# Conjugated Gradient Solver

## Shi-Lin WANG

SCIPER: 294925
URL: https://c4science.ch/diffusion/9010/

- Final Project -

2020, Spring

# Content

Consider minimizing a quadratic function,

Consider minimizing a quadratic function,

$$f(x) = \frac{1}{2}x^T A x - b^T x + c$$

Consider minimizing a quadratic function,

$$f(x) = \frac{1}{2}x^T A x - b^T x + c$$

- Conjugated gradient solver (CG)
  - Residual & $A$-orthogonal vector
  - Fewer steps during searching

Algorithm of CG:

- $x_0$

- $p_0 = r_0 = b - Ax_0$

- For $k = 0, 1, 2, \ldots n$ :
    1. $x_{k+1} = x_k + \alpha_k \, p_k$     where     $\alpha_k = (r_k^T r_k)/(p_k^T A p_k)$
    2. $r_{k+1} = b - Ax_{k+1} = r_k - Ap_k$
    3. if $r_{k+1} \leq \epsilon$ then break
    4. $p_{k+1} = r_{k+1} + \beta_k p_k$     where     $\beta_k = (r_{k+1}^T r_{k+1})/(r_k^T r_k)$
    5. $k = k + 1$

- End

Algorithm of CG:

- $x_0$
- $p_0 = r_0 = b - Ax_0$
- For $k = 0, 1, 2, \ldots n$ :
  1. $x_{k+1} = x_k + \alpha_k \, p_k$    where    $\alpha_k = (r_k^T r_k)/(p_k^T A p_k)$
  2. $r_{k+1} = b - Ax_{k+1} = r_k - Ap_k$
  3. if $r_{k+1} \leq \epsilon$ then break
  4. $p_{k+1} = r_{k+1} + \beta_k p_k$    where    $\beta_k = (r_{k+1}^T r_{k+1})/(r_k^T r_k)$
  5. $k = k + 1$
- End

Let's see how does computation/storage scales with problem size.

Consider problem size as $n$ and the number of total iterations as $\kappa$ ,

Consider problem size as $n$ and the number of total iterations as $\kappa$ ,

| Operation | Representation | Flops |
|---|---|---|
| Matrix-times-vector | $c = A \times a$ | $n^2$ |
| Vector-plus-scalar-times-vector | $c = a + \alpha \times b$ | $2n$ |
| Vector addition | $c = a + b$ | n |
| Inner product | $\alpha = b \cdot c$ | $n$ |
| Division | $\alpha = \alpha / \beta$ | 1 |
| Assuming A($n \times n$), a(n), b(n), c(n), $\alpha$(scalar), $\beta$(scalar) | | |

Table: Summary of operations vs. number of flops

Consider problem size as $n$ and the number of total iterations as $\kappa$ ,

| Operation | Representation | Flops |
|---|---|---|
| Matrix-times-vector | $c = A \times a$ | $n^2$ |
| Vector-plus-scalar-times-vector | $c = a + \alpha \times b$ | $2n$ |
| Vector addition | $c = a + b$ | n |
| Inner product | $\alpha = b \cdot c$ | $n$ |
| Division | $\alpha = \alpha/\beta$ | 1 |
| Assuming A(n×n), a(n), b(n), c(n), $\alpha$(scalar), $\beta$(scalar) | | |

Table: Summary of operations vs. number of flops

- The whole program has the total number of flops,

$$Flops = (n + n^2) + \kappa \left( (n^2) + 3(2n) + 3(n) + 2(1) \right)$$

- Thus, **time complexity** is obtained,

- Thus, **time complexity** is obtained,

$$\mathcal{O}(\kappa n^2)$$

- Thus, **time complexity** is obtained,

$$\mathcal{O}(\kappa n^2)$$

- With similar method, **space complexity** can also be calculated as,

- Thus, **time complexity** is obtained,

$$\mathcal{O}(\kappa n^2)$$

- With similar method, **space complexity** can also be calculated as,

$$\mathcal{O}(n^2)$$

- Thus, **time complexity** is obtained,

$$\mathcal{O}(\kappa n^2)$$

- With similar method, **space complexity** can also be calculated as,

$$\mathcal{O}(n^2)$$

- Based on simple observation, it shows both complexities scale in the fashion of $n^2$. It seems that the dominating parts are where matrix operation occurs.

In order to check it, a simple profiling is done by performace analyzaing tool **perf** with the problem size set to $n = 1e4$.

In order to check it, a simple profiling is done by performace analyzaing tool **perf** with the problem size set to $n = 1e4$.

Indeed,

- Around 99.75% of execution counts are spent on function cblas_dgemv.
- It does the matrix-vector-multiplication and the vector addition with double precision.

In order to check it, a simple profiling is done by performace analyzaing tool **perf** with the problem size set to $n = 1e4$.

Indeed,

- Around 99.75% of execution counts are spent on function cblas_dgemv.
- It does the matrix-vector-multiplication and the vector addition with double precision.

In the following, instead of the whole program, for simplicity, I will then focus on cblas_dgemv only.

Roofline model provides an easy way to assess the performance of program running on certain machine.

Roofline model provides an easy way to assess the performance of program running on certain machine.

$$P(AI) = min \begin{cases} \pi \\ AI \times \beta \end{cases}$$

Roofline model provides an easy way to assess the performance of program running on certain machine.

$$P(AI) = min \begin{cases} \pi \\ AI \times \beta \end{cases}$$

- Given $\pi$, $\beta$, by observing $AI$, one can estimate whether program is bounded by bandwidth or CPU on that machine.

Roofline model provides an easy way to assess the performance of program running on certain machine.

$$P(AI) = min \left\{ \begin{array}{l} \pi \\ AI \times \beta \end{array} \right.$$

- Given $\pi$, $\beta$, by observing $AI$, one can estimate whether program is bounded by bandwidth or CPU on that machine.
- I test CG on the debug node of EPFL cluster FIDIS.
  - Processor: **Xeon E5-2690 v4 processors**
  - Memory: **DDR4-2400**

- Peak band width $\beta = 71.53$ Gbytes/sec.

- Peak band width $\beta = 71.53$ Gbytes/sec.
- To estimate peak performance $\pi$, I use the equation,

$$
\begin{aligned}
\pi \quad \approx \quad & \text{Number of FP ports} \times \\
& \text{Flops/cycle} \times \\
& \text{Vector length} \times \\
& \text{Frequency} \times \\
& \text{Number of cores} \\
= \quad & 2 \times 16 \times 32 \times 2.6 \times 1 \quad \text{Gflops/sec}
\end{aligned}
$$

- Peak band width $\beta = 71.53$ Gbytes/sec.
- To estimate peak performance $\pi$, I use the equation,

$$
\begin{aligned}
\pi \quad \approx \quad & \text{Number of FP ports} \times \\
& \text{Flops/cycle} \times \\
& \text{Vector length} \times \\
& \text{Frequency} \times \\
& \text{Number of cores} \\
= \quad & 2 \times 16 \times 32 \times 2.6 \times 1 \quad \text{Gflops/sec}
\end{aligned}
$$

- To estimate $AI$, I use the equation,

$$
AI = \frac{W}{Q}
$$

- For simplicity, I consider cblas_dgemv as overall process.

- For simplicity, I consider cblas_dgemv as overall process.

$$AI = (2n^2)/(n^2 + 3n) = 1.9994 \approx 2 \quad \text{flops/bytes}$$

- Thus, the Roofline plot is obtained,

- For simplicity, I consider cblas_dgemv as overall process.

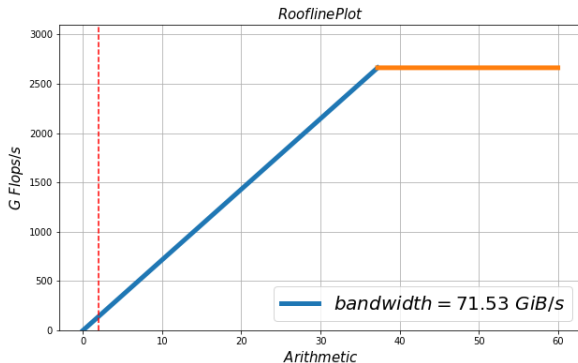$$AI = (2n^2)/(n^2 + 3n) = 1.9994 \approx 2 \quad \text{flops/bytes}$$

- Thus, the Roofline plot is obtained,

Based on Roofline plot, although the program is already bounded by bandwidth, one can still parallelize it to boost performance.

Based on Roofline plot, although the program is already bounded by bandwidth, one can still parallelize it to boost performance.

- Given a fixed problem size, Amdahl's law [1] is usually used to estimate the speed-up $\mathcal{S}$,

$$\mathcal{S}(p) \leq \frac{t_s + t_p}{t_s + t_p/p} = \frac{1}{s + a/p}$$

[1] MATH-454 PHPC lecture-01

Based on Roofline plot, although the program is already bounded by bandwidth, one can still parallelize it to boost performance.

- Given a fixed problem size, Amdahl's law [1] is usually used to estimate the speed-up $\mathcal{S}$,

$$\mathcal{S}(p) \leq \frac{t_s + t_p}{t_s + t_p/p} = \frac{1}{s + a/p}$$

- To take the problem size into consideration, Gustafson's law [1] is then used,

$$\mathcal{S}(p) \leq \frac{t_s + t_p}{t_s + t_p/p} = \frac{s + a \times n}{s + a \times \frac{n}{p}}$$

---

[1] MATH-454 PHPC lecture-01

Assuming $a \propto$ flops(cblas_dgemv), I can calculate $\quad$ perf result: 99.75%

$$(s, a) = (0.07\%,\ 99.93\%)$$

Therefore, given $n = 1e4$,

Assuming $a \propto \text{flops}(\text{cblas\_dgemv})$, I can calculate [perf result: 99.75%]

$$(s, a) = (0.07\%, \ 99.93\%)$$
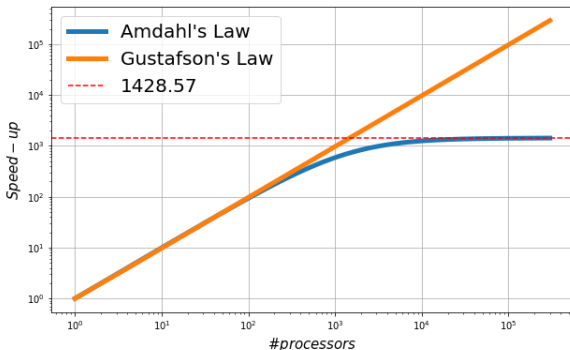
Therefore, given $n = 1e4$,



Then, how to implement parallelization for CG? MPI & CUDA

Consider a matrix-vector multiplication $y = A \times p$,

$$y = A \times p = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} \times \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{bmatrix}$$
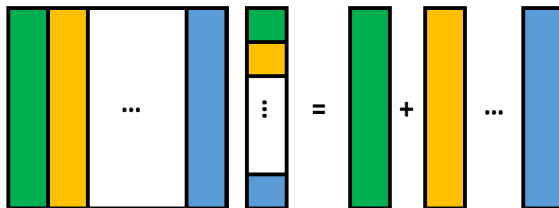
Consider a matrix-vector multiplication $y = A \times p$,

$$y = A \times p = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} \times \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{bmatrix}$$

For example, given 2 processors, one can distribute tasks as follow,

$$y = \sum_{i=1}^{2} A_i \times p_i = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \\ a_{41} & a_{42} \end{bmatrix} \times \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} + \begin{bmatrix} a_{13} & a_{14} \\ a_{23} & a_{24} \\ a_{33} & a_{34} \\ a_{43} & a_{44} \end{bmatrix} \times \begin{bmatrix} p_3 \\ p_4 \end{bmatrix}$$

A simple way to visualize such method is as follow,



Let's take a look at part of my code.

In the file *cg.c*,

```
 1   ...
 2   intv_A = m*n/size;
 3   intv_x = n/size;
 4   ...
 5
 6   while ( k < n ){
 7       /* 1. Root broadcasts p to all */
 8       /* 2. Calculate yi=Ai*pi */
 9       /* 3. Sum yi to root */
10       MPI_Bcast(p, n, MPI_DOUBLE, 0, MPI_COMM_WORLD);
11       cblas_dgemv(CblasColMajor,CblasNoTrans, m, n/size,\
12                   al , &A[intv_A*rank], m, &p[intv_x*rank],\
13                   incx, be, App, incy);
14       MPI_Reduce(App, Ap, m, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
15
16       /* 4. Root does serial task */
17       if (rank==0) {...}
18   }
19   ...
```

For CUDA, unlike MPI, the parallelization is massive due to a large number of processors of GPU (threads).

For CUDA, unlike MPI, the parallelization is massive due to a large number of processors of GPU (threads).

To parallelize cblas_dgemv,

- **NVIDIA** has provided a parallelized function for cblas_dgemv called cublas_Dgemv.
- However, the number of threads is already optimized for given matrix and vector size.
- Therefore, I will use it as benchmark and propose a naive implementation which does the same job.

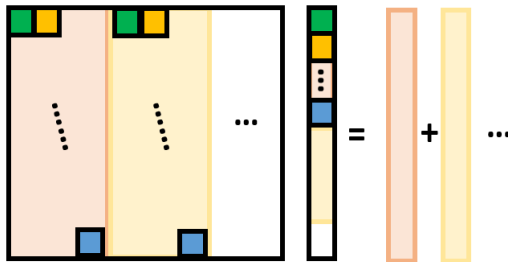Consider a matrix-vector multiplication $y = A \times p$,

$$y = A \times p = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} \times \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{bmatrix}$$

Consider a matrix-vector multiplication $y = A \times p$,

$$y = A \times p = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} \times \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{bmatrix}$$

Given 16 threads, one can distribute task to each thread as,

$$y_i = Aij \times p_j$$

To prevent interference during update, it is important to use atomic operation, which guarantees that a race condition won't occur.

A simple way to visualize such method is as follow,
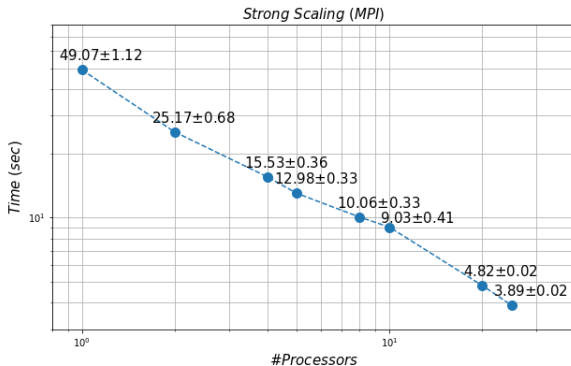


Let's take a look at part of my code.

In the file cg.cu for naive implementation,

```
1   __global__ void kernel_mv(const double * const A, const double * const X, double * const AA, int m) {
2       int i,j,k,l;
3       int t = blockDim.x*gridDim.x;
4       int id = threadIdx.x + blockIdx.x*blockDim.x;
5       double temp;
6
7       for (i=0;i<(m*m/t);i++){
8           j=(id+i*t);
9           k=(id+i*t)/m;
10          l=(id+i*t)%m;
11          if (j<m*m){
12              temp=A[j]*X[k];
13              atomic_Add(&AA[l], temp);
14          }
15      }
16  }
17  ...
18  kernel_mv<<<10000,1000>>>(d_A, d_x, d_Ap, m);
19  ...
```
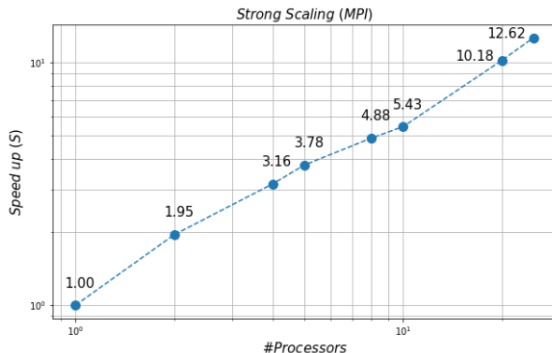
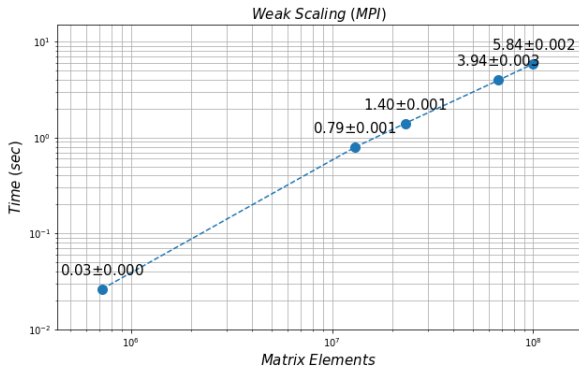- Strong scaling, $p = \{1, 2, 4, 5, 8, 10, 20, 25\}$ & $n = 10000$



Strong Scaling (MPI)

49.07±1.12

25.17±0.68

15.53±0.36
12.98±0.33

10.06±0.33
9.03±0.41

4.82±0.02
3.89±0.02

Time (sec)

$10^1$

$10^0$ $10^1$

#Processors

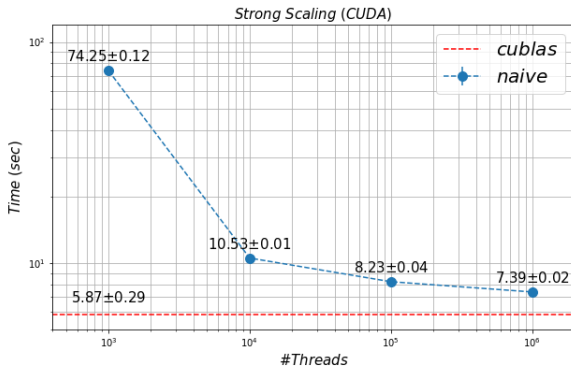- Speed up, $p = \{1, 2, 4, 5, 8, 10, 20, 25\}$ & $n = 10000$  Gustafson: 1
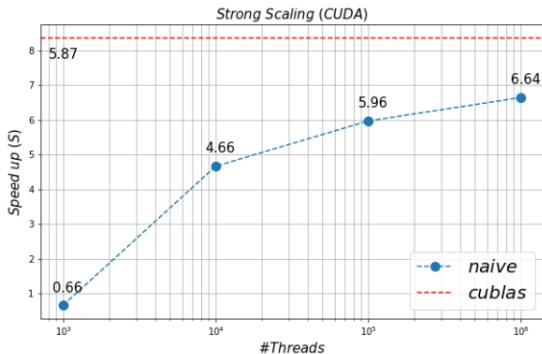


Strong Scaling (MPI)

- Weak scaling, $p = 16$ & $n = \{848, 3600, 4800, 8192, 10000\}$

- Strong scaling, $p = \{10^3, 10^4, 10^5, 10^6\}$ & $n = 10000$

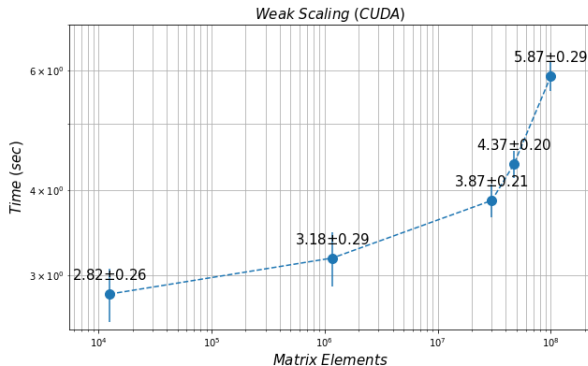- Speed up, $p = \{10^3, 10^4, 10^5, 10^6\}$ & $n = 10000$

- Weak scaling, $p = 10^6$ & $n = \{848, 3600, 4800, 8192, 10000\}$



*Weak Scaling (CUDA)*

- Indeed, the program is speed-up by paralleization. Speed-up of MPI goes linearly while speed-up of CUDA does not.

## Conclusion

- Indeed, the program is speed-up by paralleization. Speed-up of MPI goes linearly while speed-up of CUDA does not.
- The measured speed-up doesn't follow the prediction. It may arise from data transfer, processors' communication and cache size.

## Conclusion

- Indeed, the program is speed-up by parallelization. Speed-up of MPI goes linearly while speed-up of CUDA does not.

- The measured speed-up doesn't follow the prediction. It may arise from data transfer, processors' communication and cache size.

- Nevertheless, the linear behavior of MPI still provides a means of extrapolation for one to predict the budget for a bigger project.

**Thank you!**