

154B Discussion 9

March 11th, 2022

Goals

- Going over practice final.

- Q1 - Large pages (2MiB) compared to 4KiB
 Huge pages (1GiB)

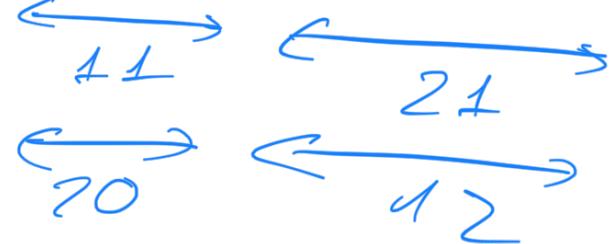
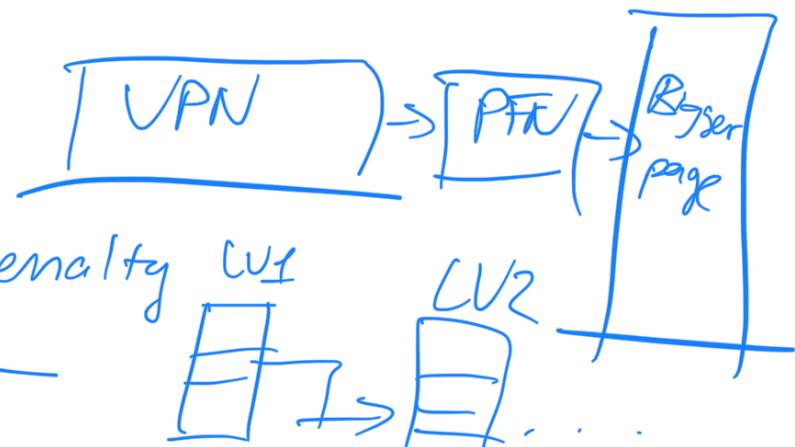
a. spatial / temporal

b. No difference in TLB miss penalty if the # of level doesn't change

c. VAddr: 32 bits (suppose this)

~~2²¹~~ bits for page offset \leftarrow 2MiB page \leftarrow 2^{21} bytes

2^{20} VPNs \leftarrow 12 bits for page offset \leftarrow 4KiB page



Assume that we have a fixed # of slots in TLR
(e.g. 96)

2^{11} VPNs \rightarrow 2^{11} VPNs competing for 96 slots in TLR
(\geq MiB page)

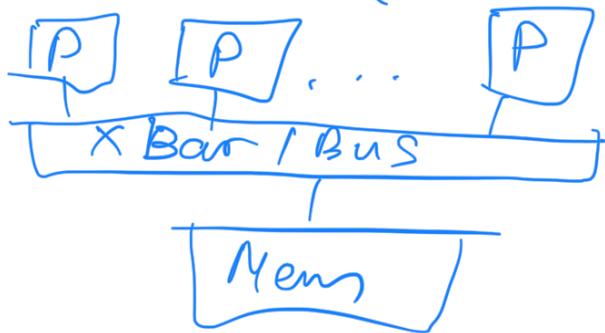
2^{20} VPNs \rightarrow 2^{20} VPNs competing for 96 slots in TLR
(LkiB pages)

2 MiB pages have lower TLR miss ratio

Q2

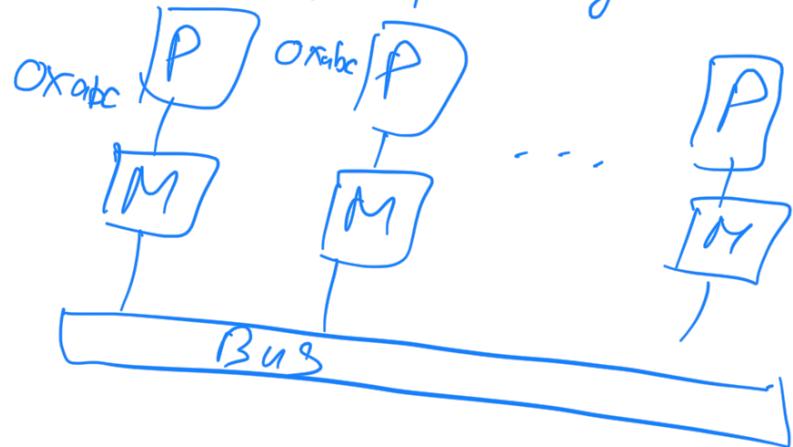
Shared memory vs. message passing

a. Shared mem



b.

message passing



c.

d. Shared memory: prog model: coherency is taken care of by HW

Message passing:

SW

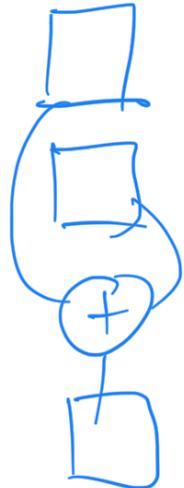
SW

→ shared memory → programming is easier.

Q3. Vector architecture

Scalar register (DINOCPU)

rs1



rs2



op

rd

~~at~~

a: int[] b: int[] c: int[]

loop: a + b = c

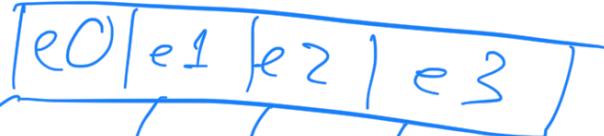
load x1, a[i]
load x2, b[i]
add x3, x1+x2
store x3

Scalar

) K times
→ K × 4 dynamic insts

Vector register

vs1



vs2



op

vd



K × 4 dynamic insts

Vector

loop:
load vs1, [a[i], a[i+1], a[i+2], a[i+3]]
load vs2, [b[i], b[i+1], b[i+2], b[i+3]]
vadd vs3, vs1, vs2
store vs3

Q4: ~~SISD | MISD~~
~~MIMD | MI MD~~

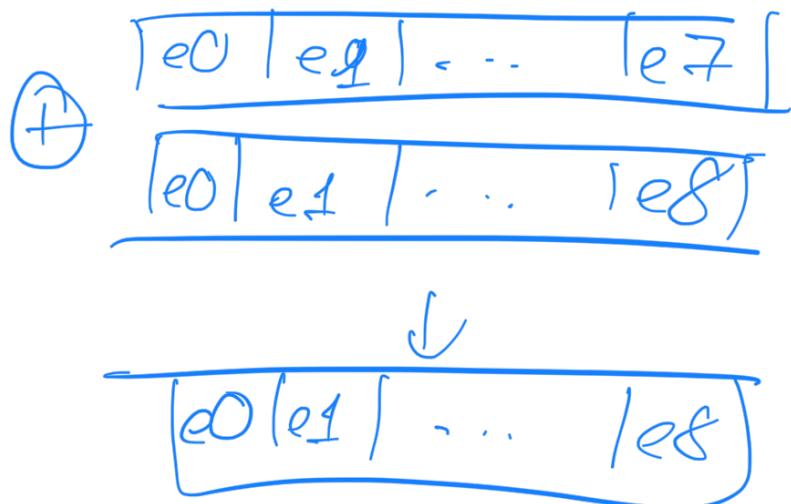
- a. SISD arch
- b. SIMD arch
- c. SISD / SIMD
Data-level Parallel (DLP)

\rightarrow ^{easier} SISD : programmers (users) don't have to find data parallelism

SIMP: /, have to do extra work
 to find and exploit DLP

Q5 . 75% of an application has data parallelism
→ 25% „ „
has to be done serially

a. 8 lanes → a vector



Amdahl's law: speed up of parallel machines over a serial one

$$\text{Time}_{\text{serial}} = 1 = (0.75 + 0.25)$$

$$\text{Time}_{\text{parallel}} = \left(\frac{0.75}{\# \text{lanes} \text{ or } \# \text{cores}} + 0.25 \right)$$

$$\text{speed up of A compared to B} = \frac{t_B}{t_A}$$

$$\frac{t_{\text{serial}}}{t_{\text{parallel}}} = \frac{4}{\frac{0.75}{8} + 0.25}$$

b. 16 lanes

c. Speed up application on 16 lanes compared to 4 lanes .

$$\frac{t_{4 \text{ lanes}}}{t_{16 \text{ lanes}}} = \frac{\frac{0.75}{4} + 0.25}{\frac{0.75}{16} + 0.25}$$

d. Maximum speed up.

$$\frac{1}{0.25 + \frac{0.75}{N}} = \frac{1}{0.25} = 4x$$

N: very big

$$3.9x \Rightarrow 117 \text{ lanes}$$

$$3.99x \Rightarrow 1197 \text{ lanes}$$

$$3.999x \Rightarrow 11997 \text{ lanes}$$

We cannot get $4x$ b/c cost / or no N could get $\frac{0.75}{N} = 0$

Q6. Efficient

→ Performance : SIMD is faster when there's a lot of DLA

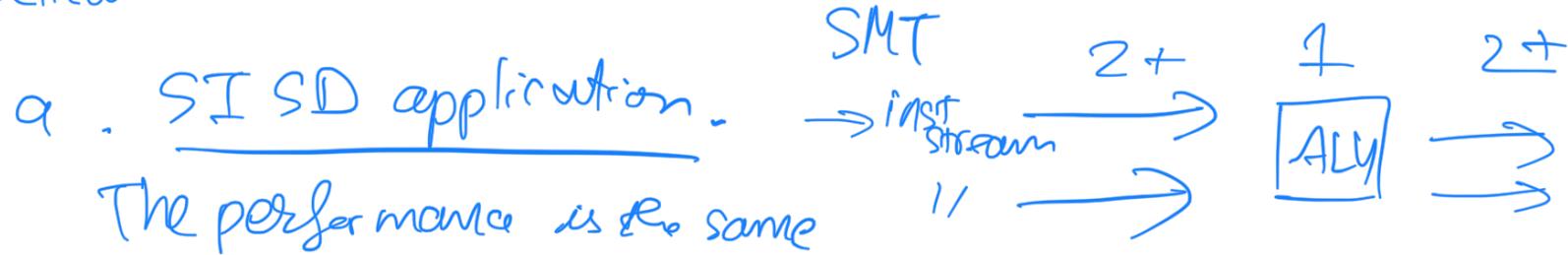
→ energy : less decoding / fetching

Q7. True .

Q8 . scalar vs. multi-threaded

HW thread \neq SW thread

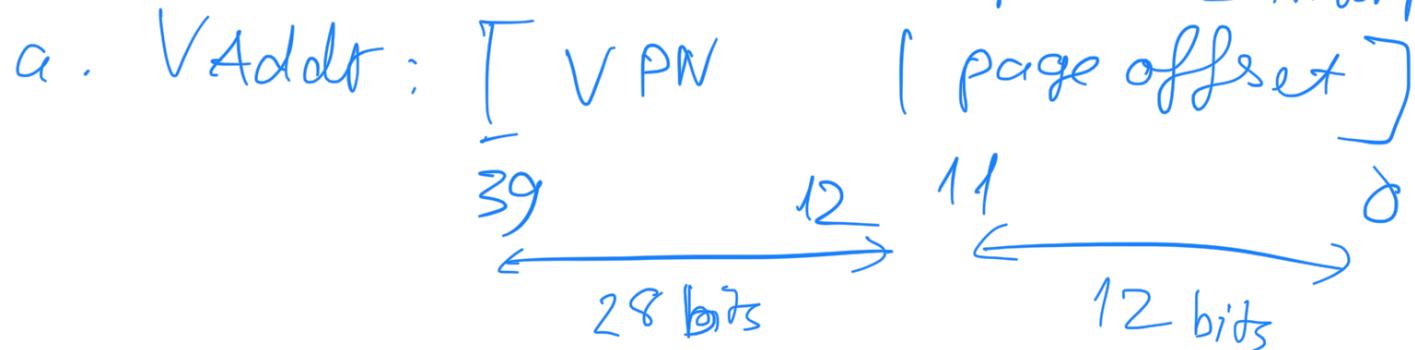
\downarrow different execution context
 \downarrow can share mem space



b. MIMD. Multi-threaded is better & Multi-threaded

c. Many different applications \rightarrow Throughput produced better throughput overall
(individual app can have higher latency)

Q9. VAddr 40 bits 1KiB page size TLR: VVA \rightarrow PFC
 $\hookrightarrow 2^{12}$ bytes \rightarrow 12 bit for page offset



b = Fully associative to TLB(cache)

UPN0	PFN0
UPN1	PFN1
UPN2	PFN2
⋮	⋮
⋮	⋮

UPN

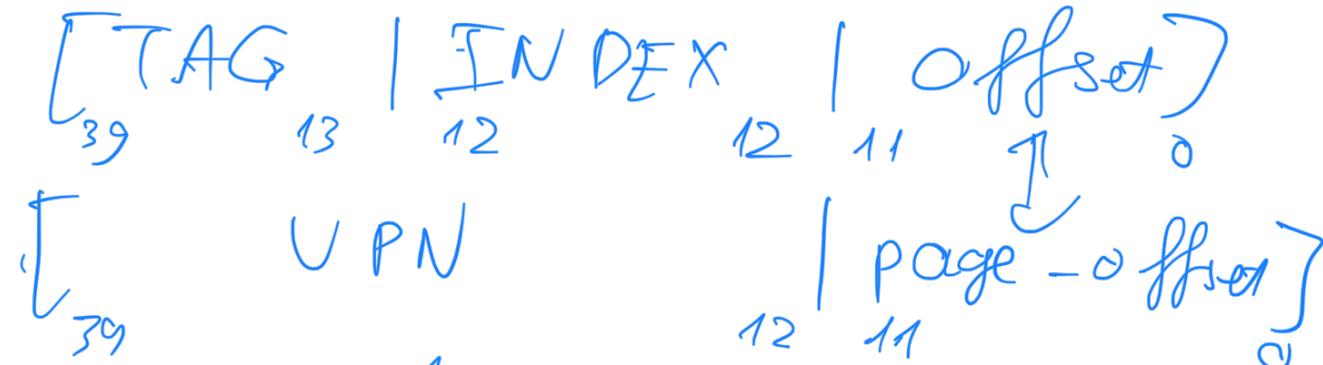
Tag: bits (39, 12)

d. 16 entries, 8-way set associative



16 entries

8 way



Index → which row?

Index bits are (12, 11)

2^4 rows → 4 bits for index