

Vad handlar kursen om?

- * Pipelinade mikroprocessorer
- * Minnes hierarki: Cache och virtuellt minne
- * Multiple-Issues processorer (Superskalär & VLIW)
- * I/O - (Skivminne, bussar, SSD)
- * Datearitmetik (float)
- * Flertrådade processorer
- * Flerkärniga processorer
- * Optimering av kod och minnesystem för en viss applikation (en Gauss-elimination)

Attande drivande ideer

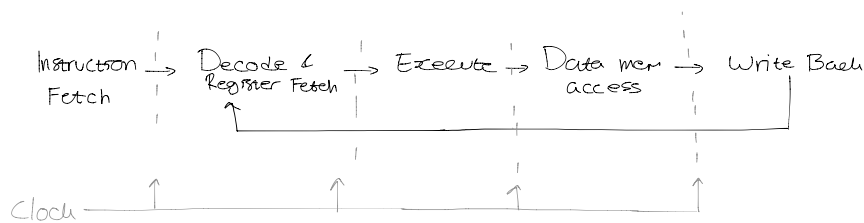
- * Konstruktionerna baserade på Moore's law
- * Abstraktion för att hantera komplexitet
- * Fokusera på common case
- * Hög prestanda via parallelism
- * Hög prestanda via pipelining
- * Hög prestanda via prediktion
- * Hierarki av minne
- * Tillförlitlighet via redundans

Minnesteknologier

- * Cacheminne görs av SRAM.
 - 0.5 ns - 2.5 ns
- * Dynamisk RAM
 - 50ns - 70ns
- * Magnetiska blivor
 - 5ms - 20ms

Idealt minne: Access-tid som SRAM
Kapacitet och kostnad/GB som hårddisk.

Prestanda via Pipelining



$$E_{\text{ffektiv}} = \text{Capacitiv last} \times \text{Matningsspänning}^2 \times \text{Frequens}$$

Multicores

- * Fler än en CPU per chip
- * Längsammare klocka → mindre värme
- * Problemet flyttas till mjukvaran

Datorprestanda

Beror av: Den underliggande algoritmen, bestämmer antalet operationer
Programmeringsförått
Processorn och minne
I/O-system (inkl. OS)

PC

Stort SW-utbud
Kostnad / prestanda

Server

Nätverksbaserade
Jobs/sekund

Superdator

Högprestanda
Top500.org
Green500.org

CPU Clocking

Klockperiod: Längden på en klockcykel
Klockfrekvens: Cykler per sekund

Bästa mått för prestanda är: Ekkveringstid av typiska program! (Kolla slides.)

$$CPU\text{tid} = \frac{\text{Instruktioner}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruktion}} \times \frac{\text{Sekunder}}{\text{Klockcykler}} = IC \times CPI \times T_c$$

Prestandan beror av: Algoritm (IC, CPI)
Prog språk (IC, CPI)
Kompilatorn
ISA

SPEC Benchmarks

Kostnad/Prestanda ökar
Pipelining & Cache = ökad prestanda
Exekveringstid.....

RISC

Reduced Instruction Set Computer

- × Enklare instruktioner och adresseringssätt
- × Gör processorerna mindre komplexa → snabbare
- × De flesta moderna processorer är av RISC
x86 är dock CISC externt

CISC

Complex Instruction Set Computer

- × Komplexa instruktioner och adressering

Mips

- × Typexempel på RISC.
- × Används i inbyggda system.
- × Aritmetiska operationer använder två källregister och ett destinationsregister. Detta ger regularitet.
- × Har 32*32 bit register file (RF)
 - 32 bit data kallas word

Massor av exempel. Fanns i kap 1 av kompendium.

Processor

CPU prestanda: $T_{tot} = IC \cdot CPI \cdot T_c$

IC avgörs av ISA, kompilatorn, algoritmen och programspråket.

CPI avgörs av HW.

Instruktionsexekvering

PC → Instruktionsminnet, hämta instruktion

Registernummer → Registerfil, läs reg

Beroende på type:

- * Använd ALU för att beräkna: Aritmetiska resultat
Minnesadresser för load/store instr
Hoppadresser och hoppvillkor
- * Läs/skriv dataminnet för load/store instr.
- * $PC = PC + 4$ eller hoppadress

Logikkonstruktion

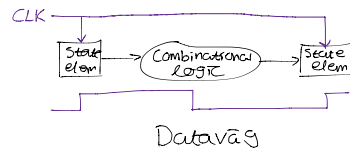
Info kodus binärt

Kombinatoriska element

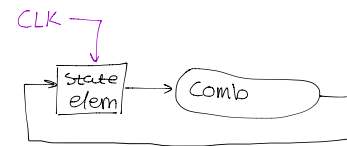
Gmindar, osv

Tillståndselement

Lagrar data



Dataväg



Tillståndsmaskin

R-Format Instruktioner

Läs två reg operand

Utför aritmetisk/logisk operation

Skriv resultatvärdet till registerfilen

I-format

Lw \$destreg, offs(\$basreg)

Sw \$källreg, offs(\$basreg)

Behöver teckenutvidga

Branch

bneq/bne \$r1, \$r2, label

Teckenutvidga

Sliffa vänster 2-bit positioner

Addera detta till PC

1cbe pipelined

En instruktion per cykel.

Slide Chapter 4 - The processor - 18

Styrenhet

R-type: O, rs, rt, rd, Shift, funct

L/S 35/43, rs, rt, offset

Jumps

Använder ord-adresser

Uppdatera PC med concat:

Högsta 4 bitar av PC

26 bitar från jump-instr

Prestanda

Den längsta signalfördröjningen avgör clockperioden

MIPS Pipeline

Fem pipesteg:

1: IF Instruction Fetch

2: ID Instruction Decode & Reg read

3: EX: Execute Instruction eller beräkna adress till DM

4: MEM: Access DM

5: WB: Write Back resultat till reg

Om alla pipesteg är balanserade: $\text{Pipeline Speedup} = \frac{T_{c,icke}}{T_{c,pipe}} \approx \text{Antal Pipesteg}$

Om pipen är obalanserad så är speedup mindre

Speedup pga genomströmning

Latenstiden för en enskild instruktion minskar dock inte

Konflikter (Hazards) i datavägen

Situationer som förhindrar start av instruktion i nästa klockcykel.

- Strukturell : Ett block är upptaget S 34
- Datakonflikt : Föregående instruktion ej klar S 35
- Styrkonflikt : Hoppvillkor och adress evalueras en bit in i pipen. S 39

Pipeline Summering

- ? Pipelining ökar prestanda via ökad instruktions- genomströmning
- ? Flera instruktioner exekverar parallellt (via överlappning av deloperationer)
- ? Pipelinekonflikter (eng. hazards) måste hanteras
- ? Strukturkonflikt, datakonflikt (RAW), styrkonflikt
- ? Tydligt beroende mellan ISA och komplexitet hos pipeline implementation

Summering

- ? Pipelining förbättrar instruktions- genomströmningen via parallellism (överlappning)
- ? Fler instruktioner slutförs per sekund
- ? Kortare signalvägar möjliggör en högre CPU klockhastighet
- ? Hazarder: strukturell, data, styr konflikt
- ? Nästa föreläsning: Detaljerad styrning av pipelinen, samt Multiple-Issue och Spekulative processorer

Styrning av Pipen

Hur löser vi RAW hazards? 4-8 §4.7

Vi kan använda forwarding. Finns några olika varianter i slides.

Double data hazard: Hazard på både ID/EX + EX/MEM och IE/EX + MEM/WB



add S1, S1, S2

add S1, S1, S3

add S1, S1, S4

Då krävs reviderad forwarding 4-14

Hur löser vi load-use? 4-17

Liknande data hazard med fm aritmetisk.

Om vi behöver ställa: sätts en NOP in.

: PC ändras inte.

Hazarder vid hopp? 4-24

Hoppstrategier - Fördröjda hopp

Kör alltid instruktionen efter hoppinstruktionen (delay slot)

Fyll den platsen med en nyttig instruktion, annars NOP

Mer avancerade metoder 4-33

Deepare pipes => prestanda förlust vid hopp större.

Använd dynamisk hopp-prediktion

Branch Target Buffer

Instruction Level Parallelism 4-39

Exekvera flera instruktioner parallellt via överlappning

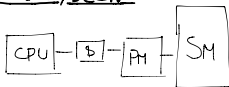
Öka ILP => Deepare pipelines => Mindre arbete p steg
Fler hazarder

Multiple Issue CPU

Mips med statisk Dual Issue 4-42

Loop unrolling 4-45

Minnesystem



- * Minnesväggen
- ✓ Minnes Hierarki
 - Cachelminne
 - Direktmappad alt associative
 - Flernivåscache
 - Hur ökar vi bandbredden mellan B och PM.
 - Sekundärminne
 - Skrivminne
 - Flash

Minnsteknologier

- SRAM
 - 0.5ns-2.5ns, \$2000-\$5000 per GiB
- DRAM
 - 50ns-70ns, \$20-\$75 per GiB
- Magnetiska Skivor
 - 5-20ms, \$0.2-\$2 per GiB

Idealt minne: Accessstid som SRAM
Kapacitet och kostnad som för hårddisk

Lokalitetsprincipen (slide 7)

Instruktioner och data accessas ganska sekvensiellt.

Spatial

Om en adress accessas är det sannolikt att en närliggande adress kommer att accessas snart.

Temporal

Om en adress accessas är det sannolikt att samma adress kommer accessas snart igen.

Cachelminne

Ett litet, snabbt, minne nära processorn och större långsammare minne längre ifrån.

h_x = Hit rate

m_x = Miss rate $AMAT = h_x T_1 + m_x T_2 = h_x T_1 + (1-h_x) T_2$

T_x = Hit time

$(T_y - T_x)$ = Miss penalty

- Fyra viktiga:
1. Var skall blocket placeras?
 2. Hur hittar vi?
 - 3.
 4. Slide 28

Cache adressering

Tag Index Offset

~slide 30

Massor av exempel, se slides.

forts.

Hur ökar vi bandbredden till PM? slide 60-62

Sekundärminne

Flash

100-1000 ggr snabbare än skivminne
Mindre och mer robust
Kostar mellan DRAM och skivminne

Skivminne

slide 67...

NOR-flash

Random read/write access
Används som instruktionsminne i inbyggda system.

NAND-flash

Tätare bits/area men access sker blockvis.
Billigare per GB.

Flash töts ut med tiden!

Inlämningsuppgift

Implementera gauss-elimination i MARS-assembly. (Flyttal)

Optimera

Välj cache parametrar

OSV OSV

for(i=0; i<N; i++) Två block med samma index
 A[i]=A[i+1] Ger miss vid varje access

Prestanberäkningar - Tema 2010-05-25

a) 1000 mätningar $\cdot 10 \text{ ms} = 10 \text{ s}$

$$T = IC \cdot CPI \cdot T_c$$

$$IC = 2 \cdot 10^9, T_c = \frac{1}{10^8} = 10^{-8}, CPI = ?$$

$$CPI = 1.3 + CPI_{\Delta D\$} = 1.3 + 0.10 \cdot 10 \cdot 0.2 = 1.5$$

Andel data ref
Missrate miss penalty

$$T = 2 \cdot 10^9 \cdot 1.5 \cdot 10^{-8} = 30 \text{ s}$$

$$T_{\text{tot}} = 30 + 10 = 40 \text{ ms}$$

c) Halvera waiting time på sensorn eller halvera miss penalty?

Halvera väntetiden $\Rightarrow 5 \text{ s}$ förbättring

Halvera inverkan av D\$ missar $\Rightarrow \frac{0.2}{2} = 0.1$

Hela beräkningstiden minskar med $\frac{0.1}{1.5} \cdot 30 = 2 \text{ s}$

Det är alltså bättre att halvera väntetiden på sensorn.

b) I\$: MP=10 cykler, MR=0.01

$$CPI_{\Delta I\$} = 10 \cdot 0.01 = 0.1$$

Beräkningstiden minskar med faktorn $\frac{1.5-0.1}{1.5} = 0.933$

$$\frac{T_{\text{tot}}}{T_{\text{tot}}} = \frac{10 + 30 \cdot 0.933}{10 + 30} = 0.95$$

Virtuellt Minne

Man har virtuella och fysiska adresser. Den virtuella är oftast större vilket ger illusionen av ett större PM.

VM använder PM som cache för SM.

Varje program har ett eget virtuellt adressrum.

Physical address: Address in PM

PAGE FAULT

Sidan måste hämtas från disk.

Tar miljoner clock cykler

Hanteras av OS kod.

Det finns sidtabeller i PM som innehåller placeringsinformation.

5.1)

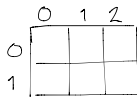
```
for (i=0; i<8; i++)
  for (j=0; j<8000; j++)
    A[i][j] = B[i][0] + A[i][j];
```

5.1.1 How many 32-bit ints in 16 byte cache block?
4

5.1.2 Which variables exhibit temporal locality?
j, i

5.1.3 Spatial locality

C codes gets stored rows first

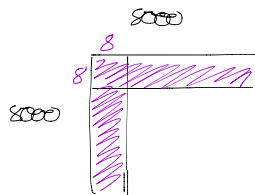


3x2 so we have spatial locality when accessing several items on the same row.

```
for i=1:8
  for j=1:8000
    A(i,j) = B(i,0) + A(j,i);
  end
end
```

5.1.4 How many 16-byte cache blocks are needed to store all 32-bit elems referenced?
 $(8 \times 8000 + 8000 + 8000 \times 8 - 8 \times 8) \cdot \frac{1}{4} = 33984$

5.1.5 Temporal
i, j



5.1.6 Spatial
A(j,i), B(j,0)

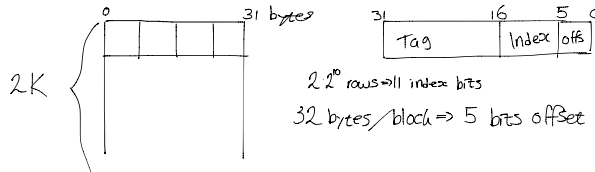
5.5,

512 KiB working set with address stream: 0 2 4 6 8 10 12 ...

5.5.1 Assume 64 KiB direct mapped cache with 32-byte block. What is the MR?
How is r sensitive to the set size or cache size?

$\frac{64}{32} = 2K$ blocks

1 bytes per elem
32 bytes per block



2^{11} rows \Rightarrow 11 index bits
32 bytes/block \Rightarrow 5 bits offset

First access \Rightarrow Miss, read 16 elems to cache (1, 3, 5, ... garbage)
17th access \Rightarrow — 11 —

Miss every 16th access. Only cold misses due to streams not having locality.

5.52 MR for 16-byte, 64 and 128-byte block size

16 byte: $\frac{1}{8}$
 64 byte: $\frac{1}{32}$
 128 byte: $\frac{1}{64}$

What kind of locality is exploited?
 Spatial. We read sequentially.

5.53 Assume a two-entry stream buffer.

First read would give a miss, then no misses.

Cache block size (B) can affect both miss rate and miss latency. Assuming CPI=1 and an avg of 1.35 references per instruction find the optimal block size given the following MR for block sizes: B:

8	16	32
4%	3%	2%

5.54 What is the optimal block size for miss latency of $20 \times B$ cycles?

B=8: $0.04(8 \cdot 20)$
 B=16: $0.03(16 \cdot 20)$ 8 wins!
 ⋮

5.55 $24+B$ cycles. 24 units of time to fetch the first elem and the consecutive elems takes 1 unit.

B=8: $0.04(24+8)$
 B=16: $0.03(24+16)$
 ⋮

5.56 For const miss latency?
 Only the missrate matters!

5.11

5.11 As described in Section 5.7, virtual memory uses a page table to track the mapping of virtual addresses to physical addresses. This exercise shows how this table must be updated as addresses are accessed. The following data constitutes a stream of virtual addresses as seen on a system. Assume 4 KiB pages, a 4-entry fully associative TLB, and true LRU replacement. If pages must be brought in from disk, increment the next largest page number.

4669, 2227, 13916, 34587, 48870, 12608, 49225

TLB

Valid	Tag	Physical Page Number
1	11	12
1	7	4
1	3	6
0	4	9

Page table

Valid	Physical Page or in Disk
1	5
0	Disk
0	Disk
1	6
1	9
1	11
0	Disk
1	4
0	Disk
0	Disk
1	3
1	12

5.11.1 Show the final stage of the system given the above info. Also show if it's a hit in TLB, Page table or page fault.

4 KiB pages $\Rightarrow 4 \cdot 2^0 \Rightarrow 12$ bits to encode
 4669:

0001	0010	0011	1101
------	------	------	------

 Tag Byte offset in page Tag 1 is not in TLB \Rightarrow TLB miss
 (Tag = Virtual Page)

5.11.2 Use 16 KiB instead of pages. Advantages/Disadvantages?

- + Fewer misses since you load larger chunks.
- Larger misspenalty!

Example: TLB och Cacheinteraktion

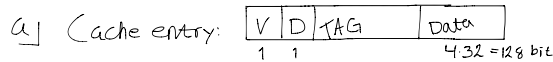
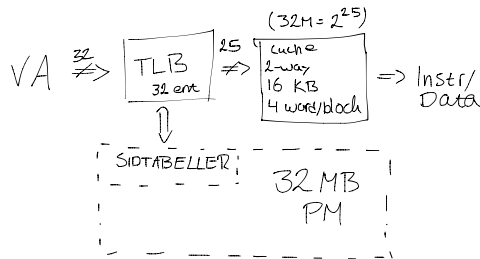
Virtuellt minne exempel (tentamensfråga 4)
2010-08-16

- Assume a computer systems with the following characteristics: The CPU uses virtual memory (unified data- and instructions) using 32-bit virtual addresses. There is a maximum of 32 MB of physical primary memory, and a two-way set associative 16 KB cache. The hit-time in this cache is 1.5 ns. A 32-entry fully associative TLB is used for virtual-to-physical address translations. The TLB hit-time is 1 ns. The page size is 16 KB, and the cache uses a block size of 4 words.
- The cache, the TLB, and the virtual memory uses a "write-back" ("copy-back") strategy. The replacement algorithm used by the cache is "Random", and for the virtual memory an approximative "LRU" using 2-bit time "time stamps" per page. Virtual memory pages that belongs to the OS are marked with a "protection bit" to protect it from user process accesses. In maximum, 16 processes are active at any given time. The Page Tables do not store secondary storage addresses when the pages are not in primary memory.

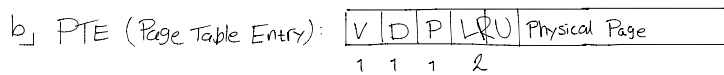
2010-08-16, uppg.4, forts

- How many bits must the cache memory be able to store? (4 p)
- How many bytes is needed to store the Page Tables? (4 p)
- How many bits must the TLB be able to store? (4 p)
- What is the maximum clock frequency the processor can have (assuming the critical timing path is constrained by TLB and cache) (6 p)?

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 27



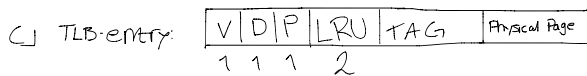
offset: 2⁴ \Rightarrow 16 byte
 Index: $\frac{16 \text{ KB}}{16 \text{ B}} = 1\text{K blocks}$
 2-way $\Rightarrow \frac{1\text{K}}{2} = 512 = 2^9 \Rightarrow 9$ bitar
 Tag: 25 - 9 - 4 = 12 bitar
 1 + 1 + 12 + 128 = 142 bitar/block $\Rightarrow 142 \cdot 1024 = 145408$ bitar



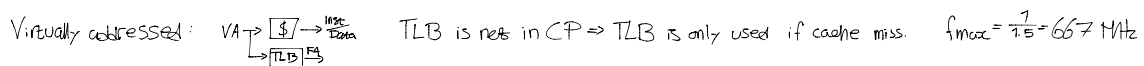
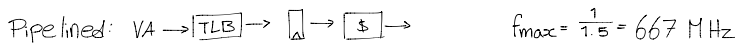
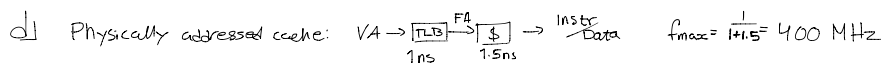
Page offset: 16 KB \Rightarrow 2¹⁴ bits
 Physical Page: 25 - 14 = 11
 PTE = 1 + 1 + 1 + 2 + 11 = 16 bitar

There's one PTE for each virtual page.
 How many pages? 2³²⁻¹⁴, 32: VA, 14: Page offset = 2¹⁸ = 256 K

256K · 2B · 16 = 8MB
 #processes



Tag in TLB = Virtuellt page number = 18
 Physical Page = 11
 Total: 1 + 1 + 1 + 2 + 18 + 11 = 34 \Rightarrow 32 entries in TLB \Rightarrow 1088 bitar i TLB



Mipsaritmetik

Multiplikation: Två 32-bitars register för produkten
HI (32 MSB)
LO (32 LSB)

Division: Använder HI/LO registren för resultat
HI: Rest
LO: Kvot

Flyttal

Representation av icke-heltal.

Binärt: $\pm 1.x_1x_2\dots x_n \cdot 2^m$ ← Normaliserat

Det gick
värdet föris.....

Standardiserad enl IEEE



$$x = (-1)^s \cdot (1 + \text{Fraction}) \cdot 2^{(\text{Exponent} - \text{Bias})}$$

Exponent: Excess representation: faktisk exponent + Bias

Sök på flyttal på
dchdwdkgo.com!

Exempel $F \rightarrow C$

C-kod:

```
float f2c(float f){
    return((5.0/9.0)*(f-32));
}
```

MIPS:

```
f2c: lwc1 $f16, const5($gp)
     lwc1 $f18, const9($gp)
     div.s $f16, $f16, $f18
     lwc1 $f18, const32($gp)
     sub.s $f18, $f12, $f18
     mul.s $f $f $f
     jr ra
```

Exempel från tenta

- Nyttjar lokaliteten bättre.
- Får bara plats med 4 block \Rightarrow Fler konflikter.
- Vad är den genomsnittliga misskostnaden?

$$4B: 0.38(4 + 4) = 1.9$$

$$16B: 0.22(4 + 1\frac{1}{4}) = 1.76$$

$$64B: 0.19(4 + \frac{1}{4}) = 3.8$$

$$256B: 0.27(4 + \frac{256}{4}) = 18.36$$

Exempel Hopp hazarder

ADDI	\$5, \$0, 2
L1: LW	\$4, 100(\$5)
ADDI	\$5, \$5, -1
ADD	\$3, \$3, \$4
BNE	\$5, \$0, L1
SW	\$3, 100, \$0

Always stall: 2 cykler extra pga hazard. Om hoppvillkor beräknas och används i ID: $10 + 2 + 4 = 16$ (4: Töm pipe) Om Ex: $10 + 2 \cdot 3 + 4 = 20$

Assume not taken: Gissar fel en gång. \Rightarrow 1 cykel extra om hoppvillkoret beräknas och används i ID-staget

$$10 + 1 + 4 = 15$$

Om vi beräknar i EX och använder i MEM \Rightarrow 3 cykler extra.

$$10 + 3 + 4 = 17$$

L1 körs 299.

Alla datahazarder kan

lösas med forwarding.

Hur många extra cykler pga hopp hazarder?

Räkna inte med försörjda hopp!

Multi-threading (MT)

Flera PC, fler registerfiler, snabb växling mellan trådar.

Fin resp grovkornig MT.

Parallelism exempel

Summera 10 tal sekventiellt och utför också en 10×10 matrisaddition. Vi kan använda 10 eller 100 CPUs

$$1 \text{ CPU: } T = (10+100) \cdot t_{\text{add}} = 110 \cdot t_{\text{add}}$$

$$100 \text{ CPU: } T = 10 \cdot t_{\text{add}} + \frac{100}{10} \dots$$

Tenta Q1 2010-08-16

RAW
 Stt (V0, a0, a1)
 beq (V0, zero, L1)
 add (V0, a0, zero)
 add (v1, a1, zero)
 beq (zero, zero, L2)
 L1:
 add (V0, a1, zero)
 add (v1, a0, zero)
 L2

- a) Sorts the values in the registers depending on the values.
 b) In data: a_x
 Out: V_x
 c) How many stalls if $a_0=5$ $a_1=1$.
 1. beq vill ha v0 i 10 men kan inte få det förrän stt skrivit tillbaka i WB. Det krävs 2 stall-cykler för detta.
 2. 3 stall-cykler
 Summa: 5
 d) För att göra till en subrutin: Lägg till Label och avsluta med jr \$ra
 anrop: addi a0, zero, 5
 addi a1, zero, 1
 jal Label

Q2 2009-05-19

50 computers @ 500 MHz
 Offer: 1 GHz samma ISA
 600 MHz bättre ISA

1 GHz = dubbelt så snabb som 500 MHz ty samma ISA.

$T_{exe} = IC \cdot CPI \cdot T_c$

$T_G = T_{mg} + T_{cg} = 0.4 \cdot IC_G \cdot 3.0 \cdot \frac{1}{500 \cdot 10^6} + 0.6 \cdot IC_G \cdot 1 \cdot \frac{1}{500 \cdot 10^6}$
 $T_N = T_{mn} + T_{cn} = 0.1 \cdot IC_N \cdot 0.5 \cdot \frac{1}{600 \cdot 10^6} + 0.9 \cdot IC_N \cdot 1.0 \cdot \frac{1}{600 \cdot 10^6}$
 $0.9 \cdot IC_N = 0.6 \cdot IC_G \Rightarrow IC_N = \frac{0.6}{0.9} IC_G$
 Uppsnabbning $\Rightarrow \frac{T_G}{T_N} = \frac{(0.4 \cdot 3.0 + 0.6) \cdot IC_G \cdot \frac{1}{500 \cdot 10^6}}{(0.1 \cdot 0.5 + 0.9) \cdot \frac{0.6}{0.9} IC_G \cdot \frac{1}{600 \cdot 10^6}} = 3.41$

1 GHz ger 2.99r 500
 600 MHz ger 3.41 89r 500

Q3 2011-01-12

a) Avg read/write time?
 Sector size: 512 byte
 7200 RPM
 Avg seek time: 8 ms
 Transfer rate: 20 MiB/s
 Controller off: 2 ms

Disk access = Seek + rotational latency + transfer + controller off =
 $8 + 0.5 \cdot \frac{60}{7200} \cdot 1000 + \frac{512}{20 \cdot 2^{20}} \cdot 1000 + 2 = 14.19 \text{ ms}$

b) 3 step process: Read 4KiB = $8 + 0.5 \cdot \frac{60}{7200} \cdot 1000 + \frac{4 \cdot 1024}{20 \cdot 2^{20}} \cdot 1000 + 2 = 14.37 \text{ ms}$
 Process = $\frac{20 \cdot 10^6}{400 \cdot 10^6} = 50 \text{ ms}$
 Write 4KiB = Read = 14.37 ms

Total time: $2 \cdot 14.37 + 50 = 78.74 \text{ ms} \Rightarrow \frac{1}{78.74} \cdot 1000 = 12.76 \text{ blocks/s}$

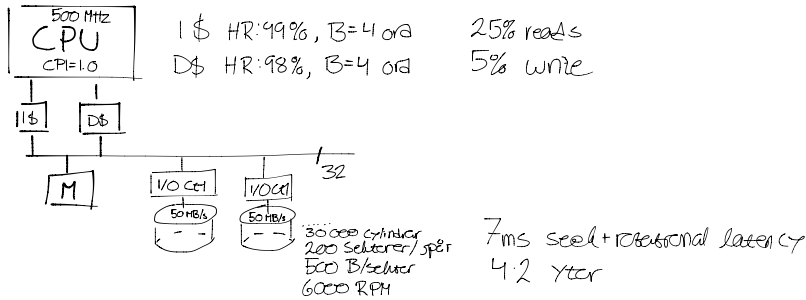
c) What is the bottle neck?

Time to process 64 Ki block $\frac{20 \cdot 10^6}{3 \cdot 10^9} = 0.67$ ms.

Transfer a 64 KiB block on bus: $\frac{64 \cdot 2^{10}}{640 \cdot 2^{20}} = 0.097$

Time for I/O $9 \cdot \frac{64 \cdot 2^{10}}{64 \cdot 2^{20}} \cdot 1000 + \frac{10^6}{3 \cdot 10^9} \cdot 1000 = 10.33$ ms

Q4 2009-05-26



Bussen används för överföring av 4 ord åt gången. (Address + Block = 1 + 4 = 5 cykler)

a) Busstid I\$Miss = $0.01 \cdot 5 \cdot 5 = 0.25$ Busstid D\$Miss, läs = $0.02 \cdot 5 \cdot 5 \cdot 0.25 = 0.125$

Missare → MP

Busstid D\$Miss, skriv = $0.05 \cdot 2 \cdot 5 = 0.5$ (Data-cyklar)

CPU 500MHz synbarare på bus #CPU cykler # Busscyklar

Systembussen är upptagen med trafik mellan cache och minne: $0.25 + 0.5 + 0.125 = 87.5\%$ av tiden.

b) Max bandbredd per disk R/W-huvud = $\frac{6000}{60} \cdot 500 \cdot 200 = 10$ MB/s
→ 8 yttre: $8 \cdot 10$ MB/s = 80 MB/s

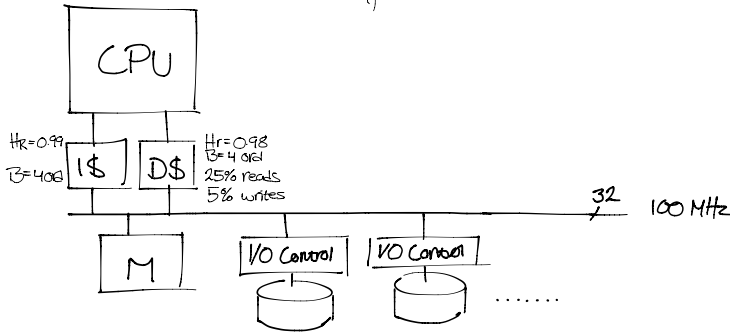
c) Max bandbredd på systembussen = $\frac{4 \text{ ord}}{5 \text{ cyklar}} = \frac{4 \cdot 4 \text{ B}}{5 \cdot 1000^6} = 320$ MB/s

d) Använd bandbredd per I/O-buss?
DMA (128 K block) = $2 \text{ ms} + 7 \text{ ms} + 10 \cdot 2 \cdot 1000$

DMA setup → SEEK + ROT

Tenta Q4 2009-05-26

1)



Misspendly: Address 4 ors
 30000 cylinderar
 200 sektorer
 500 bytes/sektor
 6000 RPM
 7ms seek och
 10ns latency
 Amount CPU cycles / bus cycles

1) $Busstid_{I\$Miss} = 0.01 \cdot 5 \cdot 5 = 0.25$

2) $Läsmiss$

$Busstid_{D\$RH} = 0.02 \cdot 5 \cdot 5 \cdot 0.25 = 0.125$
 25% reads

$Busstid_{D\$SM} = 0.05 \cdot 2 \cdot 5 = 0.5$
 5% writes
 skrimiss
 AddressData

$Busstid_{Miss} = 0.25 + 0.125 + 0.5 = 0.875 \Rightarrow$ Bussen är upptagen med missar 87.5% av tiden. 12.5% gör alltså att använda till I/O-operationer.

3) Max bandbredd per disk R/W-head: $\frac{6000}{60} \cdot 500 \cdot 200 = 10 \text{ MB/s}$. 8byter \Rightarrow 80 MB/s.

4) Bandbredd för bussen: $\frac{4}{5} \frac{\text{ord}}{\text{cykel}} = \frac{16 \text{ B}}{5 \cdot 100 \cdot 10^6} = 320 \text{ MB/s}$

5) Använd bandbredd per I/O-buss?

DMA (128 KiB) = $2 \text{ ms} + 7 \text{ ms} + \frac{128 \text{ KiB}}{10 \text{ MB/s}} \cdot 1000 = 22.1 \text{ ms}$
 DMA Setup seek/rotation transferetid

2 DMA överföringar kan ske samtidigt $\Rightarrow \frac{2 \cdot 128 \text{ KiB}}{22.1 \text{ ms}} = 11.86 \text{ MB/s}$ (effektiv databandbredd per I/O-buss)

Bandbredd för disk > Bandbredd för I/O-Buss, alltså är det I/O-bussen som flaskar.

6) $0.125 \cdot 320 = 40 \text{ MB/s}$ ledigt för I/O $\Rightarrow 3 \cdot 11.86 < 40 \Rightarrow$ Vi kan ha 3 I/O-kontroller.
 Tid ledig

Q2a. 2011-01-12

- add $R4 = R1 + R0$
- sub $R9 = R3 - R4$
- add $R4 = R5 + R6$
- lw $R2 = M[R3 + 100]$
- lw $R2 = M[R2 + 0]$
- sw $M[R4 + 0] = R9$
- and $R3 = R3 \& R1$
- beq $R9 == R1, Target$
- and $R9 = R9 \& R1$

Lila löses med ALU FWD
 R2d ger ingen hazard.
 Grön måste stallas en cykel och sedan använda MF \rightarrow Ex FWD.

Uppgiften är dock att göra ett diagram som beskriver detta.

Q4. 2012-01-11

Physical

Cache	TAG	Set Index	Byte Off
-------	-----	-----------	----------

Block size: 128 B \Rightarrow Byte offset = 7 = $\log_2(128) \Rightarrow C=7$

Cache size: 64 KiB $\Rightarrow \frac{64 \cdot 1024}{128} = 256 \Rightarrow$ Index = 8 $\Rightarrow B=8$

2-Way assoc

Physical address = 32 bits \Rightarrow Tag = 32 - 8 - 7 = 17 $\Rightarrow A=17$

128 B block $\Rightarrow 128 \cdot 8 = 1024$ bitar $\Rightarrow D=1024$

Tag/Status: 17 + (1+1) = 19 $\Rightarrow E=19$

Valid+Dirty

Virtual

TLB	Tag	Set Index	Page Offset
-----	-----	-----------	-------------

Virtual address: 64 bitar

8 KiB sidar: 13 bitar (2^{13}) sid offset $\Rightarrow H=13$

512 entries: $\frac{512}{8} = 256 \Rightarrow$ 8 bitar index $\Rightarrow G=8$

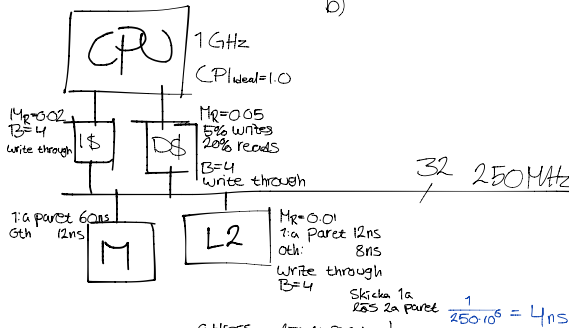
Tag: 64 - 8 - 13 = 43 bitar $\Rightarrow F=43$

PA = 32 bitar, Page = 8 KiB $\Rightarrow 32 - 13 = 19$ bitar physical page $\Rightarrow I=19$

Q4 2009-08-17

a) Beräkna CPI

b)



a) L2s misspenalty: $4 + 12 + 8 + 8 = 32 \text{ ns}$

PMs mp: $4 + 60 + 12 + 8 = 84 \text{ ns}$

Det tar 8ns att skicka men 12 att hämta.

ΔCPI

MRL₁

Inst hämt L2: $0.02 \cdot 32 = 0.64$

Inst hämt PM: $0.01 \cdot 0.02 \cdot 84 = 0.0168$

MRL₂ MRL₁ 10% reads

Data läs L2: $0.05 \cdot 0.20 \cdot 32 = 0.32$

Data läs PM: $0.05 \cdot 0.01 \cdot 0.2 \cdot 84 = 0.0084$

$\text{CPI} = 1.0 + 0.64 + 0.0168 + 0.32 + 0.0084 = 1.9852$