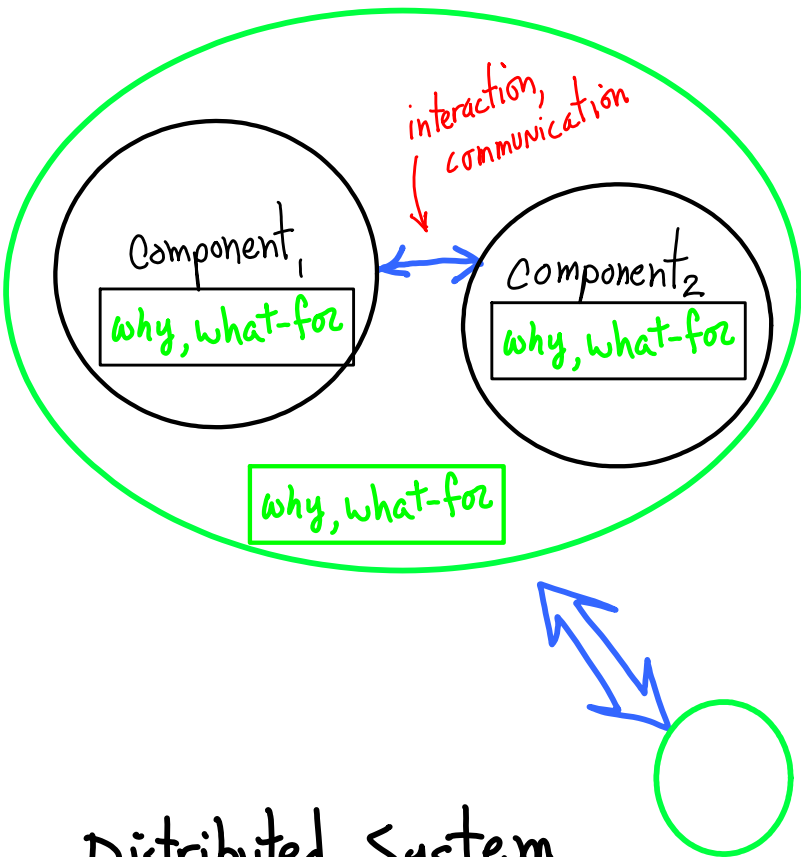
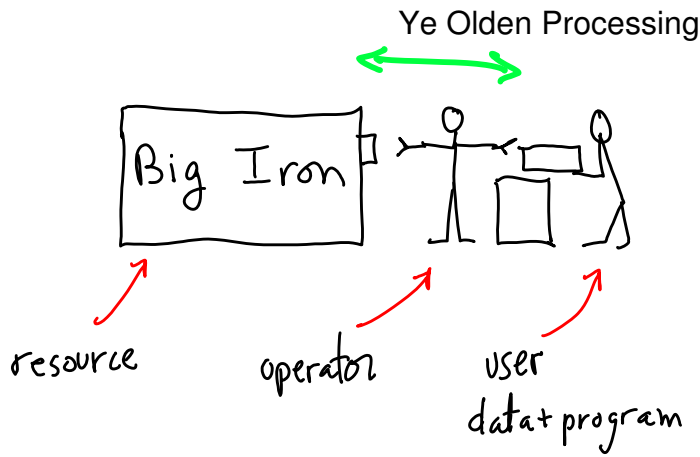


# Parallel Processing

1. Definitions
2. Goals
3. Types



a system, generally

component = internally highly connected, not easily separable, functionally distinct

communication = relatively sparse connections between components

Hierarchy of

- Structure: grouping of components
- Purpose: why, what-for

Surfaces are logical and/or time-like and/or space-like

Boundaries are drawn arbitrarily, by us, for some specific purpose.

## Distributed System

a system w/

- slow communication between components (distance [physical])
- less common dependency, more autonomy (no shared, common dependency on a single component) \*
- many components ≈ peers

- United by common/shared goals
- separated by local clocks  
local memory

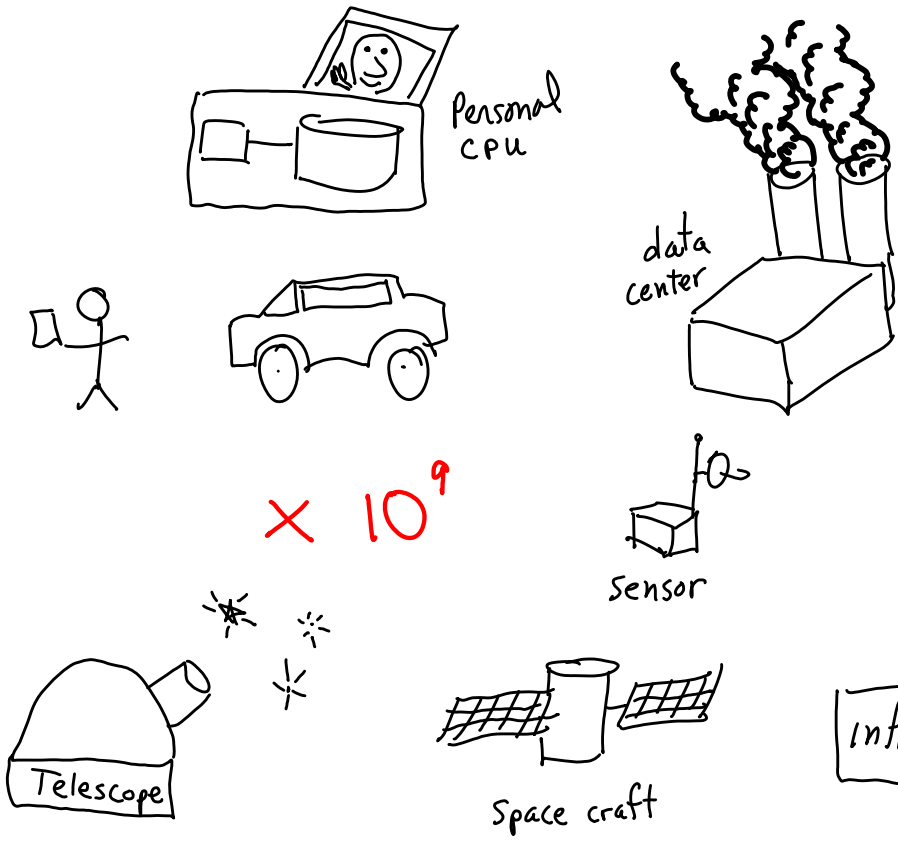
★ [more constant\* vs less constant] \*ignored

Common, but constant: Earth\*  
Common, but not constant: power switch

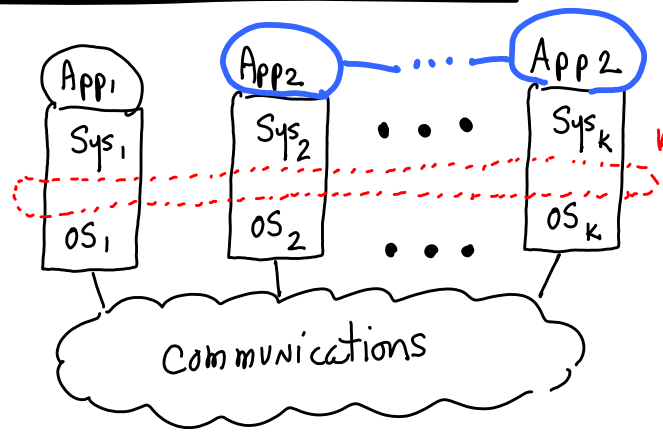
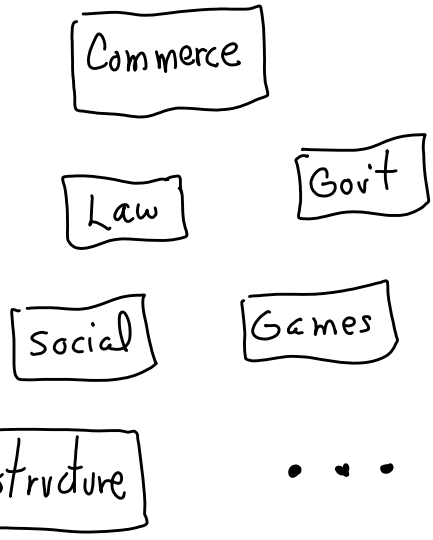
⇒ relative time scales defines the difference

# Systems Today are heterogeneous

Internally and externally



$\times 10^9$



middle ware  
logical interface: hierarchy

One version of a (slightly) larger view of our reality.

Parallelism is everywhere!

## Goals

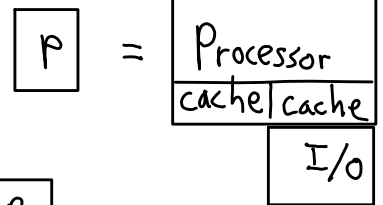
- accessible resources, sharing, efficiency, reliability
- collaboration on shared enterprise - But, with control + security
- Requires (?) - Scalability, modularity, incremental expansion, heterogeneity

# - Transparency (not visible)

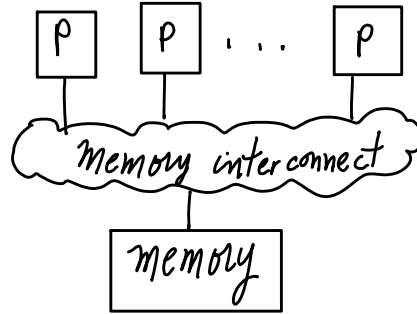
- access: hide data representation differences [types, file names, methods]
- Location: hide latencies, specific details of location/access path
- Migration: hide dynamic reconfiguration
- relocation
- replication: hide caching, manage consistency
- concurrency: hide other's usage
- failure: hide intermittent response

## Parallel Systems

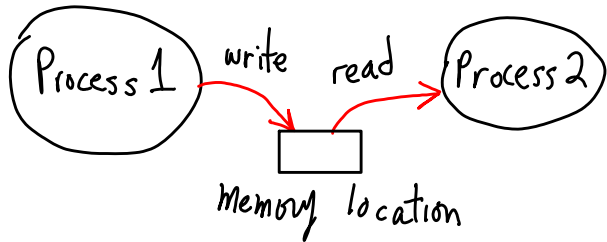
common address space



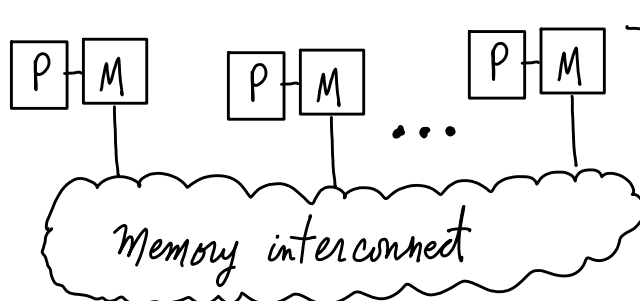
- multi-processor, shared memory (SMP)
- shared interconnect
- Local clocks/global clock



- Uniform access to memory
- UMA



coordination through memory  
(can implement messages on top)



- non-UMA: NUMA

## Classification

Hardware = serial  
Program = sequential

- e.g. single core
- Matrix mult

HW = serial  
SW = concurrent

- e.g. single core
- OS multi-processing

HW = parallel  
SW = sequential

- e.g. multi-core
- matrix mult

HW = parallel  
SW = concurrent

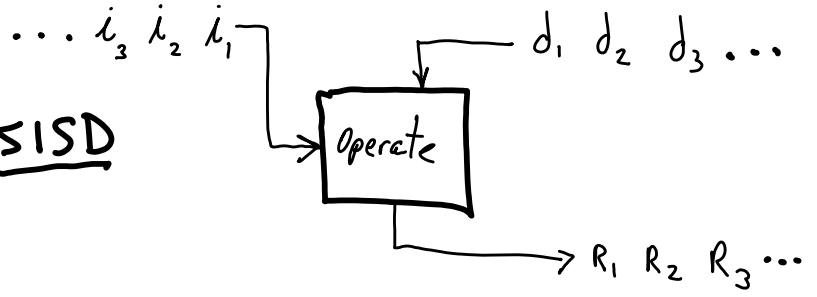
- e.g. multi-core
- OS multi-processing

# Taxonomy

Single Instruction, Single Data SISD

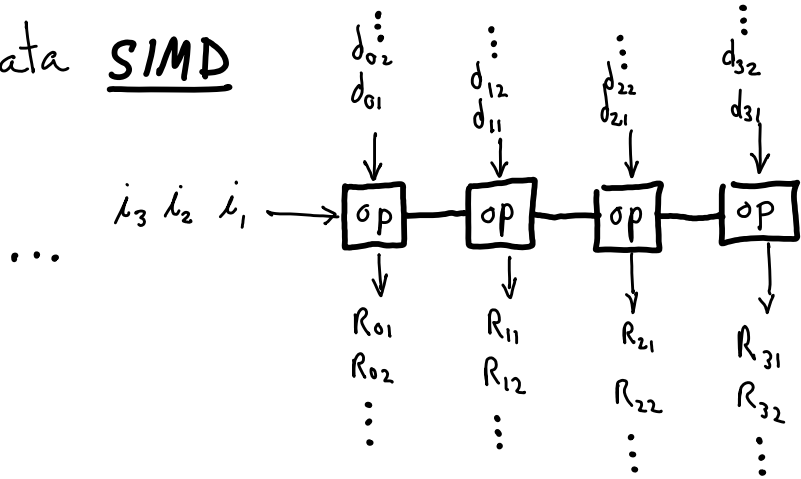
STREAM                      STREAM

(serial)



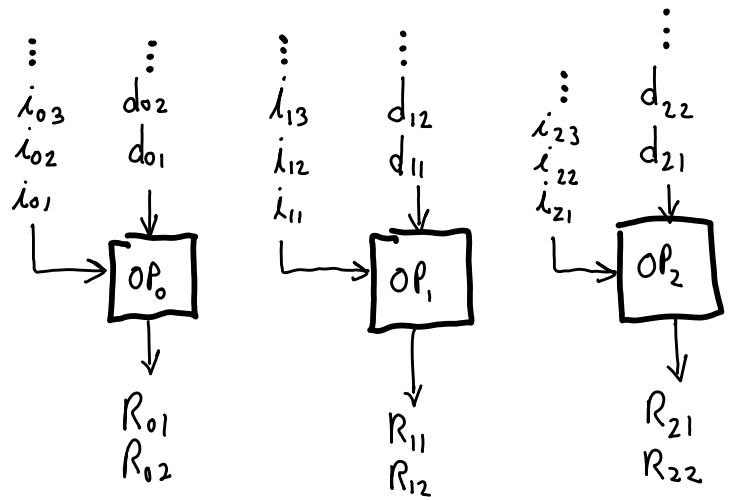
Single Instruction, Multiple Data SIMD

(data-level parallelism)



Multi-Instruction, Multi-Data MIMD

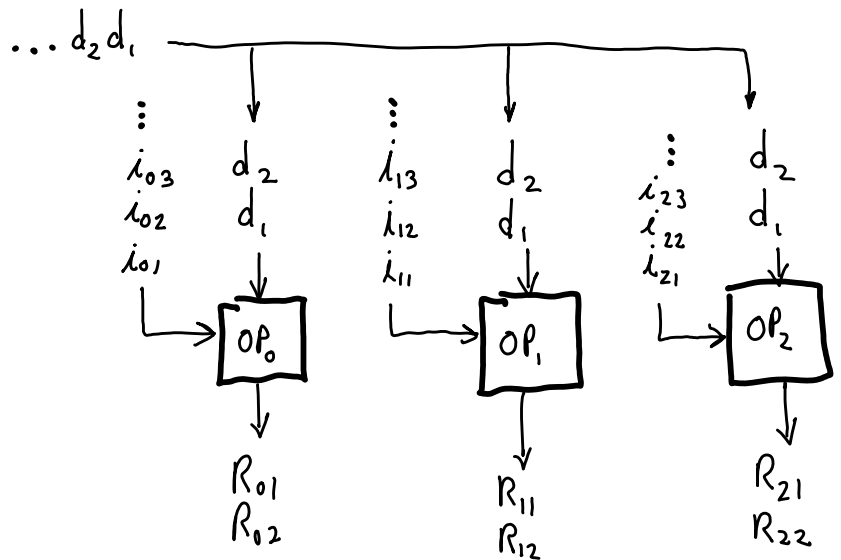
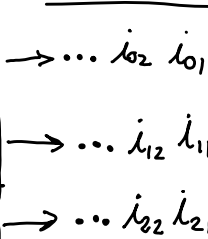
(instruction and data parallel)



SPMD



Threads on procs

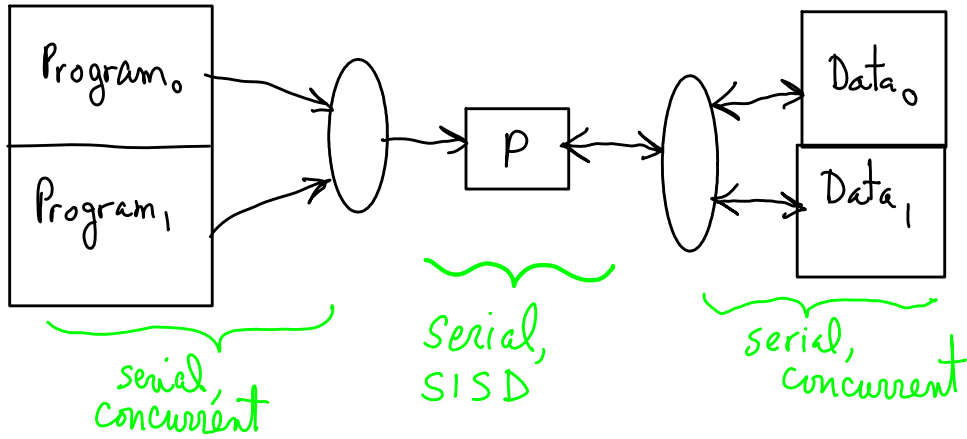
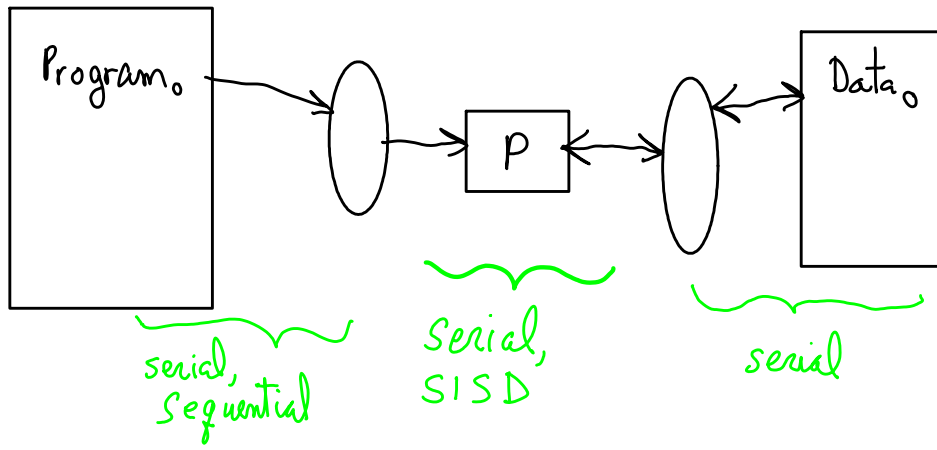


MISD

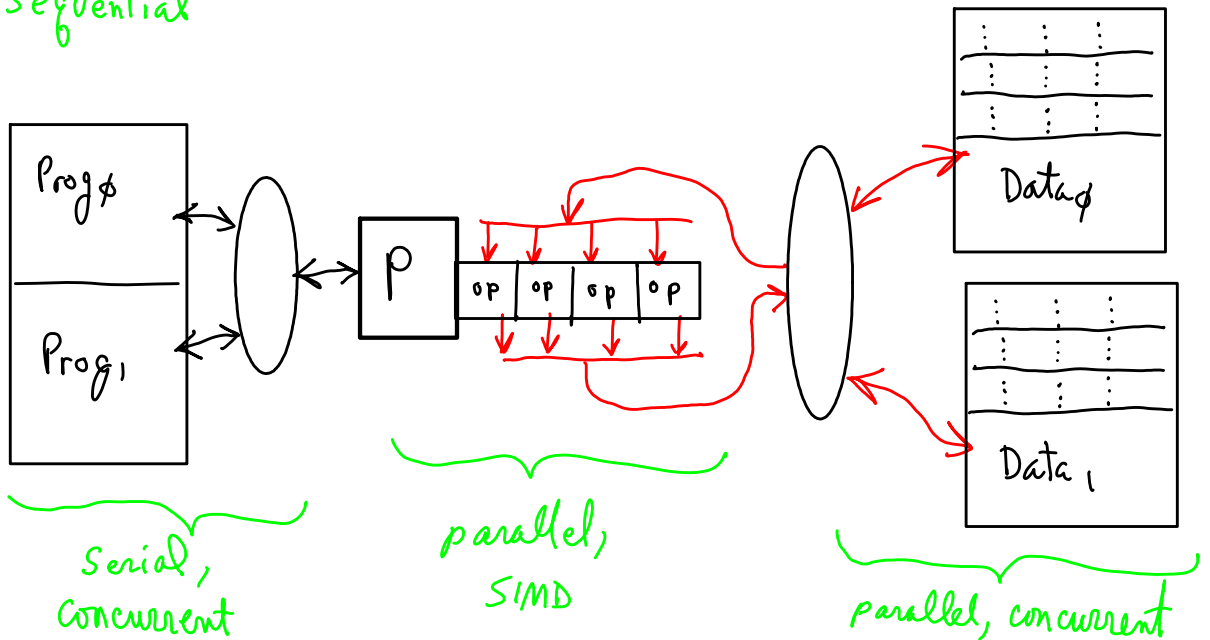
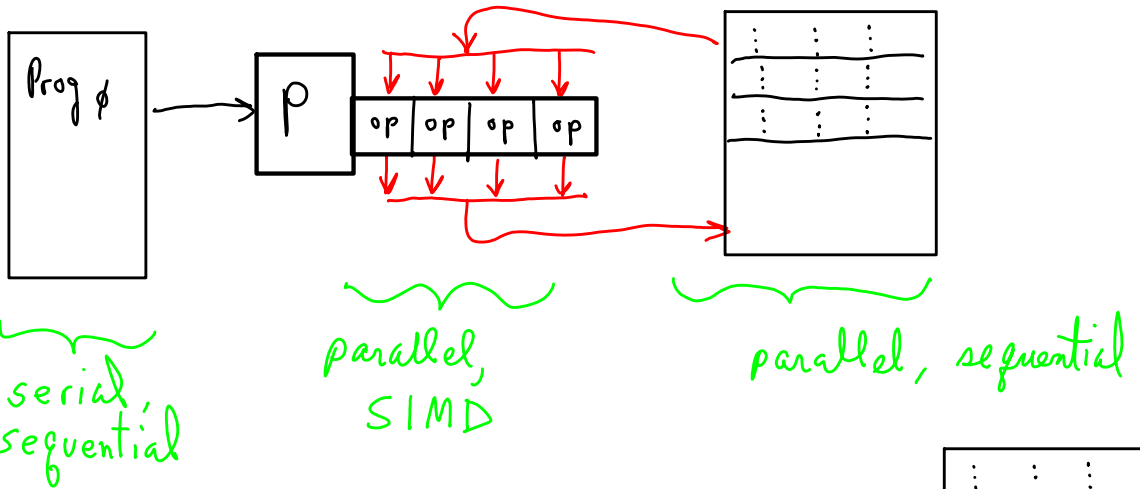
Multi-Instruction, Single Data

(instruction parallel)

# SISD



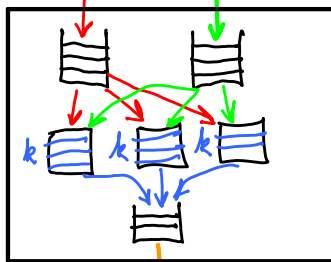
# SIMD



# ILP

$i_3 i_2 i_1$   
 $\dots i_3 i_2 i_1$

## CPU



$d_3 d_2 d_1$   
 $\dots d_3 d_2 d_1$   
 $\dots R_3 R_2 R_1$

Data matched on the fly with instructions (forwarding, etc).

Instructions queued until data is ready.

Instruction schedules reorder.

Commits (state change) are done to preserve semantics.

Nullifying (speculation)

$n$  functional units,  
 $k$ -level pipelining,  
 $\Rightarrow (n \cdot k)$  simultaneous instructions

+ speculative execution  
 + speculative prefetch

## Speed up

$$S_{2-1} = \frac{V_2}{V_1} = \frac{W_2 / T_2}{W_1 / T_1}$$

$$W_2 = W_1 \Rightarrow S_{2-1} = \frac{T_1}{T_2} = \frac{T_{\text{fixed}} + T_{\text{improvable}}}{T_{\text{fixed}} + T_{\text{improvable}} / S_{\text{improvement}}}$$

$W_{\text{real}} = \text{output, not computations}$   
 $V_{\text{real}} = W_{\text{real}} / T$

Speed up clock :  $T_{\text{fixed}} = 0 \rightarrow S = S_{\text{clock}}$

more Pipeline stages ( $k$ )

$$CR_{\text{new}} = k \cdot CR_{\text{old}}, S_{\text{clk}} = k$$

more functional units ( $n$ )

$$IPC_{\text{new}} = n \cdot IPC_{\text{old}}, S_{\text{IPC}} = n$$

$$T_{\text{fixed}} \approx (1-f)T : S \approx \frac{(1-f)T + fT}{(1-f)T + f \frac{T}{S_{\text{clk}} \cdot S_{\text{IPC}}}} = \frac{1}{(1-f) + \frac{f}{kn}} = \frac{kn}{kn \cdot g + (1-g)} \quad g = (1-f)$$

$$g \sim \frac{1}{kn} \rightarrow \frac{kn}{1+1}$$

(Speculation, forwarding, caching)  $\leftrightarrow$  (exposed latency, unused units, pipe + cache flushing/conflicts, pipe bubbles)

BUT  $\frac{1}{k}$  (W/cm<sup>2</sup>, device speed)  $\frac{1}{n}$  (usability)

# Scaling

assume:  $W_1 = W_2 = W = W_{\text{fixed}} + W_{\text{improvable}}$   
 $= (1-f)W + fW$

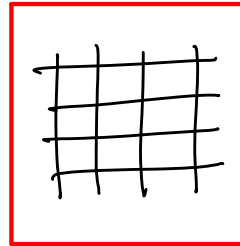
$$\begin{aligned} S_{2-1} &= \frac{V_2}{V_1} = \frac{W_2/T_2}{W_1/T_1} = \frac{T_1}{T_2} \\ &= \frac{(1-f)W/V + fW/V}{(1-f)W/V + fW/(S_f V)} \\ &= \frac{1}{(1-f) + f/n} \\ &= \frac{n}{n(1-f) + f} \xrightarrow{f \rightarrow 1} \frac{n}{0 + 1} \end{aligned}$$

$$S_f = n \quad \begin{cases} M_2 \text{ is } n \text{ times faster doing} \\ W_{\text{improvable}} = fW \end{cases}$$

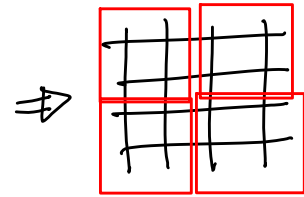
$$\begin{aligned} V &= V_{\text{fixed}} = V_{\text{improvable}} \quad \text{for } M_1 \\ V &= V_{\text{fixed}}, V_{\text{improvable}} = V/S_f \quad \text{for } M_2 \end{aligned}$$

E.G.

grid computation, 10% serial



1 processor



4 processors

$$S = \frac{4}{4(1-0.9) + 0.9} = \frac{4}{1.3} \approx \frac{4}{4/3} = 3 < 4$$

Caveat:

serial work for M1 < serial work for M4 (coordination, dividing data)

$$\Rightarrow W_1 = W_{\text{serial}} + W_{\text{parallel}} \quad \Rightarrow T_1 = \frac{W}{V}$$

$$W_2 = W_{\text{serial}} + h \cdot W_{\text{parallel}} + W_{\text{parallel}}$$

$$\Rightarrow T_2 = \frac{W_{\text{serial}} + h W_{\text{parallel}}}{V} + \frac{W_{\text{parallel}}}{V \cdot S_n}$$

$$= \frac{(1-f)W + h f W}{V} + \frac{fW}{V \cdot n} \quad \Rightarrow S_{2-1} = \frac{T_1}{T_2} = \frac{1}{(1-f) + h + f/n}$$

Here h contributes to the serial portion. Could contribute to parallel part.

$$h = h(n)$$

# STRONG scaling

But suppose  $h(n) < 0$ ?  $h \in (-1, 0]$

Less serial work needed w/ n processors

$$\sum_{2-1}^1 = \frac{1}{(1-f+h) + f/n} \xrightarrow{h \rightarrow (f-1)} n/f = \mathcal{O}(n)$$

or

$$W_{\text{parallel}} \rightarrow fW + hW$$

Less parallel work needed w/ n processors

$$\sum_{2-1}^1 = \frac{1}{(1-f) + (f+h)/n} \xrightarrow{h \rightarrow -f} \frac{1}{(1-f)}$$

# Weak scaling

More processors ==> Do a larger problem

$$W = W_{\text{fixed}} + n \cdot W_{\text{improvable}}$$

(or generally  $r(n) \cdot W_{\text{improvable}}$ )

e.g.  $r(n) = \mathcal{O}(n)$

$$T_2 = \frac{W_{\text{fixed}}}{v} + \frac{n \cdot W_{\text{improvable}}}{v}$$

$$\sum_{1-2} = \frac{\frac{W_{\text{fixed}}}{v} + \frac{n \cdot W_{\text{improvable}}}{v}}{\frac{W_{\text{fixed}}}{v} + \frac{n \cdot W_{\text{improvable}}}{n \cdot v}} = \frac{(1-f) + nf}{(1-f) + f} = 1 + f(n-1) = \mathcal{O}(n)$$

for our previous example,  
4x larger grid

$$= 1 + 0.9(3) = 3.7$$

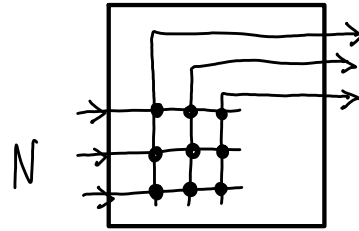
In general, all above effects can play a part (+ or -).

- serial work for coordination, initialization
- parallel work for parallel portion, problem remapping
- communication overhead and resources
- memory and cache bandwidth effects
- cache coherency



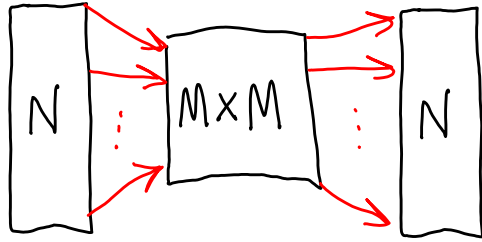
**interconnects switched**

$N \times N$   
Cross bar  
switch



- $N$   
route
- any 1-1 map
  - any  $(N-k) \times N$  \*  
fan-out
  - any  $N \times (N-k)$  \*  
fan-in
- \* Collisionless  
(buffering?)

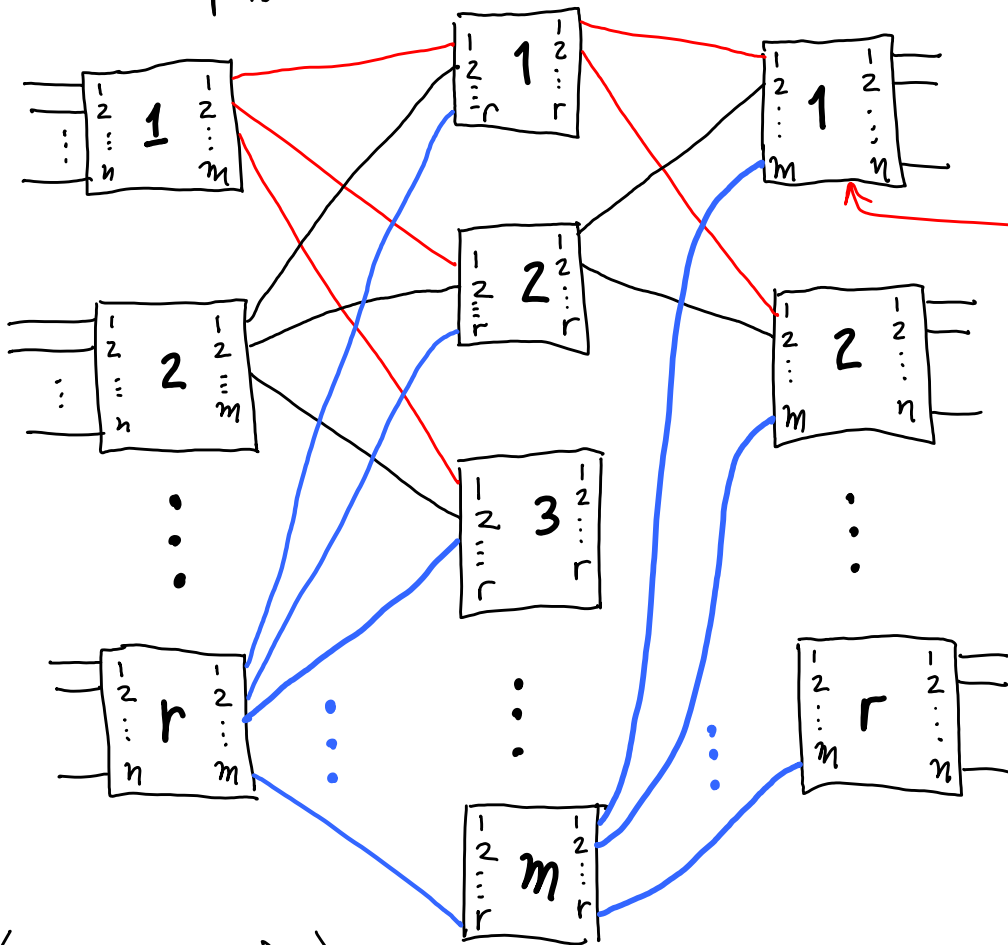
Clos



$N = n \cdot r$  inputs

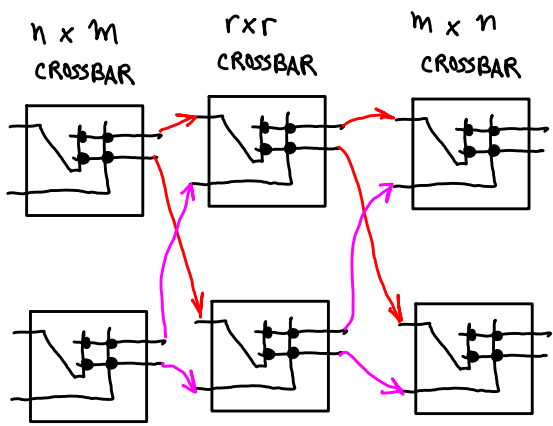
$M = r \cdot m$

$N = n \cdot r$  outputs

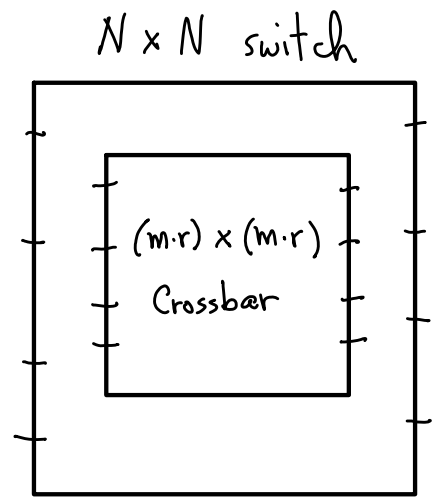


crossbar switch

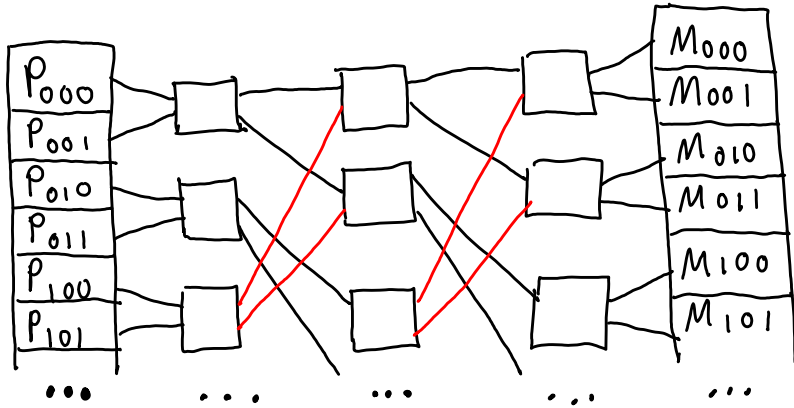
(recursive def'n)



||



# $\Omega$ 8x4 interconnect



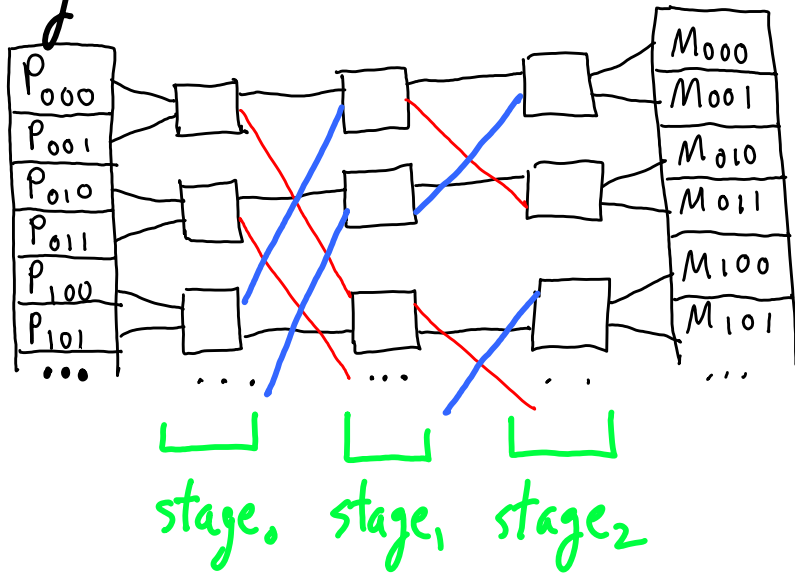
switch  $\Rightarrow$  blocking



also, can broadcast:



# Butterfly 8x4



(Banyan)

$n$   $P_s$

$n$   $M_s$

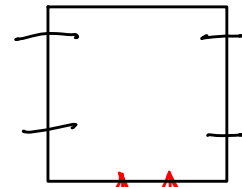
$\log n$  stages, @  $n/2$

$\log n$  address bits

Switch at  $k^{\text{th}}$  stage:

addr<sub>0</sub>  $a_{\log n} \ a_{\log n-1} \ \dots \ a_k \ \dots \ a_1 \ a_0$

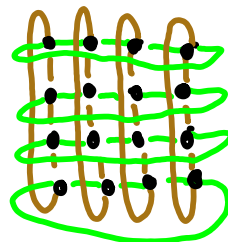
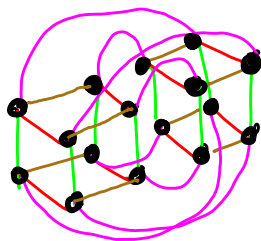
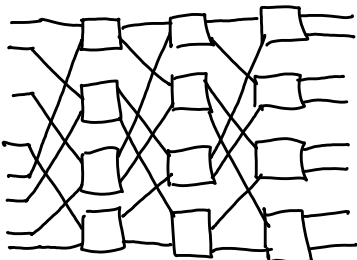
addr<sub>1</sub>  $a_{\log n} \ a_{\log n-1} \ \dots \ a_k \ \dots \ a_1 \ a_0$



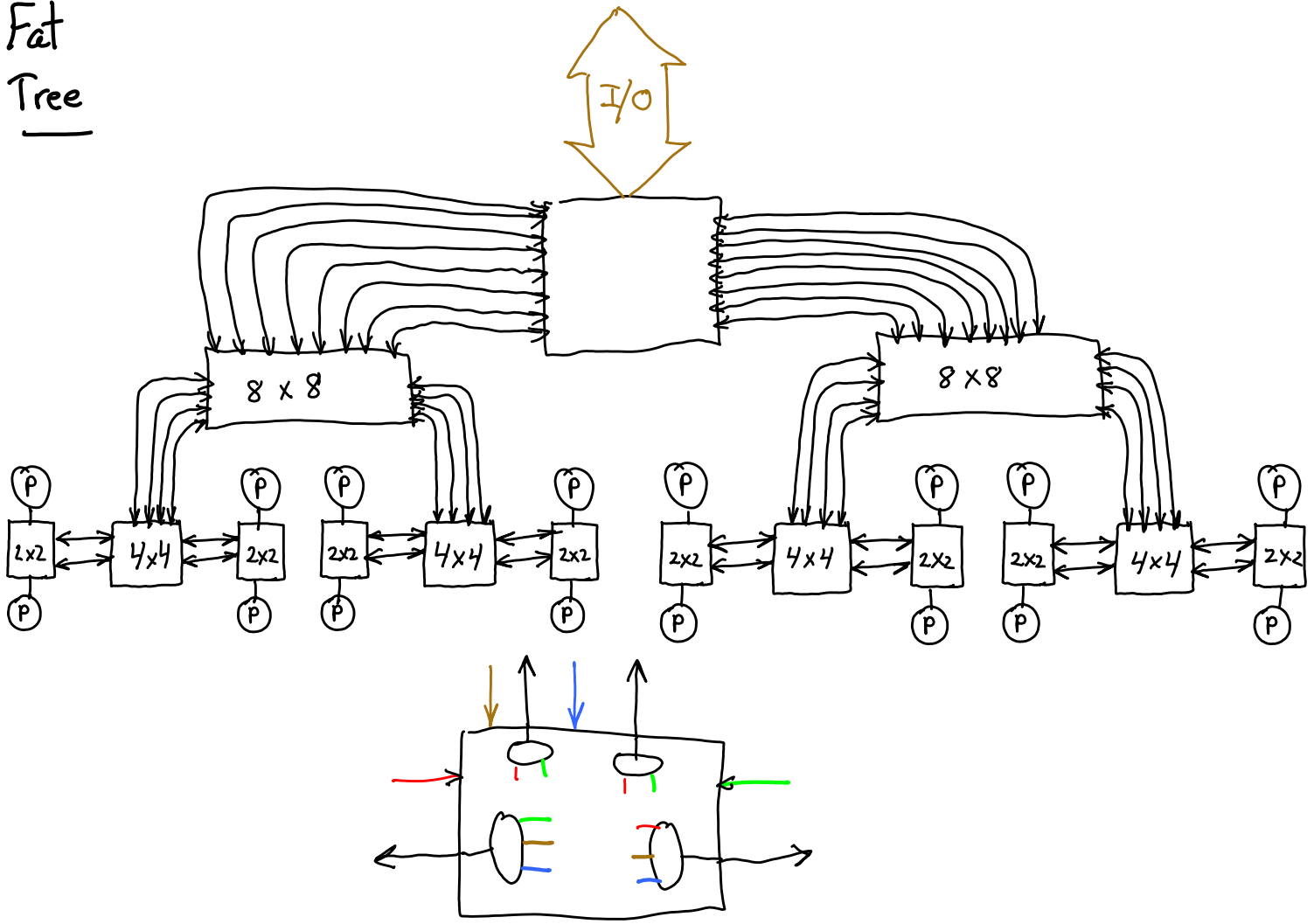
shuffle-exchange (perfect card shuffle)

hypercube

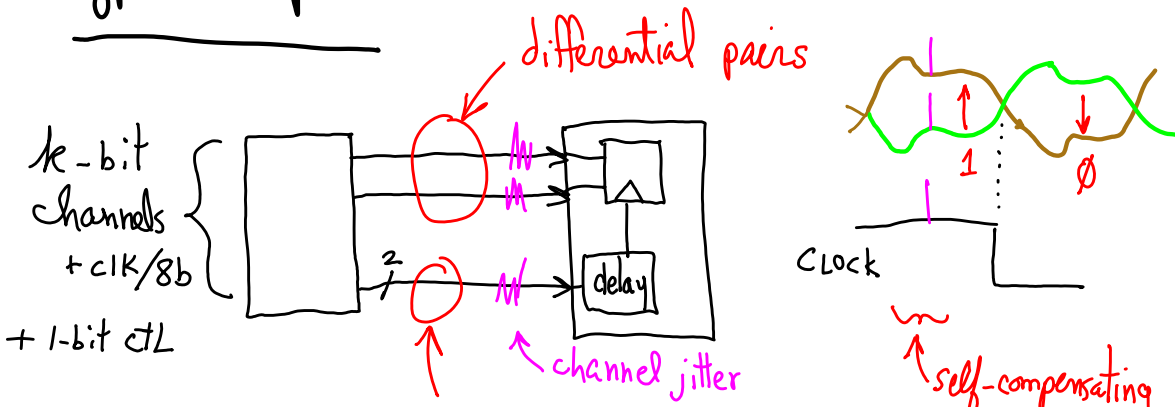
Torus



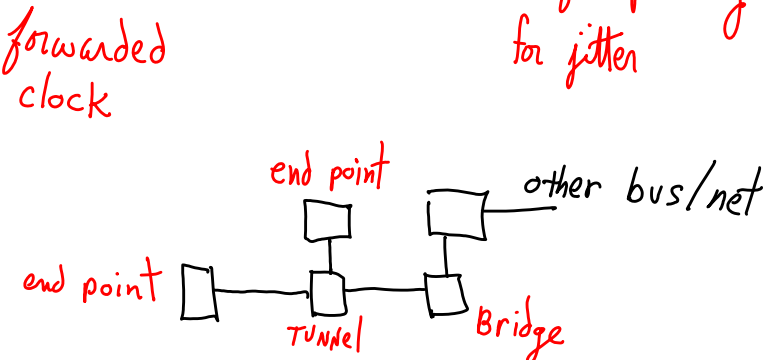
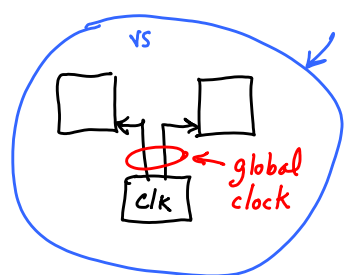
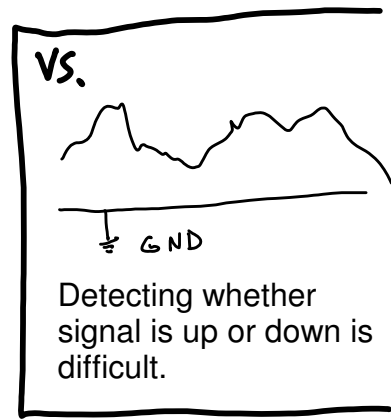
# Fat Tree



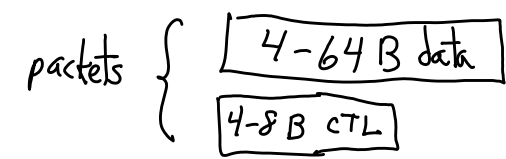
# HyperTransport



high-speed signals are bouncy (jitter), and corrupted by other signals.



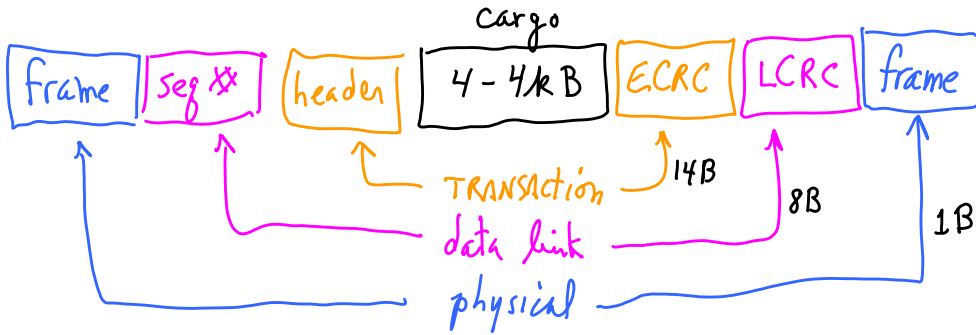
daisy chained, point-to-point



# Packet overhead

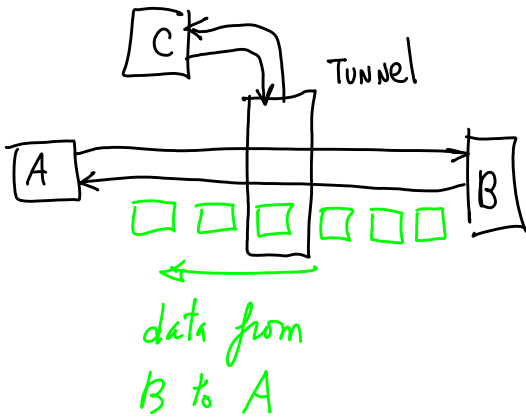
Header		cargo	HT
~8 B	4-64 B		

Transaction / data link / physical

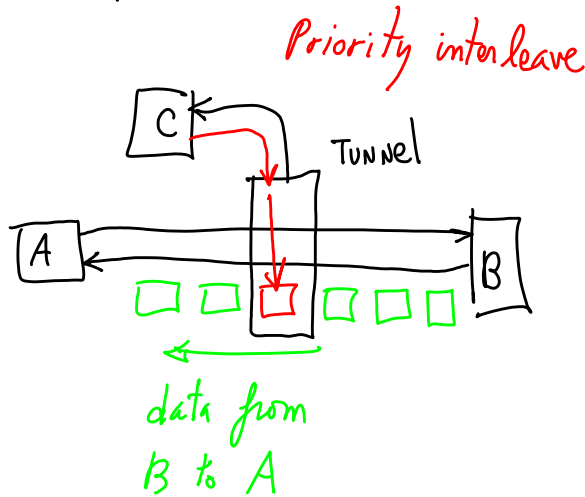


PCIe ~24 B overhead

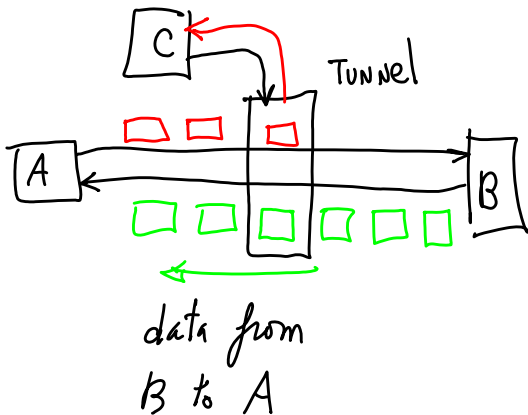
# Priority / interleaving



Links in pairs

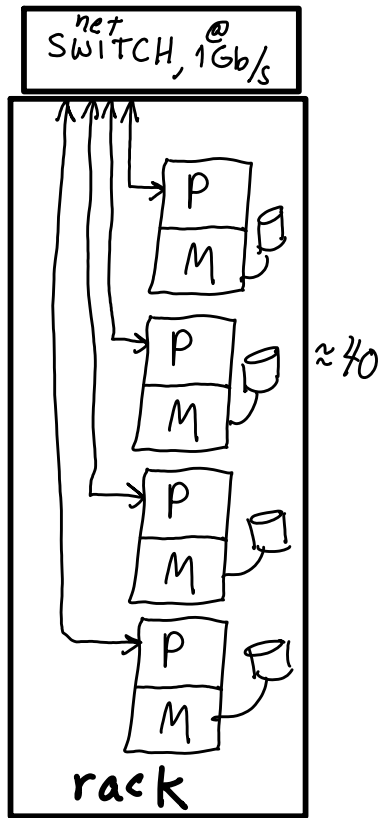


data from A to C



⇒ Low Latency response

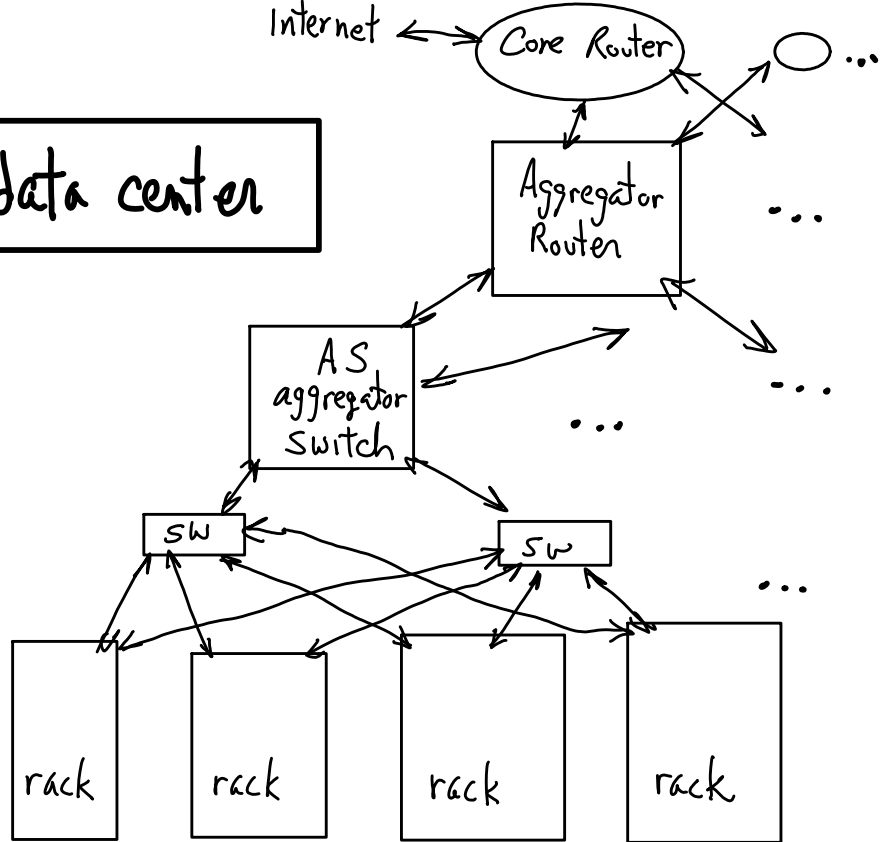
# Cluster



non-shared memory,  
message passing

Thread parallelism  
Process parallelism  
Task parallelism

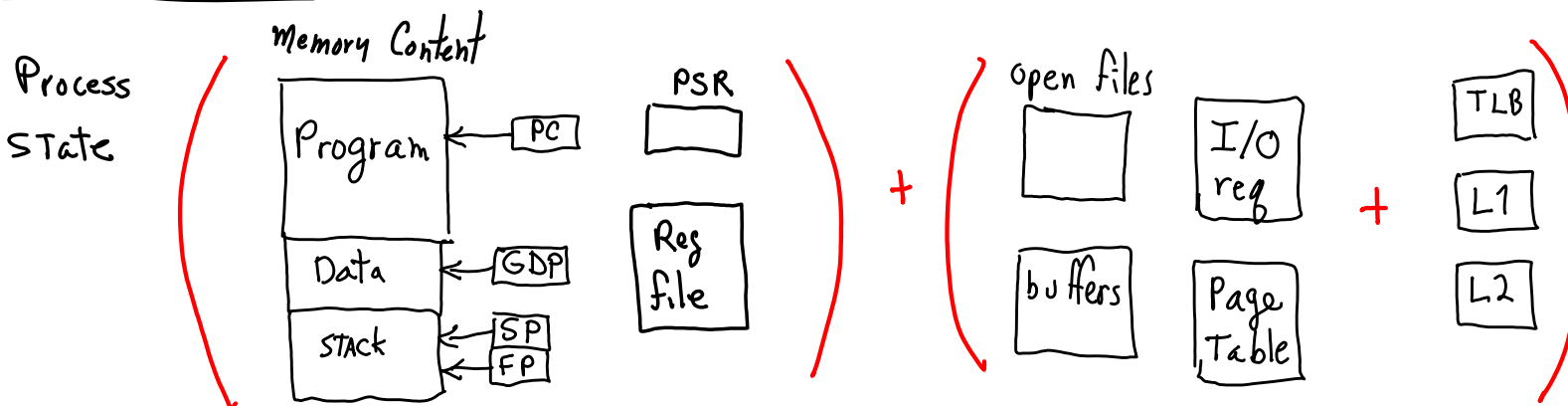
# data center



- difficult algorithm-to-program paradigm
- slow communication
- multiple copies of OS, software
- multiple administration
- middle ware overhead
- hot swappable
- expandable
- commodity components
- virtualization, migration, redundancy

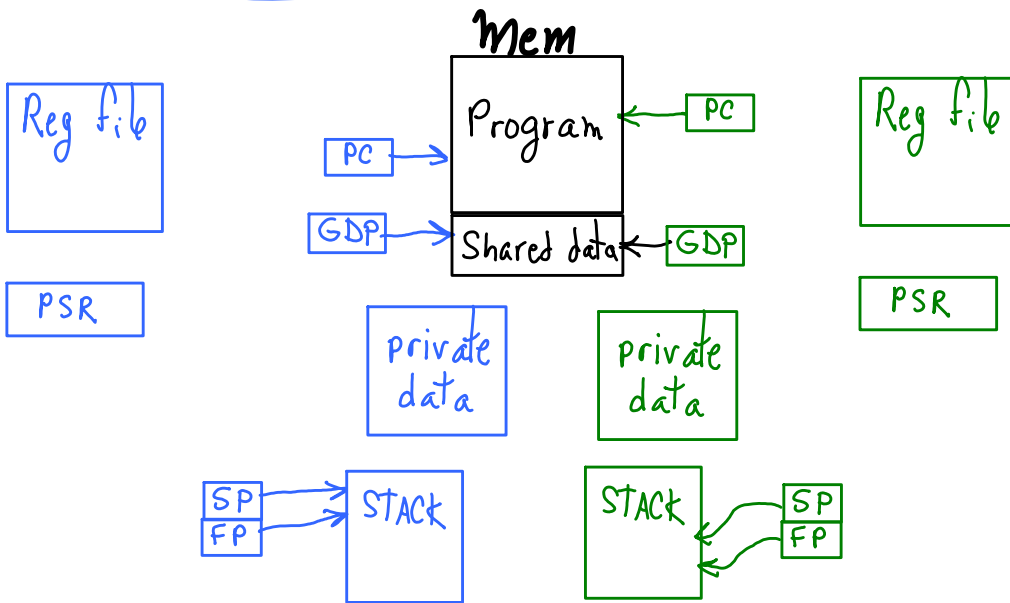
# Threads

Process context switch  $\approx 1M$  instructions



Thread 1  
STATE

Thread 2  
STATE



Switch threads

Software:

Copy state - from

Copy state - to

load PSR, PC, SP, FP

TLB, L1, L2 ok.

PT ok (mostly?)

File tables are ok?

Buffers? I/O req.?

## Software threads

### Switching thread context

In OS code, or user code:

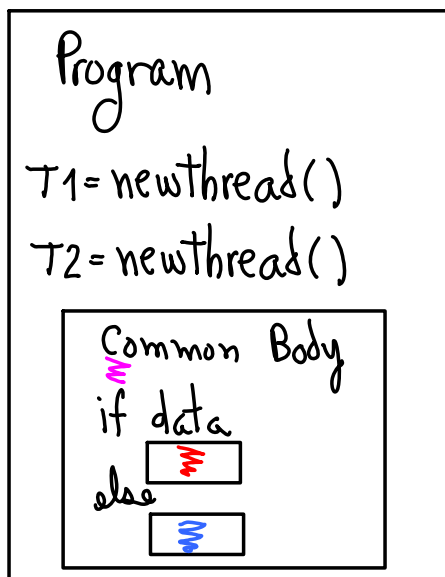
- save state T1
- load state T2

run T2 until

- HW timer
- IO event sleep
- cooperative hand-off

Switch context

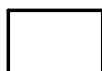
# Instruction streams



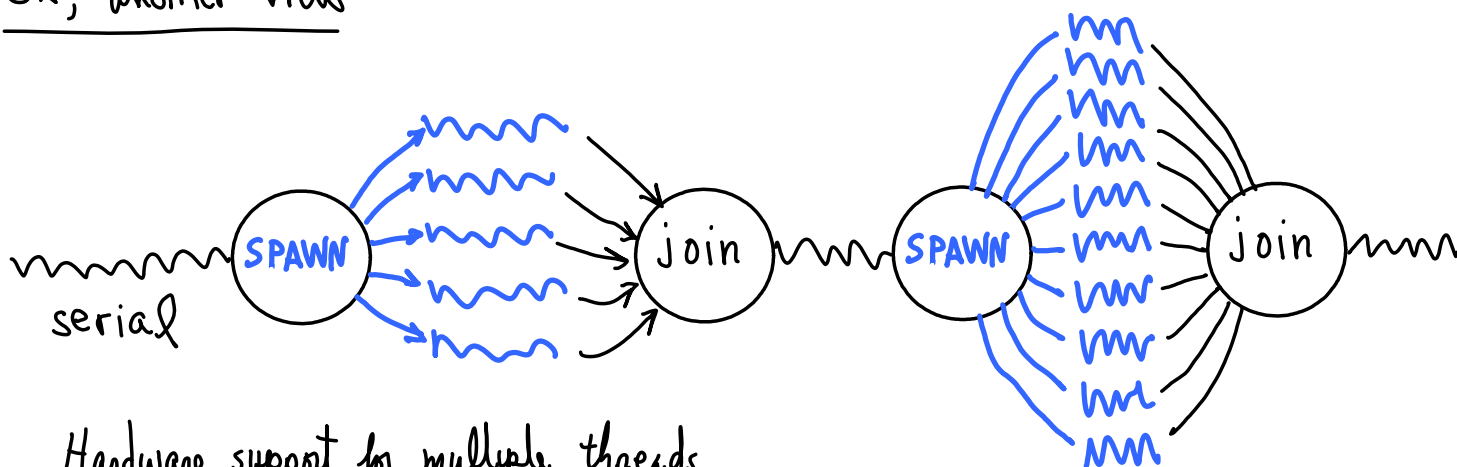
PC1



PC2



## OR, another view



## Hardware support for multiple threads

### ALL THREADS from SAME PROCESS

- Duplicate and switch: PC, PSR, RegFile, Stack, Private data
- Copy/Save/Restore state
- Shadow registers, renaming
- TLB content (separate page tables? or shared?)
  - hardware switch
  - TID, thread ID labeled

### MULTIPLE THREADS from MULTIPLE PROCESSES

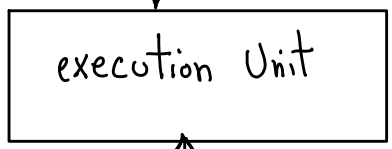
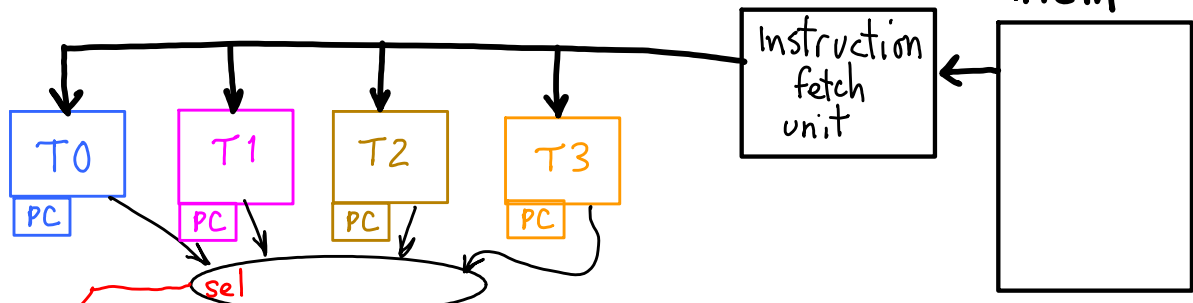
- PID + TID
- Larger state to consider (page tables, file and IO tables and buffers)

OS policy  
not known to  
HW designer?



# Fine-grained multi-threading

Instruction buffers



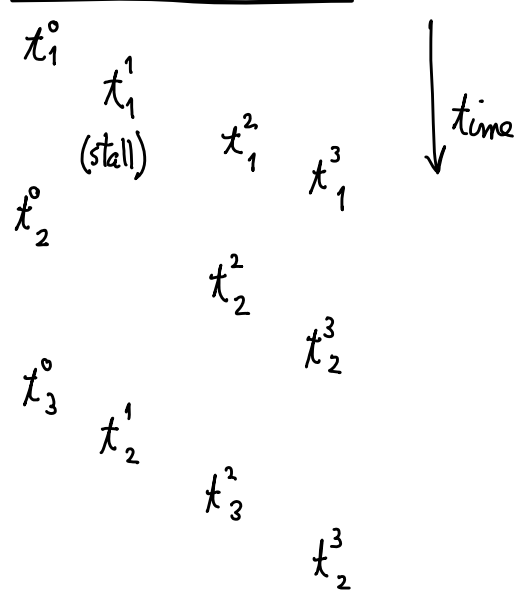
Contexts



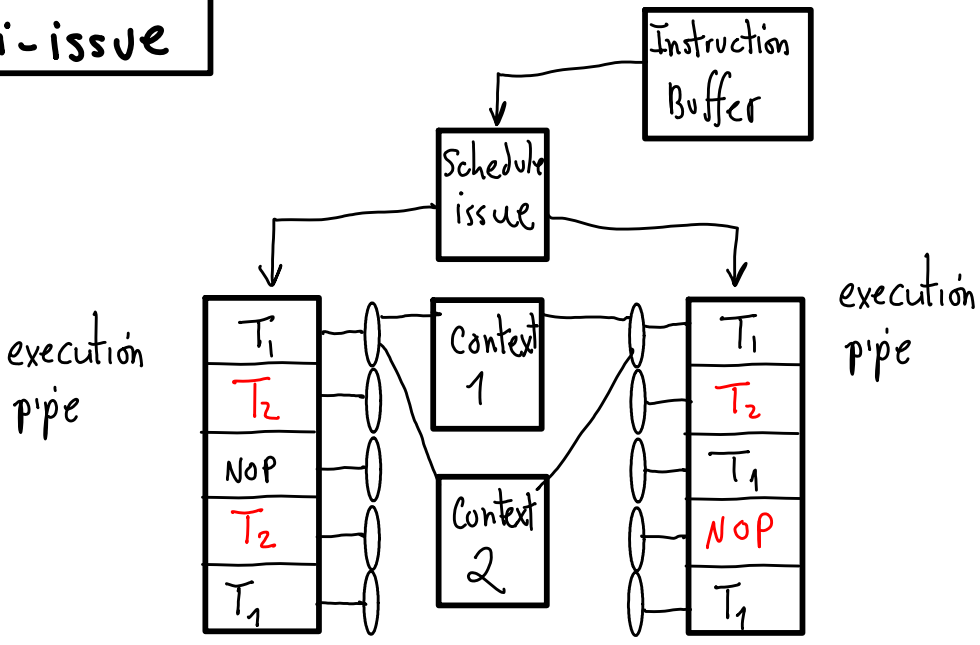
Round Robin Select

- skip stalled threads
- longer response time
- switching overhead

## Execution Trace



# Multi-issue



# Coarse Grained

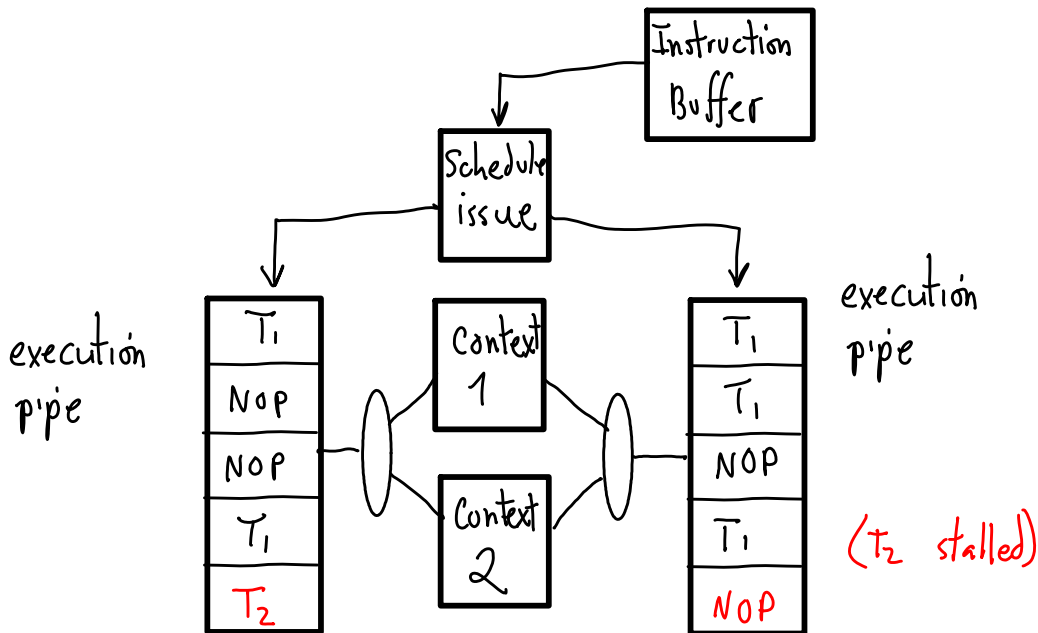
Same HW, but schedule same thread until stalled.

## Execution Trace

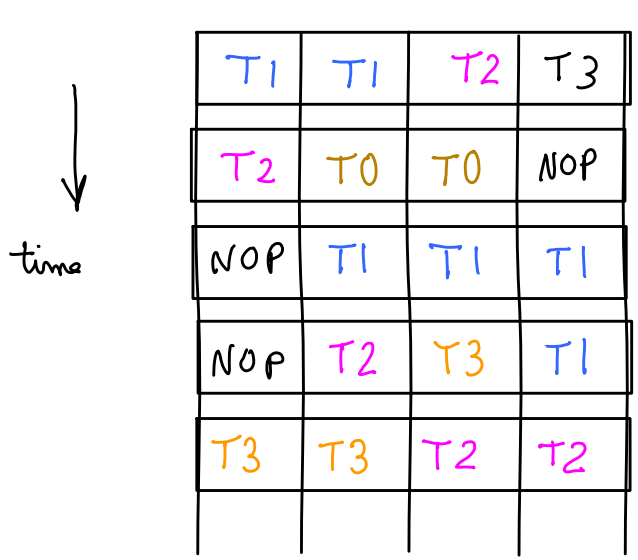
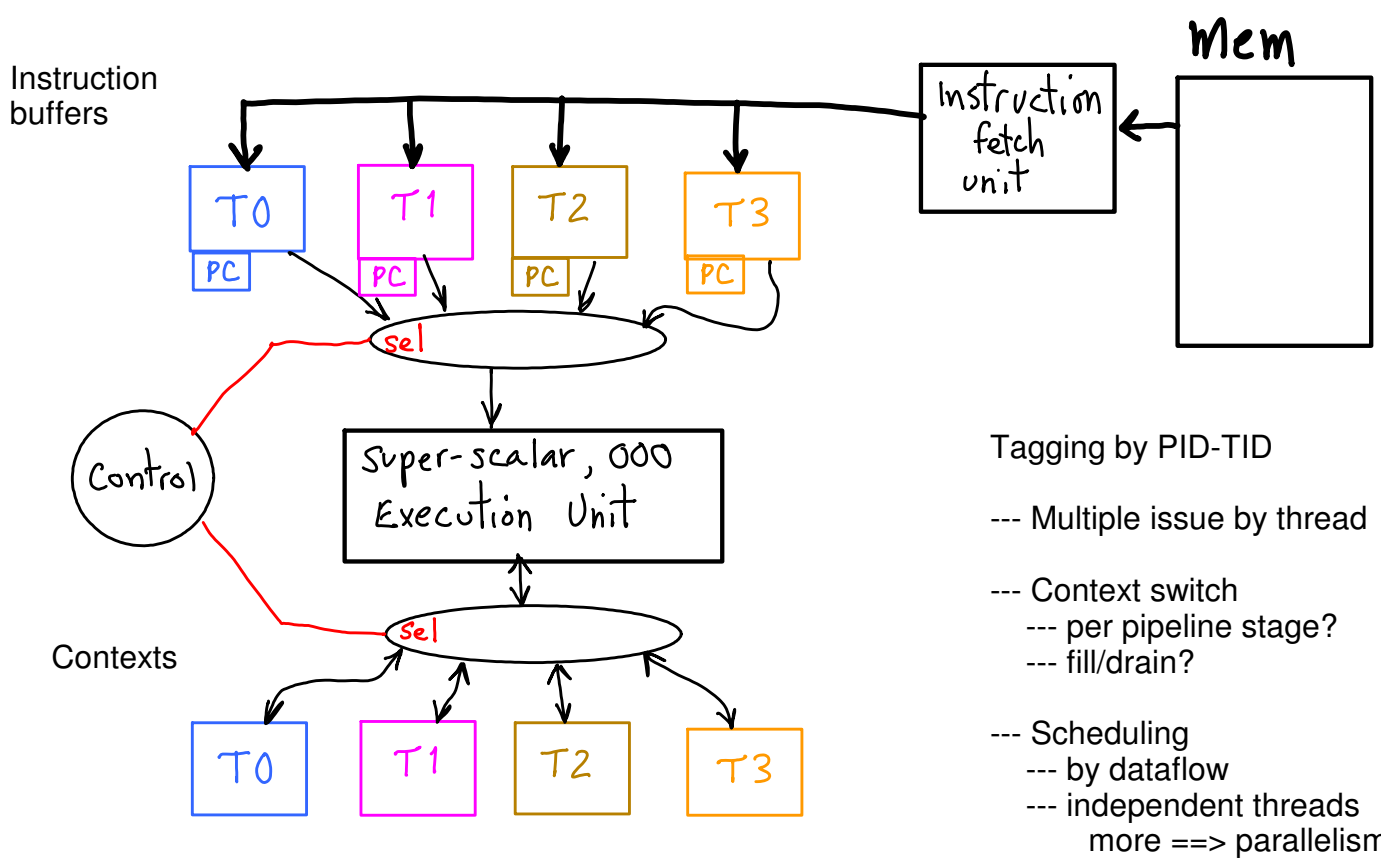
$t_1^0$   
 $t_2^0$   
 $t_3^0$   
 (stall)  $t_1^1$   
 (stall)  $t_1^2$   
 $t_2^2$   
 $t_3^2$

- Faster response time
- lower switching overhead
  - affords more complex control
  - faster execution between
- Starvation now a problem?

# Multi-issue



# (Simultaneous)(Hyper) Threading



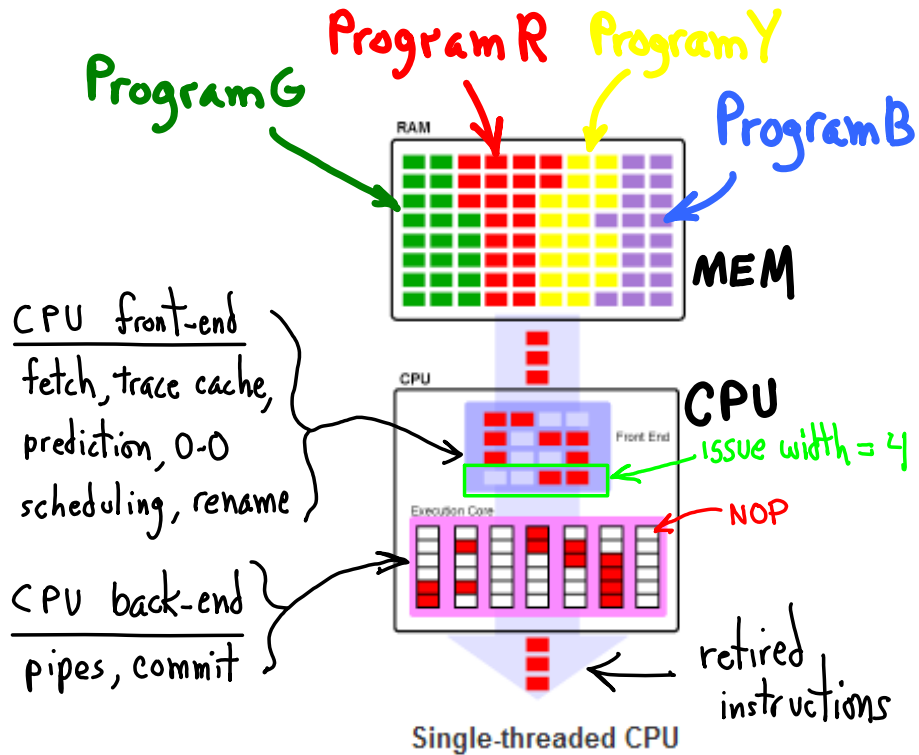
← multiple issue

- OOO scheduling across threads
- more filled slots

## Single threaded execution

### Multi-tasking

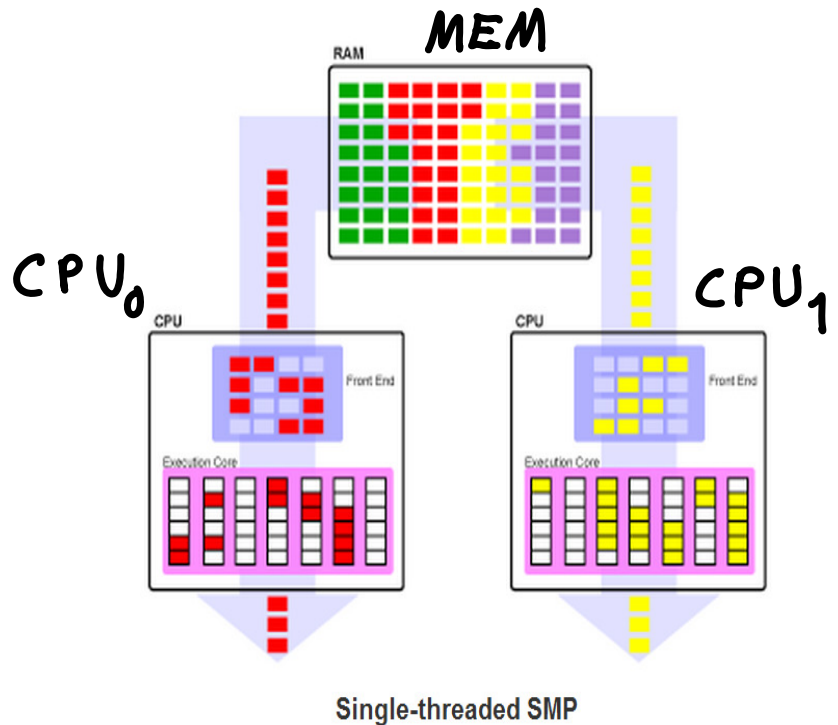
- Multiple concurrent execution (not simultaneous)
- Memory shared but separate (virtual)
- CPU time-multi-plexed
  - cooperatively, pre-emptively, IO
- Process context switching
  - drain/fill (pipes, caches, TLB, ...)
- Extract ILP from single stream
  - Unused issue slots
  - pipeline bubbles/stalls



Credits:  
*Introduction to Multithreading,  
 Superthreading and Hyperthreading*  
 By Jon Stokes

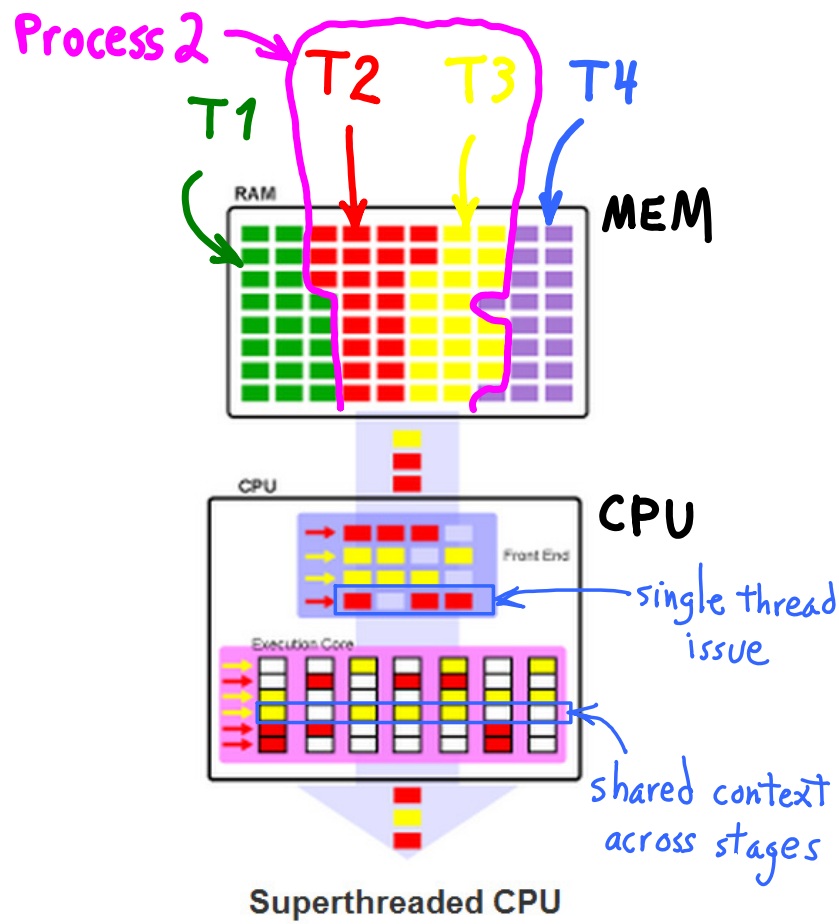
## SMP, Symmetric Multi-Processing

- Context switch per CPU
- Simultaneous execution
  - multiple programs/processes/threads
- ILP extracted per process
  - double silicon resources
  - same NOP density
  - ==> Could speedup be > 2?



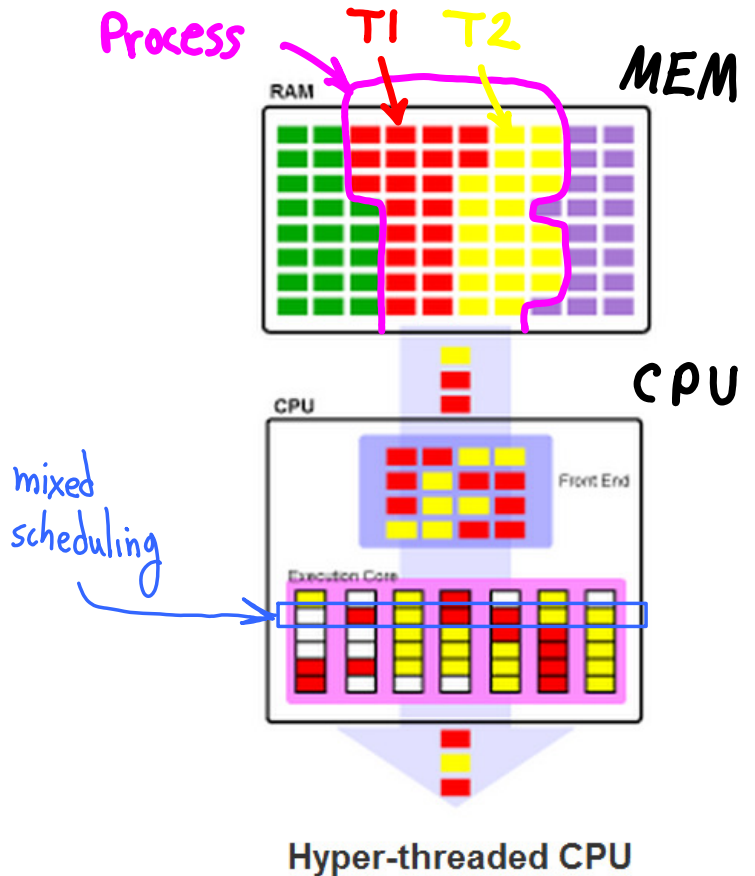
## Multi-Threaded (Superthreading)

- Concurrent process scheduling
  - process context switching
  - cooperative, pre-emptive, ...
- Single process, multiple thread execution
- Time-multiplexed thread scheduling
  - from same thread
- Instructions issue from single thread
  - thread context switch
  - in HW
  - per stage
  - across stages
- Execution slots filled
  - due to stalled threads
  - filled from non-stalled threads
  - Lower density of NOPs



## SMT, Simultaneous Multi-(Hyper)-Threading

- Concurrent Processes
  - context switching
- Thread context switching
  - independently on different pipes
  - issue from multiple threads simultaneously
- Average ILP = 2.5, empirically
  - max single-thread issue = 4 (here)
  - combined ILP ==> 4
- Logical Processors == 2
- Lowest NOP-density

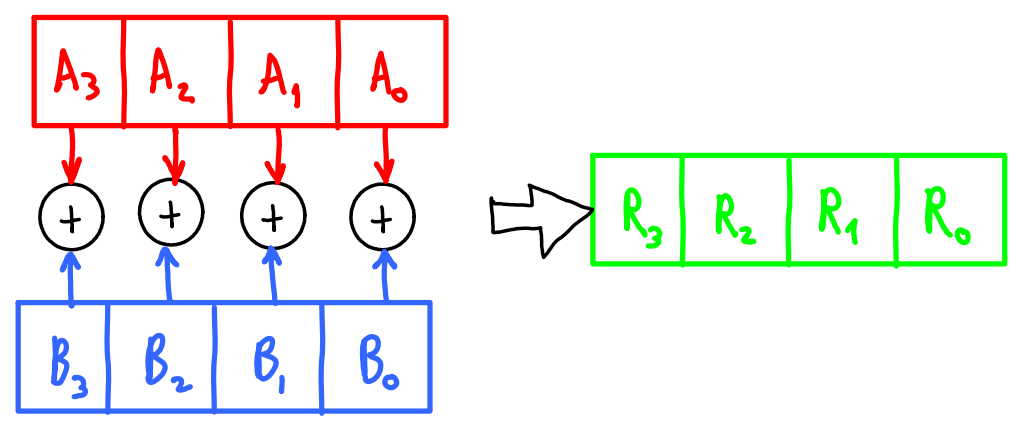


Modification of existing O-O-O CPU  
 => 10% added cost,  $\rho > 2$

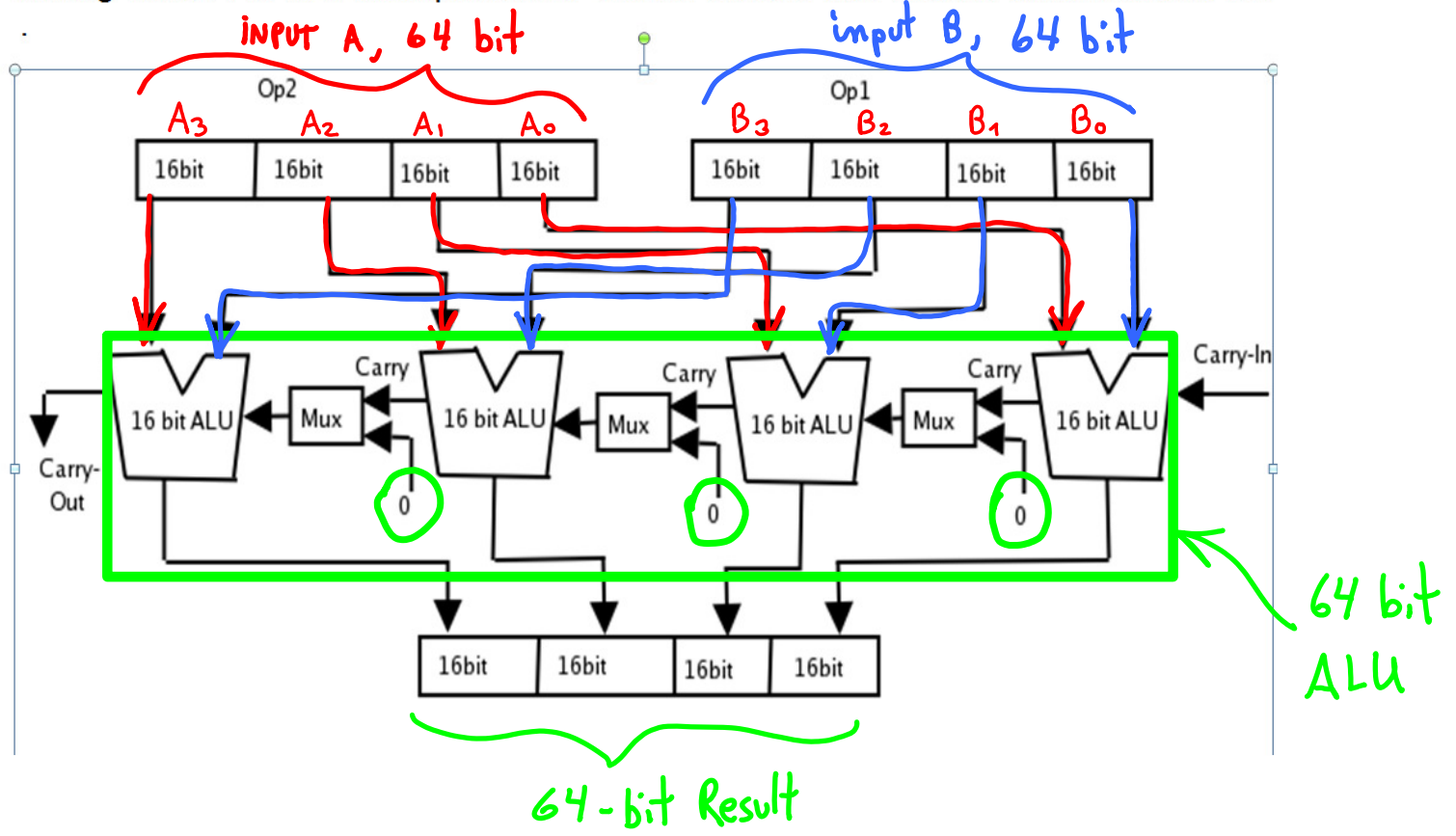
# SIMD, cheap

Why not use existing ALU as vector processor?

Vector ADD:  $R = A + B$



A long-word ALU in a microprocessor can be divided into several small-word ALU's



What for? 16-bit sound DSP?

Intel MMX  $\Rightarrow$  added to ISA: larger

- 1) Vectors (more elements)
- 2) elements (more bits)

# Predicated execution

what to do w/ IFs?

Loop (per vector)

if  $A_i > B_i$

$R_i \leftarrow A_i + B_i$

else

$R_i \leftarrow (-A_i) + B_i$

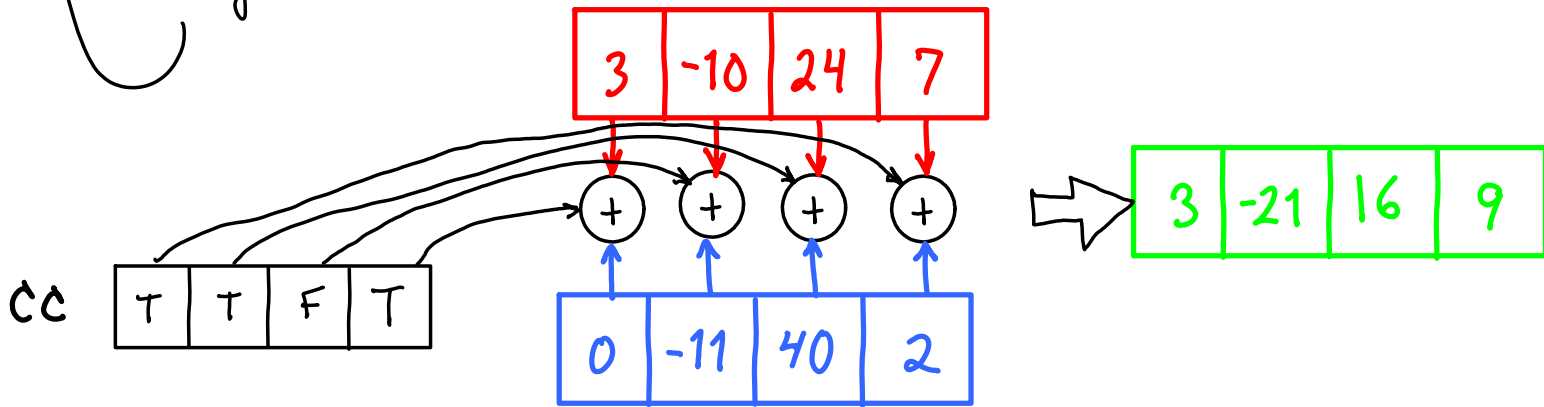
endif

$\Rightarrow$

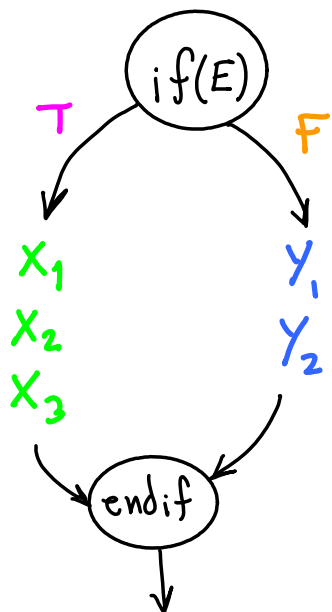
$(A - B) \Rightarrow CC$  (predicate vector)

$cc_3$	$cc_2$	$cc_1$	$cc_0$
T	T	F	T

Predicate vector



# Predicated execution for non-vector ops



Predicate register  $P$  ?

Remove Branches

