868061, 0.934856, 0.187473, 0.89147, 0.854506, 0.145506, 0.765051,

```
1 many_vec = rand(Cfloat, 2^many_log_size)
```

```
1 many_sum(x::Vector{Cfloat}; base_num_threads = 1, factor = 2, verbose = 0) =
  ccall(("sum", many_sum_lib), Cfloat, (Ptr{Cfloat}, Cint, Cint, Cint, Cint), x,
  length(x), base_num_threads, factor, verbose);
```

many_log_size =  16

base_num_threads =  2

factor =  2

Amdahl's law ⇔

Speed-up and efficiency ⇔

Def: Speed-up

$$S_p = \frac{T_1}{T_p}$$

Def: Efficiency

$$E_p = \frac{S_p}{p}$$

Let T_p bet the time with p processes

- $E_p > 1 \rightarrow$ Unlikely
- $E_p = 1 \rightarrow$ Ideal
- $E_p < 1 \rightarrow$ Realistic

Amdahl's law ⇔

- F_s : Fraction of T_1 that is sequential
- $F_p = 1 - F_s$: Fraction of T_1 that is parallelizable

$$\begin{aligned} T_p &= T_1 F_s + T_1 F_p / p \\ S_p &= \frac{1}{F_s + F_p / p} & E_p &= \frac{1}{p F_s + F_p} \\ \lim_{p \rightarrow \infty} S_p &= \frac{1}{F_s} \end{aligned}$$

Application to parallel sum ⇔

The first `sum_to` takes n/p operations. Assuming factor is 2, there is one operation for each of the $\log_2(p)$ subsequent `sum_to`.

$$\begin{aligned}
 T_1 &= n \\
 T_p &= n/p + \log_2(p) \\
 S_p &= \frac{1}{1/p + \log_2(p)/n} & E_p &= \frac{1}{1 + p \log_2(p)/n}
 \end{aligned}$$

► How to get $1/F_s = \lim_{p \rightarrow \infty} S_p$?



Activating project at '~/work/LINMA2710/LINMA2710/Lectures'



`biblio =`

► `CitationBibliography("/home/runner/work/LINMA2710/LINMA2710/Lectures/references.bib", Alp`

① Loading bibliography from '~/home/runner/work/LINMA2710/LINMA2710/Lectures/references.bib'...

① Loading completed.

0.2

1 `BenchmarkTools.DEFAULT_PARAMETERS.seconds = 0.2`