

Problems Prediction

1> Operational Intensity & Roofline Model

① # FLOP : floating point operation, double (unless noted)

$$② \text{Arithmetic (Operational Intensity)} I = \frac{W}{Q} = \frac{\# \text{Flops}}{\# \text{bytes}}$$

针对算法的，对一个 algorithm's metric，算法的本征量，与计算平台无关 [float = 4 byte, double = 8 byte]
32 bit

W: amount of work / i.e. floating point operations required

Q: memory transfer / i.e. access from DRAM to lowest level cache

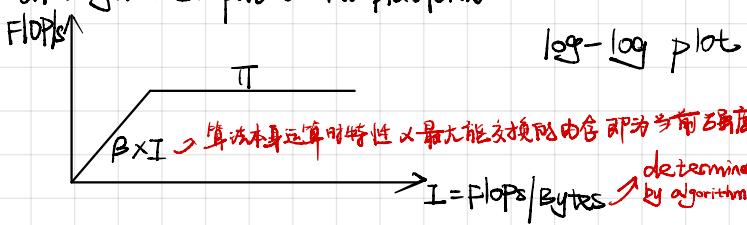
W: 直接计算访问 + 一米 / 的数且 → 配上循环

Q: Depends on cache size {No cache (reading directly from memory) | 每一次数据访问都要被 counted
perfect cache: 数据只用读取 + 写入一次}

每次循环都只读取相应 行列 → 考虑写入算不算访问 ⇒ hw01

★ 计算时注意数据类型，结果可用 $OI = O(N)$ ⇒ tutorial 02 slides

③ Roofline model: estimate how efficient is our code based on a given computational platform



④ CPU peak performance of this platform: $\Pi = \text{FLOPs}_{\text{peak}} / \text{cores}$

⑤ RAM memory bandwidth ($f = \text{bytes/sec}$) memory bound

↓? How to estimate Π and f based on the hardware specification?

$$\Pi = \text{clock speed} \times \text{arithmetic units} \times \text{Vec} \times \text{cores}$$

↓ vectorization { AVx16 bit → 8 double/8 float
AVx4 bit → 4 double/8 float }

$$f = \text{clock speed} \times 6.5 \text{ bytes/cycle} \rightarrow \text{直接用就行}$$

2> Amdahl's law

Metrics for parallel computing

→ t_{S1} : execute time on a single core

→ $t_{\text{L}}(N)$: execute time on processor out of N cores

$$\text{Speed up} : S(N_{\text{proc}}) = \frac{t_{\text{S1}}}{\min_{\text{samples}} \max_i t_{ij}(N_{\text{proc}})} \times \frac{N_{\text{proc}}}{N_{\text{L}}}$$

N_{proc} , N_{L} → problem size → { linear speedup N_{proc}
super linear speedup $S > N_{\text{proc}}$ }

Strong scaling & Weak scaling

→ 是我们对并行计算性能分析的一种手段

→ Strong scaling: 问题尺度不变，增加并行单元，效率能变化，用 Speed up 来度量

→ Weak scaling: 问题尺度随并行单元数等比变化。理想上来说时间不变。用 $\frac{t_{\text{S1}}}{t_{\text{L}}(N_{\text{proc}})}$ 度量效率衰减

Amdahl's law

$$S_s(N) = \frac{t_{\text{S1}}}{f t_{\text{S1}} + (1-f) \frac{t_{\text{S1}}}{N}} = \frac{N}{1 + (N-1)f} \Rightarrow \text{problem-size fixed}$$

⇒ partial parallel / serial

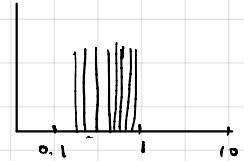
Problem: Operational Intensity & Roofline model

hw01: 把读入和写出分离开来算 注意，若 $C[n \times n]$ ，在读入时可能进行了 n^2 次，写出时则可能只有 n 次。

hw01: Be careful when calculating $\Pi \rightarrow$ should we consider vectorization and vectorization \rightarrow float/double
单位: $W \rightarrow \text{flops}$ $Q \rightarrow \text{bytes}$

Mock Exam: ① Draw Roofline model

② 计算 O.I. 时把 N 记得带上



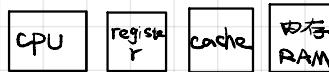
③ 给出的 iteration scheme 和最后代码不同，我们可以将一些常数的组合假设成常数，从而减少计算强度 ⇒ 在 Mock exam，我们没有将 constant 的载入算作一次 memory access。同理，要从具体代码出发，求出多少次 memory access.

Problem: Amdahl's law

HW1: 当考虑 communication cost 0.001n + D

$$S_{\text{nd}} = \frac{t_{\text{S1}}}{t_{\text{S1}} + f t_s + (1-f) \frac{t_s}{N} + 0.001n + D} \quad \begin{matrix} \text{communication is} \\ \text{related to} \\ \text{serial computation} \end{matrix}$$

3> Cache



1. Cache Optimization → Linear Algebra Operations

effects of storing and using data in different ways on the efficiency of the code. PLHWW1 HW3 ⇒ 我们在 FD 中

① Hadamard product of 2 matrix $\text{advance} \hookrightarrow$ 更新时可 row-wise 更 row-wise 间

$$\begin{matrix} * \rightarrow x \rightarrow x \\ n \rightarrow x \rightarrow x \\ C \end{matrix} = \begin{matrix} x \\ x \\ x \\ x \\ A \end{matrix} \times \begin{matrix} B \\ B \\ B \\ B \\ B \end{matrix} \quad \begin{matrix} N \times N & N \times N & N \times N \\ \text{for } i \rightarrow N, \text{for } j \rightarrow N \end{matrix}$$

row-wise: $C[i \cdot N + j] = A[i \cdot N + j] + B[i \cdot N + j]$

column-wise: $C[i \cdot N + j] =$

↳ 由于 index 不变 循环角度换顺序

• 两种不同形式的数据访问。第一种更快，因为数据在 cache 存储是 row-wise，这一次访问和下次访问都可以在 cache 中找到，不涉及 memory ⇒ cache 的频率更快

• 常用 std::swap(), avoiding copying memory

Dwarfs

<1> Diffusion Equation and Finite Difference

Logic of code

总浓度

- Define a struct Diagnostics → Store C and t

→ 之后画全局扩散曲线图用

- Define a Diffusion2D class / construct

member function ↓ Diffusion2D { \rightarrow 生成尺寸 dimension → consider ghost information }

Diffusion2D { 初始化函数, 与类同名, 初始化并赋值给 private }

↓ advance() { 依据迭代格式进行迭代, often loop, 2D
就是双loop Euler + Finite Difference [can be parallel by openMP \rightarrow collapse()] and MPI local-N] }

↓ Compute diagnostic() { 计算总浓度, 写一个循环即可 [涉及 reduction 操作] 通常也会利用 write-diagnostic
 \rightarrow do not compute ghost node }

↓ initialization() { 根据 initial values 初始值 [也可涉及到循环, 用 OpenMP collapse()] }

↓ Compute histogram() { 画出浓度分布图。 \rightarrow OpenMP 里需注意用 max and min reduce 会得到 global max and global min. Then, using MPI_Allreduce }

• In MPI

achieved in the advance()

- We need to exchange the ghost information. After that, we can use Diffusion2D.advance()

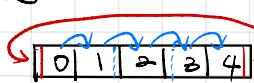
- main() 函数中实现 MPI_Init(), rank, size.

By the way D L N 这些系统参数在 Compile 时指定。

- For local-N, achieve it in the Initialization function

number of grid points of this process

- 初始化函数最后一步有 initialize(), 会根据 rank 不同进行 initialization, initialize 增加了 ghost 之和的值



advance()

信息交换时会出现经典乒乓通信, 有 sender or receiver 分开 \rightarrow 出 MPI (II) more on blocking point-to-point communication

Send-receive : send of a message to one destination.

receive a message from another process.



$W[0] \rightarrow$ Send previous rank

$W[1] \rightarrow$ receive previous rank

$W[N] \rightarrow$ Send

$W[N+1] \rightarrow$ receive

↓ mock exam coding MPI

HW5 Solution

① 固定边界, 均用的是 MPI_PROC_NIL

② 处理两种边界 { 给 pre-rank 和 next-rank 值 }

③ 同期性边界 { next = (rank + 1) % size; }

prev = (rank + size - 1) % size

size

index \rightarrow $ceil(Nt)$ + 1, N

如果是高维 ghost { ghost } ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

ghost { ghost }

<2> Linear algebra operation with A splitting MPI

核心问题是 local coordinate 与 global coordinate

$$\boxed{\quad} \quad N \times N$$

index index

```

int myN = N
size
    
```

\Rightarrow for rank 0 \rightarrow local = global
 for rank 1 \rightarrow local + myN * rank
 $(i_0 + myN * rank, j_0) = (i, j)$

之后关于 ij 操作全用 vector 处理, include

\downarrow 与 (i, j) 有关的 initialization
 \downarrow 接着, 当我们分 rank 计算完后时, 为了不影响
 对这个 vector 之后的操作, 我们可以用

\downarrow MPI - Allgather 重新回来 (就比如你先 norm
 操作)

\downarrow 每经过一次 MA 操作分一次, 再拼一次进行操作

See Mock exam

HW6 给出更加复杂的过程, 进行的不仅仅是 row, 而
 是 block 的操作, 下面给出一点分析

$$U^{n+1} = U^n + AU^n$$

3	4	5
rank 0	1	2

$$\xrightarrow{A}$$

• 由于整个过程不涉及对向量所有成分求解, 则
 我们不用 gather, 直接全分块走到底.

将 A 也根据所对应 rank 分块, A 根据所对应
 rank 离散化, 整个程序一条路走到底。最后写入
 的时候再拼回原来的 U^n .

<3> Brownian Motion Based on OpenMP

HW2

① Calculate time based on OpenMP

double GetWtime() { return omp_get_wtime(); }

② 会涉及到 Histogram 计算, 注意用之前提到两种策略

③ #include <omp.h>

<4> Dimensionality reduction with PCA

HW3 Q2

这里 PCA 通过 covariance method 得到, 其本质就是 $X^T X$
 再做 Eigen decomposition. Actually, we can do SVD
 based on X to help us achieve PCA.

Covariance Matrix : $C = \frac{1}{N-1} X^T X$ \rightarrow Here, data is
 a matrix

<1> Transpose Data \leftrightarrow compute mean (every dimension)

<2> Compute std. $s = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$ (HW3 Q2 提 N-1)

<4> Standardization / Normalization of Data

Transformation of the data to zero mean unit variance

$$\frac{Data[i] - \text{mean}}{s}$$

* 注意数据分 dimension, 每个 dimension 依次求

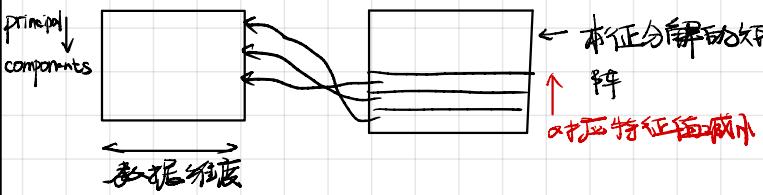
<5> 当我们 standardization 后矩阵后, 便可以对其进行求 covariance Matrix.

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \frac{1}{N-1} \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

\uparrow Due to its symmetry, 整个过程是 3 个 loop 嵌套。

<6> Compute the eigenvalues and eigen vectors.
 这里使用了 LAPACK.

<7> Extract the principal components & eigenvalues



<5> Monte Carlo

motivation: numerical Quadrature for high dimension problem.

Using Composite Trapezoidal Rule, we found

$$E = \mathcal{O}(M^{-1/d}) \quad M = n^d \Rightarrow \# \text{ point evaluation}$$

With higher dimension, the order of convergence decreases

\hookrightarrow Curse of dimensionality. $\xrightarrow{\text{possible solution}} \text{Monte Carlo}$

introduce probability

• stochastic variable \rightarrow described by PDF

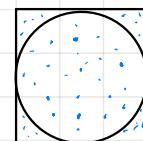


\downarrow Expected Value

即在这样一个随机可能处的范围内, 我们对其期望是多少
 \hookrightarrow variance, 可能随机的结果相对于我们期望之间是否离散

\downarrow Sampling: uniform distribution $\sim U(a, b)$

$$x = a + (b-a)\alpha$$



$$\rightarrow A = \int_{-1}^1 \int_{-1}^1 f(x,y) dx dy \quad f(x,y) = \int_0^1$$

$\mathcal{O}(M^{-1/d}) \rightarrow$ 误差相依赖点变化

\downarrow 理解, If we want to solve $\int_a^b f(x) dx$

Assuming stochastic variable $X \sim U(a, b)$ $p(x) = \frac{1}{b-a}$

$$E[f(X)] = \frac{\int_a^b f(x) dx}{b-a} \Rightarrow b-a E[f(x)] = \int_a^b f(x) dx$$

\hookrightarrow convert integral into expected value

Assuming a variable $Y = f(X)$

生成 N 个 $[a, b]$ 上随机数, 代入 $f(x)$ 得到生成的 N 个随机
 数 Y , 则 高概情况是 $Expected = \frac{\sum f(x)}{N} \Rightarrow \int_a^b f(x) dx = \frac{b-a}{N} f(x)$

HW4 (\Rightarrow) based on MPI $t_1 = MPI_time()$

MPI

Common code:

```
mpic++ -O3 xxx.cpp -o xxx
```

```
mpexec -n 4 ./xxx
```

```
#include <mpi.h>
```

```
MPI_Init argc, argv
```

```
int rank, size;
```

```
MPI_Comm_rank(MPI_COMM_WORLD, &rank);
```

```
MPI_Comm_size(MPI_COMM_WORLD, &size);
```

```
MPI_Finalize();
```

→ group communication

① Gather 这种操作的结果是 recubuf 被拼成了一个很长的结果集，从 rank 0 到 rank n-1 顺序为主

1. SPMD programming model P6

2. Point-to-Point communication P8

3. MPI_Send MPI_Recv P19 → 例子 P21 - P23

4. MPI_Send communication modes P4 通常用计语句类的 rank
车辆助力执行

↓
Example: data exchange → dead lock P25 - P28

5. MPI 常见手段 Splits

Monte Carlo
n_local_size = double(N)/procs;
n_start = rank * n_local_size;
n_end = rank == (procs - 1) ? n_start +
很远
n_local_size



Diffusion local_N = N / size;

if (rank == size - 1) local_N += N % size;
否则 C.resize (local_N+2) * (N+2), 0.0)

6. Blocking collective P31

- Reduce • AllReduce • Broadcast • Gather • AllGather
- Scatter • Barrier • AllToAll • Scan • ExScan
- Reduce_Scatter

7. Communication Cost: Bandwidth and Latency P2

8. More on blocking point to point

↳ unknown size message MPI_Probe P5

↳ Cycle communication often in diffusion P6

↳ S 交换

 | MPI_Sendrecv P7

9. Message passing protocols P8 - P12

↳ eager rendezvous

10. Non-blocking point to point communication

↳ MPI_Isend MPI_Irecv P10 { blocking non-blocking with P14
send modes P16 }

↓ 建立 MPI_Request 用下面方法管理

↳ Request management

P17 { MPI_Wait → { MPI_Waitall
MPI_Fence } MPI_Waitany
MPI_Waitsome } test multiple request P22

↳ Example: 1D Finite Difference P7 - P11

11. Non-blocking collective communication P23 - P38 (examples
+ 1D diffusing)

12. Performance metrics P40

① 提到 communication overhead 影响 ② 提到 Timing MPI program

OpenMP

```
#include <omp.h> g++ -fopenmp -o myprog myprog.cpp
```

Syntax Format:

```
#pragma omp construct [clause clause]
```

→ Structured block Unstructured block Par Opi

• the threads synchronize at a barrier at the end of the parallel region, then master thread continues

• const int id = omp_get_thread_num(); P25 parallel for

• #pragma omp critical 下面的代码每次只能一个线程执行
行, 但具体是哪个线程执行.

→ P29 Opi

→ num_threads() if (parallel: i > 2)

→ omp working sharing constructs

Boo { Loop construct P34

Opi { Section Construct P36

Single Construct P37

Both have an implicit barrier

at the end unless nowait

→ omp_get_thread_num(); Omp_get_num_threads(); P31

→ Loop scheduling (load balance) P38

Thread count P24

不同于 MPI, 其并行数是在编译时指定.

Pragma omp parallel for schedule (static, 2)

Pragma omp parallel for schedule (dynamic, 2)

→ #pragma omp for collapse (2) P43

↓ perfectly loop

→ #pragma omp for reduction (t : sum) P46

↓ Sum会有 private copy, no race condition

→ #pragma omp parallel for combined construct P47

→ ★ Data environments P49

. in OMP, data outside a parallel region are shared by default

. global/static variables are shared

→ defaults P51

→ private P52 → create a private copy of each variable in the list / avoid race / not initialized

→ initialized { copy to first private P53

update after lastprivate P54

↓ 有 race

P1 - P9

→ Recall Data environment { critical sections P38 Pigs locks P55
Atomsics P4 P6 barriels P16

↓ P18 一开始有系统自带 P33 Synchronization, P15 critical, atomsics 又讲了一遍。

Synchronization constructs ensure consistent access to memory address that is shared among a team of threads.

{ master P1 barrier P22 → implicit barrier P33 } { false sharing P55
Library routines Environment variables

13. Communication Topologies and Groups

① Communicators P5 ② Example: domain decomposition P6
communication/computation ratio

③ Create MPI cartesian Topology P10 → finding neighbors P15 → 有用的 P16

14. MPI Vector Datatype

① contiguous vs strided datatype P18 ② Strided datatype P20 → example P1

③ multidimensional arrays P23

15. MPI Struct Datatype

① riemann sums P13 P7 - P28 → 引出问题 P28 - P34

② 使用 MPI_Pack P25 - P36 → 不清楚, 难懂 P37

③ general approach MPI_Type_create_struct P28 - P40 → 例子 P1 - P43