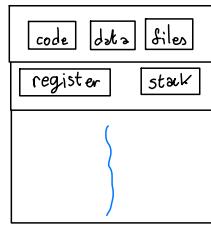


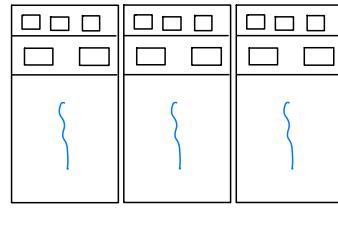
Open MP : for parallelism in a multiprocessor computer by creating multiple threads

History of OpenMP

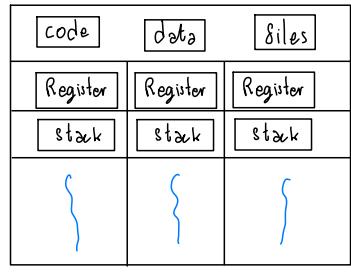
developed & accepted in 1990s as a standard which consists of compiler with library routine & environment var.
! use thread-based shared memory model



Single thread process



Multi Processing



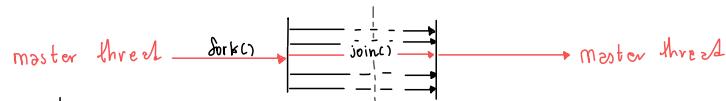
multithread Processing

fork / join

Since OpenMP uses multithread processing approach → the process can fork / join threads

Fork : master thread creates parallel child threads * run code after fork() line

Join : complete the thread process return to parent thread



#include <omp.h>

int main() {

 int var1, var2, var3;
 <sequential code>

 #pragma omp [directive] private (var1, var2) shared (var3) num_threads ()

{

 <parallel code>

} ← thread will join here & auto add barrier for synchronization

compile: gcc -fopenmp -o < > <.c>

what function to use
ex. parallel, for

if not declared = global, is declared in parallel region
= local

! '{' must in new line

int tid = omp_get_thread_num(); — return ID of each thread
int nthreads = omp_get_max_threads(); — return no. of threads created

Work can be shared Manually

```
# pragma omp parallel private (i) num_threads (CT)
{
    int part_size = N/T;
    int tid = omp_get_thread_num();
    int START = tid * part_size;
    int END = (tid + 1) * part_size - 1;
    ....
}
```

use Schedule to control work distribution

ex... for schedule(<type>)
..... for schedule(<type>, <chunk size>)

• static : split equally using round robin

• dynamic : adjust size based on thread speed (default=1)

• guided : initialize big chunk then decrease exponentially

or use Directive Functions

Pragma OMP for ... "every thread do some part"
for(i=0; ...)

! simple for loop (++, --, +=, -=)

Note if #pragma omp parallel #pragma omp for can combine to #pragma omp parallel for

Pragma OMP Sections

{
 # pragma omp section
 <code>
 # pragma omp section
 <code>
}

execute separately by threads
give compiler a hint that sections can run in parallel

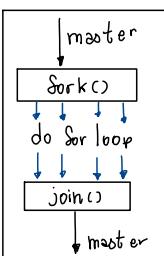
Pragma OMP parallel Sections ensures parallelism

{
 # pragma omp section
 <code>
}

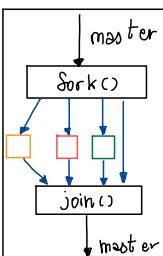
Pragma OMP single → execute only 1

Pragma OMP master → execute only master thread (0)

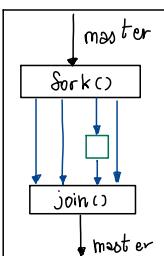
omp for



omp sections



omp single



Private Variable

even var declare before #pragma ...

private: create private non initialized copy for each thread

first private: create private initialized copy for each thread
declare before the #pragma ...

last private: set value = last value from the parallel region

Reduction Clause

use to combine local copies of variable in different thread

use with (parallel, for, sections)

ex. $\text{sum} = 0$

#pragma omp parallel for reduction(+:sum)

for ($k=0$; $k \leq 100$; $k++$) { $\text{sum} = \text{sum} + \text{func}(k);$ }

operators: $+, *, -, \&, |, ^, \&&, ||, \max, \min$

Race condition

#1 write after write

$a \leftarrow R_1$

$a \leftarrow a + 1$

$R_1 \leftarrow a$

$b \leftarrow R_1$

$b \leftarrow b + a$

$R_1 \leftarrow b$

#2 write after read

$a \leftarrow R_1$

$b \leftarrow a + 1$

$R_1 \leftarrow b$

$c \leftarrow R_1$

$d \leftarrow c + 1$

$R_1 \leftarrow d$

solve by Synchronization which is specifying a critical region to avoid race condition or using barrier to wait for slower thread

OMP Critical: only allow 1 thread to execute

#pragma omp critical [name]

<code here>

'not allow same name execute together
(optional)

Note: beware of using too many synchronization

≈ serialize

OMP barrier: wait until all thread reaches the barrier

#pragma omp barrier

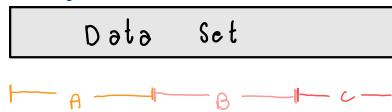
MPI Programming

what to consider when designing parallel program

- **degree of parallelism**: how many simultaneous process
- **initialization**: what is required to run
- **work distribution**: who does what
- **Data / IO access**: where to access / work part
- **Communication**: sharing information
- **Synchronization**: what to do after it is done
- **Finalization**: data collection between each thread

Work distribution

Type 1 : Data decomposition



divide into pieces for repeated operation → ex. for loop

! no overlapping → consider dependency

Type 2 : functional decomposition → divide the computational steps ex.

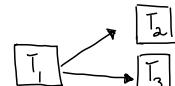
calculate A → MUL B → Add C

Dependency Analysis : find what can execute together

↳ Bernstein's Condition

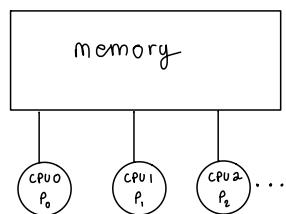
1. input of P_1 not in output of $P_2 \rightarrow I_{P_1} \cap O_{P_2} = \emptyset$
2. input of P_2 not in output of $P_1 \rightarrow I_{P_2} \cap O_{P_1} = \emptyset$
3. no overlapping in output of P_1 & $P_2 \rightarrow O_{P_1} \cap O_{P_2} = \emptyset$

↳ task graphs : use to show the dependency
; node = task, edge = dependency



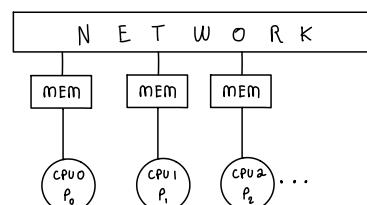
Parallel Programming Model

#1 Shared Memory



All tasks access the same memory space
(threads, openMP, fork())

#2 Message Passing



uses local memory → to share data must transport it explicitly
(MPI, MPICH)

MPI

Message Passing Interface

- + free implementation
- + support multiple languages
- + support multiple architectures (parallel pc, distributed system)

Provide 2 things

- method to create separate process
- method of sending & receiving msg.

Creating MPI Programs

Compiling → \$ mpicc -o <filename> <filename>.c

Running → \$ mpirun -np <n of process> <executable name> <arg....>

Function	Description
int MPI_Init(int *argc, char **argv)	Initialize MPI
int MPI_Finalize()	Exit MPI
int MPI_Comm_size(MPI_Comm comm, int *size)	Determine number of processes within a comm
int MPI_Comm_rank(MPI_Comm comm, int *rank)	Determine process rank within a comm
int MPI_Send(void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)	Send a message
int MPI_Recv(void *buf, int count, MPI_Datatype datatype, int src, int tag, MPI_Comm comm, MPI_Status *status)	Receive a message

- * the process execute concurrently
- * msg passing send data between process
- * msg have tag to tell the type

know who are you [0, process-1]

} Communicator

NOTE: Id start at 0

Communicator : MPI_COMM_WORLD declare the entity where communication between process can occur

initialization

```
main() {
    int rank, size;
    MPI_Init(&argc, &argv) // no need to be 1st statement
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    // code
    MPI_Finalize();
}
```

! the variable is independent

↑ "SMPD" single program multiple data

Sending & Receiving data

- | | |
|--------------------|---------------------|
| • what to send | • what to receive |
| • where to send | • where (from whom) |
| • How many to send | • How many |

MPI Send

MPI_Send (<pointer to address of what to send>, <int amount>, <MPI_Datatype data type>, <int destination rank id>, <int tag>, <MPI_Comm>, <MPI_Status *status>);

All process will be blocked until data transfer is done

MPI_Receive

MPI_Recv (<pointer to address where to store>, <int amount>, <MPI_Datatype data type>, <int source rank ID>, <int tag>, <MPI_Comm>, <MPI_Status *status>);

* status can be used in many way : error code, msg. tag, sender rank, etc.

MPI datatype

MPI_CHAR
MPI_SIGNED_CHAR
MPI_UNSIGNED_CHAR
MPI_SHORT
MPI_UNSIGNED_SHORT
MPI_INT
MPI_UNSIGNED
MPI_LONG
MPI_UNSIGNED_LONG
MPI_FLOAT
MPI_DOUBLE
MPI_LONG_DOUBLE

C datatype

signed char
signed char
unsigned char
signed short
unsigned short
signed int
unsigned int
signed long
unsigned long
float
double
long double

1D Array distribution = send & A , length

2D Array distribution = send & A[0][0] , length

(-) 1 free master process that wait for response

Collective Communication

- * 1 function but call by every process
- * no msg. tag * sometime better than send/reciv

`MPI_Bcast()` : broadcast value from root to all other process - Recv is not necessary!

`MPI_Gather()` : gather value from set of process

`MPI_Scatter()` : scatter buffer into small part and give it to group of process

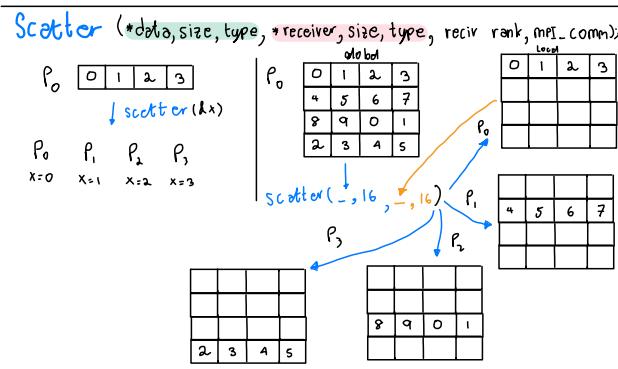
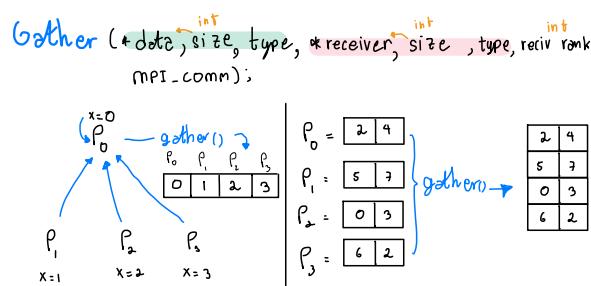
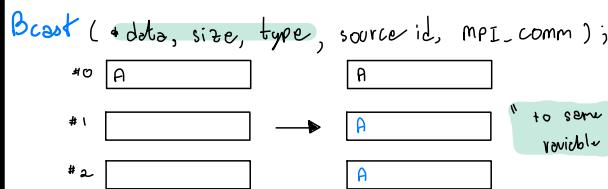
`MPI_Alltoall()` : exchange data of all process (every one have all value)

`MPI_Allgather()` : gather and give it to all process

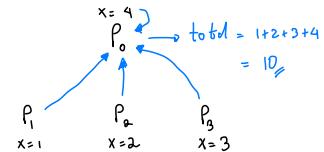
`MPI_Reduce()` : combine value from multiple process to one

`MPI_Reduce_Scatter()` : combine value and give it to all process

`MPI_Barrier()` : wait until every process reach this checkpoint



`Reduce (*data, +result, size, datatype, reduce function, destination process id, MPI_Comm);`



Ex. implement in integral function

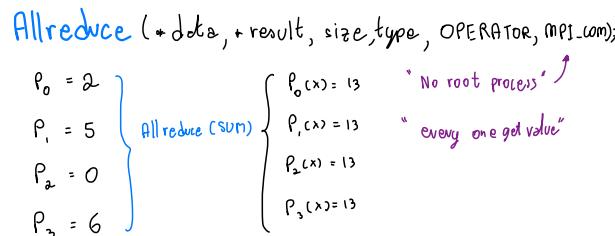
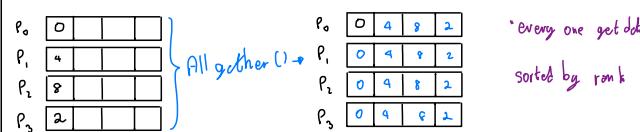
Array Reduction



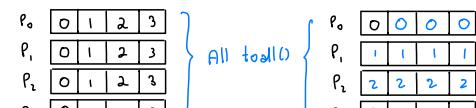
`MPI_Reduce (Min) \rightarrow [0, 2] ! compare 1 address of each process`

`MPI_Reduce (SUM) \rightarrow [13, 16] root=0`

`Allgather (*data, size, type, *receiver, size, type, MPI_Comm);`



All to All



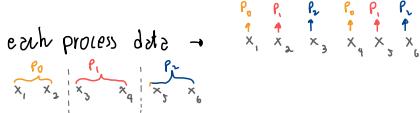
`Barrier (MPI_Comm);` wait until every process complete

Performance Evaluation

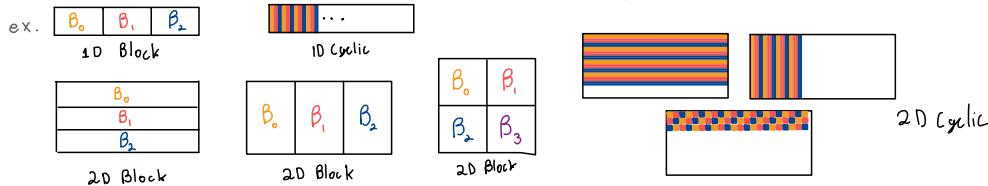
goals

- Partition the work equally → "load balance" reduce synchronization time
(wait time for another process)
- Reduce Overhead ~ penalty or setting up before execution ex. sending data

Data decomposition

Interleaved: cyclic: use mod function to provide each process data → 

Block: give process a continuous set of data



What to do if N ≠ Process

for interleaved → for (i=pid ; i < n ; i+=p)

for Block → calculate range & owner of that range

$$\left. \begin{array}{l} \text{START} \rightarrow \lfloor \frac{\text{pid} \times N}{P} \rfloor \\ \text{END} \rightarrow \lfloor \frac{(pid+1) \times N}{P} \rfloor - 1 \end{array} \right\} \text{OWNER} \rightarrow \lfloor \frac{p(\text{elemeID} + 1) - 1}{N} \rfloor$$

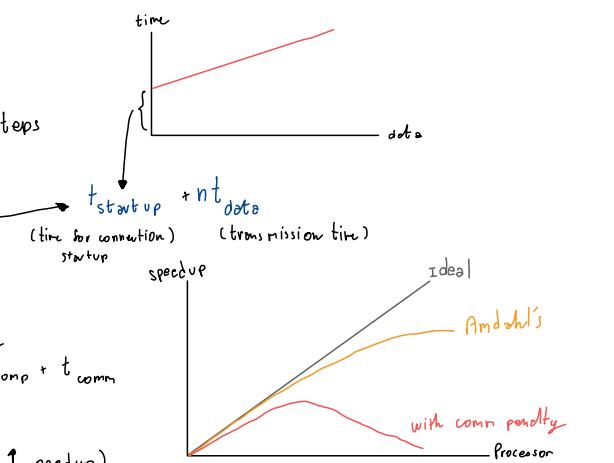
Performance Concepts

Execution time

Sequential (t_s): calculate by counting steps

Parallel (t_p): equal to slowest process

$$t_p = t_{\text{compute}} + t_{\text{communication}}$$



$$\text{Speed Up } (S(p)) = \frac{T_{\text{sequential}}}{T_{\text{parallel}}} = \frac{t_s}{t_{\text{comp}} + t_{\text{comm}}}$$

Ideally speedup should increase linearly ($\uparrow \text{Process} = \uparrow \text{speedup}$)

But program sometimes have sequential part → δ -factor "Amdahl's Law"

$$S(p) = \frac{t_s}{(f\delta t_s) + (1-\delta)t_s} + \frac{(1-\delta)t_s}{P} = \frac{1}{\delta} + \frac{(1-\delta)}{P} ; \text{ Upper limit} = \frac{1}{\delta}$$

"Amdahl's Effect" → large dataset will eliminate time penalty ≈ linear

Gustafson's Law: calculate from ideal speedup then deduct serial penalty of idle process when runs sequential code

$$\text{Scaled Speedup} \rightarrow S_g(p) = \frac{\text{Assume constant}}{st_p + p(1-\delta)t_p} = p + (1-p)\delta$$

Amdahl: sequential → parallel estimation "constant problem size scaling"
Gustafson: parallel → sequential estimation "fine constrained scaling"

Computation / Communication ratio: show the amount of communication required = $\frac{t_{\text{comp}}}{t_{\text{comm}}}$

Efficiency: show how effective of parallelism = $\frac{S(p)}{P} = \frac{t_s}{P \cdot t_p}$ } show improvement per processor

Timer in MPI

double time = MPI_Wtime() * in seconds

CUDA Programming

Parallel programming using Processor Array Architecture

GPU: graphic Processing Unit → a processor array for graphic rendering

- * calculate faster than CPU

- * have higher memory bandwidth ($\approx 10 \times$)

- * mainly consist of ALU support for simple & specific instructions

- * communicate via PCI-e bus (Peripheral Component Interconnect Express)

#1 high throughput global memory → "device memory"

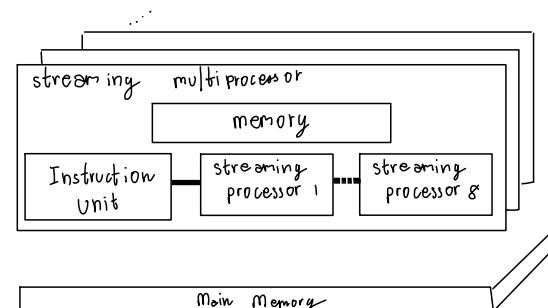
#2 set of Streaming Multiprocessor

Streaming Multiprocessor contain : work the same instruction

- instruction Unit

- group of streaming processor [SIMD]

- shared local memory



GPU Architecture

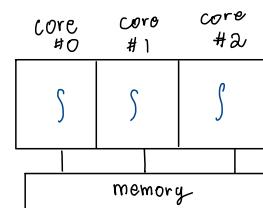
GP GPU

General Purpose graphic processing unit

- general computation using GPU and graphic API in application other than 3D graphic
- using in parallel computation (large data array, SIMD parallelism, low latency of floating point calculation)

CUDA : compute unified device architecture

- a general purpose programming model by NVIDIA
- extend ANSI C programming language
- low learning curve



Shared Memory Programming

- Process have multi thread, flow of program, which can be construct by fork() but here we focus "CUDA threads" which will run in parallel

requires

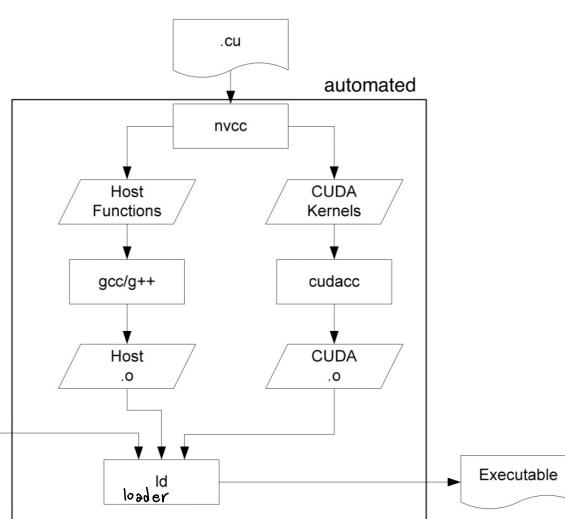
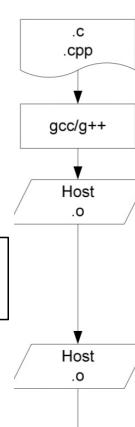
1. CUDA Driver & SDK
2. editor with CUDA source file (.cu)

Compile

```
> nvcc -o cout < code.cus
```

Run

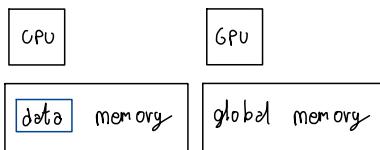
```
> ./out
```



Writing Code

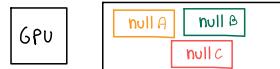
Initial

• data loaded to computer memory
 int a[N], b[N], c[N];
 //Init data of a,b,c



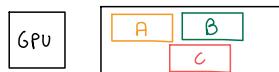
STEP #1 Allocate GPU Memory

```
int *gpuA, *gpuB, *gpuC; // pointer to GPU
                           mem location
cudaMalloc((void**) &gpuA, size),
cudaMalloc((void**) &gpuB, size),
cudaMalloc((void**) &gpuC, size),
```



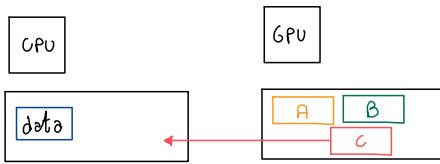
STEP #2 copy data to GPU memory

```
cudaMemcpy(gpuA, a, size, cudaMemcpyHostToDevice);
cudaMemcpy(gpuB, b, size, cudaMemcpyHostToDevice);
cudaMemcpy(gpuC, c, size, cudaMemcpyHostToDevice);
```



STEP #4 transfer data back to CPU

```
cudaMemcpy(c, gpuC, size, cudaMemcpyDeviceToHost);
```



if no. of threads < data

```
for (i = threadIdx.x; i < n; i += T) {
    C[i] = A[i] + B[i];
}
```

STEP #3 execute using kernel code

with CUDA syntax (no. of threads & physical structure)

```
myKernel << n, m >> (arg1, ...);
```

↑ function grid size n block, m thread

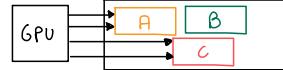
N = 256; Kernel routine CUDA specifier

```
_global_ void vecAdd(int *A, int *B, int *C) {
```

int i = threadIdx.x; ast threads to add

```
    C[i] = A[i] + B[i]; // 1 index simultaneously
}
```

```
vecAdd << 1, N >> (gpuA, gpuB, gpuC);
```



STEP #5 free GPU memory

```
cudaFree(gpuA);
cudaFree(gpuB);
cudaFree(gpuC);
```

Finding Prime Number → Sieve of Eratosthenes

- 1st prime is 2
- delete all number of stored prime tuples
- next number will be prime & repeat

sequential

List = 2, 3, ..., n

K = 2

while $K^2 < n$

- all no. of k tuples between $K^2 - n$ is marked
- k = smallest un marked

unmarked = prime

Make it Parallel

② marking the non prime from known prime

↓
data decomposition by break to sub array

$$\text{first} = \left\lfloor \frac{\text{pid}(n)}{P} \right\rfloor, \text{last} = \left\lfloor \frac{(P\text{id}+1)n}{P} \right\rfloor - 1$$

$$\text{owner} = \left\lfloor \frac{(P(i+1) - 1)n}{P} \right\rfloor$$

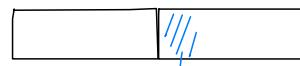
Note

- we sieve to \sqrt{n}
- P should smaller than \sqrt{n}
- so it contain all prime

↪ no need to reduce data to P_0

Parallel the marking

P_0 broadcast prime to other process to mark



Local = 0
global = 3

for (i = first ; i ≤ last ; i++) {

↓ if $i \cdot K = 0 \rightarrow \text{mark } i$

! slow
* is faster

if $\text{first} \cdot K = 0 \rightarrow \text{mark(first, first + nK)}$

until it reach last value of its block

Final Pseudo code

1. create list [2, n]
2. K=2
3. repeat until $K^2 > n$

• mark all K between $K^2 - n$ ✓ multiple threads

• K = smallest un marked → P_0 only

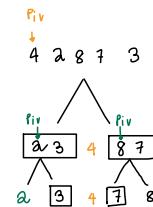
• P_0 broadcast K

4. unmarked is prime

5. reduce the result to prime only

Parallel Sorting

Optimal speed $O(n \log n)$ of sequential
 \therefore ideal for parallel = $O(\log n)$



Sequential Quick Sort

1. choose Pivot value
2. split the data into 2 group
 Lower than Pivot, Higher than pivot
 ! pivot value is in the right position
- repeat

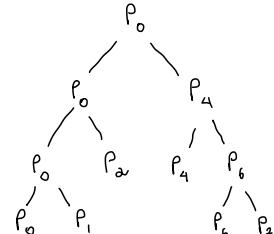
Average: $O(n \log n)$

Worst: $O(n^2)$

unbalanced at every step

Parallel quick sort

Use Process tree to handle the splitted data section
 simple \rightarrow use $P=n$ \rightarrow impossible



but $n \log n$ is not guarantee because splitted array might unbalance

\hookrightarrow choose better pivot "Median of 3": pick 3 numbers from data and choose the middle one

Computation time (t_{comp})

$$t_{comp} = n + \frac{n}{2} + \frac{n}{4} + \frac{n}{8} + \dots \approx 2n \quad \leftarrow \text{measure using computation step}$$

Communication time (t_{comm})

! no. of process = 2^k

$$\begin{array}{l} \stackrel{1^{st}}{P_0 \rightarrow P_4} = t_{start} + (n/2) t_{data} \\ \stackrel{2^{nd}}{P_0 \rightarrow P_2, P_4 \rightarrow P_6} = t_{start} + (n/4) t_{data} \\ \vdots \end{array} \quad \left\{ \text{total} = (\log n) t_{start} + n t_{data} \right.$$

! to begin next step must wait all process to finish current step

! picked pivot might not be the "True Median"

ex. initial : data split to all process

Pivot	75	91	15	64	21	P ₀
	50	12	47	72	65	P ₁
	83	66	13	0	88	P ₂
	20	40	89	86	85	P ₃

P₀ broadcast Pivot

STEP 1: pair P₀-P₂, P₁-P₃ & exchange low/high

≤ 75	75, 15, 64, 21, 66, 19, 0	P ₀
> 75	50, 12, 47, 72, 65, 80, 40	P ₁
	83 91 88	P ₂
	89 86 85	P ₃

! execution time
 = last process finish
 ! pivot effect balancing

STEP 2: pick pivot of P₀, P₂ & exchange

≤ 21	15, 21, 13, 12, 0, 20	P ₀
> 21	50, 47, 72, 65, 40, 75, 64, 66	P ₁
≤ 88	83 88 86 85	P ₂
> 88	89 91	P ₃

step 3: each processor sort its own value

0, 12, 13, 15, 20, 21	P ₀
40 47 50 64 65 66 72 75	P ₁
83 85 86 88	P ₂
89 91	P ₃

Merge Sort

Sequential: $O(n \log n)$

* Divide & Conquer

- ; split unsorted list into half until only 1 then merge into sorted by each branch
- Parallel: $O(p \log n) \rightarrow \log n$ for splitting
 $\rightarrow \log n$ for merging

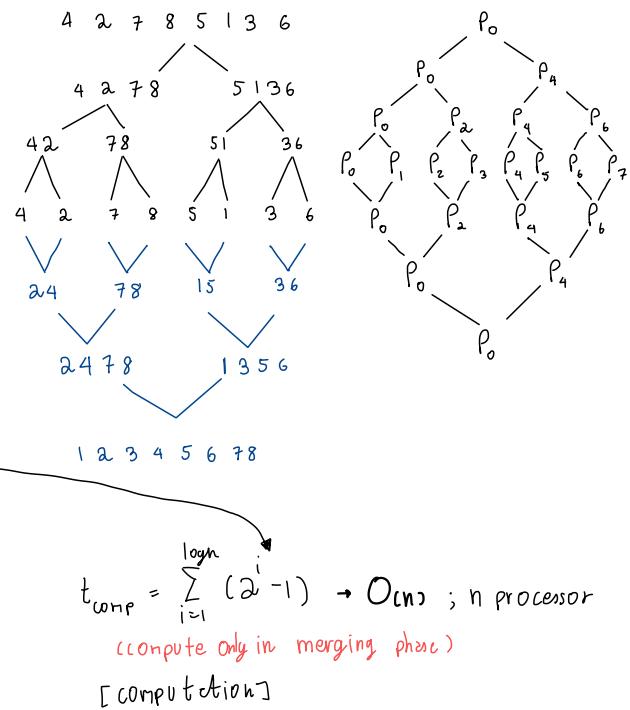
Merging Step

- use 2 index to point data in each list
- put the lower data in another list
- total = $M+N-1$ comparison

$$C[i] \leftarrow A[i], B[i]$$

COMPLEXITY $\sim 2 \log n$ steps

$$\begin{aligned} &\hookrightarrow 2 \log n (t_{\text{startup}} + n t_{\text{data}}) \\ &\quad [\text{communication}] \end{aligned}$$



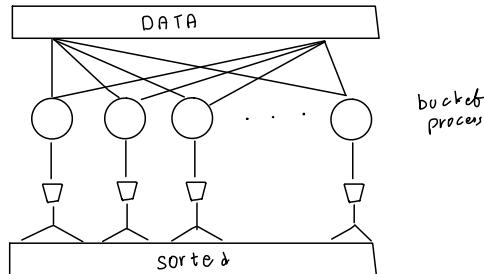
Bucket Sort

each "bucket" have a n acceptance range

- 1st Broadcast data to every bucket process
- 2nd Pick only value in range of the bucket
- 3rd Sort locally

COMPLEXITY

$$\begin{aligned} \text{broadcast} \rightarrow t_{\text{comm}} &= t_{\text{start}} + n t_{\text{data}} \\ \text{merge } (1-1) \rightarrow (p-1)(t_{\text{start}} + n_p t_{\text{data}}) \end{aligned}$$



$$\text{compute} = n + \frac{(n_p)}{\text{bucket filter}} \log(n_p) \quad \text{assume balanced}$$

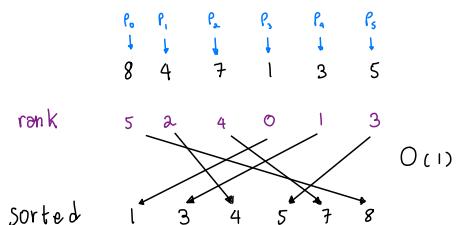
` sorting inside = $(n \log n)$

Rank Sort

calculate rank from counting lower data in data set

rank = position on the sorted array

- 1st broadcast data to every process
- 2nd each process pick value based on process id
- 3rd compare all data rank + 1 if there is data lower
- 4th insert it to result [rank]



Sequential = $O(n^2)$; every one compare

Parallel = $O(n)$

COMPLEXITY

$$\begin{aligned} \text{broadcast} \rightarrow t_{\text{comm}} &= t_{\text{start}} + n t_{\text{data}} \\ \text{merge } (1-1) \rightarrow (p-1)(t_{\text{start}} + t_{\text{data}}) \end{aligned}$$

compute = $n-1$ comparison