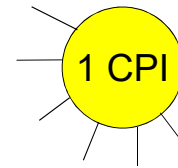
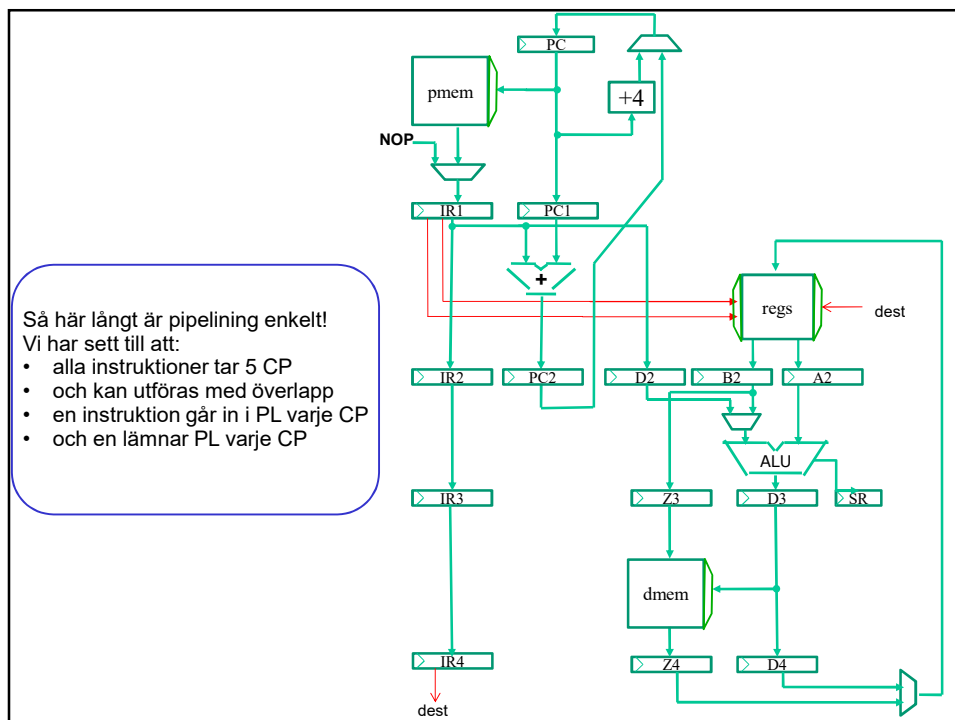


Föreläsning 5

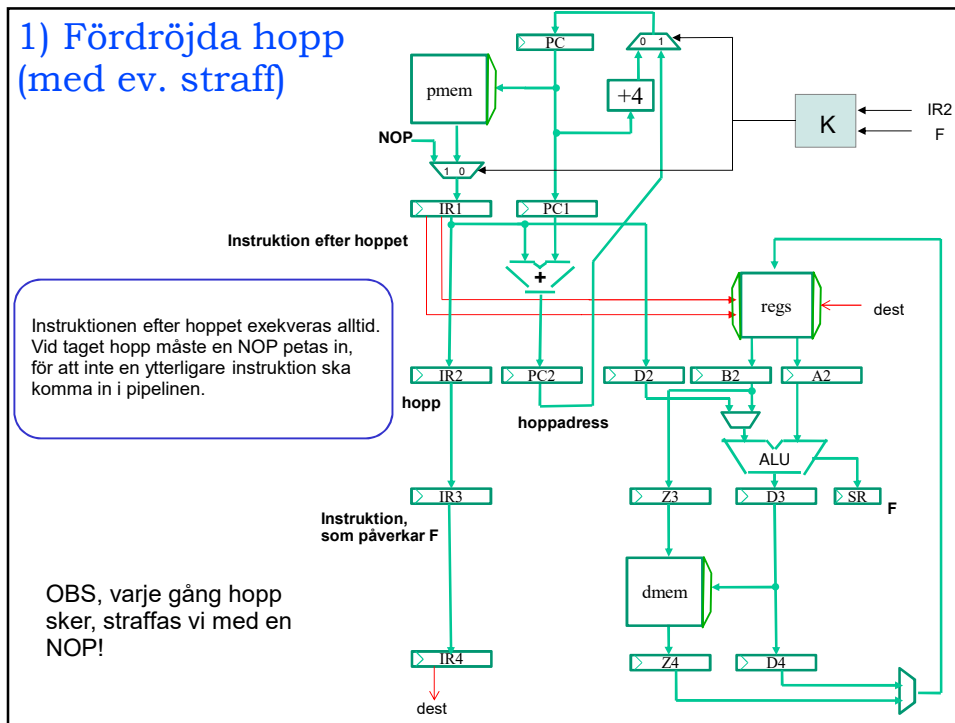


- Sammanfattning pipelining
- Cacheminnen
 - associativt minne som cache
 - associativt minne som BPT
 - direkt-mappad cache
 - flervägscache (2,4)
 - I/D-cache

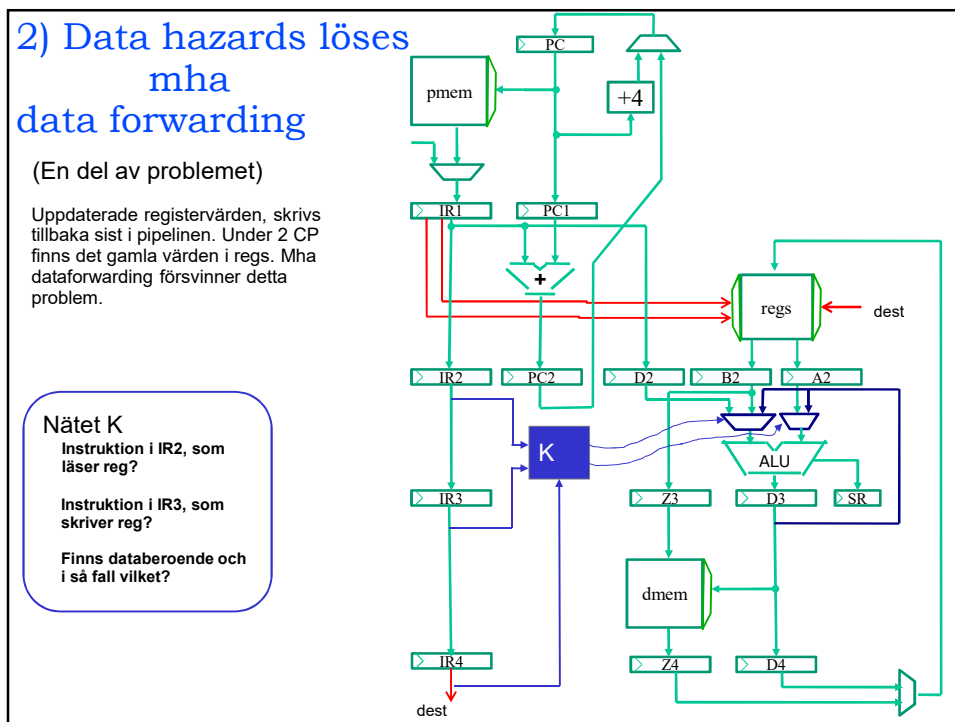
1



2



3



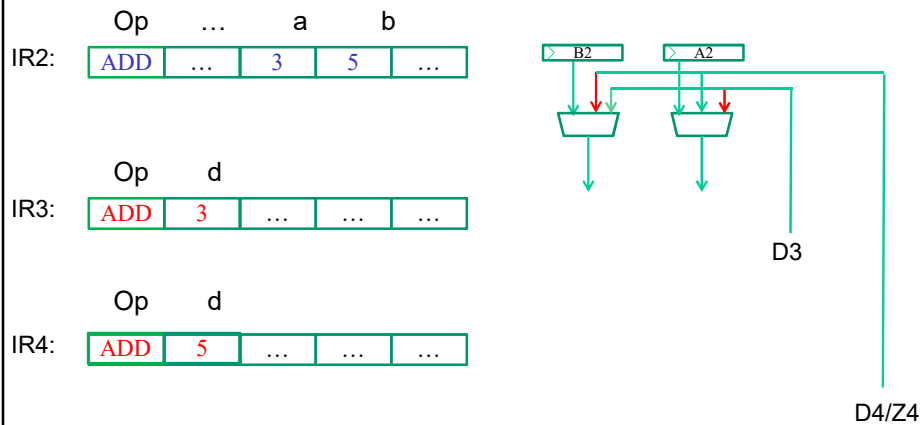
4

Databeroende

Regs kan innehålla gamla data!

Var kan nya resultat finnas? D3, D4/Z4

Hur ska muxarna styras?



5

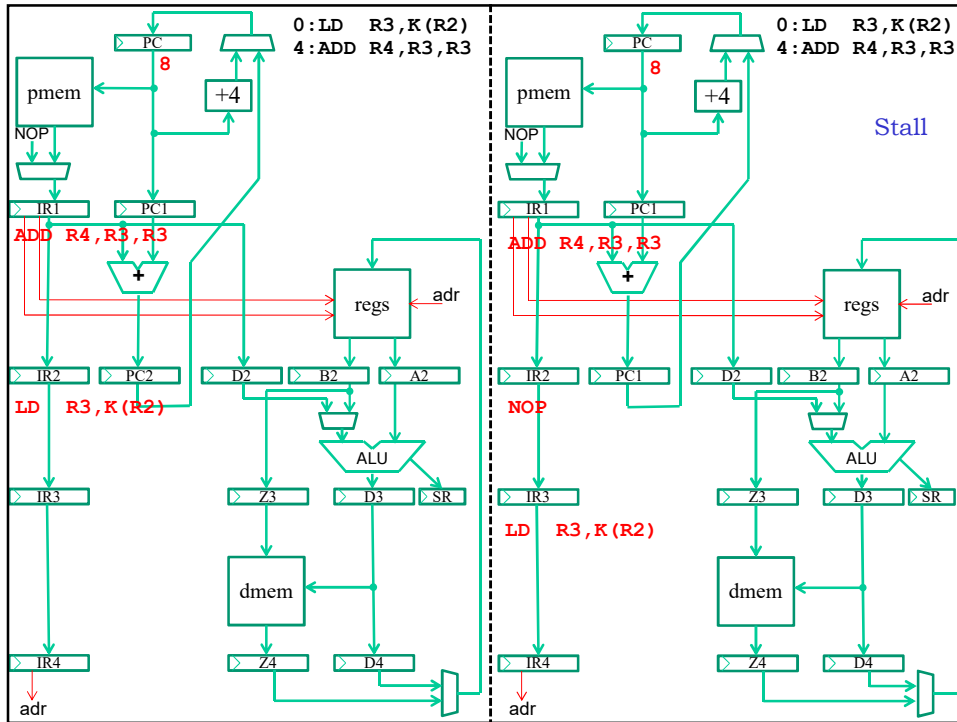
Problem 3: Stall

0:LD R3,K(R2) ; läs från minnet

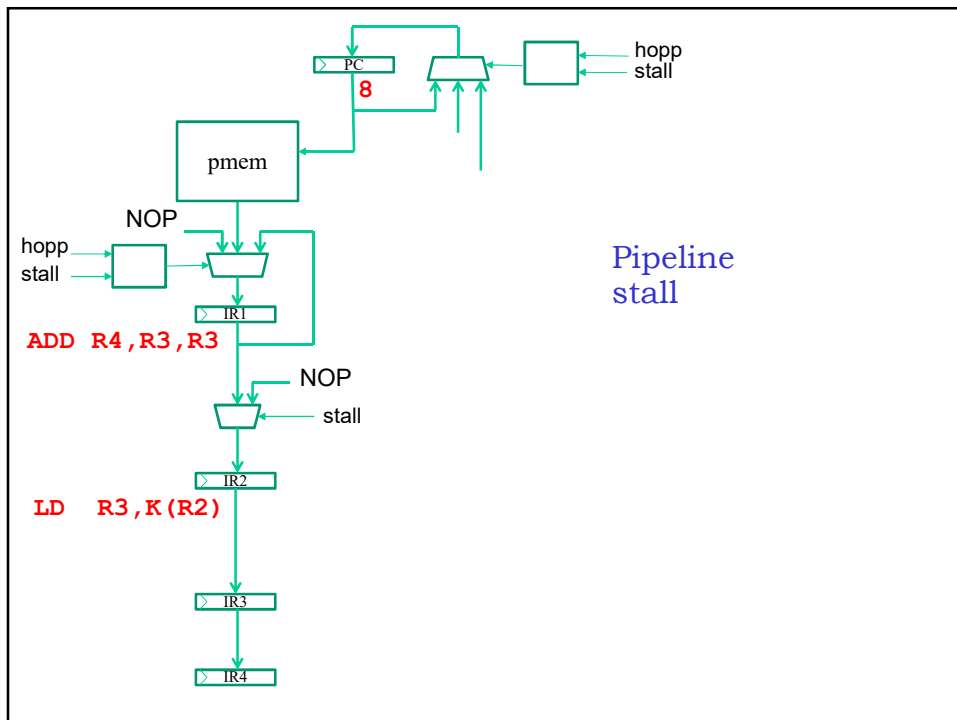
4:ADD R4,R3,R3 ;

Problemet beror på att minnet sitter 1 klockcykel efter ALU-n.

6



7



9

Observera att stall innebär prestandaförlust!

```
LD R2, 0 (R0)
ADD R3, R2, R1
```

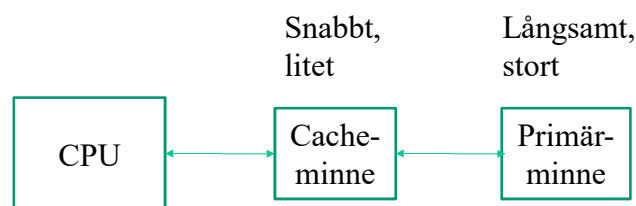
Vi kan lika gärna skriva

```
LD R2, 0 (R0)
NOP
ADD R3, R2, R1
```

Det går kanske att flytta en annan instruktion till NOP:ens plats?
Ovanstående bor kanske i **heta loopen**, där varje klockcykel är viktig.

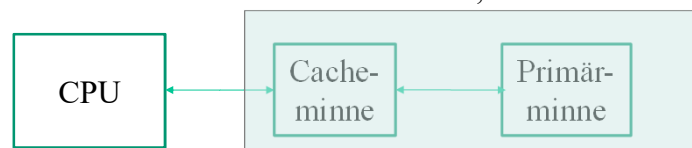
10

Cacheminnen



- CPU:n är normalt mycket snabbare än primärminnet. Det tar kanske 10-50 CP för att läsa ett data från minnet.
- Cacheminnet innehåller kopior av de instruktioner/data som används mest
- Kopieringen sker helt automatiskt, utan att programmeraren behöver tänka på det
- CM finns i CPU-chipet
- CM innehåller ett minne och en styrenhet.

Snabbt, stort



11

Lönar sig cacheminne i vår FPGA?

I FPGA-n finns sammanlagt 64 kB sk block RAM.

Det räcker inte med dessa, så **nej!**

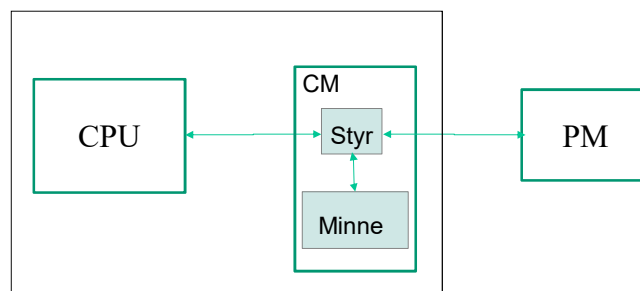
Utanför FPGA-n finns externt minne (flera MB)
men med accesstiden 7 CK.

Det är stort nog, så förvisso **ja!**

Ett CM måste i FPGA:n byggas med block RAM för att bli tillräckligt snabbt.

13

Cacheminnen



- När datorn startar är cacheminnet tomt
- Varje minnesaccess ger en **cachemiss**, kopiering till PM->CM
- Vanligt är att vid en cachemiss kopiera flera (t ex 4) instruktioner/data till CM. Kallas en **cacheline**.
- Sannolikheten är nu stor (t ex 95%) för **cacheträff**.

14

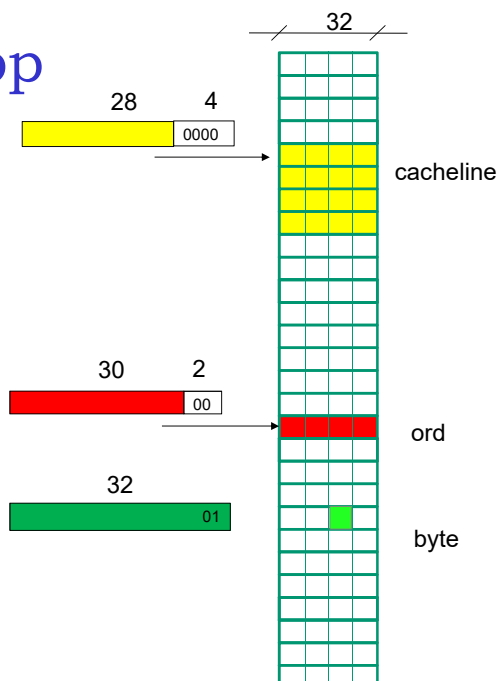
Cacheminnen

- ⤴ Program innehåller loopar => sannolikheten är stor att samma instruktion accessas igen (**temporal lokalitet**)
- ⤴ Samma sak gäller ofta för data
- ⤴ Program är sekvensiella => sannolikheten är stor att också nästa instruktion accessas (**spatial lokalitet**)
- ⤴ Dataaccesser brukar också vara lokala

15

Minnesbegrepp

- Minnesadress: 32 bitar
- Minnets bredd: 32 bitar, dvs 4 bytes
- Cacheline: 4 ord, dvs 16 bytes
- Möjligt att adressera bytes
- Ord finns på en multipel adress av 4
- Cachelines på en multipel adress av 16



16

Enkelt cacheminne

En styrenhet gör följande:

Vid **cachemiss**: Läs en hel cacheline.

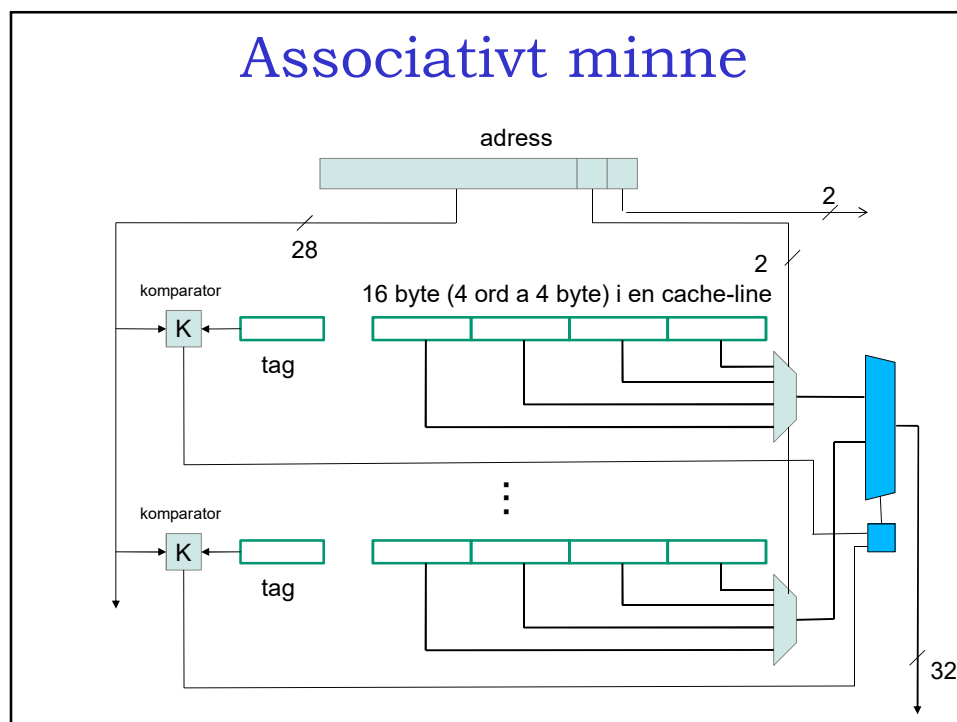
Placera de 4 orden på lämplig rad i CM

Placera de 28 MSB i adressen i taggen
(annars vet man ju inte var cacheline
kommer ifrån)

Vid access: Jämför de 28 MSB i adressen samtidigt
med alla taggar. Finns sökt data i cachem
har vi en **cacheträff** annars cachemiss.

17

Associativt minne

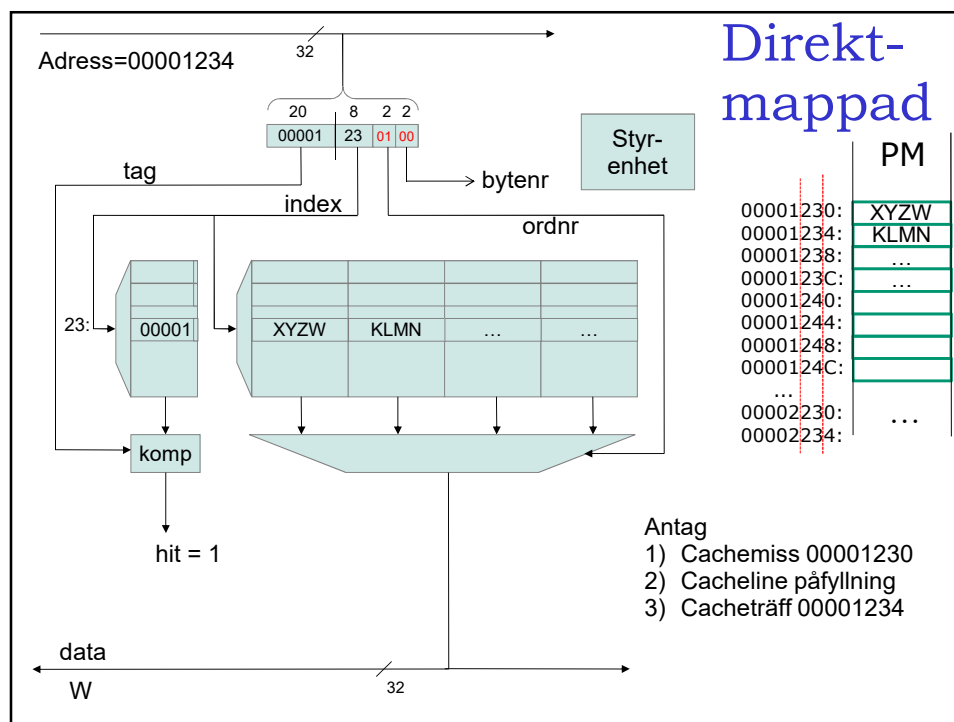


18

Cacheminnen

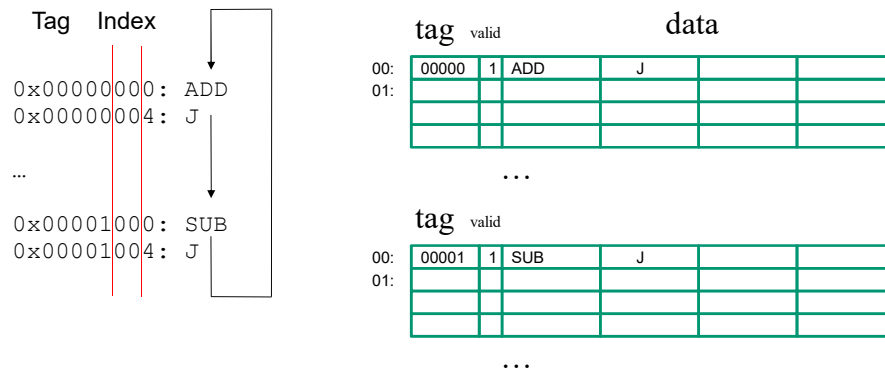
- Associativt minne (AM) är det mest generella sättet att bygga ett CM.
 - Data kan placeras var som helst i AM.
 - Vi hoppar över diverse problem:
Vilka rader är lediga?
När AM är fullt, vilken rad ska skrivas över?
 - AM är dyrt => en komparator per rad!
 - AM används inte till cache för instruktioner/data
-
- En sk **direct mapped cache** har endast 1 komparator
 - Data kan **inte** placeras var som helst!

19



24

Ett enkelt exempel?



Fastän cachen för övrigt är tom, får vi cachemiss hela tiden! (för ADD och SUB)
 Detta är nackdelen med en 1-vägs direktmappad cache.

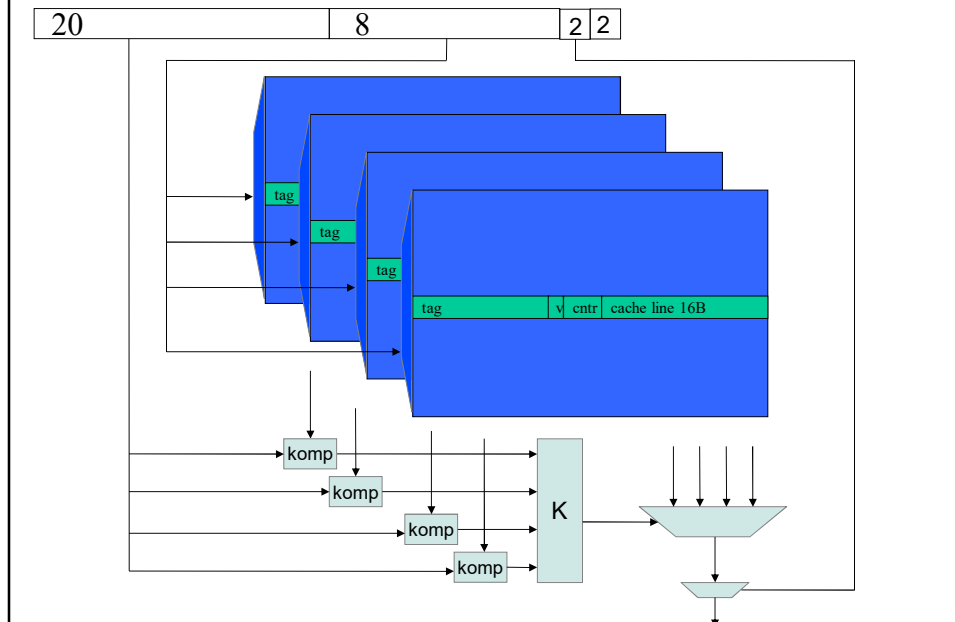
25

Förbättring av direktmappad

Vi gör en riktig cache.
 Vi gör cachen 4 gånger större, associativitet=4,
 genom att parallellkoppla 4 st 4kB cachar,
 alltså 16kB totalt.

26

4-vägs cache, 16 kB



27

Cacheminnen

- När man kört en stund, så är cachen full
- Vid en cachemiss, så ska ju cachen uppdateras, och ett data skrivs över
- En flervägs-cache behöver en utbytesalgoritm, vem blir offret? LRU är en vanlig metod

LRU = least recently used.

Till varje cacheline hör en räknare:

Hit: max->cntr

Miss: välj cacheline med minst cntr. Fyll på. max->cntr

Inget: räkna ner alla cntr (aging, inte på varje cp)

28

Cacheminnen

Mycket vanligt är att ha separata cacheminnen för instruktioner/data

För datacachen tillkommer då vad som ska ske vid skrivning:

Write-thru: skriv till CM och PM (PrimärMinne)
alltså CM,PM uppdateras samtidigt

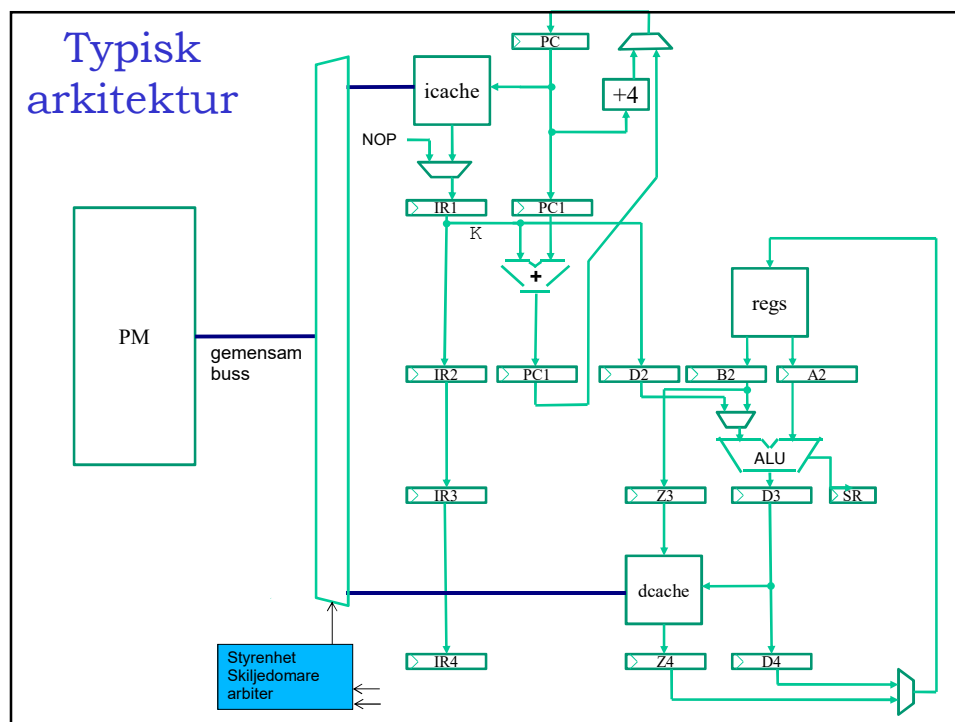
Write-back: PM uppdateras endast vid cachemiss

ytterligare en bit per cacheline: **dirty bit.** (CM förändrat)

Hit: läs/skriv CM

Miss: uppdatera PM om dirty=1. Skriv över.

29



30

Typisk arkitektur

- Vid cachemiss i I-cache eller D-cache måste pipelinen stoppas.
- Båda cacharna använder en gemensam buss mot minnet. Samtidig miss i I-cache och D-cache, innebär påfyllning av en cache i taget. Pipelinen fryst hela tiden.
- Observera att write-thru belastar den gemensamma bussen.

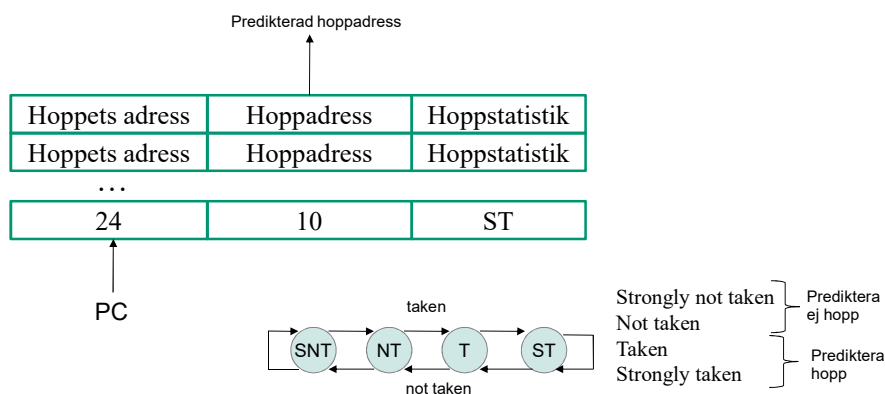
31

BPT = branch prediction table

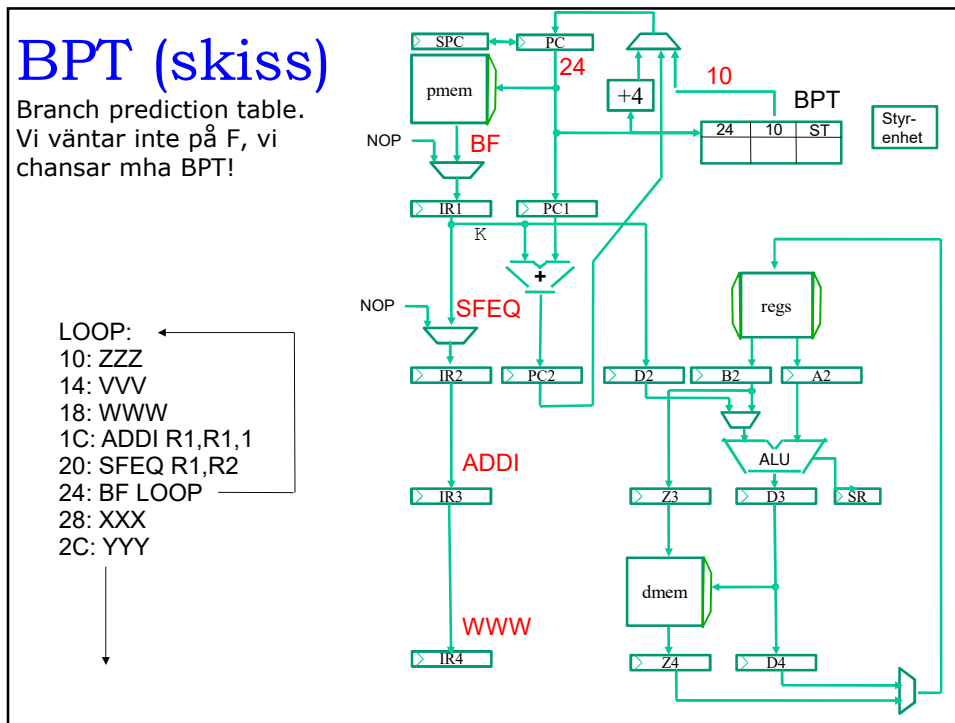
Ett försök att få bort den där NOP-en vid hopp

Skiss: Varje gång ett hopp påträffas

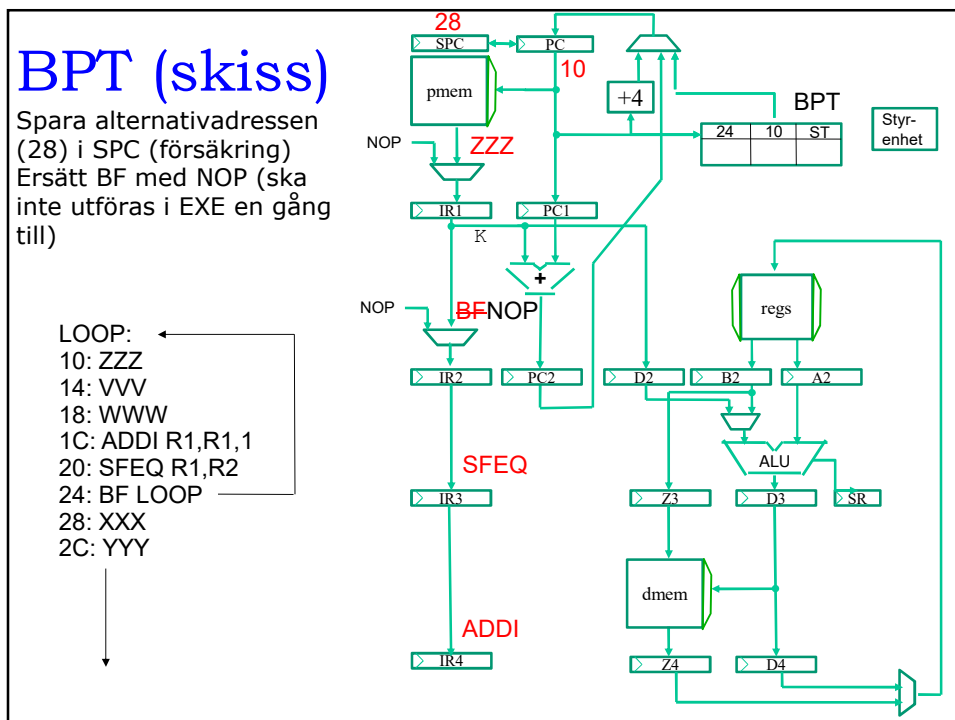
- 1) Finns inte i BPT: stoppa in det och uppdatera hoppstatistiken
- 2) Finns i BPT: läs ut hoppadressen



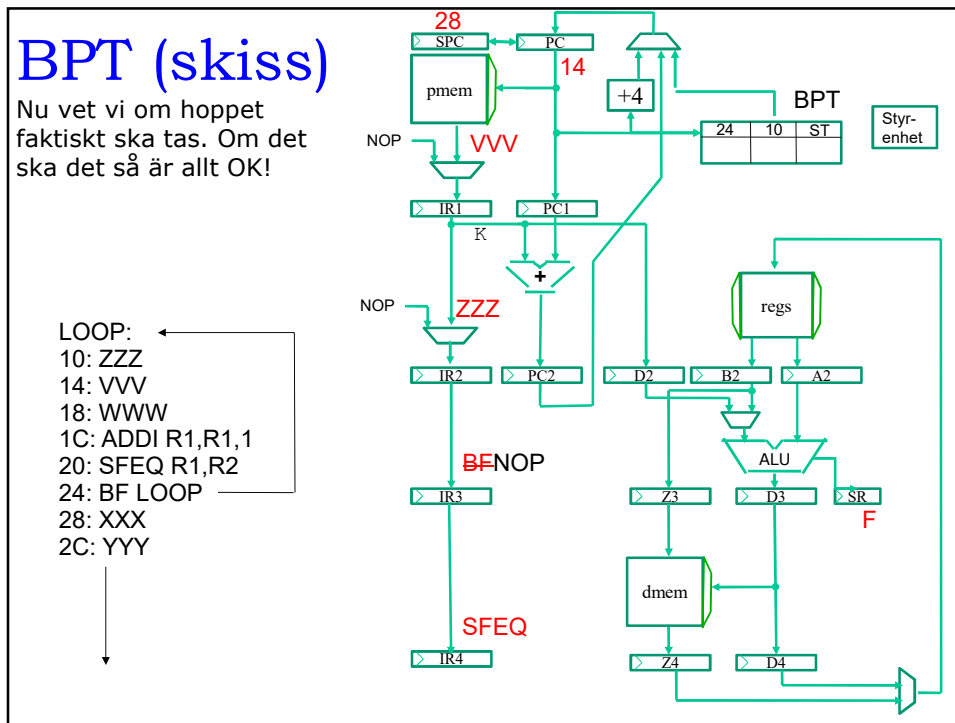
37



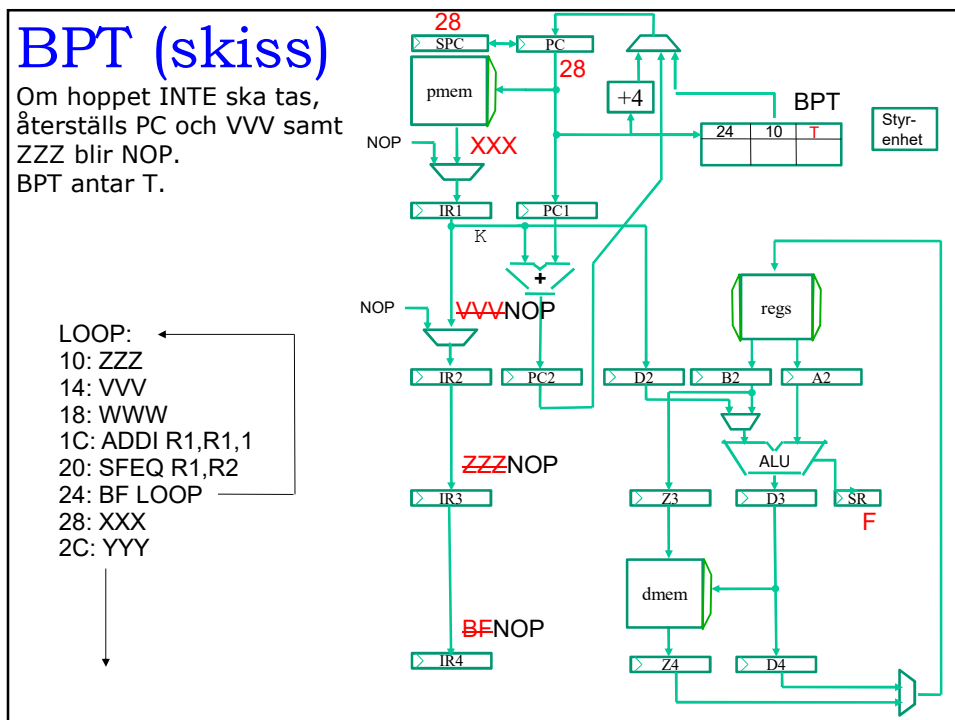
38



39



40



41

```

// Andreas Ehliar
...
int buf[4096]
...
unsigned int countones(void)
{
  unsigned int i;
  unsigned int numones = 0;
  unsigned int numzeroes = 0;
  volatile unsigned int dummy;

  for(i=0; i < BUFSIZE; i++){
    if(buf[i] == 1){
      numones++;
    }else{
      numzeroes++;
    }
  }

  dummy = numzeroes;
  return numones;
}

int main(int argc, char **argv)
{
  unsigned int i, numones;
  fillmem();

  for(i=0; i < 1000000; i++){
    numones = countones();
  }
  printf("Number of ones: %d\n", numones);
}

```

Benchmark

Mönster	Tid	
Bara 1:or	17 s	3.5 s
Varannan 1:a, 0:a	58 s	3.5 s
Slumpföljd av 1:or och 0:or	58 s	21 s

Kört på pandaboard
(ARM Cortex A-9
Smartphone-CPU)

Laptop core i5

44

Design of a baseband processor

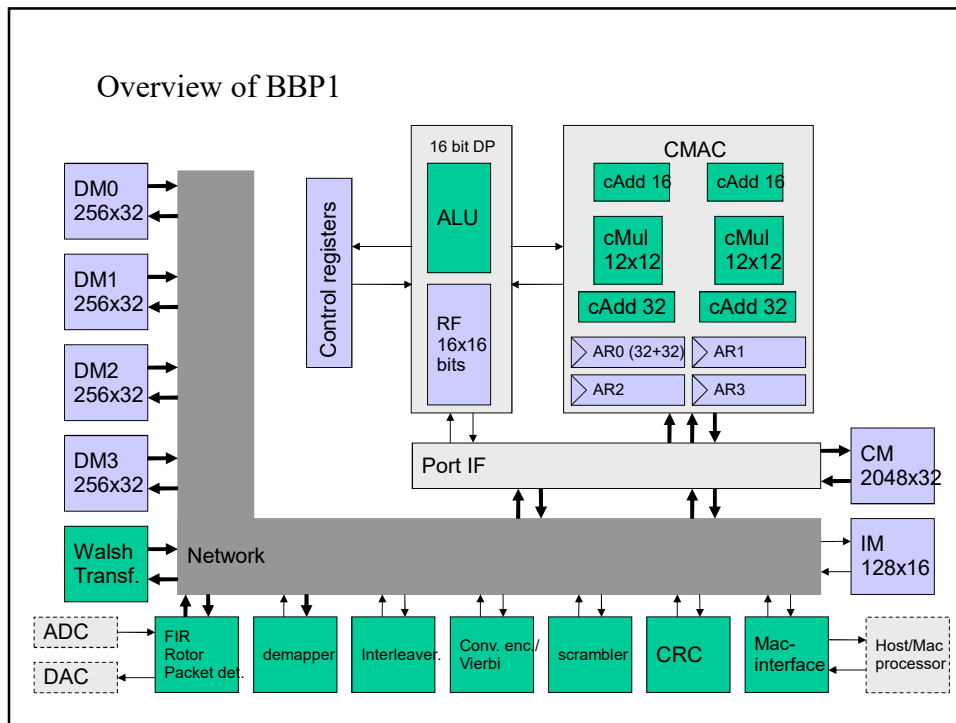
Eric Tell & Anders Nilsson

Coresonic (uppköpt av Mediatek)

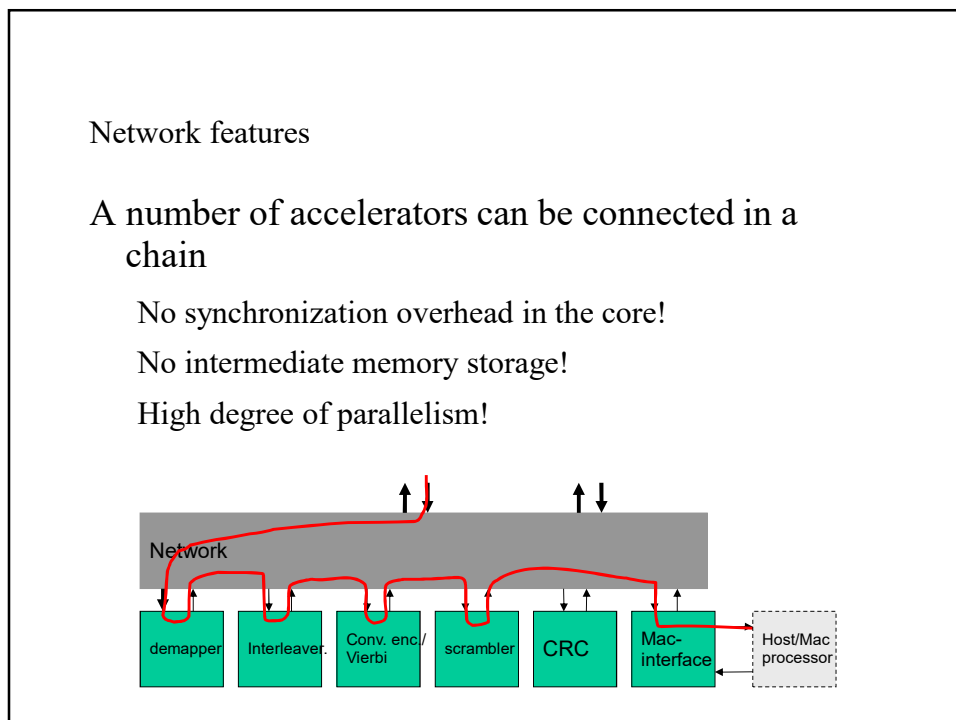
Ska klara så många standarder som möjligt
WLAN, 3G, DVB-T, ...

Använder en kombination av
mjukvara och hårdvara (acceleratorer)

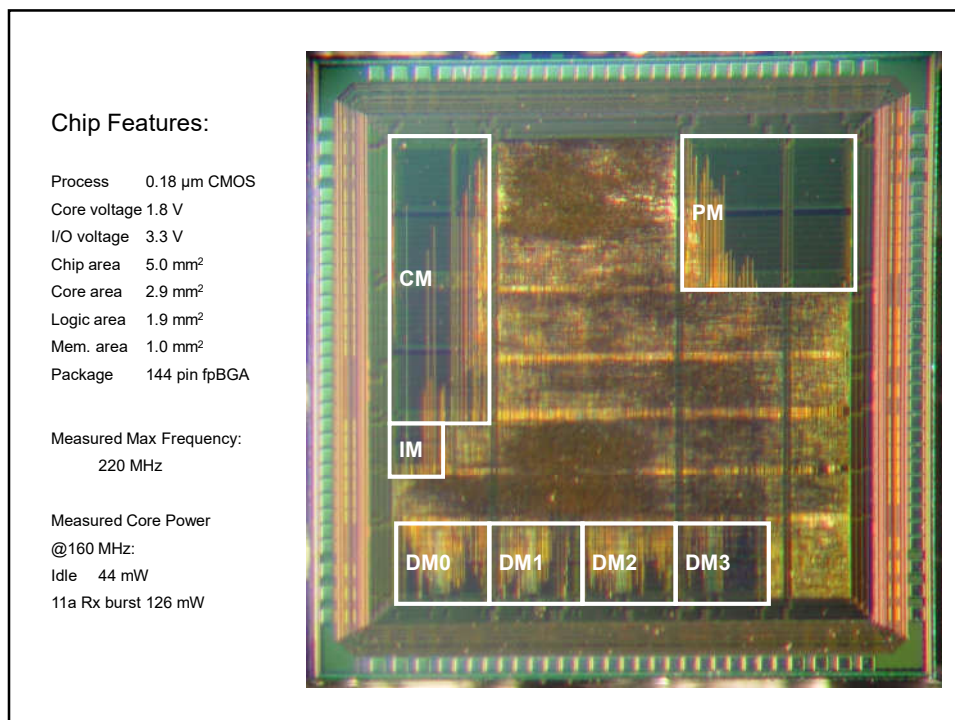
45



46



47



48