

FKEF.02.143

Arvuti arhitektuur

Computer Architecture and Organization




dotsent
Toomas Plank

©Toomas Plank, 2008

FKEF.02.143

Vahemälu

Cache



11. loeng,
9. mai 2008

©Toomas Plank, 2008

Vajadus kiire mälu järele

- Võrreldes protsessori kiirusega on mälu aeglane seade
- Vahemälu võimaldab mädul näida kiiremana kui ta tegelikult on
- Vahemälu efektiivsus põhineb järgmistel tõdemusel:
 - Enamus programmide tööajast kulub tsükklites ja samade protseduuride sagedases täitmises
 - Hiljuti täidetud käsud tulevad suure tõenäosusega varsti uuesti täitmisele
 - Viimati täidetud käsule lähedal (mälu paiknemise mõttes) paiknevad käsud tulevad samuti suure tõenäosusega peagi täitmisele

Arvuti arhitektuur FKEF.02.143

Vahemälu (Cache)

- Paigutades programmi aktiivse osa kiirsesse vahemälusse, saame suure võidu programmi tööajas
- Eelmisel slaidil kirjapandud **ajalise läheduse tõdemus** soovitab lugeda kõik vajaminevad käsud vahemällu
 - tõenäoliselt läheb neid varsti uuesti tarvis
- Eelmisel slaidil kirjapandud **ruumilise läheduse tõdemus** soovitab lugeda vahemällu mitte üksikud käsud vaid ka nende naabruses asuvad käsud
 - s.t loeme andmed vahemällu plokkide kaupa (*block, cache line*)
 - tõenäoliselt läheb ka neid käskke varsti tarvis

Arvuti arhitektuur FKKE02.143

Põhimälu ja vahemälu



- Kui vajatakse mälust andmeid, siis loetakse need vahemällu
- Kui vajalikud andmed on juba vahemälus olemas, siis loetaksegi kohe sealt
- Tagatakse vastavus põhimälu ja vahemälus resideeruva koopia vahel (*mapping function*)
- Kui vahemälu on täis, tuleb uute andmete lisandumisel miskit olemasolevat vahemälust välja visata (*replacement algorithm*)

Arvuti arhitektuur FKKE02.143

Vahemälu töökorraldus

- Protsessor ei pea olema teadlik vahemälu olemasolust,
 - tema annab ainult mälust lugemise ja mälusse kirjutamise käskke
 - Vahemälu kontrollahelad selgitavad välja, kas vajalik info on juba vahemälus (tabamus)
 - või tuleb info alles vahemälusse lugeda (möödalask)
 - Kui andmed on juba vahemälus, siis **põhimälust lugeda** pole tarvis
- Kirjutamiseks on kaks varianti:
- *write-through* – nii vahemälu kui ka mälu uuendatakse üheaegselt
 - *write-back, copy-back* – uuendatakse ainult vahemälu
 - Vahemälust ploki eemaldamisel uuendatakse info lõpuks ka põhimälus
 - Info uuendamise vajaduse kohta kajastub muudetud lipubitits (*modified, dirty bit*)

Arvuti arhitektuur FKKE02.143

Vahemälu töökorraldus

- Lugemine. Kui infot pole vahemälus, siis
 - loetakse kogu plokk põhimälust vahemällu ja seejärel edastatakse vajalik sõna protsessorile
 - loetakse info põhimälust ja edastatakse vajalik sõna protsessorile kohe kui ta on olemas (veidi kiirem, aga nõuab keerukamaid juhtimisahelaid)
- Kirjutamine. Kui infot pole vahemälus, siis
 - Kirjutatakse info põhimällu ilma vahemälu vahenduseta (*write-through*)
 - Loetakse esmalt kogu plokk vahemällu ja siis muudetakse vastava sõna vahemälus (*write-back*)

Arvuti arhitektuur FKFE.02.143

Info paigutus vahemälus

Paigutusviisid

- Otsene paigutus (*Direct Mapping*)
- Assotsiatiivne paigutus (*Associative Mapping*)
- Assotsiatiivse rühmana paigutus (*Set-Associative Mapping*)

Vaatleme neid paigutusviise lähemalt, kasutades näitena

- 4096 plokist koosnevat põhimälu
- 128 plokist koosnevat vahemälu
- Igas plokis olgu meil 16 sõna
- Seega põhimälu $16 \times 4096 = 65536$ sõna
- ja vahemälu $16 \times 128 = 2048$ sõna

Arvuti arhitektuur FKFE.02.143

Otsene paigutus (*Direct Mapping*)

- Mälu plokk j pannakse C plokist koosnevas vahemälu plokki

$$j \bmod C$$

Silt

Plokk nr 0

Silt

Plokk nr 1

Silt

Plokk nr 2

⋮

Silt

Plokk nr 127

- Mälu aadress jagatakse kolmeks osaks

- Käesolevas näites on meil 16-bitised mälu aadressid

Plokk nr 0

Plokk nr 1

Plokk nr 2

⋮

Plokk nr 127

Plokk nr 128

Plokk nr 129

Plokk nr 130

⋮

Plokk nr 4095

Arvuti arhitektuur FKFE.02.143

Otsene paigutus (Direct Mapping)

Silt Plokk Sõna

b₁₅ b₁₄ b₁₃ b₁₂ b₁₁ b₁₀ b₉ b₈ b₇ b₆ b₅ b₄ b₃ b₂ b₁ b₀

16 bitti

Silt Plokk nr 0

Silt Plokk nr 1

Silt Plokk nr 2

⋮

Silt Plokk nr 127

- Sõna-väli määrab ära sõna asukoha plokkis
- Ploki-väli määrab ära ploki asukoha vahemälus
- Sildi-väli määrab ära milline põhimälu plokkidest vahemälusse kirjutati

Arvuti arhitektuur FKFE02.143

Assotsiatiivne paigutus (Associative Mapping)

- Mälu ploki *j* saab panna suvalisse kohta vahemälus
- Mälu aadress jagatakse kaheks osaks
 - Käesolevas näites on vaja 12-bitti ploki identifitseerimiseks (sildi-väli)
 - Vahemälu kasutus efektiivsem
- Miinus: vahemälu kasutamisel tuleb läbi vaadata kõik sildi-väljad

Plokk nr 0

Plokk nr 1

Plokk nr 2

⋮

Plokk nr 127

Plokk nr 128

Plokk nr 129

Plokk nr 130

⋮

Plokk nr 4095

Silt Plokk nr 0

Silt Plokk nr 1

Silt Plokk nr 2

⋮

Silt Plokk nr 127

Arvuti arhitektuur FKFE02.143

Assotsiatiivne paigutus (Associative Mapping)

Silt Sõna

b₁₅ b₁₄ b₁₃ b₁₂ b₁₁ b₁₀ b₉ b₈ b₇ b₆ b₅ b₄ b₃ b₂ b₁ b₀

16 bitti

Silt Plokk nr 0

Silt Plokk nr 1

Silt Plokk nr 2

⋮

Silt Plokk nr 127

- Sõna-väli määrab ära sõna asukoha plokkis
- Ploki-välja pole, asukoht vahemälus vabalt valitav
- Sildi-väli määrab ära milline põhimälu plokkidest vahemälusse kirjutati

Arvuti arhitektuur FKFE02.143

Assotsiatiivne rühm (Set-Associative Mapping)

- Kahe eelmise variandi kombinatsioon
- Mälu plokki j saab panna vahemälu ettemääratud alamhulga suvalisse plokki

Silt	Plokk nr 0	Mälu aadress jagatakse kolmeks osaks	Plokk nr 0
Silt	Plokk nr 1		Plokk nr 1
Silt	Plokk nr 2		Plokk nr 2
Silt	Plokk nr 3		Plokk nr 3
...
Silt	Plokk nr 126	Käesolevas näites on vaja 6-bititi alamhulga identifitseerimiseks (alamhulga-väli) – praegu kaks plokki alamhulgas	Plokk nr 63
Silt	Plokk nr 127		Plokk nr 64
Arvuti		ja 6-bititi plokki identifitseerimiseks (sildi-väli)	Plokk nr 65
			Plokk nr 66
			...
			Plokk nr 4095

Assotsiatiivne rühm (Set-Associative Mapping)

Silt	Alamhulk	Sõna
$b_{15} b_{14} b_{13} b_{12} b_{11} b_{10}$	$b_9 b_8 b_7 b_6 b_5 b_4$	$b_3 b_2 b_1 b_0$

16 bititi

- Sõna-väli määrab ära sõna asukoha plokis
- Alamhulga-väli määrab ära plokki võimalikud asukohad vahemälu
- Sildi-väli määrab ära milline põhimälu plokkidest vahemälusse kirjutati

Silt	Plokk nr 0
Silt	Plokk nr 1
Silt	Plokk nr 2
...	...
Silt	Plokk nr 127

Arvuti arhitektuur FKEE.02.143

Asjakohasuse bitt Valid bit

- See bitt ütleb, kas andmed vahemälu on asjakohased
- Arvuti sisselülitamisel kõik need bitid väärtusega **0**
- Info laekumisel seatakse selle biti väärtuseks **1**
- Kui põhimälu uuendatakse vahemälust mööda minnes, tuleb kontrollida selle plokki olemasolu/puudumist vahemälu
 - kui see plokk vahemälu, siis seatakse selle biti väärtuseks **0**
- Samalaadsed mured DMA ülekannete puhul mälust kõvakettale
 - Siin on võimalik olukord, kus põhimälu ei ole adekvaatne info (write-back protocol)
 - Üks võimalus andmed jõuga kaasajastada ja alles siis DMA ülekannet teha
 - Seda probleemi nimetatakse vahemälu koherentsuse küsimuseks

Arvuti arhitektuur FKEE.02.143

Asendusalgoritmid

- Otsese paigutuse korral on põhimälu ploki paigutuskoht vahemälus üheselt teada ja pole vaja pead murda selle üle, millist plokki vahemälus asendada
- Assotsiatiivse paigutuse ja assotsiatiivse rühma paigutuse korral on valik suurem
- **NB! Asendusalgoritmi valik mõjutab programmi jõudlust väga palju!**
- LRU (*least recently used*) *replacement algorithm*
Vahetame välja ploki, mille viimatisest kasutamisest on kõige rohkem aega möödas
 - kasutusinfo jaoks peab olema spetsiaalne loendur
- Vanima elemendi väljavahetamine
 - miinus: ei arvesta hiljutisi kasutusmustrid
- Juhusliku ploki asendamine
 - Vaatamata lihtsusele üllatavalt head tulemused!

Arvuti arhitektuur FKEE.02.143



Näide

- Oletame, et mõõtsime kolme patarei pinget, igat patareid 9 korda
 - Mõõdised M panime arvuti mälli 3x9 maatriksina aadressidele 2700...271A
- | | |
|---------------------|--------|
| 0010 0111 0000 0000 | M(0,0) |
| 0010 0111 0000 0001 | M(1,0) |
| 0010 0111 0000 0010 | M(2,0) |
| 0010 0111 0000 0011 | M(0,1) |
| 0010 0111 0000 0100 | M(1,1) |
| 0010 0111 0000 0101 | M(2,1) |
| ... | ... |
| 0010 0111 0001 1000 | M(0,8) |
| 0010 0111 0001 1001 | M(1,8) |
| 0010 0111 0001 1010 | M(2,8) |

Arvuti arhitektuur FKEE.02.143

Sildi väli:

Assotsiatiivne
Arvuti



Näide: arvutusülesanne

- Oletame, et meil on vaja leida iga mõõdise hälve aritmeetilisest keskmisest ja kirjutada see mõõdise asemel samasse mälupeksasse
- C# programm selle rehkenduse tegemiseks võiks olla alljärgnevalt:

```
for (int j = 0; j < 3; j++)
{
    //Leiame rea summa
    double Rea_Summa = 0;
    for (int i = 0; i < 9; i++)
        Rea_Summa += M[j, i];
    //Leiame keskmise
    double Keskmise = Rea_Summa / 9;
    //Rehkendame hälbed
    for (int i = 0; i < 9; i++)
        M[j, i] = M[j, i] - Keskmise;
}
```

Arvuti a



Otsese paigutusega vahemälu

0	...	M(0,0)
1	...	M(0,1)
2	...	M(0,2)
3	...	M(0,3)
4	...	M(0,4)
5	...	M(0,5)
6	...	M(0,6)
7	...	M(0,7)
		M(0,8)

- Vahemälu pesa aadressi saame mälu aadressist tehtega mod 8
- Sisuliselt määravad asukoha ära mälu aadressi kolm viimast bitti

Arvuti arhitektuur FKEEF.02.143

Assotsiatiivse paigutusega vahemälu

0	...	M(0,0)
1	...	M(0,1)
2	...	M(0,2)
3	...	M(0,3)
4	...	M(0,4)
5	...	M(0,5)
6	...	M(0,6)
7	...	M(0,7)
		M(0,8)

- Esimene tsükkel $i = 0 \dots 8$
- Teine tsükkel $i = 0 \dots 8$

Arvuti arhitektuur FKEEF.02.143

Assotsiatiivse paigutusega vahemälu

0	...	M(0,8)
1	...	M(0,7)
2	...	M(0,6)
3	...	M(0,5)
4	...	M(0,4)
5	...	M(0,3)
6	...	M(0,2)
7	...	M(0,1)
		M(0,0)

- Esimene tsükkel $i = 0 \dots 8$
- Teine tsükkel $i = 8 \dots 0$

Arvuti arhitektuur FKEEF.02.143

Vahemälu assotsiatiivse rühmana

0	...	Alamhulk 0	• Olgu meil kaks alamhulka	M(0,0)
1	...			
2	...			
3	...			
4	...	Alamhulk 1	• Igas võimalik hoida neli sõna	M(0,4)
5	...			
6	...			
7	...			

• Viimase biti aadress määrab siis ära alamhulga
 - 0 → alamhulk 0
 - 1 → alamhulk 1

Arvuti arhitektuur FKEF.02.143

Vahemälu assotsiatiivse rühmana

0	...	Alamhulk 0	• Olgu meil kaks alamhulka	M(0,8)
1	...			
2	...			
3	...			
4	...	Alamhulk 1	• Igas võimalik hoida neli sõna	M(0,4)
5	...			
6	...			
7	...			


• Viimase biti aadress määrab siis ära alamhulga
 - 0 → alamhulk 0
 - 1 → alamhulk 1

Arvuti arhitektuur FKEF.02.143

FKEF.02.143

Jõudlus

Performance



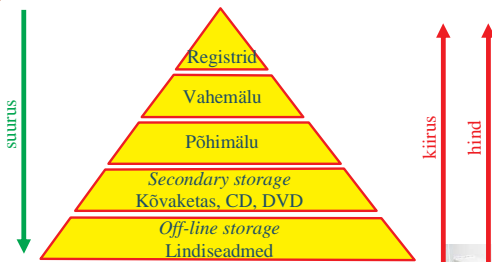
©Toomas Plank, 2008

Jõudlust mõjutavad tegurid

- Kaks peamist tegurit, mis mõjutavad konkreetse arvuti ärilist edukust on
 - hind
 - jõudlus
- Soovime saada parimat jõudlust võimalikult madala hinnaga
- Mälu puhul on tähtis,
 - kui kiiresti saame käsud/andmed mälust protessorisse
 - kui suur on mälu maht

Arvuti arhitektuur FKEE.02.143

Mälu hierarhia



Arvuti arhitektuur FKEE.02.143

Tabavus (hit rate)

$$\text{Tabavus} = \frac{\text{Sõna vahemälust leidmiste arv}}{\text{Mälupöördumiste koguarv}}$$

- Näide:
 - Oletame, et sõnade mälust lugemisel leiti 952 juhul sõna vahemälust
 - ja 47 juhul sõna vahemälus ei olnud
- siis tabavus

$$\text{Tabavus} = \frac{952}{952 + 47} = \frac{952}{999} = 95,3\%$$

Arvuti arhitektuur FKEE.02.143

Efektívne pöördumisaeg

$$\text{Efektívne pöördumisaeg} = \frac{\text{tabamusi} \times t_{\text{tabamus}} + \text{möödalaske} \times t_{\text{möödalask}}}{\text{Mälupöördumiste koguarv}}$$

Näide:

- Oletame, et sõnade mälust lugemisel leiti 952 juhul sõna vahemälust
- ja 47 juhul sõna vahemälus ei olnud
- vahemälust lugemisele kulus 25 ns
- põhimälust lugemisele koos vahemällu kirjutamisega kulus 2000 ns

$$\text{Efektívne pöördumisaeg} = \frac{952 \times 25 \text{ ns} + 47 \times 2000 \text{ ns}}{952 + 47}$$

$$= \frac{23800 \text{ ns} + 94000 \text{ ns}}{999} = 118 \text{ ns}$$

Arvuti arhitektuur FKEE.02.143

Mitmetasemeline vahemälu

$$\text{Tabavus}_{L1} = \frac{\text{Sõna vahemälust L1 leidmiste arv}}{\text{Mälupöördumiste koguarv}}$$

$$\text{Tabavus}_{L2} = \frac{\text{Sõna vahemälust L2 leidmiste arv}}{\text{Sõna vahemälust L1 mitteleidmiste arv}}$$

$$\text{Efektívne pöördumisaeg} = \frac{L1 \text{ tabamusi} \times t_{L1} + L2 \text{ tabamusi} \times t_{L2} + L2 \text{ möödalaske} \times t_{\text{möödalask}}}{\text{Mälupöördumiste koguarv}}$$

Arvuti arhitektuur FKEE.02.143



Näide

- Oletame, et sõnade mälust lugemisel leiti 9 520 juhul sõna vahemälust L1
- 470 juhul sõna vahemälus L1 ei olnud
- ja 10 juhul polnud seda sõna ka vahemälus L2
- vahemälust L1 lugemisele kulus 5 ns
- vahemälust L2 lugemisele kulus 25 ns
- põhimälust lugemisele koos vahemällu kirjutamisega kulus 200 ns

• siis efektívne pöördumisaeg

$$\text{Efektívne pöördumisaeg} = \frac{9520 \times 5 \text{ ns} + 460 \times 25 \text{ ns} + 10 \times 200 \text{ ns}}{9520 + 470}$$

$$= \frac{47600 \text{ ns} + 11500 \text{ ns} + 2000 \text{ ns}}{9990} = \frac{61100 \text{ ns}}{9990} = 6,12 \text{ ns}$$

Arvuti arhitektuur FKEE.02.143

Kasutatud kirjandus

- Carl Hamacher, Zvonko Vranesic, Safwat Zaky, Computer organization 5th edition (2002) 805 p.
- Miles R. R. Murdocca, Vincent Heuring, Computer Architecture and Organization: An Integrated Approach (2007) 544 p.

