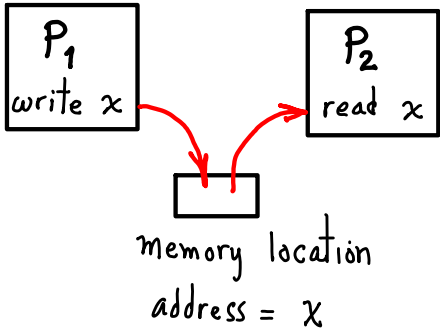


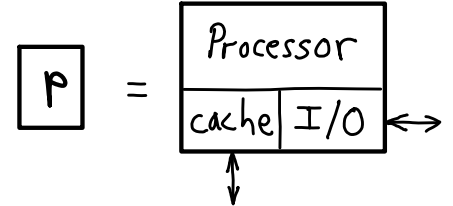
Parallel Systems

common address space PRAM model

multi-processor, shared memory (SMP)

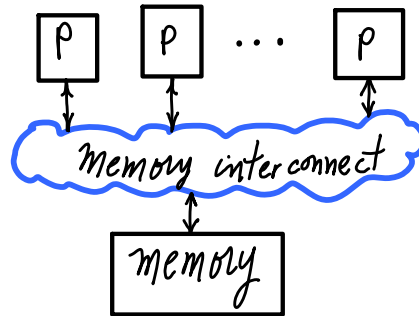


coordination through memory
(can implement messages on top)



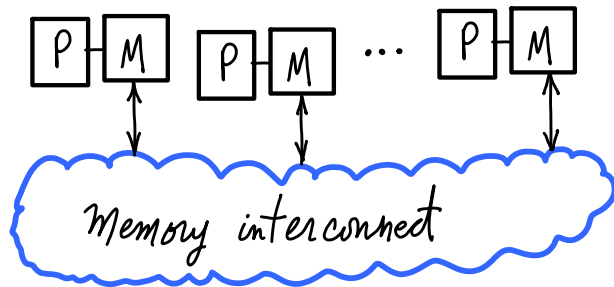
- Uniform access to memory: UMA

- Global clock
- ALL Procs communicate at same time
- ALL data sent/received in one step



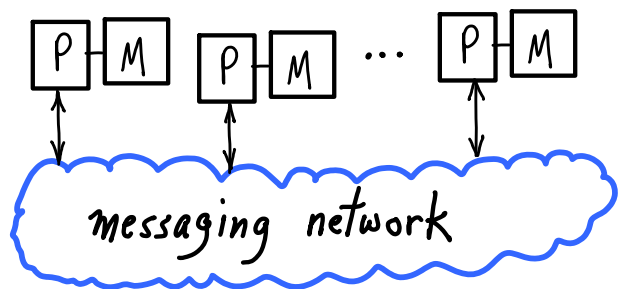
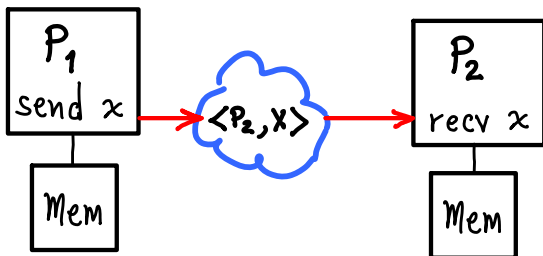
- non-UMA: NUMA

- Local clocks
- Memory delays vary
 - Local: fast
 - Remote: slow
- interconnect hierarchy



Non-shared Memory: message passing

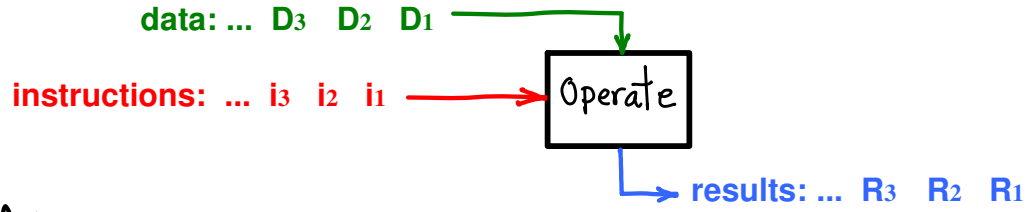
coordination through messages



Taxonomy

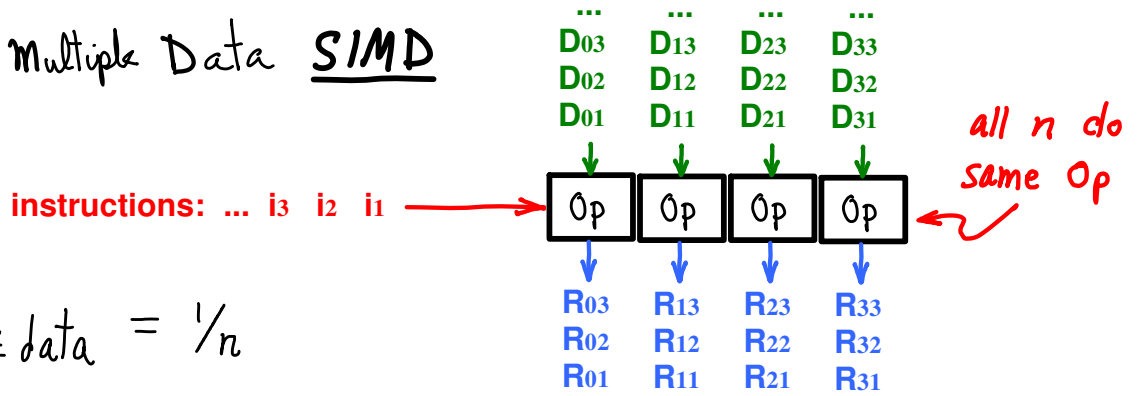
Single Instruction, Single Data SISD
STREAM STREAM

serial



data-level parallelism

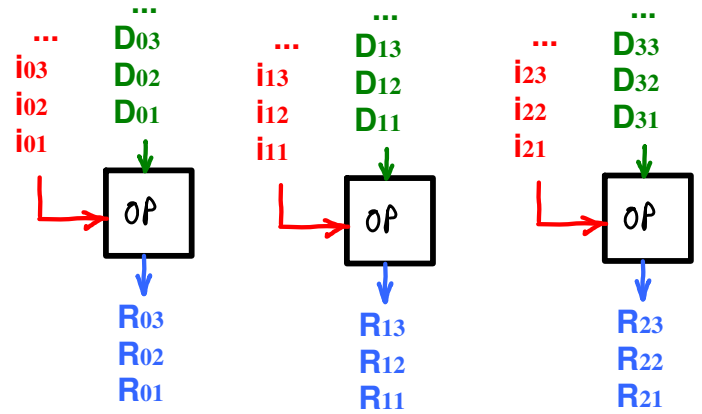
Single Instruction, Multiple Data SIMD



$$\# \text{ instructions} / \# \text{ data} = 1/n$$

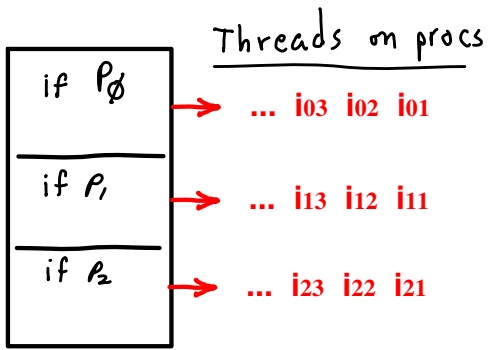
instruction and data parallel

Multi-Instruction, Multi-Data MIMD



SPMD

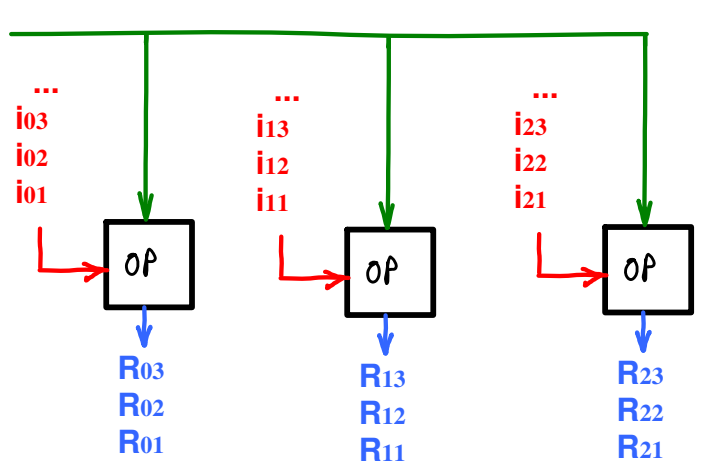
SPMT

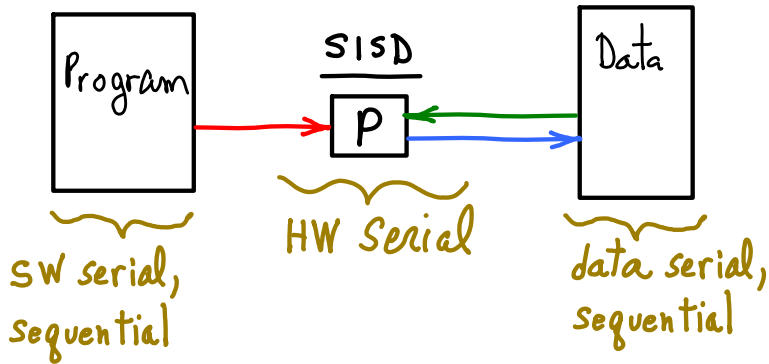


instruction parallel

Multi-Instruction, Single Data

MISD

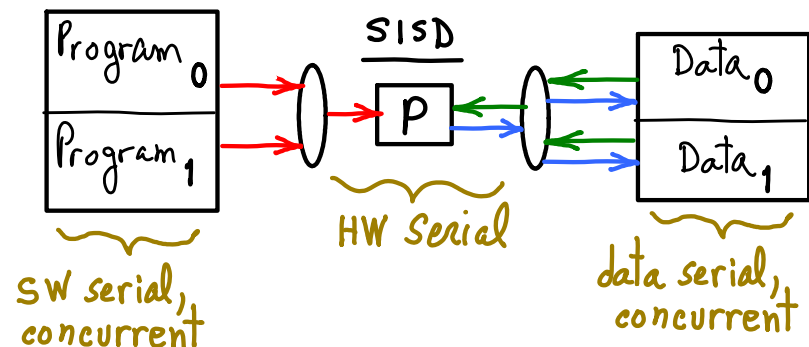




SW serial: **single instruction stream**
 SW sequential: **single-process**

HW serial: **non-parallel execution**

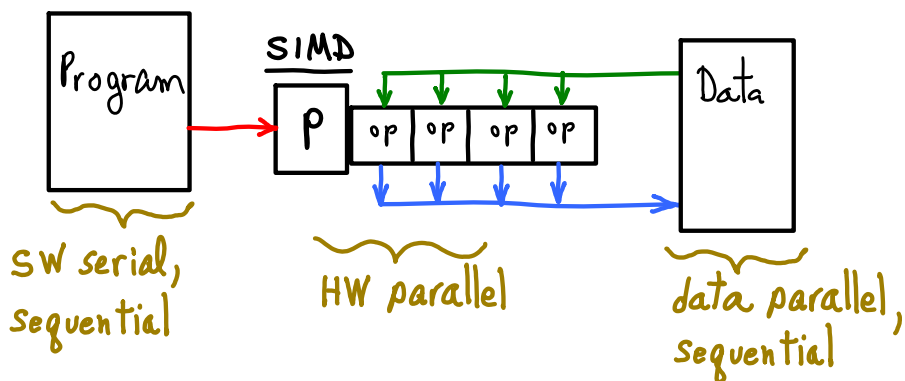
Data serial: **single data stream**
 Data sequential: **single-source**



SW serial: **single instruction stream**
 SW concurrent: **multi-process**

HW serial: **non-parallel execution**

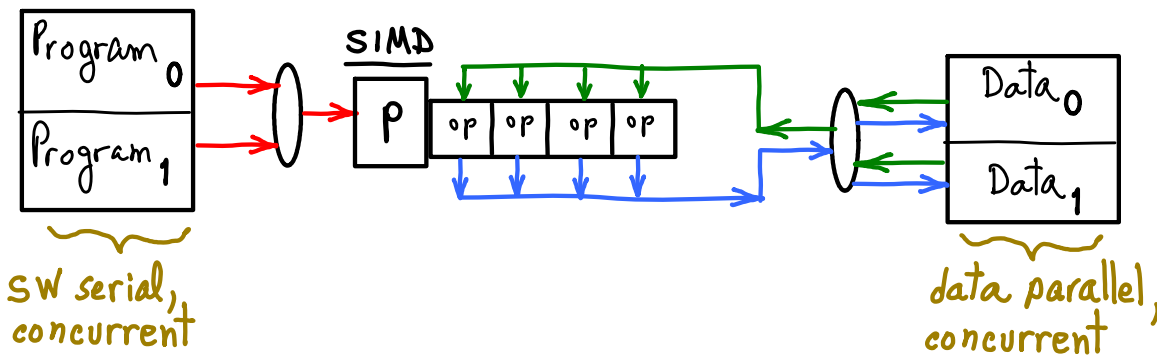
Data serial: **single data stream**
 Data sequential: **multi-source**



SW serial: **single instruction stream**
 SW sequential: **single-process**

HW: **parallel execution**

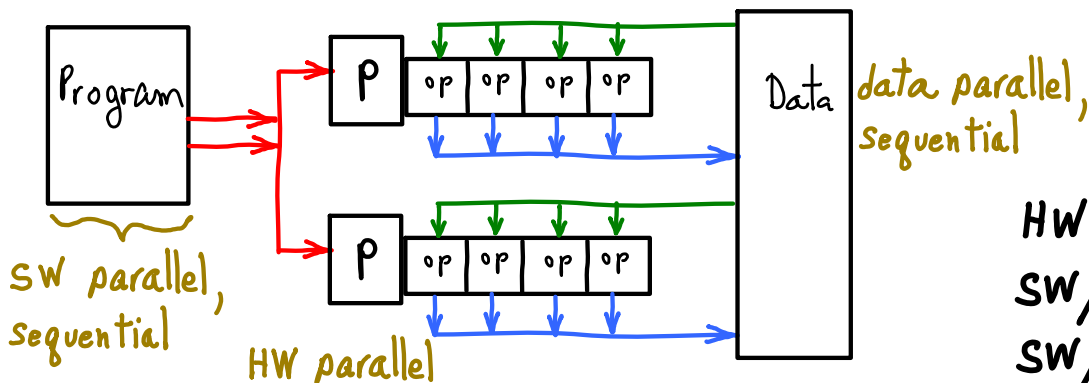
Data parallel: **multi-data stream**
 Data sequential: **single-source**



SW serial: **single stream**
 SW concurrent: **multi-process**

HW: **parallel**

Data parallel: **multi-data**
 Data concurrent: **multi-source**



etc.

HW parallel: **SIMD, MIMD**
 SW/data: **parallel/serial**
 SW/data: **concurrent/seq.**

ILP

instruction streams

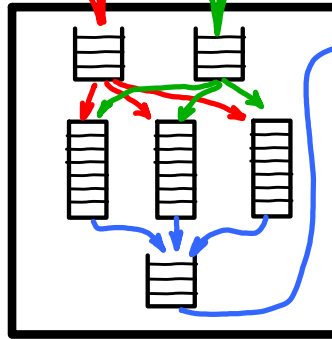
$i_3 \ i_2 \ i_1$
 $\dots \ i_3 \ i_2 \ i_1$

CPU

data streams

$d_3 \ d_2 \ d_1$
 $\dots \ d_3 \ d_2 \ d_1$
 $\dots \ r_3 \ r_2 \ r_1$

n functional units,
 k -level pipelining,
 $(n \cdot k)$ simultaneous instructions
 + speculative execution
 + speculative prefetch



Instructions queued for data.
 Data matched (forwarding, etc).

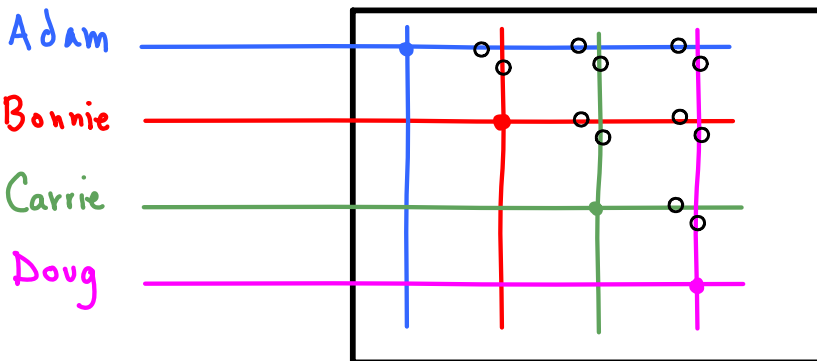
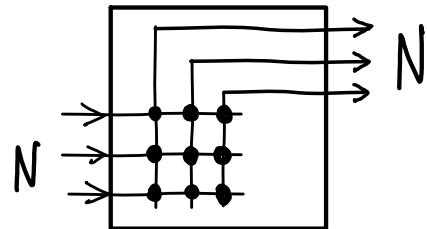
Instructions reordered.

Commits (state change) preserve semantics.

Speculation, nullifying

interconnects, switched

$N \times N$
 Crossbar
 switch



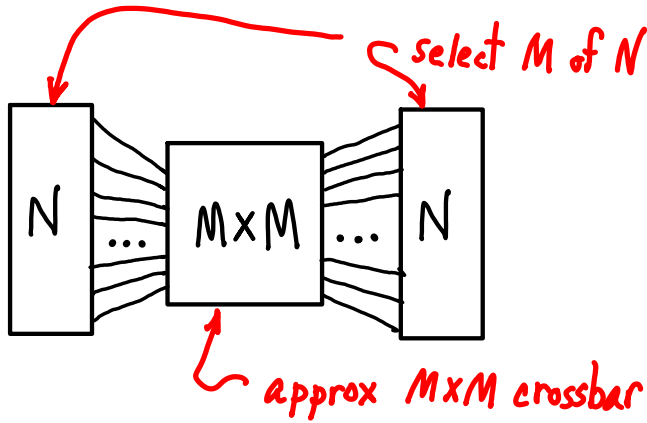
route

- any 1-1 map
- any $(N-k) \times N$ *
fan-out
- any $N \times (N-k)$ *
fan-in

* Collisionless
 (buffering?)

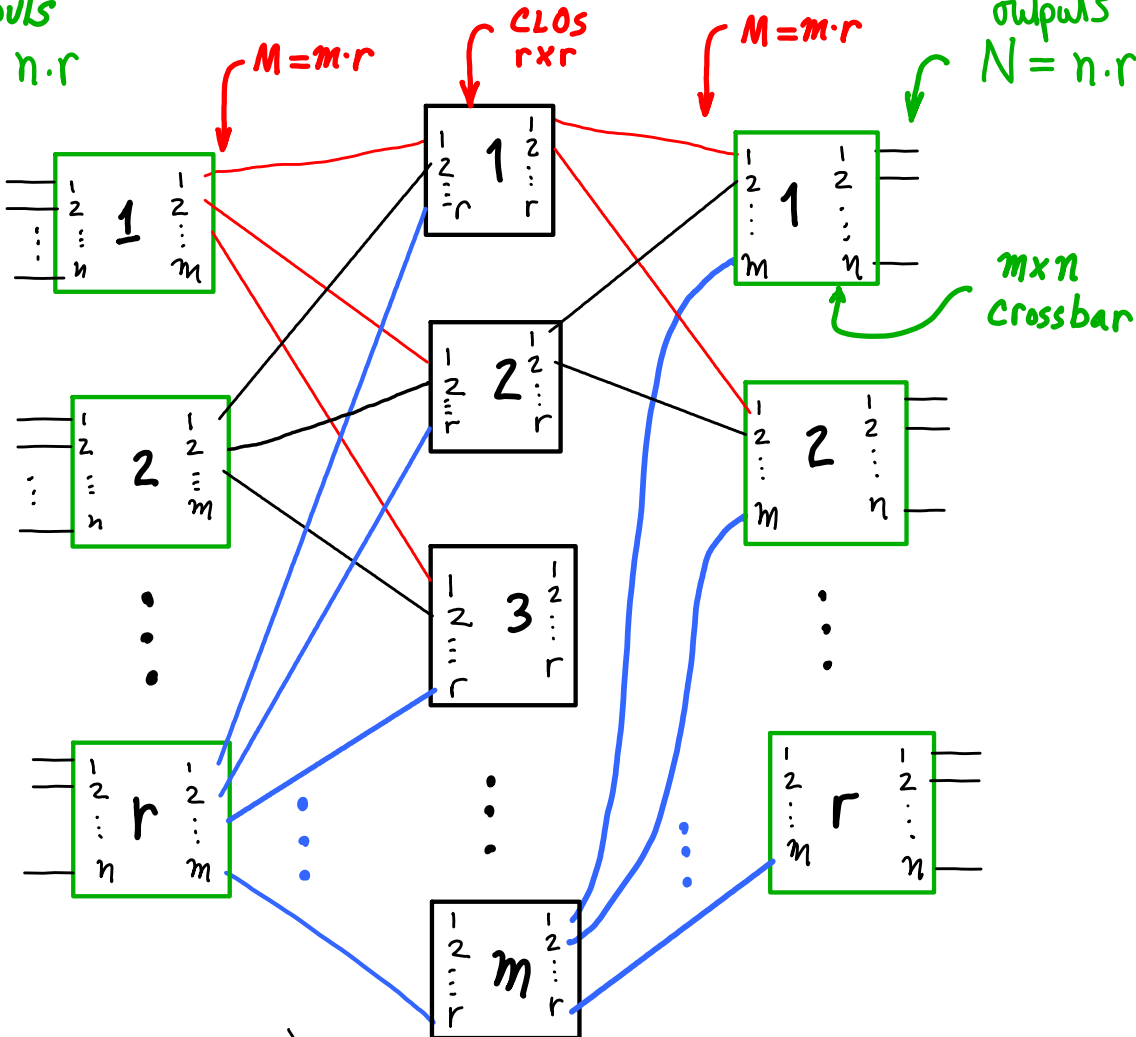
Clos

Approx
 $N \times N$

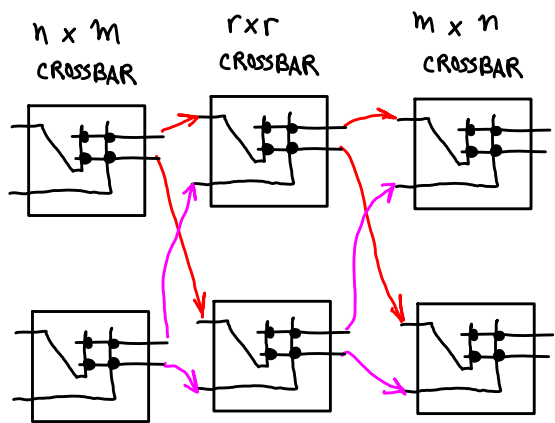


inputs
 $N = n \cdot r$

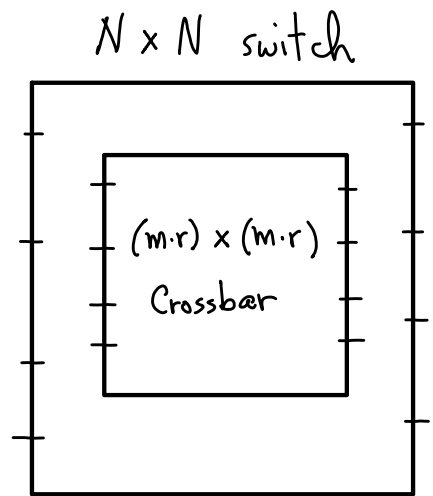
outputs
 $N = n \cdot r$



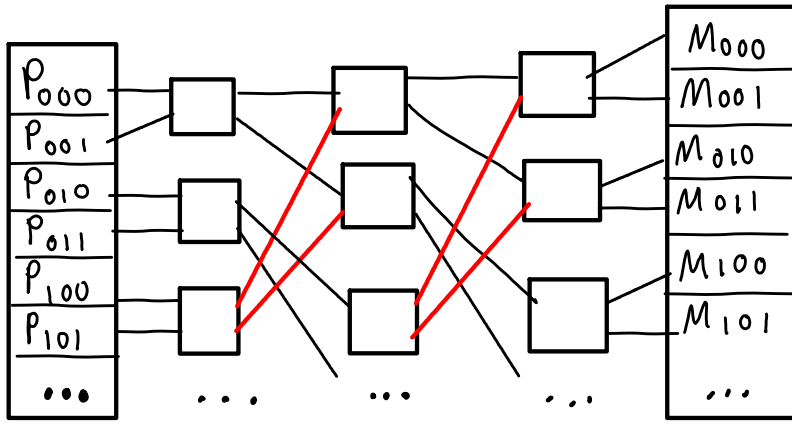
(recursive def'n)



=



Ω 8x4 interconnect



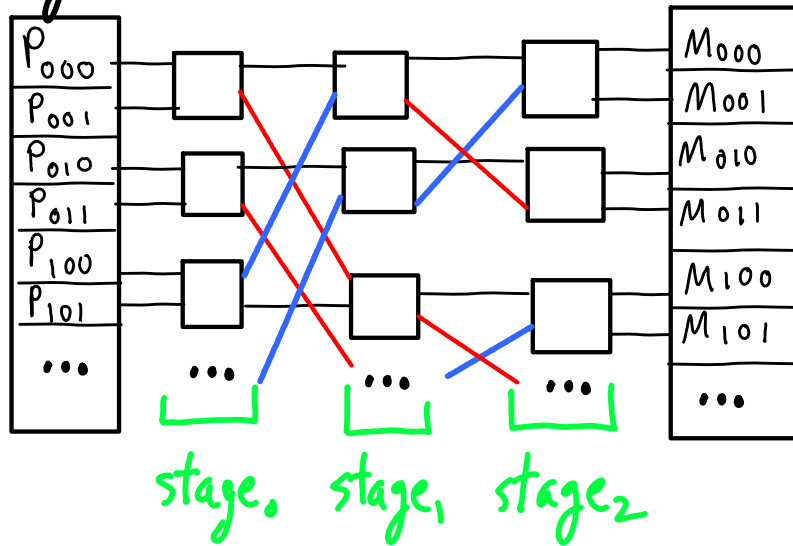
switch \Rightarrow blocking



also, can broadcast:



Butterfly 8x4



(Banyan)

n P_s

n M_s

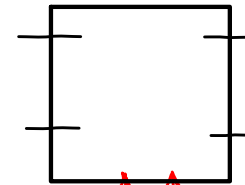
$\log n$ stages, @ $n/2$

$\log n$ address bits

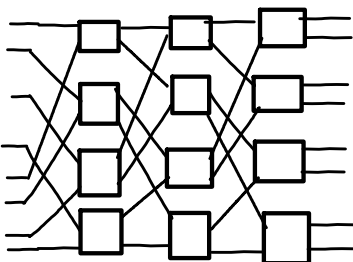
Switch at k^{th} stage:

addr₀ $a_{\log n} a_{\log n-1} \dots a_k \dots a_1 a_0$

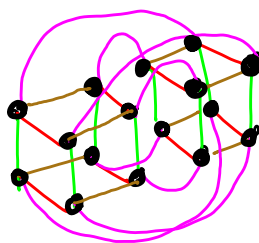
addr₁ $a_{\log n} a_{\log n-1} \dots a_k \dots a_1 a_0$



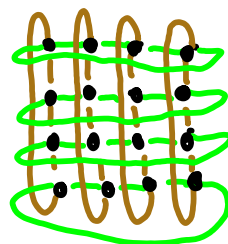
shuffle-exchange (perfect card shuffle)



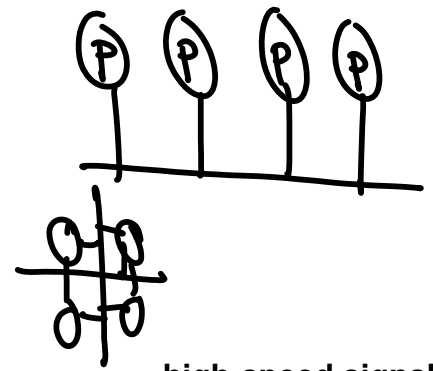
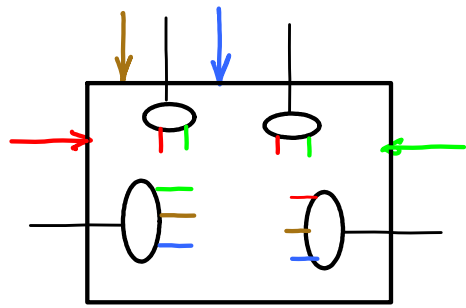
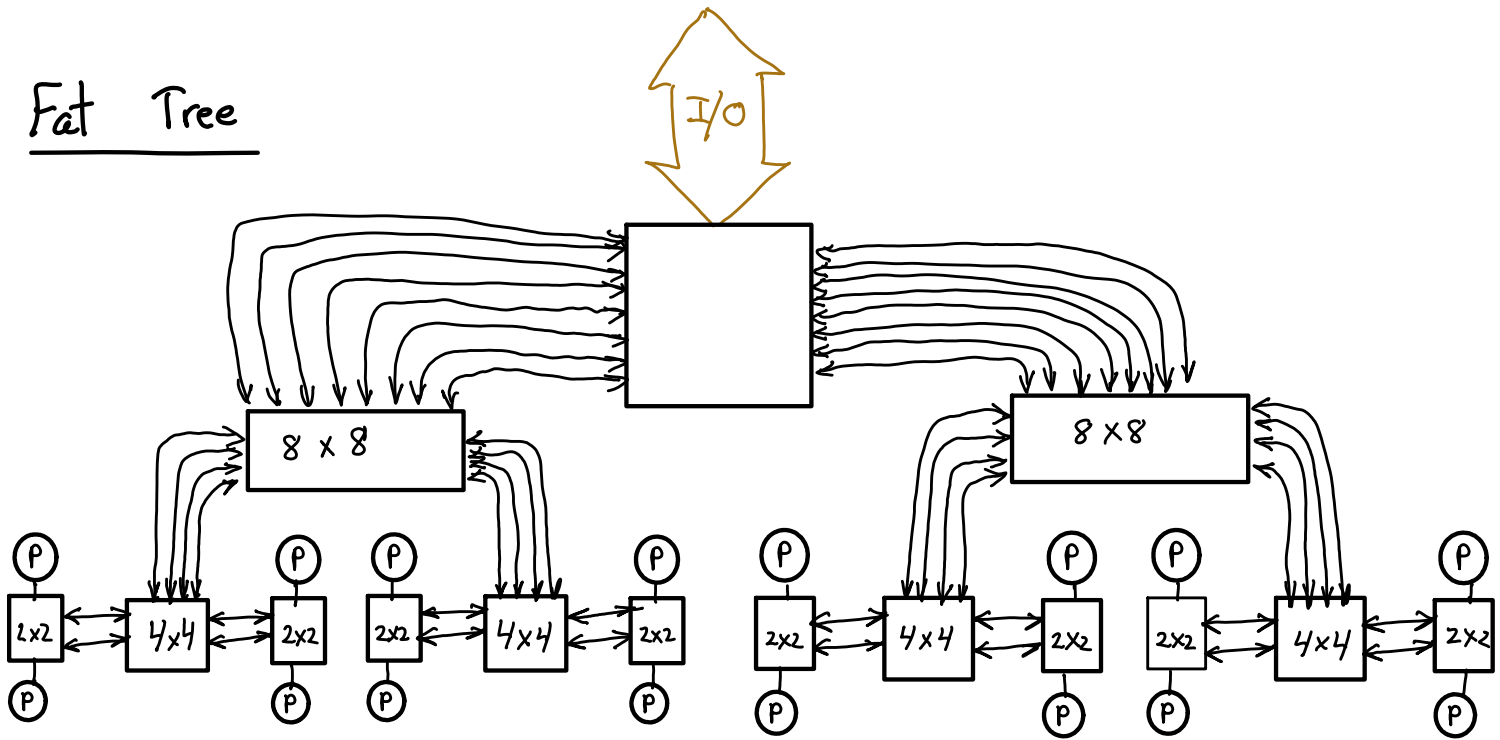
hypercube



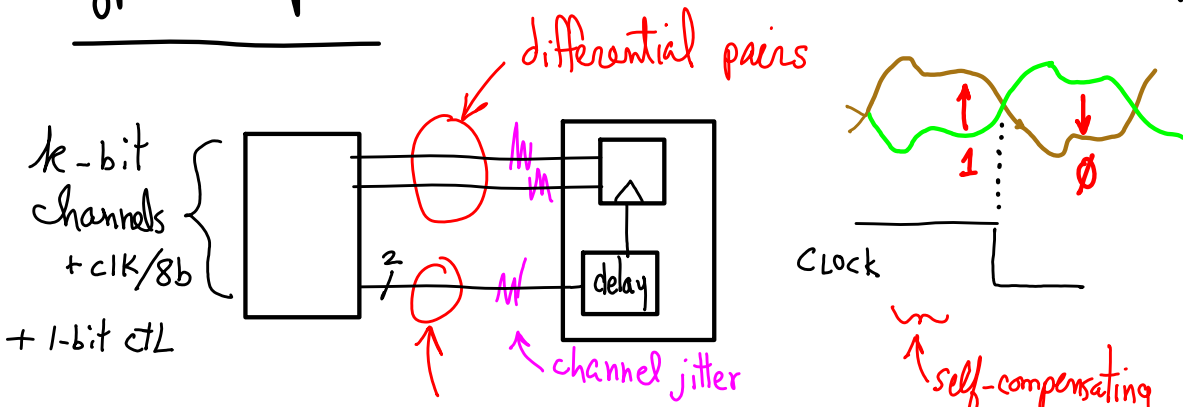
Torus



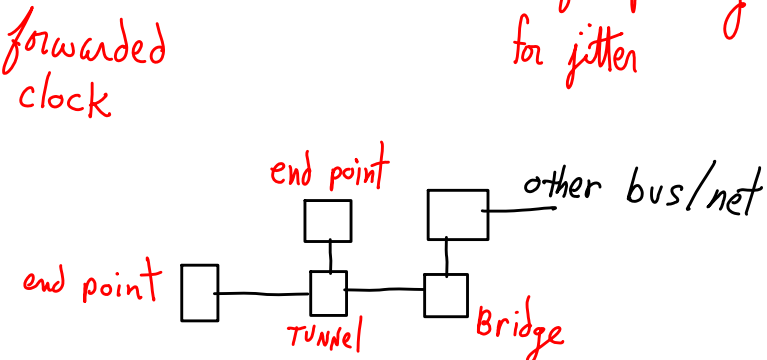
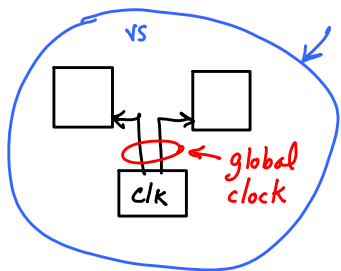
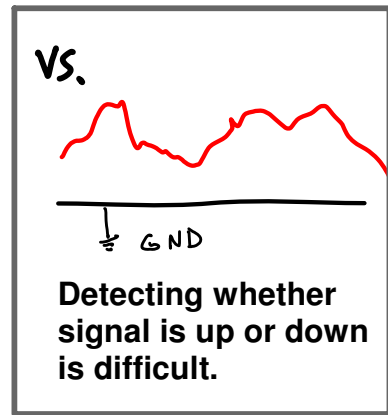
Fat Tree



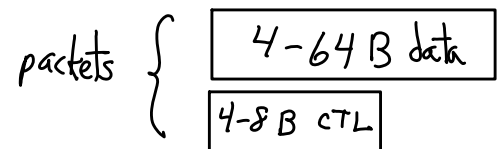
HyperTransport



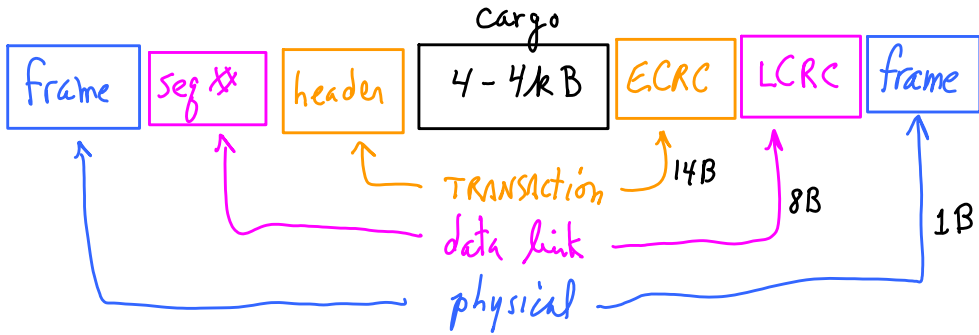
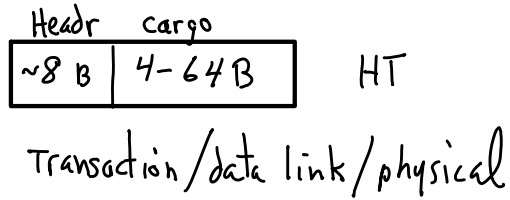
high-speed signals are bouncy (jitter), and corrupted by other signals.



daisy chained, point-to-point

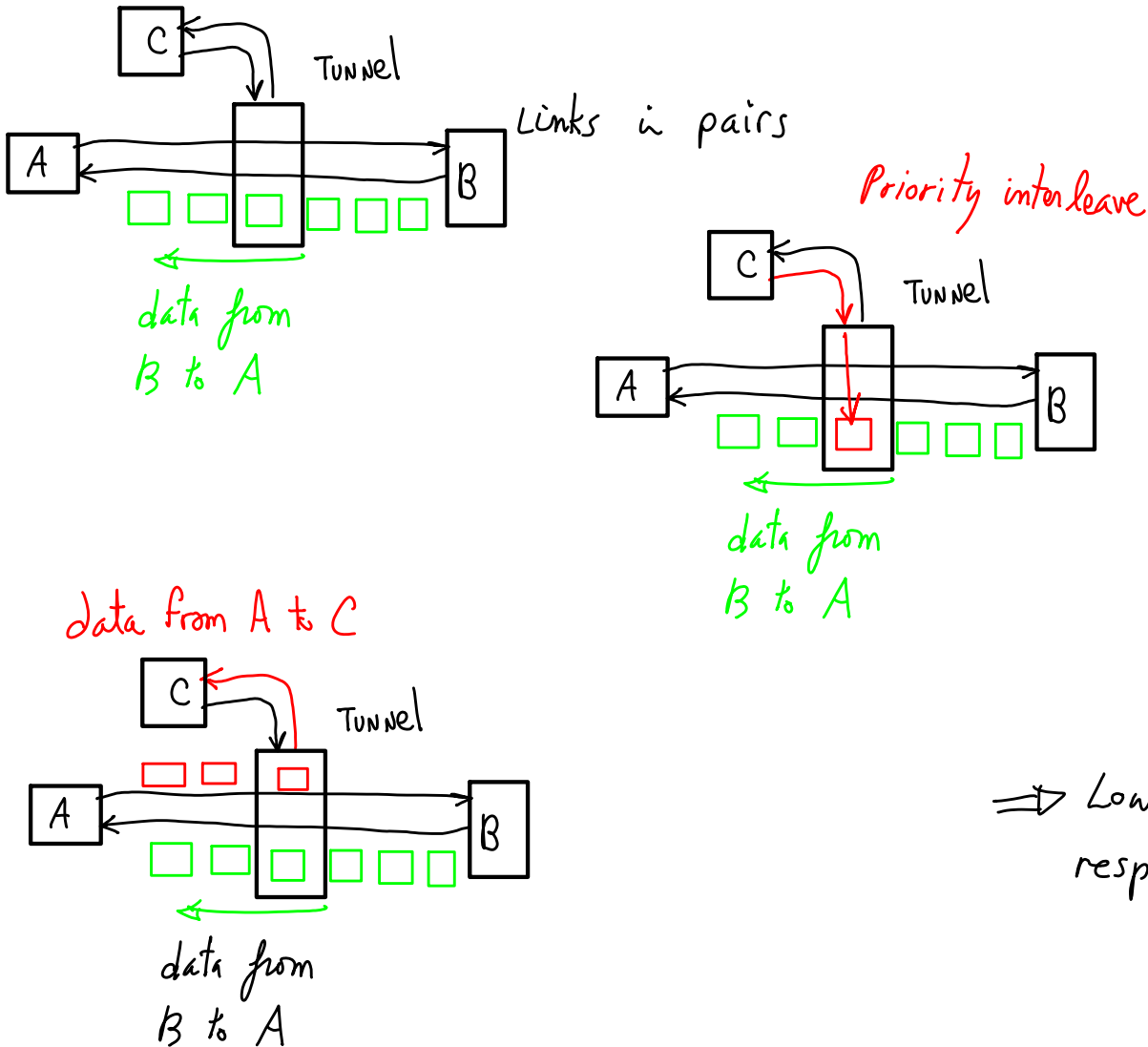


Packet overhead

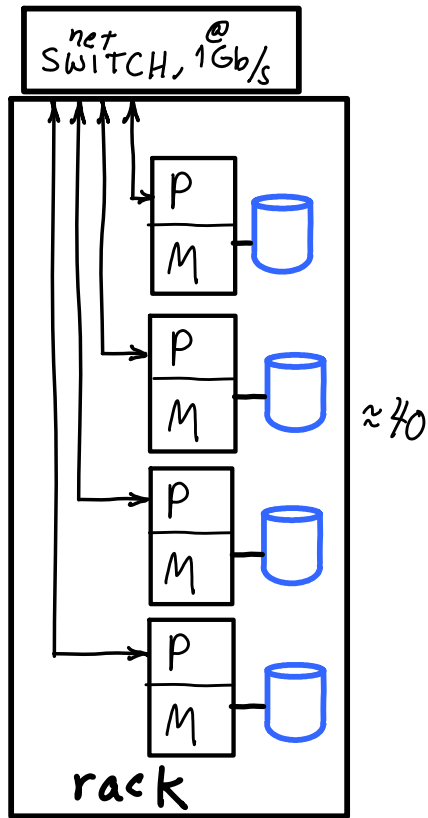


PCIe ~24 B overhead

Priority / interleaving



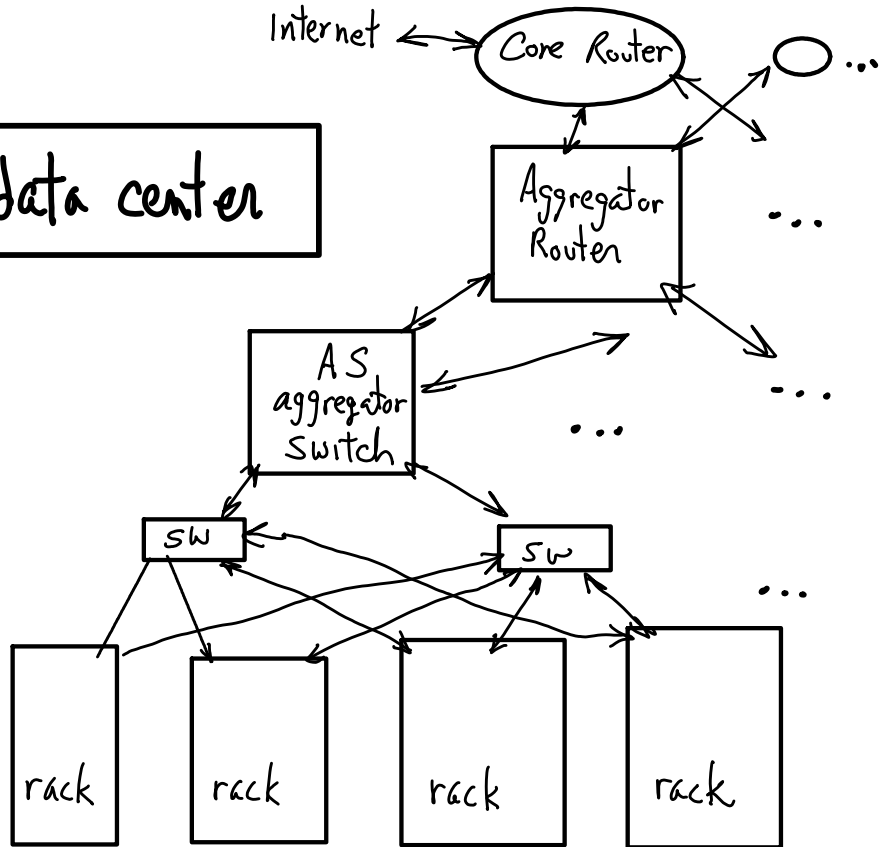
Cluster



non-shared memory,
message passing

Thread parallelism
Process parallelism
Task parallelism

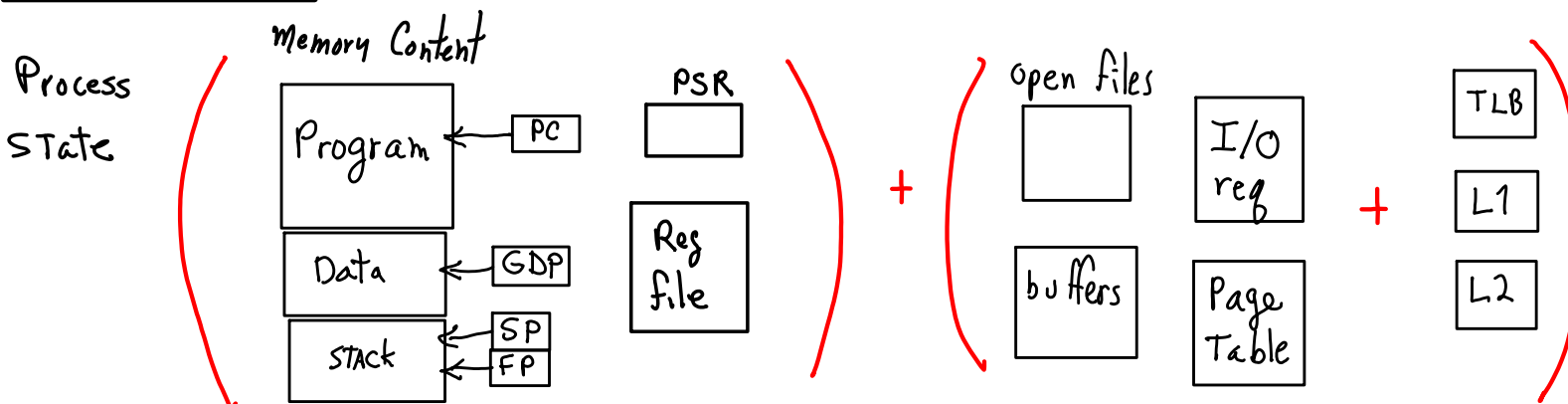
data center



- 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

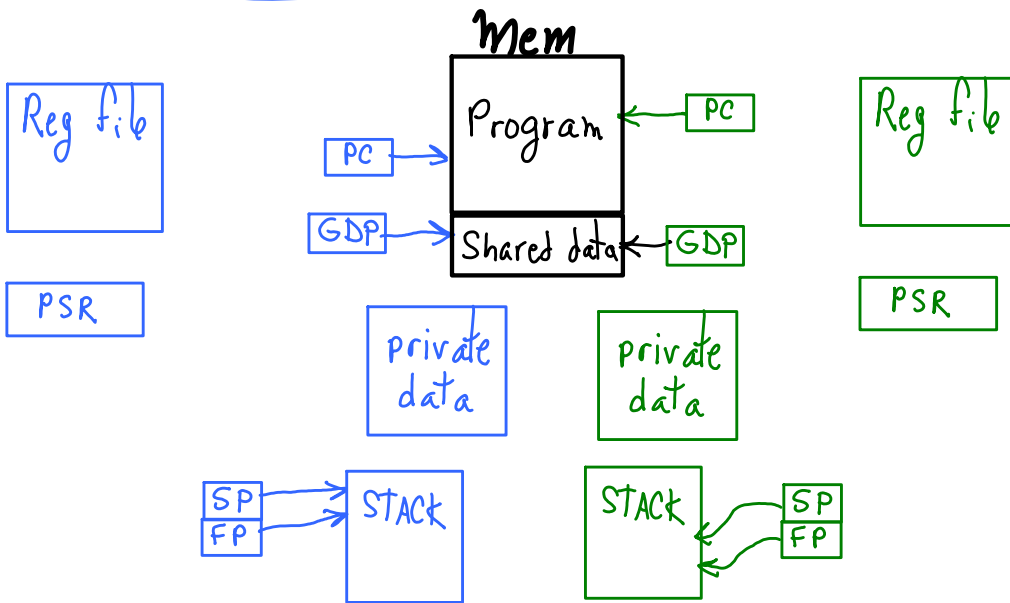
Threads

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:

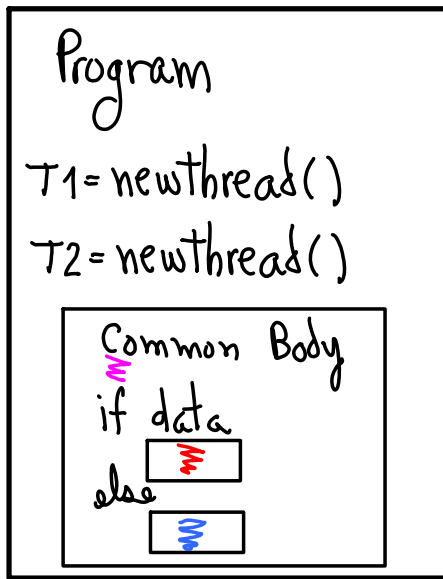
- save state T1
- load state T2

run T2 until

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

Switch context

Instruction streams



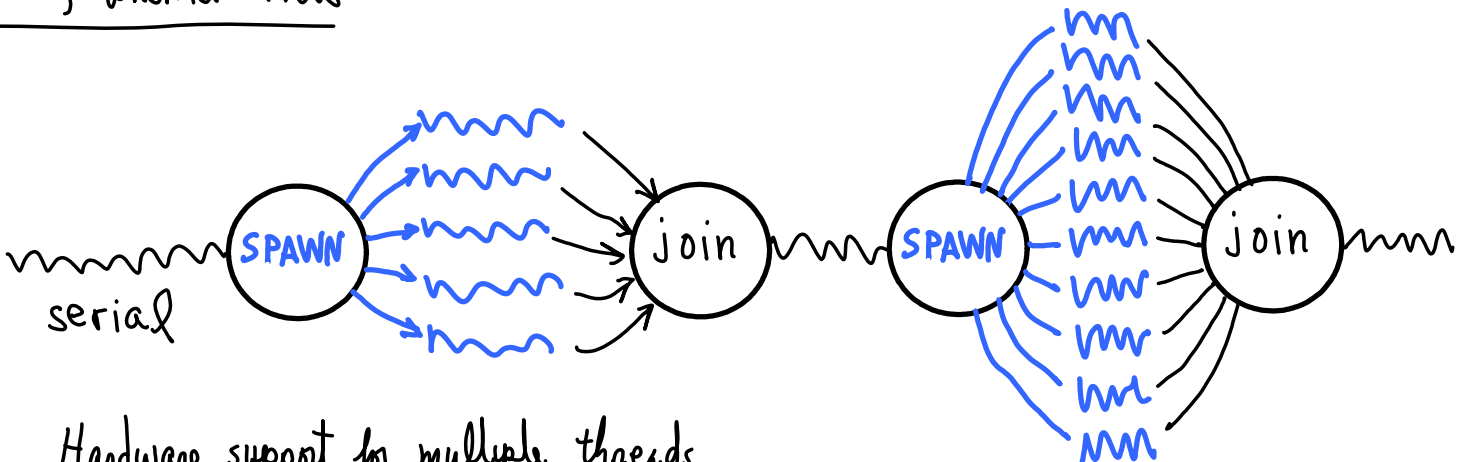
PC1



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: thread ID labeled (TID)

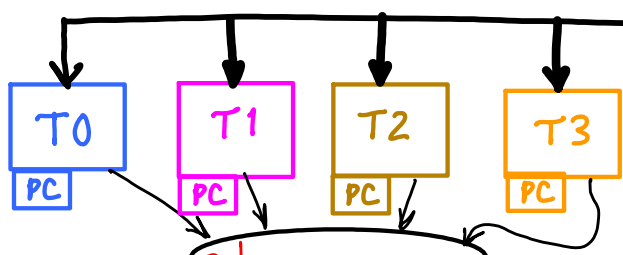
--- MULTIPLE THREADS from MULTIPLE PROCESSES: PID + TID

- Larger state to consider (page tables, file and IO tables and buffers)
 - Page Tables
 - File and IO tables and buffers
 - TLB and Cache content

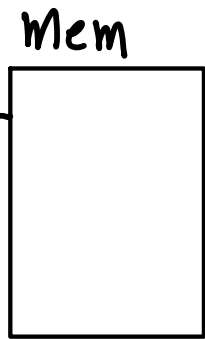
OS policy
not known to
HW designer?

Fine-grained multi-threading

Instruction buffers



Instruction fetch unit



Control

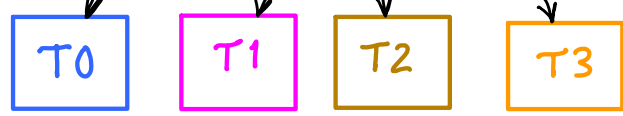
execution Unit

Round Robin Select

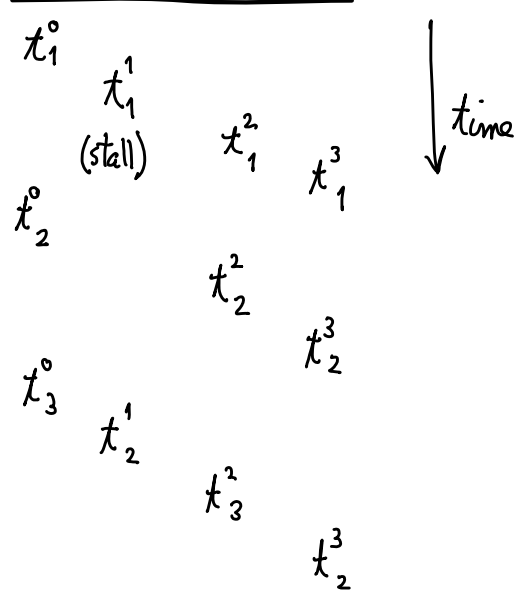
- skip stalled threads
- longer response time
- switching overhead

thread context switch per instruction

Contexts



Execution Trace

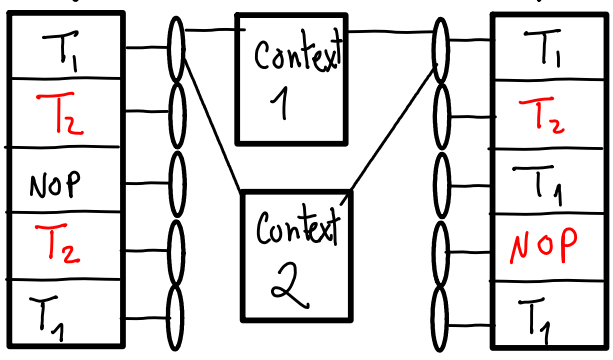


Multi-issue

Instruction Buffer

Schedule issue

execution pipe



execution pipe

Thread context switch per stage.

Coarse Grained

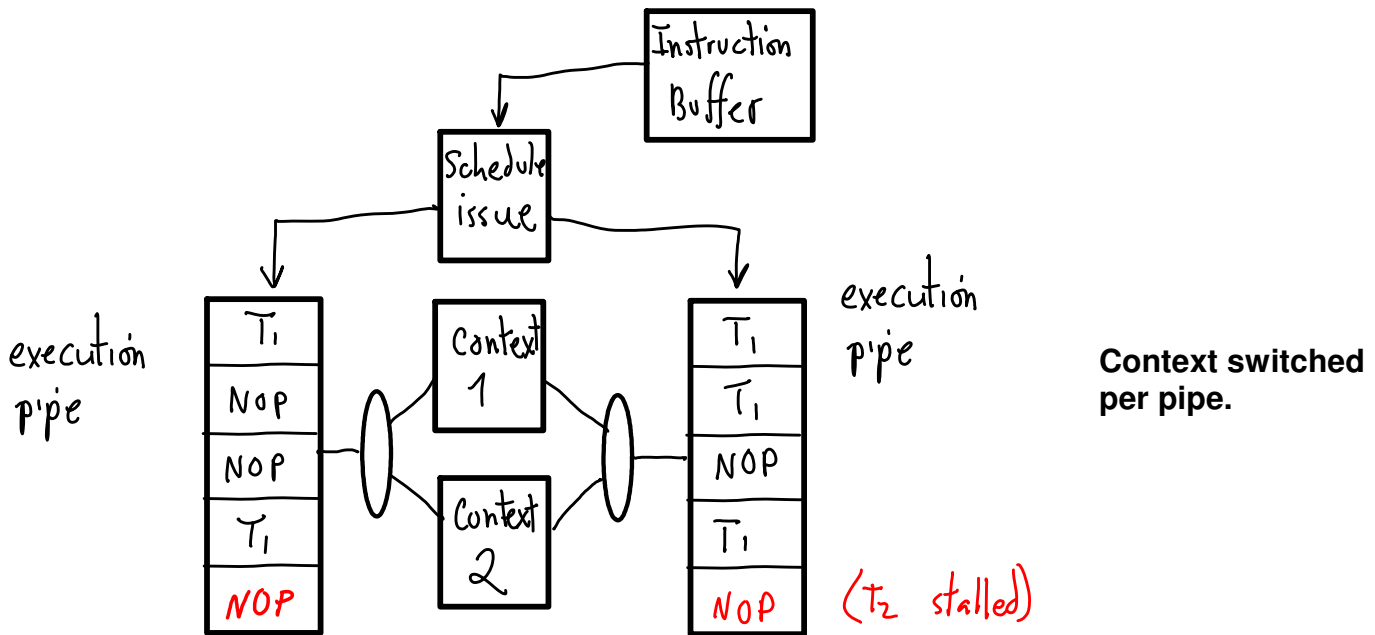
Same HW, but schedule same thread until stalled.

Execution Trace

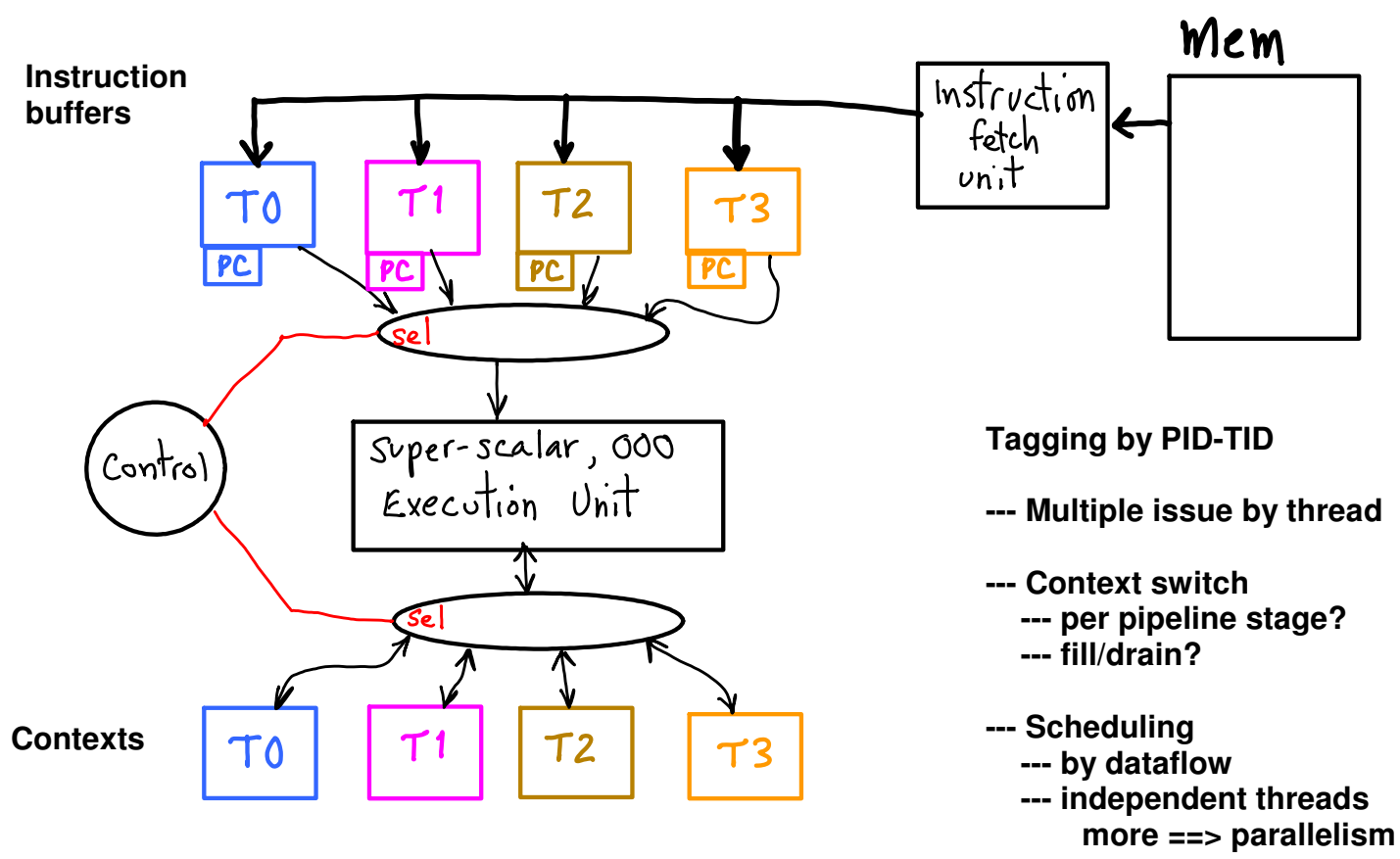
t_1^0
 t_2^0
 t_3^0
 (stall) t_1^1
 (stall) t_1^2
 t_2^2
 t_3^2

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

Multi-issue



(Simultaneous)(Hyper) Threading



time ↓

T1	T1	T2	T3
T2	T0	T0	NOP
NOP	T1	T1	T1
NOP	T2	T3	T1
T3	T3	T2	T2

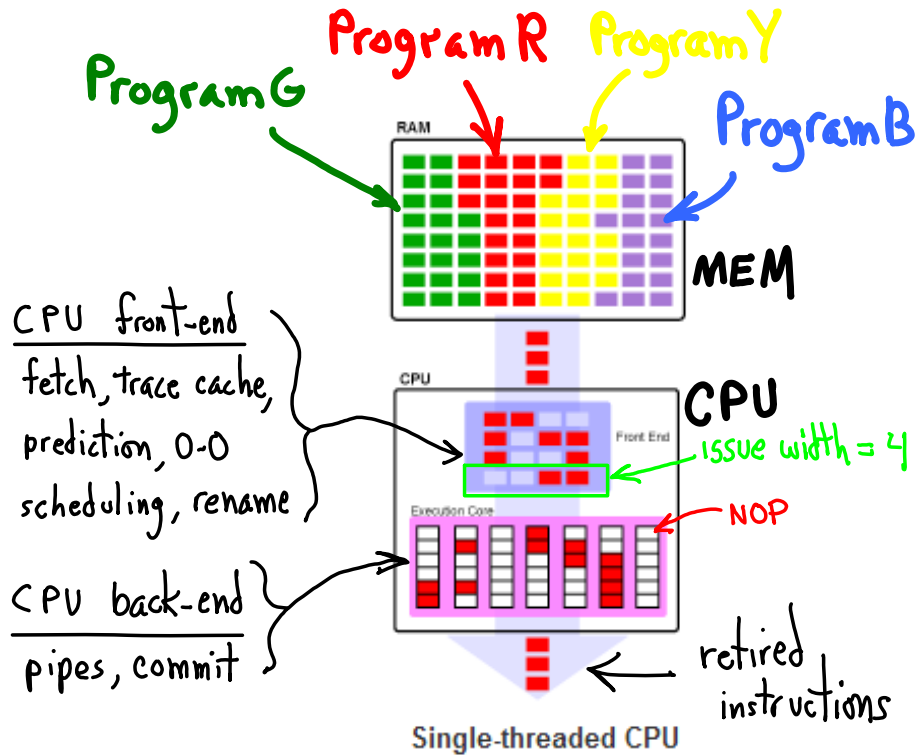
← multiple issue

- OOO scheduling across threads
- more filled slots

Single threaded execution

Multi-tasking

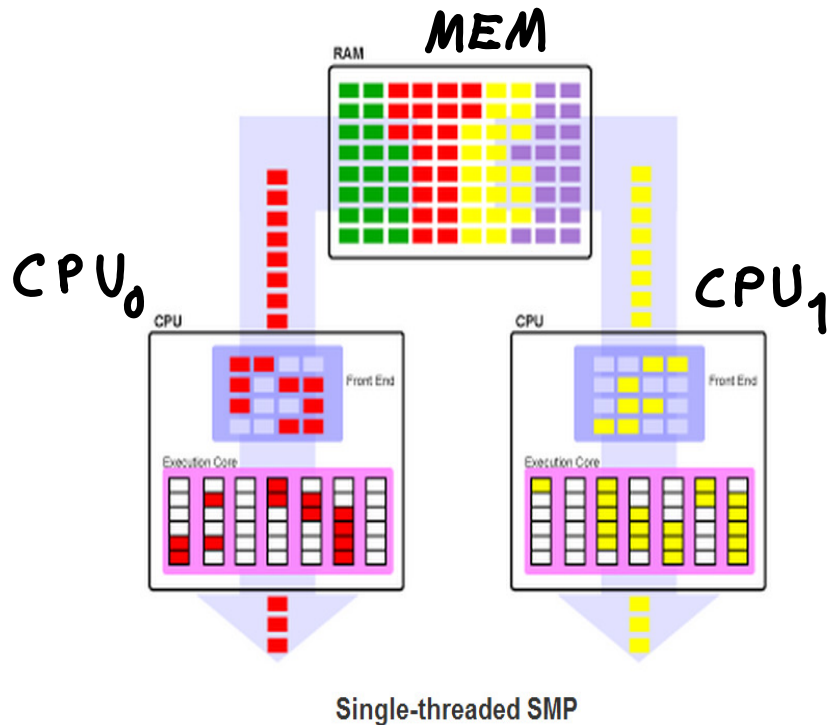
- Multiple concurrent execution (not simultaneous)
- Memory shared but separate (virtual)
- CPU time-multi-plexed
 - cooperatively, pre-emptively, IO
- Process context switching
 - drain/fill (pipes, caches, TLB, ...)
- Extract ILP from single stream
 - Unused issue slots
 - pipeline bubbles/stalls



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

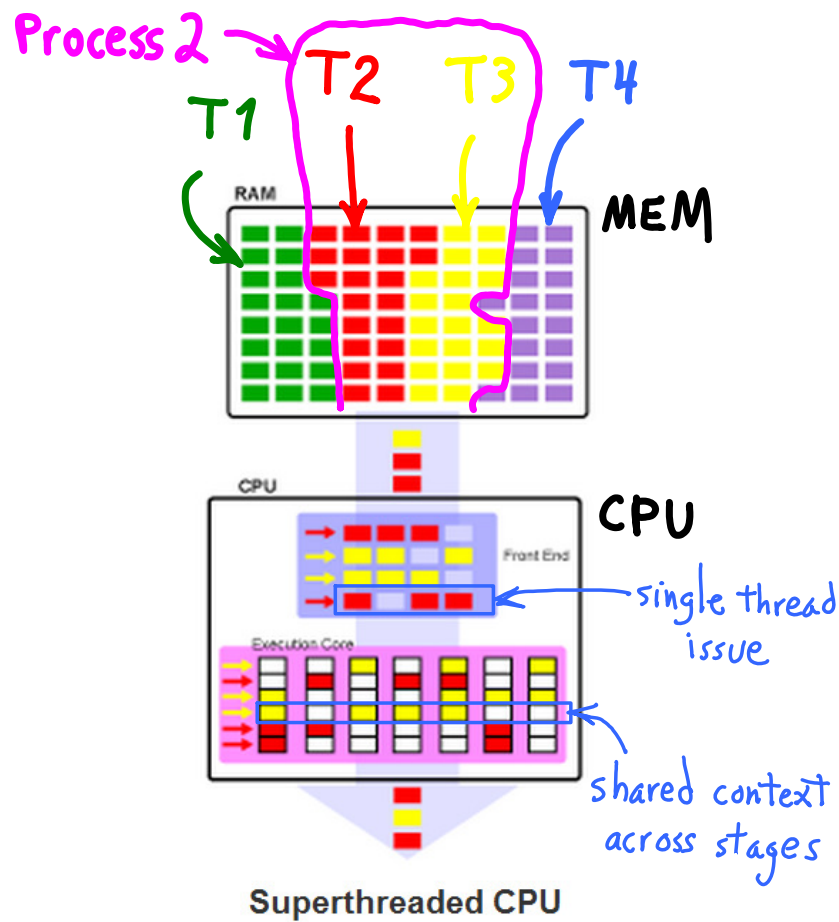
SMP, Symmetric Multi-Processing

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



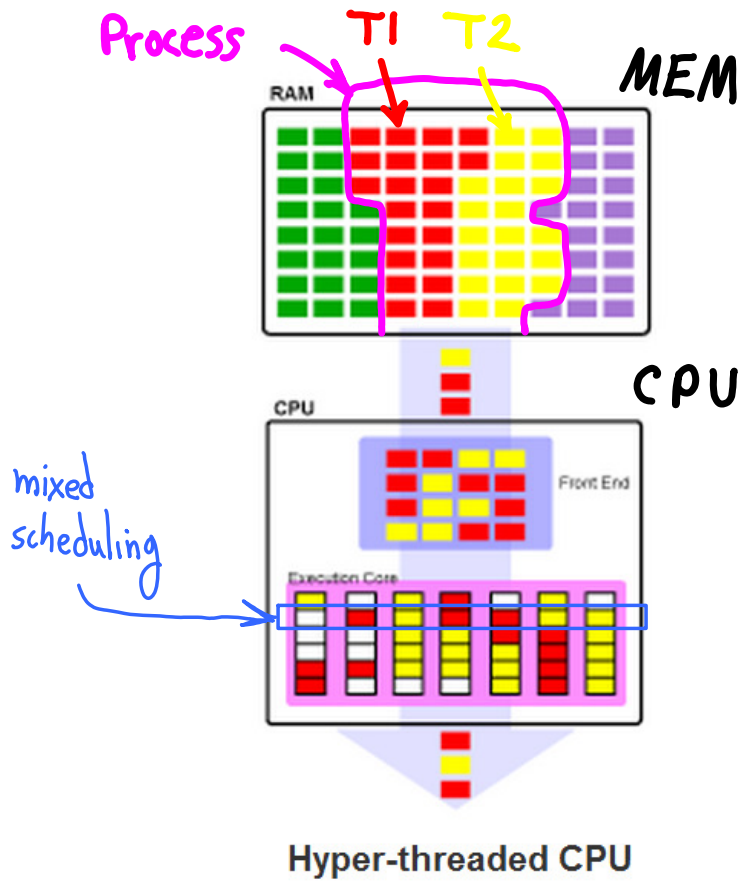
Multi-Threaded (Superthreading)

- 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

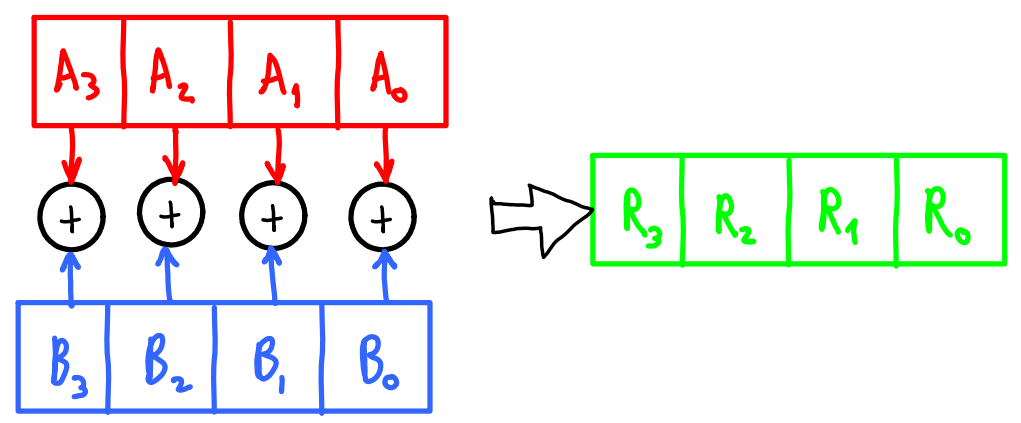


Modification of existing O-O-O CPU
 => 10% added cost, $\rho > 2$

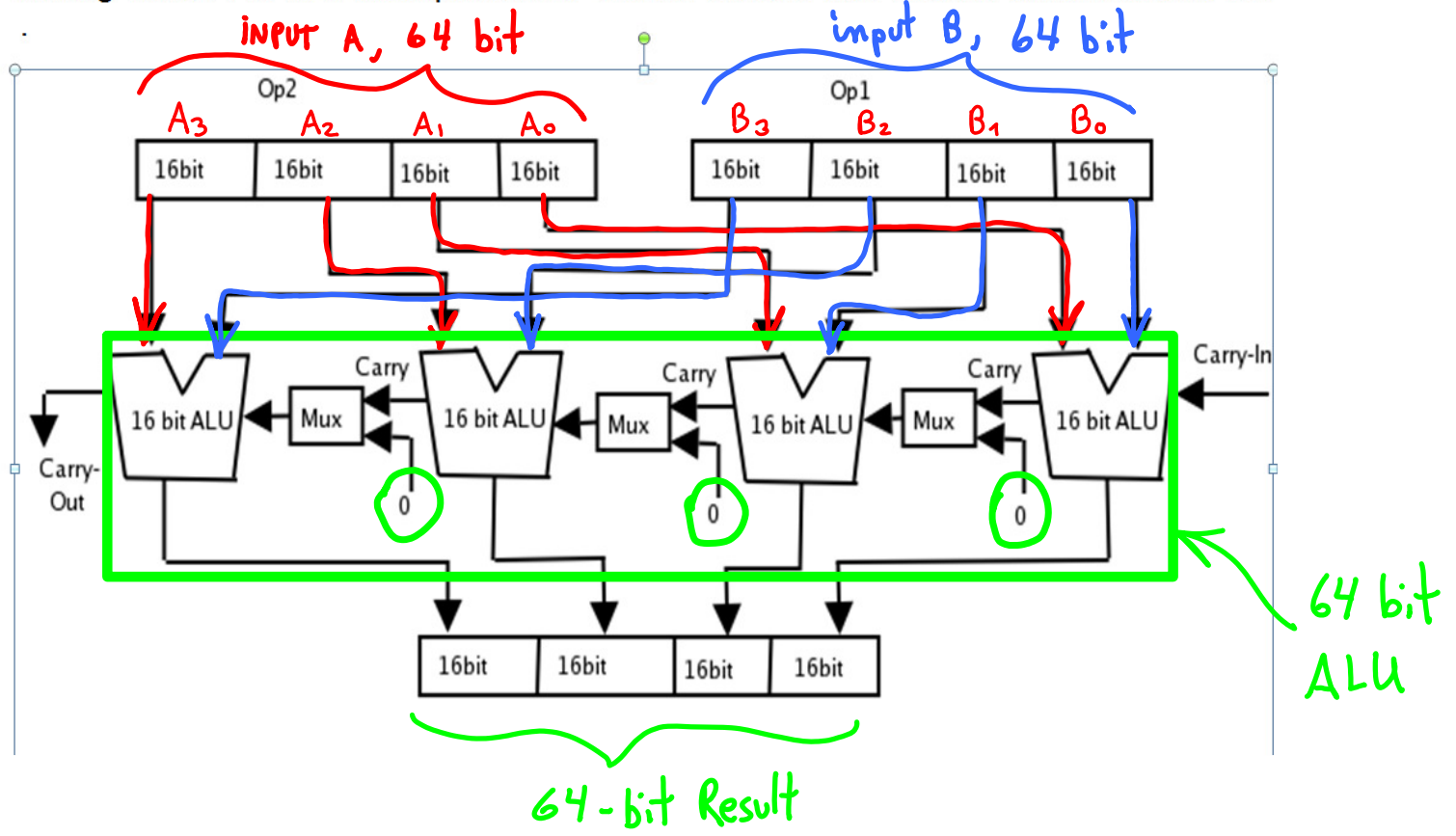
SIMD, cheap

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

$R_i \leftarrow (-A_i) + B_i$

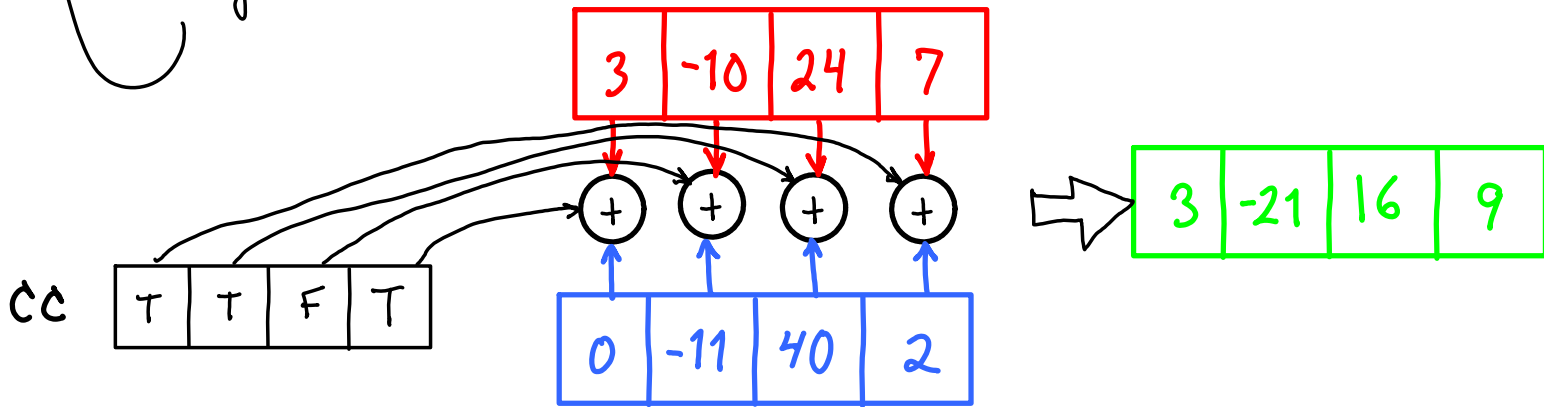
endif

\Rightarrow

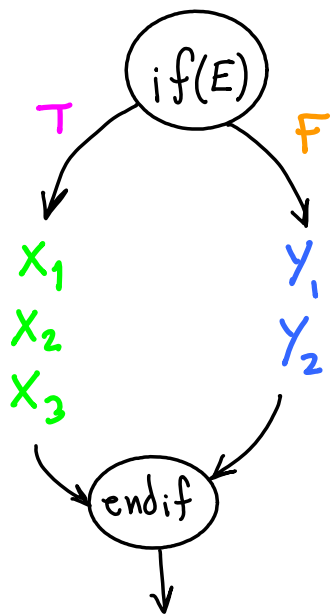
$(A - B) \Rightarrow CC$ (predicate vector)

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

