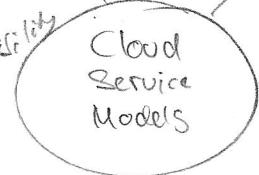


PaaS: Platform as a Service

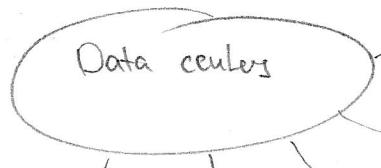
- provide SW platform
- e.g. Google's App-Engine
- avoid worrying about platt. scalability

SaaS: Software as a service

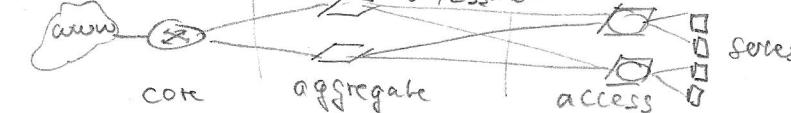
- Provide licensed apps
- no cost of maintenance, installation

IaaS: Infrastructure as a service

- raw compute, storage, network → AWS

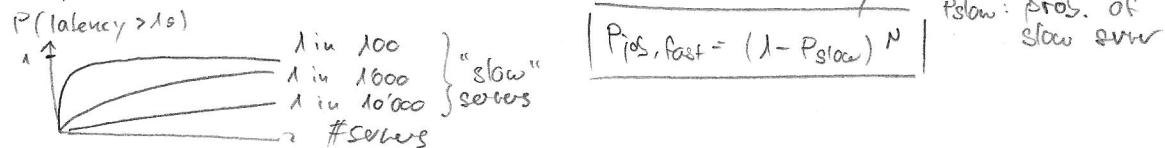
Specialized

ML cloud HW (Google TPU)

Topology:

$$\begin{aligned} \text{TCO} &= \text{Total Cost of Ownership} \\ &= \text{CapEx} + \text{OpEx} \\ &\quad \text{HW} \quad \text{electricity} \end{aligned}$$

Latency Op depends on all servers finishing. A few slow servers  
large scale → more prone to high total latency



Reliability FIT Failure in time, failures per billion h of op  
MTTF mean time to fail, time to first error  
MTTR mean time to repair, detect & fix fault

$$\text{Steady state availability} = \frac{\text{MTTF}}{\text{MTTF} + \text{MTTR}}$$

2010 → 2015: Storage &amp; network ×10, CPU not: why?

Performance

$$P_{\text{dyn}} \propto C \cdot V^2 \cdot f$$

Why multicore?

Replicate smaller cores with only little less performance increases total performance and keeps power down.

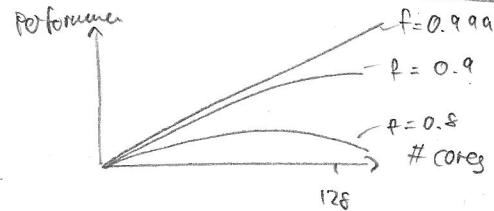
	1 core baseline	2 core	4 core
Core A	A	~A/2	~A/4
Core P	W	~W/2	~W/4
Chip P	W+O	W+O'	W+O''
Core Perform.	P	0.9 · P	0.8 · P
Chip Perform.	P	1.8 · P	3.2 · P

Amdahl's law

f: fraction running in parallel  
1-f: non-parallelizable functions

$$\text{Speedup} = \frac{1}{(1-f) + \frac{f}{n}} \xrightarrow{n \rightarrow \infty} \frac{1}{1-f}$$

↳ ignores (power) cost of n cores  
↳ more cores → each is slower  
↳ parallel speedup < n



Figures  
response time  
throughput

elapsed time  
- total, incl. I/O, OS, idle

$$\frac{\text{CPU time}}{\text{Spent processing}} = \frac{\text{CPU clk cycles}}{\text{clock rate}}$$

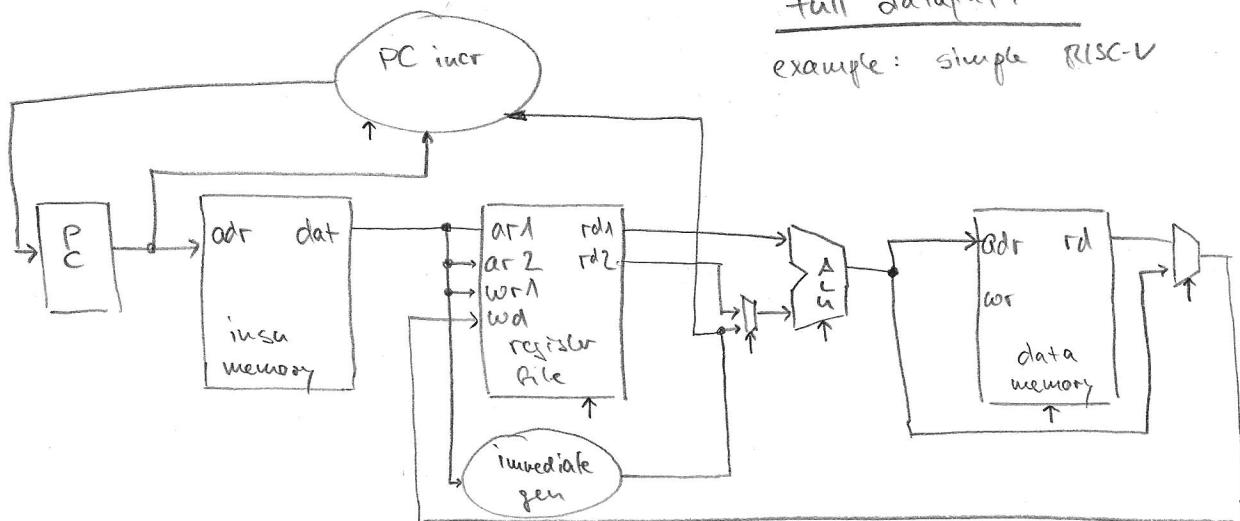
cycles per instruction  
↳ determined by program,  
ISA, compiler

$$\text{CPI} = \frac{\text{CPU time} \times \text{fclk}}{\text{Instr. count}}$$

avg. cycles per instruction  
↳ CPU hardware

Performance depends on

- ✓ Algorithm
- ✓ Programming
- ✓ Compiler
- ✓ ISA



## Full datapath

example: simple RISC-V

PC Program counter instr. memory  $\rightarrow$  fetch insc.

RF register file contains temporary data

ALU arithmetic logical unit calculates stuff used for - result for data  
- new adr for d/st  
- branch compare

$\uparrow$  control signals

you Neumann architect can't fetch & run at same time  $\rightarrow$  you Neumann bottleneck

parallel improves performance

- run different tasks at same time

example

IF: insc fetch  
ID: insc decode and RF read  
EX: exec op or calc address  
MEM: access memory  
WB: write back to RF

speedup if balanced all stages take same time

$$\text{speedup} = \sqrt{N \text{ stages}}$$

non-ideal due to wait & hazards

Id x1, 100(x4)  
Id x2, 200(x4)  
Id x3, 400(x4)

IF	ID	EX	MEM	WB
IF	ID	EX	MEM	WB
IF	ID	EX	MEM	WB
IF	ID	EX	MEM	WB

RISC-V designed for pipelining

- All insc 32-bit
- $\rightarrow$  easy IF/ID in logic
- $\rightarrow$  x86 has 1...17 byte insc
- Few & regular insc fmt
- Id/SD addressing

pipeline

- conflicting use of resource

- example simultaneous Id/st  
- causes "bubble" - pipelined arch need separate insc/dat mem  
or separate caches

HAZARDS situation that prevents starting next insc in cycle.

control fetch depends on last insc (branch)

data insc depends on completion of data access  
of prev. insc

add x1, x0, 11  
sub x2, x1, 83

IF ID EX MEM WB  
stall s4 2 IF ID EX MEM WB

solve sy forwarding (Sypos)

- requires extra connection in datapath
- doesn't always work: Id x1, 0(x2)  $\rightarrow$  sub x4, x1, x5

code scheduling

- compiler re-arranges insc to mitigate data haz.

exceptions & interrupt unexpected events  
o save PC & prog state, jump to handler

Stall on branch wait for b insc complete before reading next insc

branch predict pred. outcome, only stall if wrong

static: based on typical branch behav.

dynamic: HW measures actual branch behav.

$\rightarrow$  used in deeper pipelines

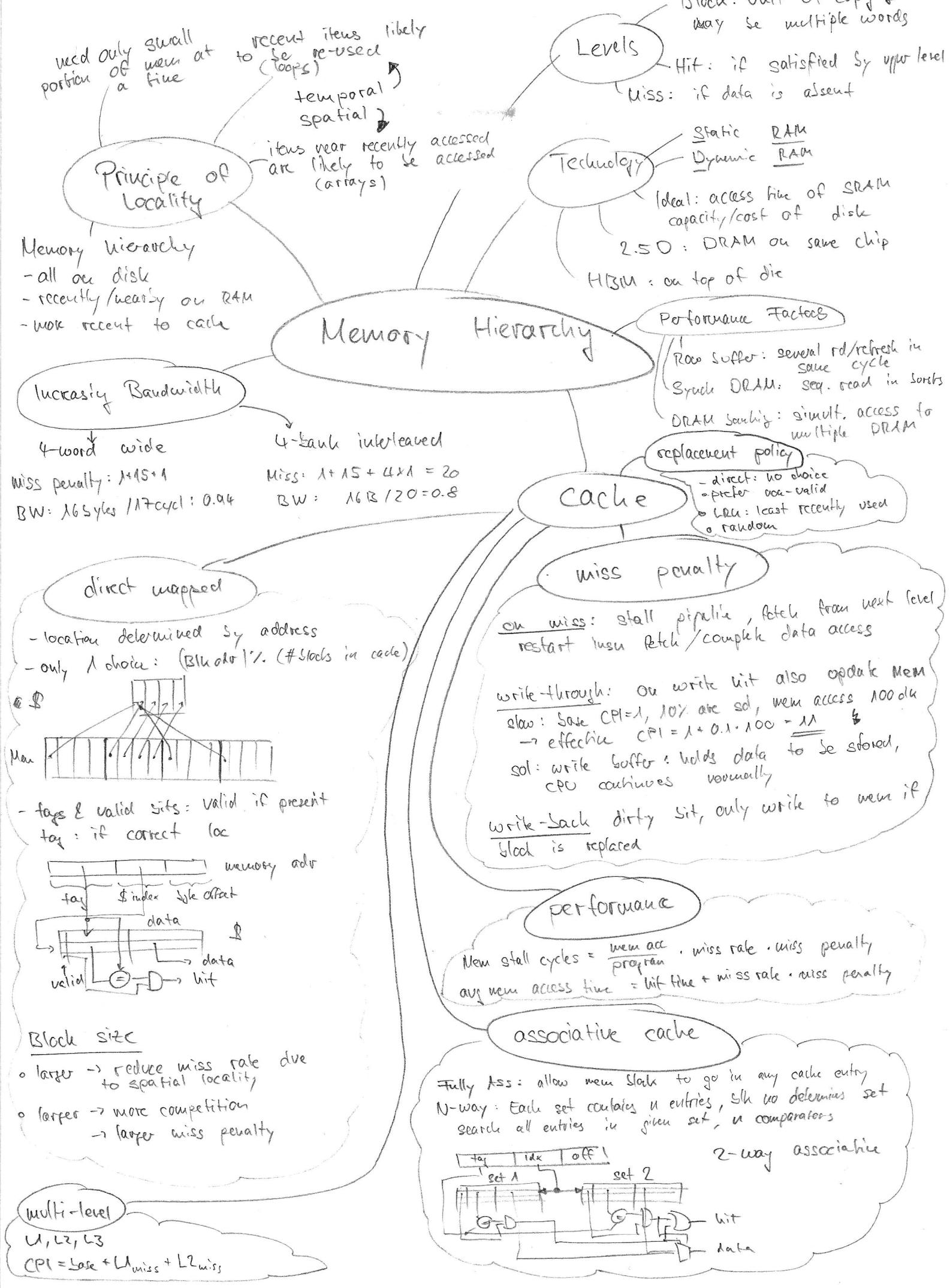
1/2 bit predictor, only change on 2 mispred.

$\rightarrow$  cache branch targets

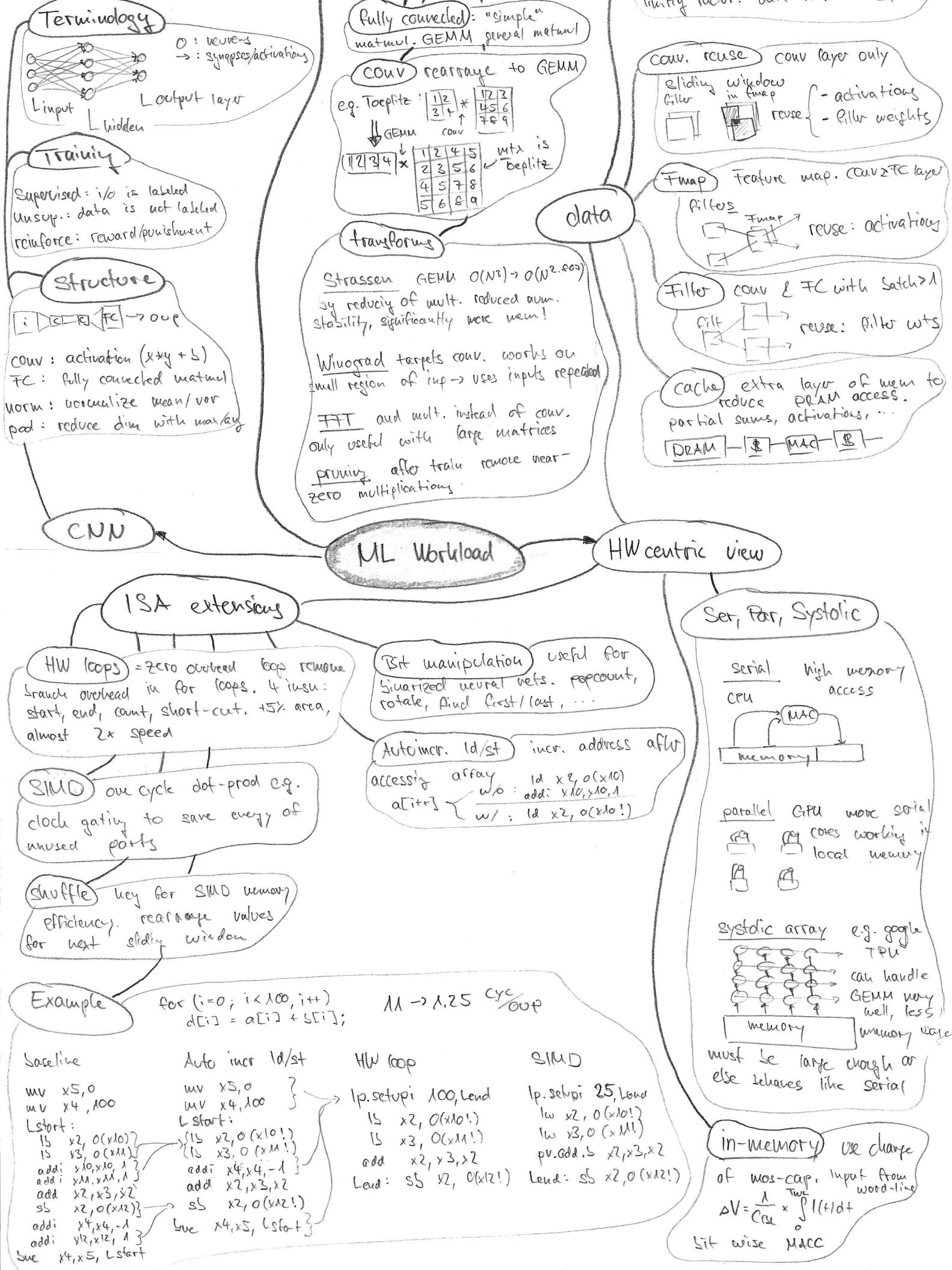
## 1: Memory Architecture

1 soc

2



## L2: ML workloads



# L3: Out of Order Processors

SOC

3

## Terms

Flynn Bottleneck: Scalar pipeline is  
by desired upper bound to IPC ≤ 1

ILP: Instruction level parallelism

IPC: Instructions per cycle

Superscalar: IF DE EX WB

D pipeline depth  
IP: Intra Parallelism

## Out of Order Processors

## Processors with Superscalar

Forwarding is no longer effective

	1	2	3	4	5	6	7	8	9	10	11	12
ADD F0, F1, F2	F	D	E+	E+	E+	W						
MUL F2, F3, F2	F	D	I*	I*	I*	E*	E*	E*	E*	E*	E*	W
SQRT F0, F1, F4	F	P*	P*	D	(E+)	W						

E+ is busy so suff most stall.

Solution: ILP

$$\begin{cases} r1 \leftarrow r2 + 1 \\ r3 \leftarrow r1 / r7 \\ r4 \leftarrow r0 - r3 \\ r11 \leftarrow r12 + 1 \\ r13 \leftarrow r12 + 1 \\ r14 \leftarrow r13 - r20 \end{cases}$$

ILP=1

$$\begin{cases} r1 \leftarrow r2 + 1 \\ r3 \leftarrow r1 / r7 \\ r4 \leftarrow r0 - r3 \end{cases}$$

ILP=2

OoO exposes more ILP reordering free  
insu results in ILP=2. → register renaming  
required for reading dependencies.

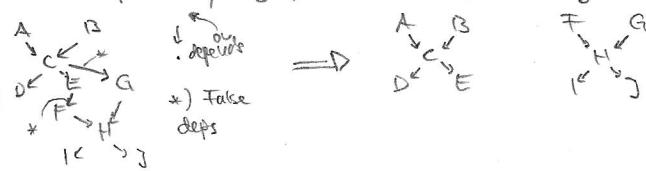
## Register Renaming

Dependency exists indep. of HW. It's only a hazard if  
HW has to deal with it

True data dependency: RTW  $r1 = r2 + r3$  prevent reordered  
 $r4 = r1 + r5$

False dep: WAW  $r1 = r2 + r3$  WAR  $r2 = r1 + r5$  WAW, WAR can  
 $r1 = r4 + r5$  be fixed by renaming

Data dependency graph



Reg names are arbitrary and only need to be  
consistent between writes

$r1 = \dots$   
 $\dots = r1$   
 $\dots = \dots r1$

If reg is W, assign it - until it is written  
to physical reg

$\Rightarrow$  again, translate it

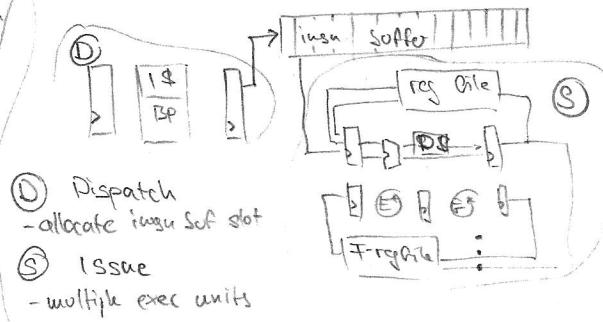
$\Rightarrow$  false dep. resolved!

Add r2, r3, r1  
SUS r2, r1, r3  
mul r2, r3, r3  
div r1, 4, r1

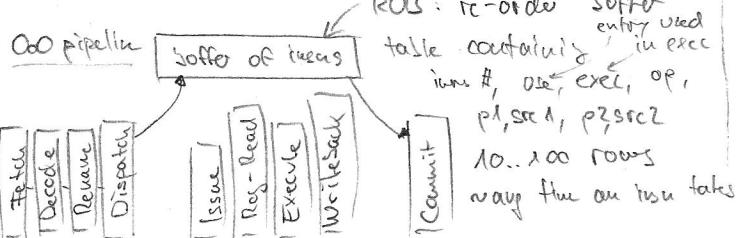
Keep track in map table  
and free list.

## OoO & dynamic scheduling

- 100% in HW at runtime
- Fetch many ins in insu window
- reorder to avoid false dependencies (concurrent > 100%)
- exec insu ASAP
- commit insu in order



Issue policy: If multiple insu are waiting, which to choose? → HW intensive



## Speculative Execution

Branch prediction: Target address speculation → BTB  
- Condition resolution  
- Smith predictor, HW complex

Trace-Cache: Hold out taken br by tracing ins as  
they commit to fill buffer

Loop-Scheduler: cache target itself instead of address

Mispredict: roll-back ROB - recover register renaming

Speculative ID/ST: - store buffer

- exec all ID/ST in program order

st x1, (x2)  
Id x3, (x4)

guess x2=x4  
guess x2=x4

## Recap

3 ideas: renaming, speculative exec, in-order commit  
complex, enables: OoO exec, multi-issue, loop control, ~  
Exec ASAP, commit in order

exception before commit? → flush pipeline

## L3: Virtual Memory

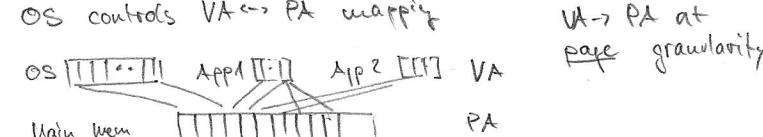
How do apps share main memory?

- ↳ Goal: Each app thinks it has infinite mem
- Solution: treat memory as "cache", store overflowed as "swap" on disk, add lvl of indirection

### Virtual Memory (VM)

VA: virtual addr PA: physical addr

OS controls VA  $\leftrightarrow$  PA mapping



→ Protection: Apps can't access each other's memory

### Address Translation

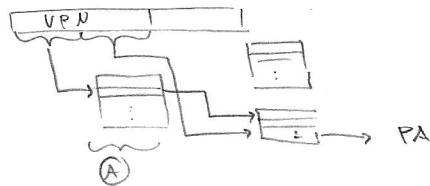
VA:	<table border="1"> <tr> <td>VPN</td><td>POFS</td></tr> <tr> <td>translate ↓</td><td>↓ don't clamp</td></tr> <tr> <td>PPN</td><td>POFS</td></tr> </table>	VPN	POFS	translate ↓	↓ don't clamp	PPN	POFS	VPN: Virtual page number POFS: page offset PPN: physical page number
VPN	POFS							
translate ↓	↓ don't clamp							
PPN	POFS							

Each process has a PT (page table) that maps VP  $\rightarrow$  PP or to disk (swap)

### Multi Level PT

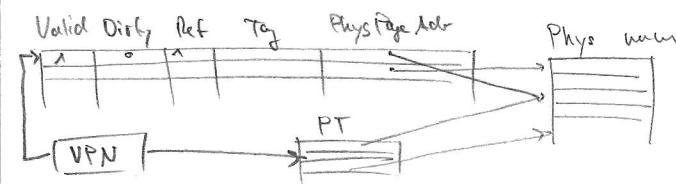
Problem: PT gets very large  
Solution: Hierarchy

large parts of ① are unused  $\rightarrow$  2nd level entries need not exist and entry in ② is null.



### Fast translation using TLB

Use cache to store PTE (PT entries) in CPU, these accesses have good locality. They are called Translation Look-aside Buffer (TLB). Typical 16-512 PTEs, 0.5-1 cycle, 10-100 nsec, 0.01-1% miss rate.



- TLB is 4+ fully associative between CPU & L1\$ / D\$
- Can perform TLB/\$ access at same time  $\approx$  10-100 nsec & needs tag only at very end
- PT "walker": on TLB miss HW searches/walks table

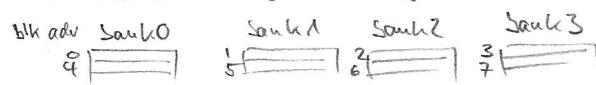
### Memory Optimizations

#### 6 basic cache optimizations

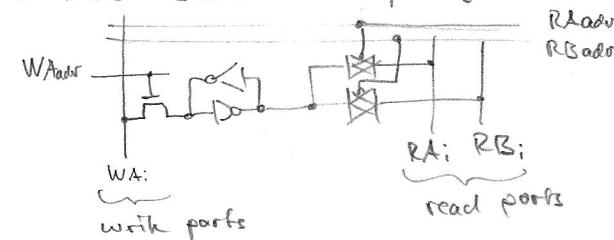
- larger block size  $\rightarrow$  larger total \$ capacity  
 $\hookrightarrow$  reduce compulsory miss  $\hookrightarrow$  reduce miss rate
- Higher associativity  $\hookrightarrow$  reduce conflicts
- Prioritize reads over writes  $\hookrightarrow$  reduce miss penalty
- Cache levels  $\hookrightarrow$  reduce overall access time
- Avoid addr translat. in \$ index  $\hookrightarrow$  reduce hot fine

#### Advanced Optimizations

- small & simple 1st lvl \$  
  - critical timing path
  - lower associativity reduces power
- way prediction
- pipelining cache  
  - to improve bandwidth
  - makes it easier for associativity increase
- non-blocking cache
- multibanked cache  
  - interleave banks according to block address



#### True SRAM multiplexing



- incr. area - slower access (with delay)

#### Compiler optimization

- loop interchange
- blocking

#### HW prefetch

- fetch two blocks on miss (incl. next 1th)

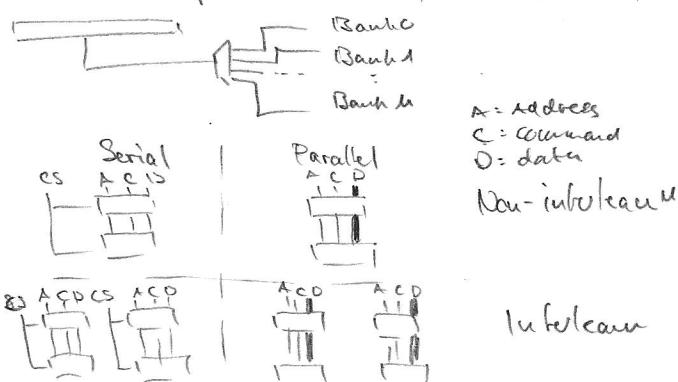
### Dependability

Introduce error correcting codes to increase reliability.

Example: Hamming SEC code. Parity bits

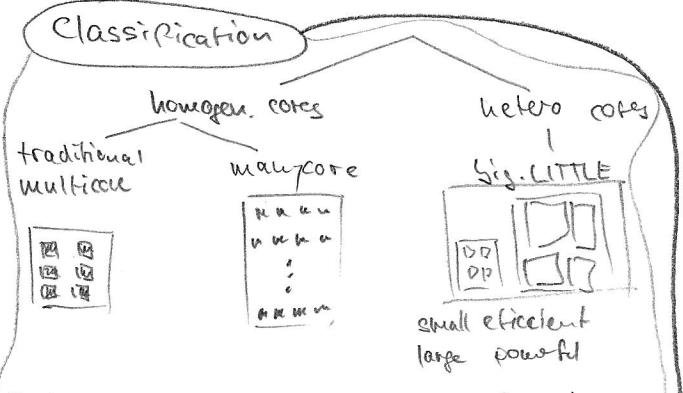
### Main memory

Interleaved: - Break info M Banks: Word A is in bank =  $A \bmod M$  at offset =  $A \div M$   
- Banks can operate concurrently & independently



Parallel: increase adr width

Goal: Maximize requests to an open row



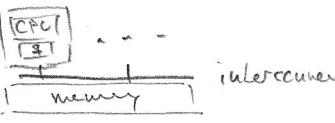
Challenge: With more cores, memory & core-to-core communication gets difficult

## parallel architecture

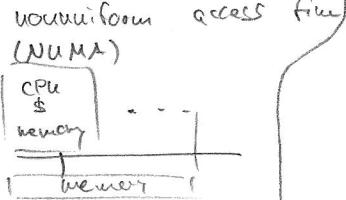
Two models: Symmetric multi processor interconnect and distributed shared memory

### SMP

centralized view w/ uniform access time (UMA)



### OSM



Non-shared memory: message passing

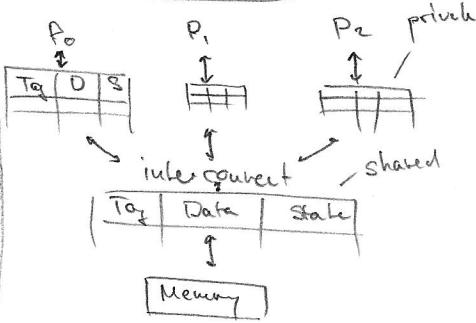
- simpler HW
- explicit comm.
- implicit sync

Shared memory: Open MP

- sharing large data
- supports HW caching

## Shared memory issues

### Cache Coherence



### Multicore

Need to regulate access to shared data. → use of "locks"

lock: low-level semaphore: high

region between acquire(lock) and release(lock) is a critical section

An interlocking acquire() will lock

A barrier is where all threads wait for all others to complete

spin lock acq(x): while(x!=0) { x=x-1; }  
↳ doesn't work ↳ acq is multi-op.

Solution: ISA provides atomic compare-and-swap CAS

cas r3 ← r1, r2, 0 (lock)  
↳ atomically executes  
ld r3 ← 0 (lock)  
if r3 == r2:  
  st r1 → 0 (lock)

Now acquire: A0: cas r3, r1, 0, 0 (lock)  
A1: bnez r3, A0

test-and-test-and-set esp repetitively  
↳ heavy interconnect use. Better load to local cache and st if lock is free  
coarse lock: easy but slow  
fine lock: difficult but faster

Always acquire multiple locks in same order

### omp others

reduction, master, static, ...

### Memory consistency

All processors "see" all I/O in the same order

???

### Open MP

Directives & libraries → generates calls to thread library sa. pthread

#### omp parallel

code executed in parallel  
id = omp\_get\_threadnum();

shared(a) declare a as shared var

private(id) id is priv. to each thread

firstprivate copy in, private storage

lastprivate copy out,

#### omp for

splits/unrolls a loop

schedule(static) for simple regular loops with equal duration

schedule(dynamic) iteration space to work queue

↳ fine/coarse, parallelization overhead ↳

#### omp critical

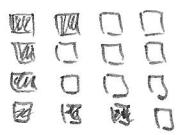
section executed by only 1 thread at a time

# LS: SMT, Vector, SIMD

## Categories

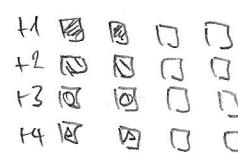
1+1 2+2 3+3 4+4 ...

### Superscalar



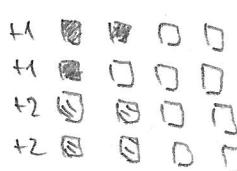
- single thread
- multi ops per cycle possible
- OOO

### Fine-Grained



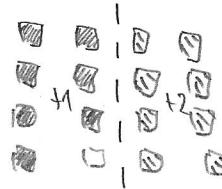
- fast context switches between threads

### Coarse-Grained



- slow switches

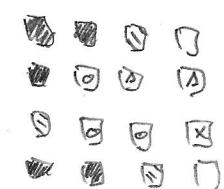
### Multiprocessing



- two superscalars
- each working on a thread

### SMT

Simultaneous multithreading



- multi threads at single cycle

### Simultaneous Multithreading:

- performing multiple threads of execution in parallel
- replicate PC, RF but keep EX the same

Fine/Course grained: - can hide large/short stalls but not both

Example Adding SMT in Power5: - locr. U1 associativity, - Add per thread id/st queue  
- add separate ins prefetch & buffering per thread - virtual reg 152 to 240 → An increase of 24%.

## Instruction and Data Streaming

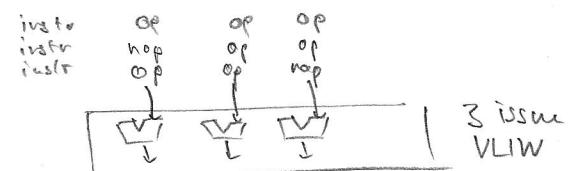
		data stream	
		single	multi
func Stream	single	SISD	SIMD
multi	MISD	MIMD	

### Multimedia Extensions

- very short vectors added to existing ISAs for memory
- limited issue set: no vector len control, no stride, ...

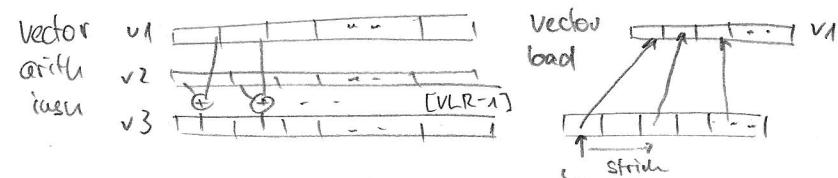
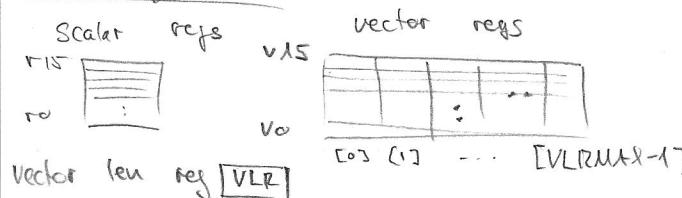
Vector - highly pipelined functional unit - simplify data-parallel program  
- good with regular access patterns  
- more general than ad-hoc media extensions  
SIMD - operate elementwise on vectors of data  
- simplifies synchronization

VLIW ISA  
vector HW without vector



## Vector Processor

### Programming model



### Advantages

- 1 issue encodes N operations
- expressive:
  - o index o save func unit o ...
- scalable

### Mem layout

- perfect for strides 1, 2

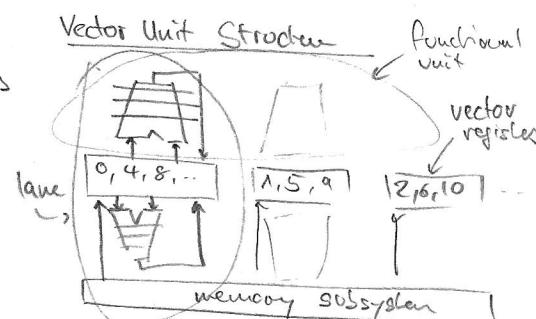
- shitty for strides 4, 8, ...

- mod 4 = 0 04
- mod 4 = 1 15
- mod 4 = 2 26
- mod 4 = 3 37

### Conditional exec

for (i=0; i<10; i++)  
if (a[i] == 1)  
then ...

- maskable
- vector instructions



Ariane is a vector processor for Ariane

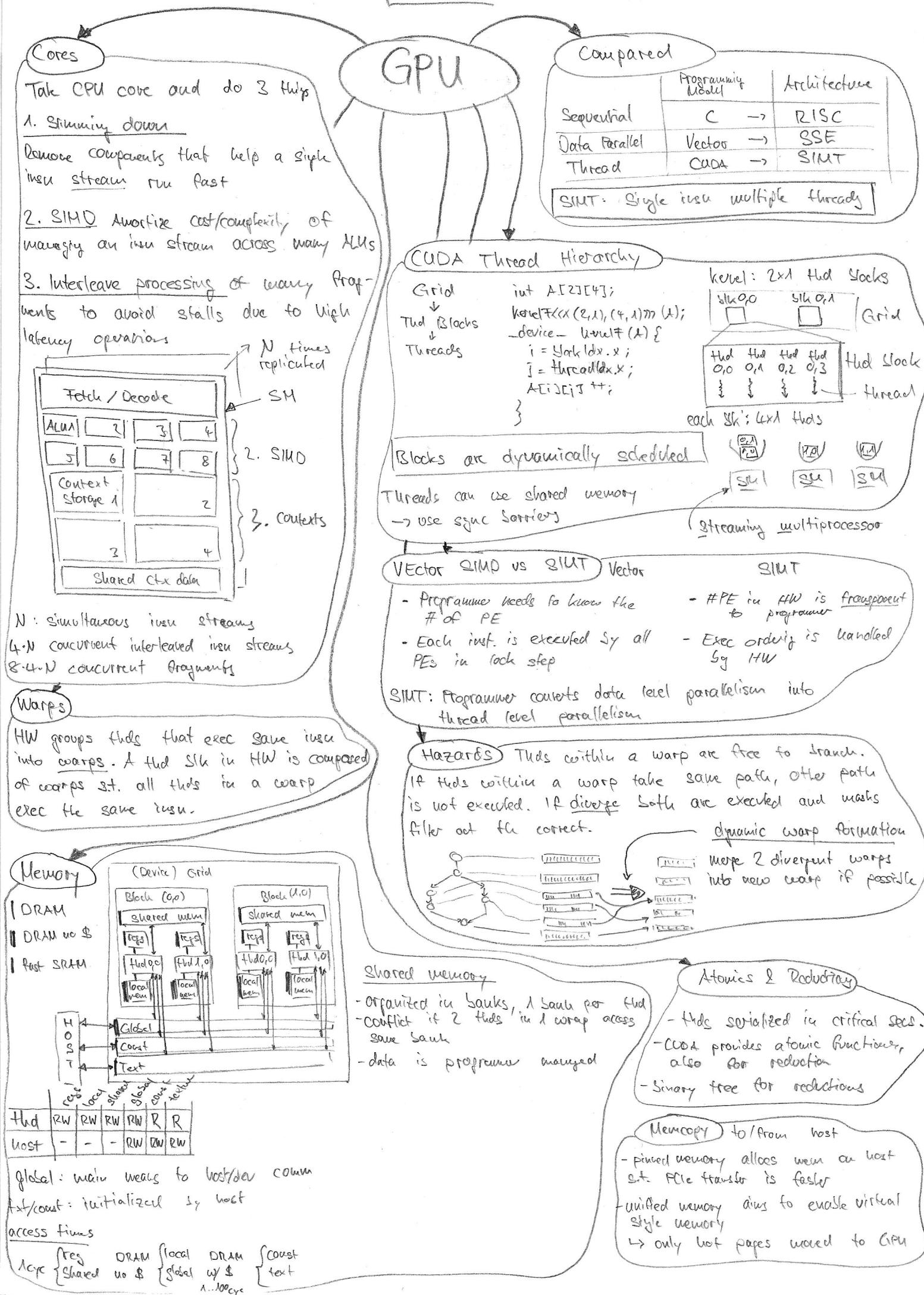
reductions are slow. Use trick to do vector MAC instead

### Strip mining

Vector regs have finite len.  
→ Break loops into fitting pieces and a remainder

### Mem Op

- unit stride: cont block in mem
- stided: prime number of data banks best supports different stride
- indexed: scatter-gather sparse



# L8: Heterogeneous Architectures

## Andahl's Law (again)

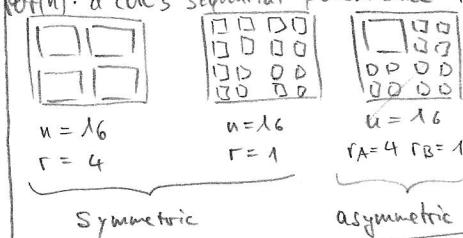
$$\text{Speedup}(f_{\text{unif}}) = \frac{1}{1-f} + \frac{f}{\text{perf}(r)} \cdot \frac{n}{r}$$

$f$ : parallelizable fraction

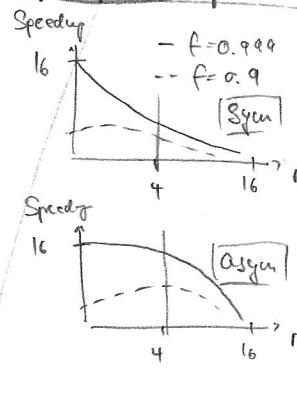
$n$ : total resources (e.g. transistor count)

$r$ : resources dedicated to each core

$\text{perf}(r)$ : a core's sequential performance  $\approx \sqrt{r}$



## Heterogeneous Architectures



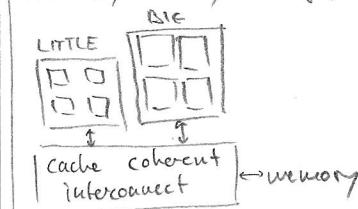
Observation: Most real-world loads are complex

- parallelizable & non-parallelizable loads
- amenable to SIMD & divergent
- predictable data access & random

Idea: the most efficient processor is a heterogeneous mixture of resources.

## big.LITTLE

Aims at power reduction. Smartphones have dynamically changing usage patterns.

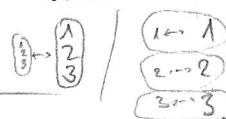


## Exclusive Cluster Switching

OS runs workload on cluster that is powerful enough but least power hungry.

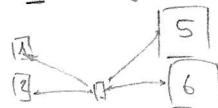
If demand grows, cluster is switched. Can be done on cluster or core base.

Unused forced off.



## Inclusive core migration

OS scheduler has all cores at its disposal. Can also activate all at same time. Heterogeneous Multiprocessing (HMP).



## Accelerators

On- and off-chip dedicated accelerators. CAPI: coherent Accelerator Processor Interface. Examples: Mobile phone incl. Nucleus, GPU, DSP, Neural accelerator

## CPU + FPGA

Zynq Dual ARM + FPGA.

UltraScale+ 4x A53, 2x RS, Mali GPU, FPGA

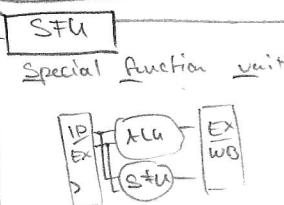
Important PC-to-PL interconnects

## Motivation

GP processor too slow for some applications. ASIP designed for specific application. Compromises performance & flexibility adds custom logic to CPU commercial tools for fast devl.

## ASIP

Application Specific Processor



## can be:

- combinational 1h ALU
- with internal state
- with special register file (SPR)

## implementation

- possible in embedded FPGA

E1: Single-Core RV MCU

RISCV: RV32IMCXCpulp

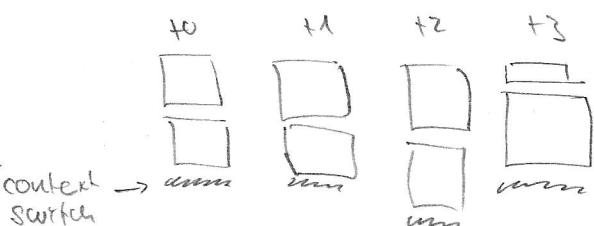
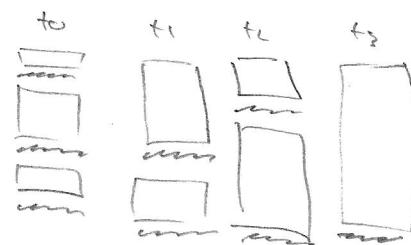
I: Integer M: Mult/Div F: float C: compressed Xpulp:

- Fast inst 1d/ef
- HW loop
- Packet SIMD
- Fixed point
- Bit Manipulation

Loop unrolling exploits inner-level parallelism

E2: OpenMP Programming

OMP parallel for num-threads (N) schedule (dynamic, C)

c=1  
=>

Fine granularity can increase total performance but overhead grows as well.

SPMD: single prog, multiple data → OMP parallel for

Shared memory programming

NPMID: multi prog, multiple data → OMP sections / tasks

execute diff streams of insns over possibly different data

omp nowait: removes the implicit sync barrier after a parallel/for/section

E4: Cuda1. Data

Need to malloc on device and copy data to it.

2. Funcglobal \_\_<sup>pid</sup> add (int\* a, ..) - parameters must be pointers

call with: add&lt;&lt; N, M &gt;&gt;(..)

N: number of threads to create, can be 2D  
M: number of threads per thread block, up to 3Duse `blockIdx.x` to access block index, `blockDim.x` for block size (#thds in block)use `threadIdx.x` to access thread index

threads can efficiently communicate &amp; synchronize.

## E5: Heterogeneous SoC

Zynq - Contains PS & PL and interconnects

SDSOC: Software Designed System on Chip: From C to:

on all-in-one tool.

- PL IP
- PL Datapath
- PS Driver
- PS Application

- Do all in SW, if reqs not met, evaluate HW solution
- Trade-off between data movement cost & Acceleration benefits

Zynq Ports: ACP (Accelerator Coherency Port)

- Has access to CPU cache hierarchy
- DDR/CCM access through ARM L1 L2 cache
- Useful for data sizes that fit in cache L1: 32K L2: 256K

HLS : - Main tasks: schedule & binding

- Schedule: what op to do in which clock cycle
- Binding: Map operators to HW blocks, share resources or duplicate

## E6: HERO Heterogeneous SoC

HERO: Heterogeneous Research Platform

- single source, single binary cross compile toolchain
- OpenMP 4.5
- shared virtual memory for host & pmlc
- HERO includes the first non-commercial heterogeneous cross-compilation toolchain

