

# OPEN MP

MODEL

- More threads than processes possible
- Each thread **unique id**
- **Sync. constructs** to prevent **DATARACE**
- **Shared & private** vars poss
- **SPMD**
- **Fork-join thread model**  
 M thread forks & creates threads  
 when done: threads join M thread



# MSA PASSING MODEL

- **Finite processes with LM**
- **Expl. Comm**
- **No implicit sync.** → only **comm.**
- **Nothing shared**

# MPI

- **Finite processes** in **comm. domain** (more than 1 procs)
- Each process **identified by rank** in domain
- **Ranks successive** (0, ..., p-1)
- **LOCAL DATA**
- **EXPLICIT COMM** → **Reliable & ordered**
- Structure of **com. data** ⊥ to **node/mode**
- **Comm. d.** reflect **physical topology**
- **No comm. cost model**
- 3 basic **comm** models
  - → **Point-to-Point** (diff. modes, local & <sup>compl.</sup> ~~best~~ semantics)
  - → **1-sided** (diff. sync., local <sup>compl.</sup> semantics)
  - **collective ops** ( ~~local~~ <sup>compl.</sup> semantics)

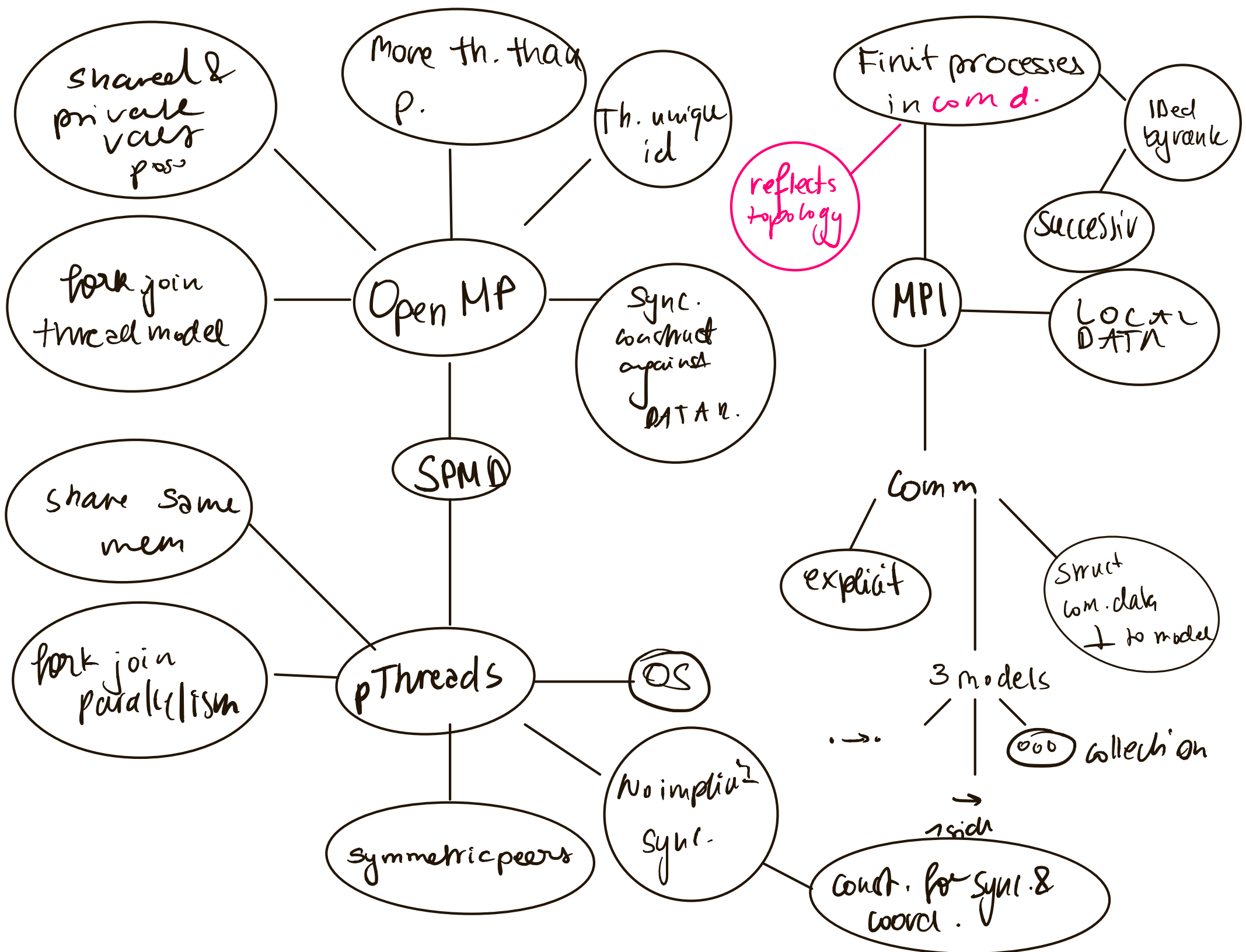
# pThreads

- **Fork-join parallelism**
- → can spawn new threads
- **symmetric peers**
- → threads **can wait** for others to complete
- **SPMD**  
 scheduled by OS
- **No implicit sync.**
- → threads are independent
- **Share same memory**
- **Construct for sync. & coordination**
  - → **MUTEX**
  - → **R/W locks**
  - → **condition vars**

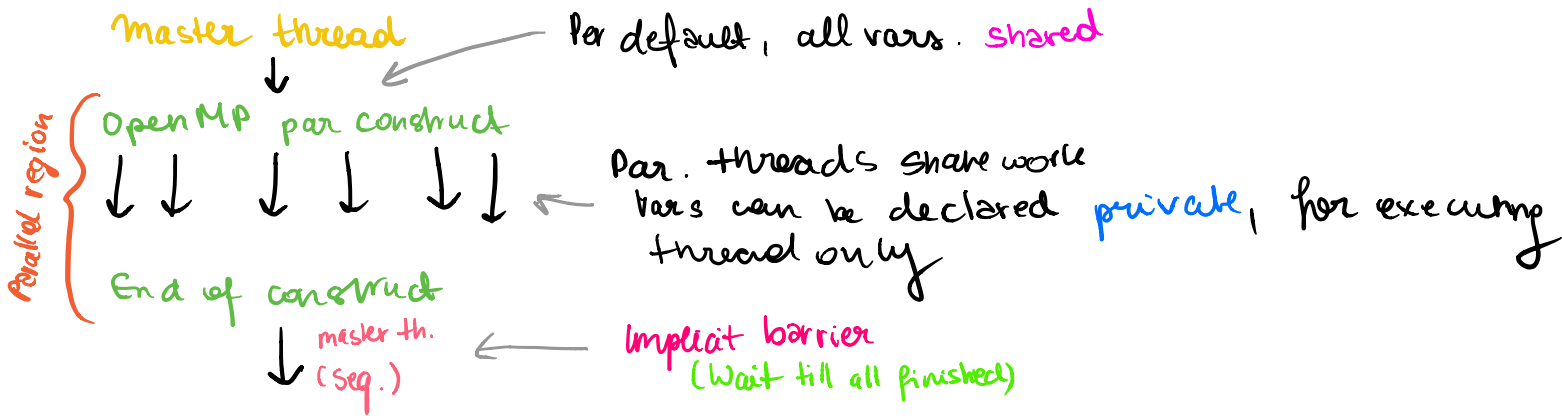
- → **Point-to-Point**
  - **MPI-Send**
  - **MPI\_Recv**
  - **MPI\_Recv**
  - **MPI-Send**
- } MPI sendrecv  
 deadlock free  
 reliable (ordered)  
 non-overtaking

- → **1-sided**
- **MPI-get (MPI-put)** → **passiv**
- **origin** → **target**
- **active**

- **Collective**
- **MPI-Barrier**
- **Bcast (root)**
- **Scatter(v)**
- **(All) gather(v)**
- **(All) Reduce**
- **MPI\_AlltoAll (ex) Scan**



# OPENMP



Data transfer between s. & p. copies var. is transparent, impl. \* can be controlled

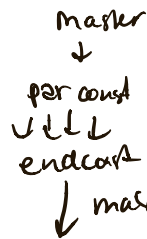
## Parallel region

```
#pragma omp parallel
```

```
{
```

```
...
```

```
}
```



"memory model" ← state

→ memory is consistent by end of const.

↳ updates to shared var visible

Threads asynchron

# MPI Equivalents

All collective func. blocking & not sync

MPI_Send	} = MPI_Sendrecv
MPI_Recv	
MPI_request	
MPI_Irecv	
MPI_Send	
MPI_Wait	

MPI\_Isend } = MPI\_Send  
MPI\_Wait }

MPI\_Recv } = MPI\_Gather  
MPI\_Sendrecv }  
MPI\_Send }

MPI_gather	} = MPI_Allgather (= P <u>Bcast</u> )
MPI_Bcast	
MPI_Bcast	
memcpy	

blocking  
notsync

Can do:

Bcast i  
Bcast j

i ... root

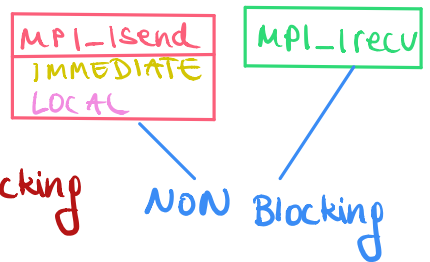
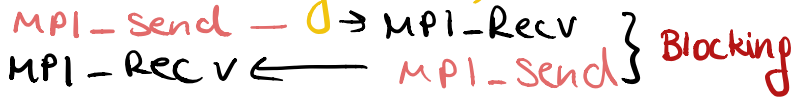
Bcast } i      Bcast } k      gather } j  
gather }

Bcast } i      ISend(i) } j  
recv(j) }      Bcast }



# POINT-TO-POINT

Reliable  
 Deadlock-free  
 Non-overtaking (ordered)

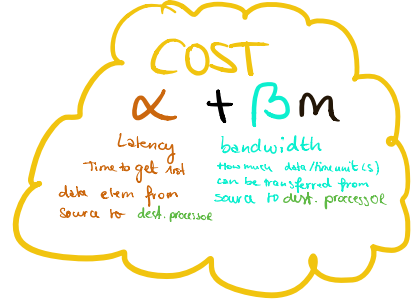


- Successful:  $0 \leq \text{rank} < \text{size}$
- Sender specifies valid rank within comm. & valid (allocated) buffer
  - Receive with matching source rank & tag is eventually posted on same comm.
  - Sent data  $< | =$  to received amount
  - Type sig of sent & received data matches

ERROR: Caught by MPI-Send: error  
 ↳ IF NOT: **DEADLOCK**  
 ↳ OTHERWISE: **MPI\_ERR\_TRUNCATE**  
 || memory corruption

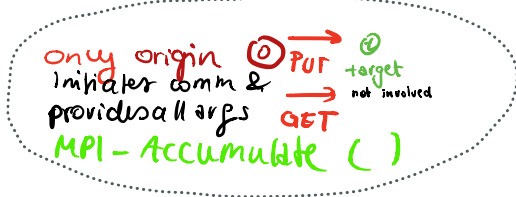
**MPI-Send**

- Short msg: Buffered at Sender, processed later. MPI-Send can return IMMEDIATELY.
- Medium msg: May be buffered locally. Sent when rec. has been posted. MPI-Send can return fast.
- Large msg: Participation of receiving process needed. MPI-Send cannot return before MPI-Recv becomes active.



# → 1-sided

Comm. between 2 processes  
**ONLY 1** Explicitly involved with comm



origin puts/gets data from stand. MPI-buffer (buf, count, type)

local completion semantics

No guarantee that data have arrived before Sync. operation

Comm. window

# Collective Operations

All processes involved

Same root → If diff. roots **DEADLOCK**

#sent = #received

matching type signatures (pairwise same size)

Difficult to impl: perf. correct. safety

# COLLECTIVE OPERATIONS

MPI\_Bcast: Data from root  $\xrightarrow{a}$  all  
not evenly  $\rightarrow$  MPI\_Scatter: Indiv. data root  $\xrightarrow{a}$   
MPI\_Gather: Indiv. data all  $\rightarrow$  root

MPI\_Allgather: Indiv. all  $\rightarrow$  all

MPI\_Reduce: Apply ass. function( $t, T, \dots$ )  
to data from each process  
res at root  
MPI\_Allreduce: res to all  
MPI\_Reduce\_Scatter: res  
scattered/distr.

$P_i$  (non-root)  
 $r$  (sbuf, rbuf,  
... root (comm));  
 $P_j$  (root);  
 $r$  (sbuf, rbuf,  
root, comm);  
only sig.  
for root

MPI\_Scan: Prefix sum (reduction)

MPI\_Exscan: value of process  
not used in  
reduction

MPI\_Barrier: Sync. (waits until  
have reached routine)

△ Caveat = Processor's Point of view

# Performance of Processors

= Nominal processor performance

↳ often measured in FLOPS

(= Floating Point Operations per second)

→ max. num of FLOPS

Depends on 2 factors:  
 - Micro-arch. (how well utilize processor cores)  
 - MEMORY (can be used to supply data needed to keep processor busy)

Performance of 1-core  
 num of instr. per clock cycle  
 \* clock frequency

Performance of multi-core

num of cores \* perf. of 1-core

determined by architecture (SIMD, num of pipelines, p. depth, ...)

c.f. ... num of ticks per sec (GHz)

# PRAM model

= Par. Random Access Machine

→ unrealistic model but useful

→ LOCK STEP (= strictly synchronized)

example:  $a[i] = c[i] + b[i]$   
 $a[i+k] = a[i]$

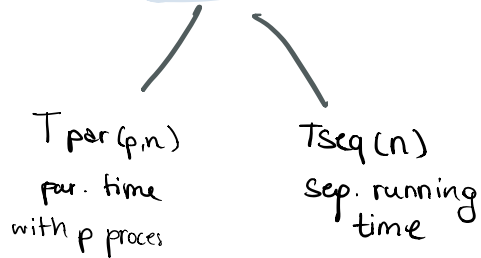
(practically not possible)

easiest model

strongest model

EREW	can't ACCESS memory simult.
CREW	can't WRITE in memory simult.
CRCW	Arbitrary CRCW: 1 write succeeds
	Priority CRCW: priority scheme determining pr.
Common CRCW: conflicting p. must write same value	

# TIME



# SPEED UP

Relative	Absolute	Linear
$\frac{T_{par}(1,n)}{T_{par}(p,n)}$	$\frac{T_{seq}(n)}{T_{par}(p,n)}$	$T_{par}(n) \geq \frac{T_{seq}(n)}{p}$
Compares par alg. against ITSECF	Compares par alg. (seq. alg) against BASELINE	Best possible speed up Argument: par alg. running on p cores can be simulated on a single core no worse than $p * T_{par}(p,n)$
Expresses to WHAT EXTENT p are well EXPLOITED	Expresses HOW well par alg. over the BASELINE IMPROVES	

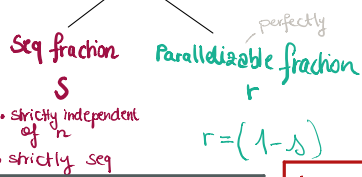
not-linear	Simulation can STILL run faster ⇒ new best known
Super-linear	$W_{par} \neq W_{seq}$ → random. alg search alg.

# PARALLELISM

$$\frac{T_{par}(1,n)}{T_{\infty}(n)}$$

# AMDAHL'S LAW

Assuming an alg. can be divided into



Max speedup by parallelization

$$\frac{1}{s + \frac{1-s}{p}} = \frac{1}{s} \text{ if } p \rightarrow \infty$$

$$T_{par}(n,p) = T_{seq}(n) \cdot \left( s + \frac{r}{p} \right)$$

$$= T_{seq}(n) \cdot \left( s + \frac{1-s}{p} \right)$$

→ Devastating ⇒ Seq. alg. that satisfy A.D  
 ↳ CANNOT be used as a basis of a good paralg.

- ⇒ VICTIMS (of par alg.)
- IN-/OUTPUT
  - Seq. preprocessing
  - Seq. data structures
  - Seq. operations

Good alg.: Don't fall victim of A.L.  
 → seq. part seq. ≠ const. fraction of total work  
 ⇒ BUT depends & decreases with  $n$

# TIME

Linear Speedup	Non-linear Speedup
$T_{par}(p,n) = \frac{O(T(n))}{p + t(n)}$ $t(n) \dots T_{oo}(n)$ non-par-term $T(n) \dots$ parallelizable term	$T_{par}(p,n) = \frac{O(T(n))}{f(p) + t(n)}$ $f(p) \dots f(p) \leq p$ $p(f(p)) \in O(p)$

# SCALED SPEED-UP

→ The faster  $\frac{T(n)}{t(n)}$  converges ⇒ the faster the speedup converges

Speedup as  $n$  increases

$$S_p(n) = \frac{T_{seq}(n)}{T_{par}(n)} = O\left(\frac{T(n)}{T(n)}\right) = O\left(\frac{1}{\frac{1}{p} + \frac{t(n)}{T(n)}}\right) \rightarrow O(p)$$

# SCALABILITY ANALYSIS

Strong Scaling A.	Weak Scaling A.
$n$ is constant IF run time ↓ proportionally to $p$ (Linear speedup) ⇒ Alg. strongly scalable	Work / processor is constant (through bigger $n$ ), IF par. running time is also constant alg. is weakly scalable

# EFFICIENCY

= comparison of par against best possible paralleliz. of Seq

Parallel efficiency ( $E_p$ )

$$E(p,n) = \frac{T_{seq}(n)}{p T_{par}(p,n)} = \frac{S_p(n)}{p}$$

Following holds

- \*  $E_p \leq 1$
- \* If  $E_p = e$  (= const) ⇒ LINEAR SPEEDUP
- \* Cost-opt. alg. have constant  $E_p$

Iso-efficiency Does not change with  $n$

IF alg. does NOT have const efficiency & speedup is independent of  $n$

THEN: calculate  $n$  as function of  $p$

→ Done by calculating  $e$  as a function of  $n$  &  $p$   
 ⇒ Then changing  $f$  to  $n$

# SCALABILITY

$f(p)$

→ Tells how  $n$  should grow as a function of  $p$  to maintain constant efficiency

→ Should not grow faster than the input size

scaling function  $g(p)$   
 ↳ Tells how much  $n$  can at most grow if  $par=const$

→  $f(p)$  should be in  $O(g(p))$

Weak Scaling

IF for (desired) constant efficiency  $e$  there is a slowly growing function  $f(p)$  such that efficiency  $E(p,n) = e$  for  $n \in \Omega(f(p))$

$f(p) \dots$  iso efficiency function

IF average work per processor  $\frac{T_{seq}(n)}{p}$  is const at w the running time of the alg. is const at  $T_{par}(p,n)$

$g(p) = T_{seq}^{-1}(pw)$

Efficiency varies with the num of processors if the problem size remains constant  
 → may be more efficiently solved with fewer processors

# FLYNN'S TAXONOMY

= differentiates machines based on inst

SISD	Seq. Computer
SIMD	Single inst can work on big am. of data ⇒ vector instructions
MIMD	= general PRAM Each processor has own inst. & data
MISD	= Pipeline like system Single data stream passes through dif. stages
SPMD	All processors execute same program

	Found	Processors	Operations
Common CRW	$O(1)$	$n^2$	$O(n^2)$
CRCW PRAM	$O(\log n)$	$\frac{n}{2}$	$O(n)$
		$n \times n \rightarrow 1 \text{ matrix}$	operations
CRCW Pram		$O(n)$	$O(n^3)$

} Max of n num stored in array

# OVERHEAD & LOAD BALANCE

→ per alg often perform more work than BEST seq  
 ⇒ OVERHEAD

- ↳ communication
- ↳ Synchronization
- ↳ Preprocessing

If overheads in bounds <sup>( $O(T_{seq}(n))$ )</sup>  
 → THEN: Alg CAN be WORK-OPT

Granularity (of alg)  
 = communication intervals

coarse	fine
RARE	FREQUENT
communication	communication

= Problem of making sure that time of all processors is about the SAME

$V_{i,j}:$   
 $T_{par}(i,n) \approx T_{par}(j,n)$

## ↳ IMBALANCE

If  $T_{par}(i,n)$  = time processor  $i$  takes ( $0 \leq i \leq p$ )

⇒ THEN: **LOAD IMBALANCE**

$$\max_{0 \leq i \leq p, 0 \leq j \leq p} T_{par}(i,n) - T_{par}(j,n)$$

static	dynamic				
= when amount of work is divided BETWEEN processors	= exchange of work DURING execution of alg				
<table border="1"> <tr> <th>oblivious</th> <th>Adaptive</th> </tr> <tr> <td>= SUBDIVISION of prob ONLY by Input SIZE &amp; STRUCTURE REGARDLESS of input</td> <td>= prob.-indp. INPUT itself taken into consideration</td> </tr> </table>	oblivious	Adaptive	= SUBDIVISION of prob ONLY by Input SIZE & STRUCTURE REGARDLESS of input	= prob.-indp. INPUT itself taken into consideration	
oblivious	Adaptive				
= SUBDIVISION of prob ONLY by Input SIZE & STRUCTURE REGARDLESS of input	= prob.-indp. INPUT itself taken into consideration				

Embarrassingly / Trivially / Pleasantly Parallel

bs where input can be statically distributed & NO interaction is needed

# COST & WORK

Cost	Work				
$p * T_{par}(p, n)$ = what we "pay" for → time of $p$ processor cores occupied by $P$	<table border="1"> <thead> <tr> <th><math>W_{seq}(n)</math></th> <th><math>W_{par}(p, n)</math></th> </tr> </thead> <tbody> <tr> <td>= num of ops carried out by alg. (= "work time")</td> <td>= total work carried out by <math>p</math> cores <small>(excluding anything not related to the work itself)</small></td> </tr> </tbody> </table>	$W_{seq}(n)$	$W_{par}(p, n)$	= num of ops carried out by alg. (= "work time")	= total work carried out by $p$ cores <small>(excluding anything not related to the work itself)</small>
$W_{seq}(n)$	$W_{par}(p, n)$				
= num of ops carried out by alg. (= "work time")	= total work carried out by $p$ cores <small>(excluding anything not related to the work itself)</small>				
<b>Cost opt</b> $\underline{IF}$ cost $p T_{par}(p, n)$ $= O(T_{seq}(n))$ for best known seq alg	<b>Work opt</b> $\underline{IF}$ $W_{par}(p, n)$ $= O(T_{seq}(n))$ for best known seq alg				

alg = ~~cost-opt~~ work-opt  $\Rightarrow$  can have **LINEAR SPEED** for small range of proce



# Par vs Dist vs Concurrent

**Par** M Efficiently utilizing dedicated par. resources to solve **comp** **probs**

- ⊕ Problem solving efficiency
- Deals with: alg., impl. & structure of comp. arch.
- ↳ Issues: many different parad, lang. & inter.
- Interesting problems ⇒ significant interaction between cores
- = Theoretical + Practical + Experimental C.S.

**Dist.** M of making indep., non-det. resources AVAILABLE & COOPERATE toward solving specific prob. complexes

- ⊕ Availability of resources (that aren't ready at hand)

**Concurrent** M of MANAGING & REASONING about **INTERACTING** processes that may / may not progress simult.

- ⊕ Activities that may / may not happen at the same time
- & Reasoning about
- & Establishing correctness

Some prob can be viewed from all perspectives