

1. Cache

Resymé: Leksjonen beskriver cache-mekanismen. Tre forskjellige mapping-funksjoner diskuteres. Forskjellige utskiftingsalgoritmer beskrives også. Strategier for skriving via cache beskrives, og en del fundamentale problemer med cache diskuteres.

Innhold

1. CACHE.....	1
1.1. HVA ER CACHE?.....	2
1.1.1. Innledning	2
1.2. PRINSIPPET OM LOKALITET	2
1.2.1. Romlig og temporal lokalitet.....	3
1.2.2. Cache-størrelse	3
1.3. CACHENS VIRKEMÅTE	4
1.3.1. Cachens plassering	4
1.4. TREFFRATE	5
1.5. EFFEKTIV AKSESSTID VED BRUK AV CACHE	6
1.6. MER OM CACHENS VIRKEMÅTE.....	7
1.7. BLOKKSTØRRELSE	8
1.8. MAPPING-FUNKSJON	8
1.8.1. Full-assosiativ cache.....	8
1.8.2. Ikke-assosiativ cache (direkte mapping)	9
1.8.3. Sett-assosiativ cache	9
1.9. UTSKIFTINGSALGORITMER	10
1.10. SKRIVING TIL MINNET	10
1.11. CACHE-KOHERENS.....	10
1.12. HASTIGHETSFORBEDRING VED BRUK AV CACHE	12
2. MODERNE CACHE-DESIGN	12
2.1.1. Fler-nivå caching	12
2.1.2. Splittet cache kontra enhetlig cache.....	13
2.1.3. Egen kontra felles (delt) cache.....	13
2.1.4. Sett-størrelse	13
2.1.5. Utviklingstrender innen cache.....	14
2.2. SAMMENKOBLING AV DELENE	15
VEDLEGG A. MER OM MAPPING-FUNKSJONER.	16
Full-assosiativ cache	16
Ikke-assosiativ cache (direkte mapping).....	18
Men hvordan sjekkes det om en blokk ligger i cache?.....	19
Sett-assosiativ cache.....	20

1.1. Hva er cache?

I en tidligere leksjon s  vi at moderne CPUer har et problem fordi CPU er s  kjapp at prim rminnet basert p  DRAM ikke greier   levere instruksjoner og data med tilstrekkelig hastighet. For   b te p  dette brukes en mekanisme som vi skal beskrive i denne leksjonen, nemlig Cache.

1.1.1. Innledning

Ting en bruker ofte er det kj kt   ha liggende like ved seg, og slik er det ogs  i datamaskinen. Vi kan oversette cache med hurtigbuffer. Cache er en mekanisme som fors ker   holde den informasjonen som brukes ofte inne i et kj ppere minne enn prim rminnet.

P  k kkenet mitt har jeg vekt, hurtigmikser og sleiv liggende innerst i skapet. N r jeg skal bake tar jeg dette fram. Mens jeg baker lar jeg det ligge p  k kkenbenken. P  den m ten slipper jeg   bruke tid p    grave det fram fra skapet for hver gang jeg skal bruke det. N r jeg er ferdig med   bake trenger jeg plassen p  k kkenbenken til andre ting, s  da m  jeg legge bakeutstyret tilbake i skapet.

Jeg bruker ogs  k kkenbenken min som hurtigbuffer.

Legg merke til at:

- 1) N r jeg f rst er inne i skapet og graver tar det ikke s rlig lenger tid   hente alt bakeutstyret. Istedenfor   hente bare en og en ting, henter jeg flere ting i samme slengen.
- 2) K kkenbenken min virker som et hurtigbuffer fordi "aksesstiden" til k kkenbenken er kortere enn "aksesstiden" til skuffene.

Cache er ogs  et minne som har kortere aksessid enn prim rminnet - s  kort aksessid at det holder tritt med CPUen

Dette f rer til en liten endring i virkem ten til datamaskinen. Tidligere sa vi at prosessoren hentet inn en instruksjon og utf rte den - deretter hentet den neste instruksjon og utf rte den.

Med cache hentes det isteden flere instruksjoner fra prim rminnet. Alle disse legges i cachen. Neste gang det skal utf res en instruksjon sj kker prosessoren f rst om instruksjonen allerede ligger i cachen. Hvis den gj r det, ja s  f r CPUen tilgang til den straks. Dermed slipper prosessoren   vente p  det trege prim rminnet. Dersom neste instruksjon ikke ligger i cachen, m  CPU vente p  at den hentes fra prim rminnet. I praksis er det faktisk slik at CPU nesten alltid finner neste instruksjon i cachen. Derfor er cache en meget effektivt virkemiddel til    ke ytelsen til en prosessor.

For bakeutstyret i rammen ovenfor var det enkelt   se hvorfor det l nte seg   hente alt bakeutstyret fra skapet med en gang. N r vi f rste begynner   bake, skal vi jo gj re oss ferdig - vi begynner for eksempel ikke   lage middag midt oppi all bakingen. S  lenge bakingen foreg r, er det bakeutstyret vi har bruk for. Sp rsm let er: hvorfor er det tilsvarende p  en datamaskin? Hvorfor er det sannsynlig at neste instruksjon ligger like i n rheten av de foreg ende? For   forklare det m  vi se p  noe som kalles «prinsippet om lokalitet».

1.2. Prinsippet om lokalitet

Prinsippet om lokalitet er et sv rt viktig prinsipp i datateknikk.

Prinsippet sier at dersom en minnelokasjon er benyttet en gang, er det sv rt sannsynlig at den - eller en lokasjon like ved siden av - snart vil bli benyttet en gang til. I en cache utnyttes dette p  den m ten at cachen tar vare p  en kopi av de sist aksesserte minnelokasjonene og deres nabolokasjoner. Det er nemlig sannsynlig at det snart vil bli bruk for dem igjen.

Prinsippet om lokalitet medf rer alts  at minnereferansene ikke er spredt tilfeldig rundt i minnet, men har en tendens til   "klumpe seg sammen". Det betyr at det til en hver tid er bare en liten del av minnet som blir brukt intenst. I l pet av denne stunden ligger resten av minnet nesten ur rt. Over tid blir hele minnet brukt - men over relativt lange tidsrom (i forhold til CPU-hastigheten) aksesseres en liten del av minnet gjentatte ganger.

Noen grunner til dette er:

- Som regel ligger instruksjonene sekvensielt. Det betyr at neste instruksjon ligger i minnelokasjonen etter den forrige.
- Et program inneholder som regel mange l kker som gjentas flere ganger. Et enkelt eksempel er en while-l kke. Instruksjonene i en slik l kke gjentas mange ganger, noe som betyr at de samme minne-lokasjonene stadig refereres p  nytt.
- Eksempler p  instruksjoner der prinsippet **ikke** gjelder, er blant annet hopp-instruksjoner og funksjonskall (kall av subrutiner/prosedyrer). Begge deler er eksempler p  at programutf ringen flyttes til et nytt sted i minnet. Slike instruksjoner utgj r imidlertid bare en liten del av instruksjonene i et typisk program, og utf res derfor sjelden.
- Det er sjelden slik at en funksjon straks kaller en ny funksjon som straks kaller en ny, som straks kaller en ny, osv ... Det er alts  sjelden at funksjoner n stes s rlig dypt. Da gjelder igjen at de samme instruksjoner har en tendens til   bli brukt p  nytt, relativt kort tid etter at de ble brukt forrige gang.
- Mange programkonstruksjoner er iterative. Det vil si at et lite antall instruksjoner gjentas mange gang.
- De aller fleste program benytter datastrukturer som er sammenhengende. Eksempler er arrays og records. Prinsippet om lokalitet gjelder alts  ikke bare for instruksjoner, men ogs  for data.

1.2.1. Romlig og temporal lokalitet

Vi skiller gjerne mellom *romlig lokalitet* (engelsk: spatial locality) og *temporal* (eller tidsmessig) *lokalitet* (temporal locality p  engelsk).

Romlig lokalitet betyr at minnereferansene ligger n rt hverandre. Sekvensiell utf ring er et eksempel p  slik romlig lokalitet.

Temporal lokalitet betyr at minnereferansene gjerne brukes gjentatte ganger. L kker er et eksempel p  at temporal lokalitet finner sted.

1.2.2. Cache-st rrelse

Cache kan alts  sees p  som en arbeidskopi av det som i  yeblikket er den mest brukte delen av prim rminnet. Hvor stor m  denne arbeidskopien v re? I praksis viser det seg at selv en sv rt liten cache har god effekt.

Vi kan f rst ser p  noen gamle maskiner. P  80-tallet hadde NORD 10 maskinen fra Norsk Data en cache p  2 KB. Dette utgjorde   promille av prim rminnet. P  en Intel 486-prosessor var cachen 8 KB, eller rundt en promille av prim rminnet.

Moderne PCer har gjerne prim rminne p  4-16 GB, og en total cache-st rrelse som m les i omtrent like mange MB (alts  4-16 MB). Hva er forskjellen p  GB og MB? Jo, siden 1 GB er 1000 ganger s  stort som 1 MB, s  utgj r cachen fortsatt rundt en promille av st rrelsen til prim rminnet.

Dette viser hvor sterkt prinsippet om lokalitet er! Over relativt lange tidsrom (i forhold til prosessorens klokkefrekvens) brukes bare omkring en promille av prim rminnet. N r vi kopierer denne promillen til cache  ker ytelsen til prosessoren betydelig.

Men: Det vi trenger er en automatisk mekanisme som bestemmer hvilke deler av minne som til enhver tid skal ligge i cachen.

1.3. Cachens virkem te

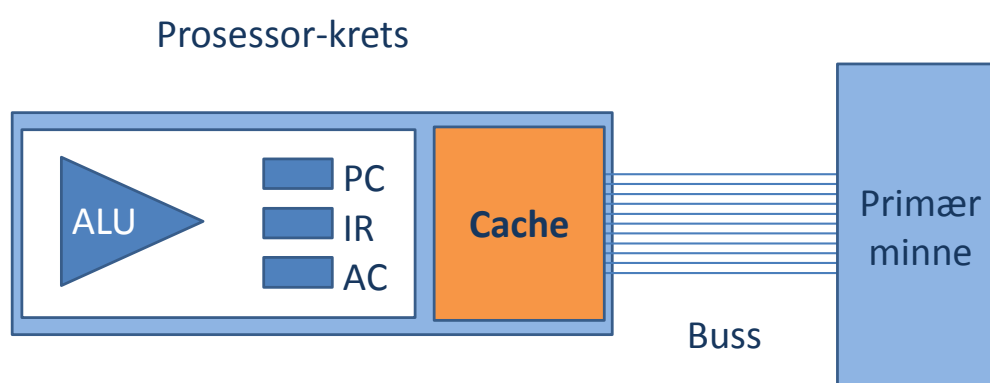
Virkem ten til cache er kort forklart at n r prosessoren skal lese en lokasjon fra prim rminnet, s  sjekker den f rst om denne lokasjonen allerede er kopiert til cachen. Hvis den er det s  vil prosessoren lese derfra - og den sparer mye tid siden cache-minnet er mye hurtigere enn prim rminnet.

Dersom lokasjonen **ikke** er kopiert til cache vil prosessoren lese fra prim rminnet. Da leses imidlertid ikke bare denne ene lokasjonen. N r prosessoren f rst m  ut i prim rminnet, s  leses det like godt inn en st rre *blokk* av data best ende av et fast antall lokasjoner - deriblant den som prosessoren vil aksessere. Denne blokken kopieres inn i cachen. Siden prinsippet om lokalitet gjelder, er det stor sannsynlighet for at prosessorens neste lagerreferansen vil ligge i den samme blokken.

Cachen inneholder alts  et antall *blokker* med fast st rrelse.

1.3.1. Cachens plassering

Hvor er cachen fysisk plassert? Det er viktig at prosessoren har hurtig tilgang til cachen. Derfor ligger cachen i umiddelbar n rhet av prosessoren. Ja, n  til dags er den en integrert del av prosessorkretsen slik Figur 1 viser.

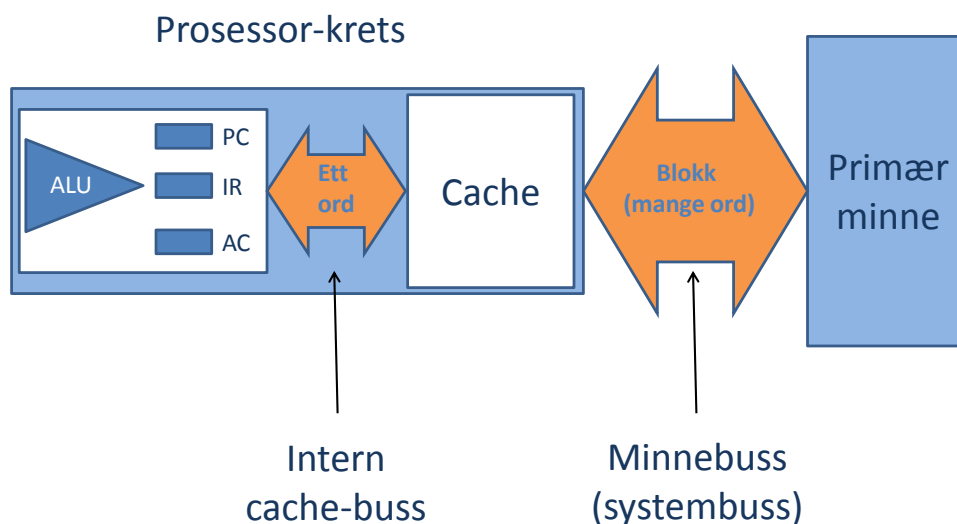


Figur 1 Cache m  plasseres n rt CPU slik at data mellom cache og CPU kan overf res hurtigst mulig. P  moderne prosessorer er cachen innebygget i prosessor-kretsen.

Som nevnt tidligere overf res det blokker mellom prim rminnet og cachen. En slik blokk inneholder mange ord. P  en PC brukes for eksempel en blokkst rrelse p  64 bytes. Siden

prosessorens registre fortsatt inneholder ett ord, s  overf res det selvf lgelig enkelt-ord mellom cache og registre. Dette er fremstilt i Figur 2.

Cachen er bygget opp av statisk RAM (SRAM). Som det ble nevnt i forrige leksjon har SRAM kortere aksessetid enn DRAM. Typiske aksessider kan v re 40-70 ns for prim rminne (DRAM) og 1-20 ns for cachens minne (SRAM).



Figur 2 Mellom CPU og cache overf res det ord. Mellom minne og cache overf res det b kker som inneholder mange ord.

1.4. Treffrate

Hensikten med cache er at prosessoren skal slippe   vente p  det langsomme prim rminnet. Hvis vi er s  heldig at alt som CPU skal lese befinner seg i cachen, s  vil CPU jobbe uten forsinkelser fordi cachen holder tritt med CPU. Hvis vi derimot aldri finner det vi trenger i cachen, s  m  prosessoren hele tiden ut i prim rminnet for   hente instruksjoner og data.

I praksis vil vi alltid ha en situasjon mellom disse ytterpunktene. En stund finner vi det som vi trenger i cachen, men s  ettersp r CPU noe som ikke ligger kopiert til cachen. Da m  vi ut i prim rminne og hente en b kk derfra.

Over litt tid vil vi derfor alltid ha noen aksesser der vi finner det vi  nsker i cachen (dette kaller vi treff), og noen aksesser der vi m  ut i prim rminnet (og det kaller vi bom).

N r prosessoren finner det den  nsker i cachen, s  er aksessiden lik cachens aksessid. Men n r vi m  ut i prim rminnet, s  m  rett b kk f rst hentes fra prim rminnet og deretter m  vi hente  nsket informasjon i denne b kka.

Den andelen av aksessene som finner data i cachen kalles *treffrate* (eller *hitrate* p  engelsk). Treffraten varierer, men ligger ofte mellom 80%-98%. N r et nytt program startes er den lik null - fordi prosessoren ikke har aksessert programmets instruksjoner tidligere - men den  ker raskt.

1.5. Effektiv aksesstid ved bruk av cache

Enkelte ganger er vi heldige og finner det vi leter etter i cachen, mens andre ganger er vi uheldige og bommer. N  skal vi finne en gjennomsnittlig aksesstid ved en gitt treffrate. Da m  vi vite aksesstiden ved treff, aksesstiden ved bom, og vi m  kjenne treffraten.

Aksesstid ved treff: N r vi f r et treff vil aksesstiden v re lik aksesstiden til cachen. Vi kaller cachens aksesstid for T_c .

Aksesstid ved bom: N r vi f r en bom vil vi f rst m tte lese data fra prim rminnet og inn i cachen. Deretter leser vi data fra cachen. Vi kaller aksesstiden til prim rminnet for T_p . Da blir aksesstiden ved bom lik $(T_p + T_c)$.

Treffraten: Vi kaller treffraten for H . Da vil H v re et tall mellom 0 og 1. Hvis treffraten er 0 har vi bare bom. Med en treffrate p  1 har vi bare treff. I alle andre tilfeller er H et sted mellom 0 og 1. La oss se p  et eksempel: Hvis vi treffer i 70% av tilfellene vil treffraten, H , v re 0,7. Legg merke til at da vil vi bomme i 30% av tilfellene. Vi bommer alts  med andelen $(1-H)$, som er 0,3 i dette tilfellet.

Den gjennomsnittlig aksesstid over mange aksesser kaller vi *effektiv aksesstid*, eller T_e . N  kan vi sette opp et matematisk uttrykk for denne. Vi antar en treffrate p  H . Det betyr at H -delen av aksessene vil v re treff. De  vrige aksesser, alts  $(1-H)$ delen, gir bom. Den delen av aksessene som gir treff vil ha en aksesstid p  T_c . Den delen som gir bom, vil ha aksesstiden $(T_p + T_c)$. Den effektive (eller gjennomsnittlige) aksesstiden blir:

$$\begin{aligned} T_e &= (\text{Andel treff}) \times (\text{aksesstid ved treff}) + (\text{Andel bom}) \times (\text{aksesstid ved bom}) \\ &= H \times T_c + (1-H) \times (T_p + T_c) \\ &= H T_c + (1-H) (T_p + T_c) \\ &= H T_c + T_p + T_c - H T_p - H T_c \\ &= T_c + T_p - H T_p \\ &= T_c + (1-H)T_p \end{aligned} \quad (1)$$

Her er

T_e - er effektiv aksesstid.

H - er treffraten.

T_c - er cachens aksesstid.

T_p - er aksesstid til prim rminnet.

Fra ligning (1) ovenfor ser vi f lgende. Hvis vi har en h y treffrate vil H nesten v re lik 1. Da vil $(1-H)$ bli nesten null, og f lgelig blir $(1-H)T_p$ ogs  nesten null. Det betyr at T_e nesten bare er avhengig av cachens aksesstid, T_c . Ved h y treffrate er alts  den effektive aksesstiden nesten lik cachens aksesstid.

Hvis treffraten er lav, alts  at H nesten er null, vil $(1-H)$ n rme seg 1. Det betyr at leddet $(1-H)T_p$ begynner   gj re seg sv rt gjeldende. Fordi T_p er mye st rre enn T_c vil dette bety at T_e nesten bare avhenger av T_p . Ved lav treffrate vil den effektive aksesstiden alts  bli tiln rmet lik aksesstiden til prim rminnet.

Eksempel p  beregning av effektiv aksesstid ved bruk av cache:

Anta at vi har et prim rminne med aksesstid 60 ns. Mellom CPU og prim rminnet har vi en

cache med aksesstid p  10 ns. Treffraten er 95%, det vil si at $H=0,95$. Finn effektiv aksesstid.

Vi bruker formelen v r:

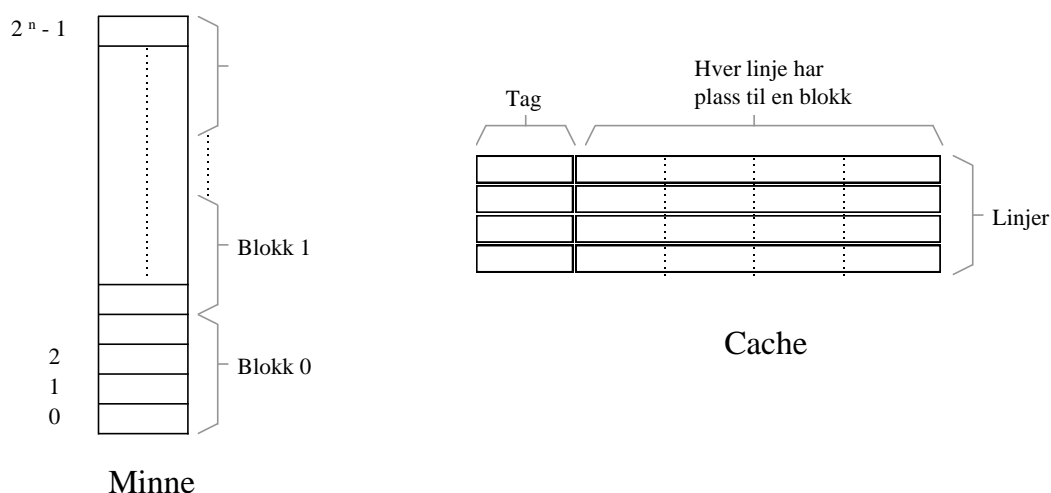
$$T_e = T_c + (1-H)T = 10\text{ns} + (1-0,95)*60\text{ns} = 10\text{ns} + 0,05*60\text{ns} = 10\text{ns} + 3\text{ns} = 13\text{ns}$$

Den aksesstiden som CPU "opplever" er alt  13 ns. Det vil si en effektiv aksesstid som er omtrent fem ganger s  hurtig som prim rminnet.

1.6. Mer om cachens virkem te

Vi skal n  se i mer detalj p  hvordan cachen virker. Det er tre ting som er viktig   huske n r man leser resten av denne leksjonen:

1. Tidligere har vi sett p  minnet som en samling av adresserbare ord, som regel tenker vi p  minnet som bygget opp av bytes. Vi har sagt at n r vi leser fra minnet, s  henter vi ett og ett ord. N r vi bruker cache henter vi imidlertid aldri enkelt-ord eller enkelt-bytes. Isteden henter vi en hel blokk. Derfor g r vi over til   se p  minnet som en samling av slike blokker. Vi nummererer blokkene fra 0 og oppover. N r CPU ettersp r ett bestemt ord, henter cache inn hele den blokken som inneholder dette ordet.
2. Cachen er mye mindre enn prim rminnet, men har likevel plass til mer enn bare en blokk. Vi sier at cachen best r av s kalte *linjer* hvor hver *linje* har plass til en blokk. Vi sier at cachen best r av s kalte *linjer* hvor hver *linje* har plass til en blokk.
3. Siden cachen er s  sv rt mye mindre enn prim rminnet s  vil mange blokker fra minnet m tte bytte p    bruke samme *linje* i cache. CPU m  notere seg *hvilken* blokk som i  yeblikket er kopiert til hver enkelt *linje* i cachen. Hver cache-linje m  derfor utstyres med en identifikator som forteller hvilken blokk av minnet som ligger der. Denne identifikasjonen lagres sammen med blokken og kalles en *tag* («merkelapp»). Se Figur 3



Figur 3 Minnet er organisert som adresserbare ord, men det kan ogs  sees p  som en rekke med blokker der hver blokk best r av flere ord. Cache best r av et antall "linjer", hvor hver linje kan lagre n yaktig  n blokk. N r vi leser fra minnet leser vi en hel blokk - isteden for et ord - og legger blokken i en linje. "Tagen" til linjen forteller hvilken blokk som ligger der. En moderne PC har en blokkst rrelse p  64 bytes.

1.7. Blokkst rrelse

En cache kan enten deles opp i et lite antall store blokker, eller i et stort antall sm  blokker. Blokkst rrelsen er alltid en avveining mellom flere forhold, og den optimale blokkst rrelsen vil variere med hva slags oppgaver prosessoren brukes til   l se.

Programmer med stor grad av romlig lokalitet vil jobbe lenge innenfor en sv rt avgrenset del av minnet. Da b r blokkst rrelsen v re stor for   holde blokka i cachen s  lenge som mulig.

N r man f rst f r en bom s  straffes man imidlertid hardere ved store blokker enn ved sm  blokker. Dette er fordi en ny stor blokk vil ta lengre tid   overf re fra prim rminnet via bussen og over til cachen. (Med en gitt bussbredde vil en stor blokk kreve flere buss-sykluser enn en liten blokk.)

Et moment som taler for   ha sm  blokker er at dersom man har lav romlig lokalitet (man hopper litt frem og tilbake i minnet), men stor temporal lokalitet (man bruker likevel disse spredte lokasjonene om og om igjen), s  vil en cache med mange sm  blokker ha en st rre sannsynlighet for   inneholde de rette blokkene.

Kompromisset som moderne cache-design har landet p  er gjerne en blokkst rrelse p  mellom 32 og 64 bytes.

Eksempel p  blokkst rrelse

P  en eldre PC brukes en blokkst rrelse p  32 bytes. Da sees alts  minnet p  som en samling blokker hvor hver blokk er 32 bytes. Blokk **nr 0** består av adressene 0..31, **blokk nr 1** består av adressene 32..63, **blokk nr 2** av adressene 64..95 og s  videre.

Dersom CPU ettersp r adresse nr 35 i prim rminnet, henter alts  cachen inn blokk nr 1 siden adresse 35 befinner seg i denne blokken. Cachen vil plukke ut en av linjene i cachen, og legge blokk 1 i denne linjen.

Nyere PCer bruker en blokkst rrelse p  64 bytes.

1.8. Mapping-funksjon

N r en blokk kopieres til cachen, m  det bestemmes i hvilken linje i cachen blokka skal ligge. Dette bestemmes med en mappingfunksjon.

Valg av mapping-funksjon er en avveining mellom enkelhet p  den ene siden og et  nske om   ikke un dig kaste blokker ut av cachen p  den andre siden. Vi skal se p  tre ulike mapping-funksjoner.

1.8.1. Full-assosiativ cache

Her kan en blokk legges i en hvilken som helst linje. N r CPU skal hente data, m  den sjekke hver eneste tag i cachen. Dette krever komplisert elektronikk.

Det er viktig   huske at hvis det skal v re noen gevinst i   bruke cache s  har maskinen bare noen f  nanosekund til   finne ut om data ligger i cachen. Selv med den hurtigste elektronikk som finnes i dag er det alt for tidkrevende   sjekke linjene etter tur. Samtlige tag-verdier m  derfor sammenlignes i parallell mot tagen til de data CPU vil hente. Dette krever like mange sammenligningskretser som det er linjer i cachen - noe som betyr at full-assosiative cacher er

komplisert   implementere. Dette er grunnene til at full-assosiative cacher ikke brukes i praksis.

1.8.2. Ikke-assosiativ cache (direkte mapping)

For   slippe   ha sammenligningslogikk for hver eneste linje kan vi isteden la hver blokk i minnet ha en *fast* plass i cache. Dette kalles *ikke-assosiativ cache*, eller *direkte-mappet cache*. Da kan en blokk bare legges p  en bestemt linje i cachen, men det er mange blokker som "tilh rer" hver linje.

N r CPU skal sjekke om en bestemt blokk ligger i cachen, trenger den bare   sjekke en eneste linje. CPU vet p  forh nd hvilken linje den skal sjekke – den vet ikke om den aktuelle blokken faktisk ligger i linjen, men den vet hvilken linje den skal sjekke.

Eksempel p  ikke-assosiativ cache (direkte mapping)

Dersom vi har 4 linjer i cachen, vil
blokk nr 0, 4, 8, 12, konkurrere om   ligge i linje nr 0
blokk nr 1, 5, 9, 13, konkurrere om   ligge i linje nr 1
blokk nr 2, 6, 10, 14, konkurrere om   ligge i linje nr 2
blokk nr 3, 7, 11, 15, konkurrere om   ligge i linje nr 3

Dersom CPU ettersp r ett ord i blokk nr 14, trenger ikke cachen   sjekke alle linjene. Den vet at dersom blokk nr 14 befinner seg i cachen, s  ligger den i linje 2. Derfor sjekker den taggen kun til denne ene linjen.

I praksis har en cache mye mer enn fire linjer, men eksemplet illustrerer virkem ten.

Trashing

Ulempen er at to ulike ord som brukes samtidig kan komme til   havne opp i samme linje. For eksempel kan det hende at rett etter at en instruksjon er hentet inn m  den vike plassen igjen for en data-verdi som havner i samme linje. Hvis dette skjer ofte kan verdien av caching reduseres betydelig.

Denne situasjonen hvor to blokker gjensidig tvinger hverandre ut av cache hele tiden kalles *trashing*.

1.8.3. Sett-assosiativ cache

For   unng  at blokker gjensidig tvinger hverandre ut av cache hele tiden, kan vi isteden organisere linjene i *sett*: Blokkene som sloss om samme linje f r flere linjer   dele p , s  det blir plass til flere av dem. Vi sl r sammen for eksempel 2 eller 4 linjer til *et sett*, der flere «konkurrenter» kan ligge side om side.

En cache med 2 linjer pr sett kalles «2-veis», og kan ha to konkurrerende blokker inne samtidig; har vi 4 linjer pr sett kan fire konkurrenter v re inne samtidig, og vi kaller det 4-veis cache.

I praksis er det nesten alltid sett-assosiative cacher som brukes. Sett-assosiative cacher gir en rimelig avveining mellom enkelhet og fare for trashing, og vil som regel gi den h yeste treffraten. I de aller fleste tilfeller brukes 2-veis, 4-veis, 8-veis eller 16-veis cacher.

1.9. Utskiftingsalgoritmer

For fullt assosiativ og sett-assosiativ cache m  vi velge mellom flere mulige linjer der data fra minnet kan kopieres til cache. (For direkte-mappet cache har man intet valg.) Dette valget m  gj res i hardware, p  noen f  nanosekunder, om det ikke skal forsinke CPU. Metodene m  derfor v re enkle.

- **LRU** : Least Recently Used. Bytter ut den blokka det er lengst siden ble aksessert. Hver cache-linje har en teller som  kes for hver minne-aksess der denne blokka ikke ble etterspurt, og nullstilles hvis blokka ble referert. (For sett-assosiativ cache endres telleren bare for referanser til samme sett.) Linjen med h yest telle-verdi er den som velges som offer f rst.
- **FIFO** : First In First Out. Bytter ut den blokka som har v rt lengst i cachen. Som LRU, men telleren  kes bare n r en ny blokk hentes fra minnet; den nye blokkas teller settes til 0. (For sett-assosiativ cache bare n r det hentes inn til samme sett.)
- **LFU** : Least Frequently Used. Bytter ut den blokka som har v rt aksessert f rrest antall ganger. Telleren  kes for referanser til dette ordet, ellers ikke. Den linjen som har lavest telleverdi velges som offer f rst. Med jamne mellomrom telles alle tellerne ned med samme verdi for   «glemme» referanser som skjedde for en stund siden.
- **Random** : Tilfeldig utskifting. En tilfeldig av blokkene velges som offer.

Fordi valget m  gj res p  noen f  nanosekunder, brukes ofte vesentlig enklere utgaver av disse strategiene enn hva vi kan akseptere i andre sammenhenger. En form for tilfeldig utskifting er s v rt vanlig.

1.10. Skrivning til minnet

S  langt har vi konsentrert oss om lesing fra minnet. Men av og til skal jo CPU skrive til minnet ogs . Hvordan skal dette gj res n r vi har cache?

Problem: CPU vil endre verdien i en lokasjon i minnet. Hva skal skje hvis en kopi av lokasjonene ligger i cachen?

Her brukes to forskjellige l sninger:

Write through:

Skriving skjer samtidig b de p  cachen og i prim rminnet. Fordelen med denne l sningen er at prim rminnet alltid er oppdatert. Ulempen er at bruken av cache ikke reduserer belastningen p  bussen ved skrivning til minnet; dette kan sinke b de CPU (fordi den m  vente p  at skrivningen fullf res f r neste skrivning godtas) og andre enheter som arbeider i parallell p  samme buss.

Write back:

Da skrives det kun i cachen, prim rminnet oppdateres f rst n r innholdet av en cache-linje m  vike plassen. Fordelen er at vi f r redusert buss-trafikk, og dermed maksimal hastighet. Ulempe er at utskifting av innholdet i en linje tar lengre tid, og CPUen m  st    vente i mellomtiden. Dessuten f r vi problem med I/O-utstyr (se neste punkt)

1.11. Cache-koherens

I et datasystem har man ofte flere enheter som kan endre innholdet i minnet. For eksempel kan IO-utstyr legge inn nytt innhold i deler av minnet med DMA (Direkte Minne-Aksess). Et

annet eksempel er fler-prosessor-systemer der flere prosessorer deler p  ett og samme prim rminne.

Siden en cache inneholder en kopi av deler av prim rminnet, kan vi fort risikere at lokasjoner som er kopiert til en cache blir endret. Hva skal vi da gj re. Da samsvarer ikke innholdet i cachen med prim rminnet. Vi sier at cache og prim rminnet ikke er konsistent.

Slike praktiske problem er  rsaken til at cache-mekanismen p  enkelte CPUer er omtrent like komplisert som resten av CPUen. Det at det er samsvar mellom prim rminnet og cache kalles cache-coherence.

Problem med at I/O-kontrollere kan endre minnet (DMA)

Problem 1: Hva skjer hvis I/O-utstyr endrer prim rminnet, og en kopi av blokken ligger i cache?

Alternativ 1: Hvis ikke programmet skal lese foreldede verdier m  cache-innholdet markeres som ugyldig. Som regel er det ikke mulig   gj re dette for hver enkelt linje. Isteden gj res det for hele cache-lageret under ett.

Alternativ 2: Hvis datamaskinen er et flerprosess-system kan operativsystemet overlate CPU (og cache) til et annet program enn det som ba om I/O-operasjon. Da kan I/O-utstyret endre prim rminnet f r programmet som startet I/O-operasjonen slipper til igjen. P  denne m ten unng s det at problemet oppst r.

Alternativ 3: Alle I/O-operasjoner kan i prinsippet g  via cachen. Dette gj r system-design og hardware-modularisering s  problematisk at det sjelden benyttes.

Problem 2: Anta at maskinen har write-back cache. Hva skjer n r et program endrer data i cache, og s  ber om at I/O-utstyr (som leser direkte fra prim rminne) skriver disse dataene ut til disk. Da er de nye data enda ikke skrevet tilbake til prim rminnet, de nye verdiene ligger bare i cache?

Alternativ 1: De nye verdiene m  skrives ut til prim rminne f r operativsystemet tillater at I/O-operasjonen settes igang. Alle write-back-cacher har en kontroll-inngang som «t mmer» endrede verdier til prim rminne («cache flush»).

Alternativ 2: I/O-operasjoner kan g  gjennom cachen - se ovenfor.

Problem 3: Flere prosessorer (i et multi-CPU-system) har hver sin kopi av de samme data, i hver sin cache. Hva skjer n r en av prosessorene endrer innholdet?

P  fler-prosessorssystemer finnes det et eget system for   sikre samsvar mellom cachene, eller *cache-coherence* som det kalles. Det g r egne kontroll-linjer der cachene rapporterer til hverandre hvilke data som blir endret hos dem. Andre cache-kontrollere med kopier av samme data vil da merke sine kopier som ugyldig, og neste gang de refereres hentes det ny, oppdatert kopi fra prim rminne.

Senere skal vi l re om fler-kjerne-prosessorer. Dette er prosessorkretser som inneholder flere prosessorer som jobber i parallell. En av fordelene med   samle kjernene i en krets i stedet for   ha flere separate prosessor-kretser er at mekanismen for   sikre cache-koherense er samlet i en prosessor-krets.

1.12. Hastighetsforbedring ved bruk av cache

Vi har n  sett at bruk av cache kompliserer CPUens arbeid i stor grad. Det kan v re naturlig   sp rre om det er bryet verdt.

Gevinsten med cache varierer b de med datamaskinens  vrige oppbygging, hva slags type program som kj res og hva slags oppgave datamaskinen brukes til. Hvis det brukes langsom RAM til prim rminnet vil cache gi st rre gevinst enn n r det brukes hurtig RAM. Mye venting p  I/O - som skjermfunksjoner o.l - gir mindre forbedring. Helt tilfeldig data-aksess gir mindre forbedring. Og til sist: sv rt tette l kker og/eller sm  datasett gir st rre forbedring.

Fors k og forskning viser at optimal cachest rrelse varierer fra programmeringsoppgave til programmeringsoppgave, men selv sv rt sm  cacher (noen f  ti-talls KB) gir god effekt.

2. Moderne cache-design

Alle moderne CPUer har cache. For eksempel ble cache introdusert p  Intel sine prosessorer i 1984. Det var p  486-prosessoren. Etter hvert som tiden har g tt har cache-mekanismen blitt mer og mer avansert. Den har ogs  blitt viktigere og viktigere for prosessorens ytelse.

Typisk for de nye prosessorer er at de har s kalt *fler-niv  cache*. Andre begreper vi skal se p  er *splittet cache* kontra *enhetlig cache*. I tillegg skal vi innom *delt cache* (eller *felles cache*). Alle disse begrepene er sentrale i moderne cache, og vi skal bruke litt tid p  dem.

2.1.1. Fler-niv  caching

Forskjellen p  prosessorens ytelse og minnets aksessetid blir bare st rre og st rre, og behovet for cache bare  ker. Derfor kan man tenkte seg    ke st rrelsen p  cachen.

Cache-st rrelse er en avveining mellom flere forhold. Vi  nsker at cache-st rrelsen skal v re s  stor at vi oppn r en h y treffrate. Da vil jo den effektive aksesstiden i all hovedsak bli bestemt av cache-minnets aksessetid. Dessverre er det samtidig slik at dess st rre en cache er, dess mer elektronikk trengs for   kontrollere den. Elektronikkens hastighet er avhengig av st rrelsen slik at en stor cache vil ha en tendens til   v re langsommere enn en liten, selv n r de er bygget opp med samme teknologi.

En godt alternativ kan da v re   bruke to cacher: en liten og sv rt hurtig cache for de aller mest brukte lokasjonene, og en st rre cache for lokasjoner som er litt mindre brukt. Den store cachen er langsommere enn den lille, men fortsatt mye hurtigere enn prim rminnet. P  denne m ten blir ytelsen bedre enn om man brukte en eneste stor cache.

Dette kalles to-niv  caching, og har i mange  r v rt vanlig p  generelle prosessorer. Den lille cachen kalles *niv  1 cache* (p  engelsk: *Level 1 cache* eller bare *L1-cache*), mens den store cachen kalles *niv  2 cache* (*Level 2 cache* eller *L2-cache*).

Etter som prosessor-ytelsen  ker mye hurtigere enn ytelsen til DRAM blir cache bare viktigere, og de siste  rene har mange prosessorer kommet med tre cache-niv er. *Niv  3 cachen* (L3-cachen) er enda st rre enn L2-cachen. Den er ogs  litt tregere. Men den er mye hurtigere enn prim rminnet, og avlaster prim rminnet i stor grad.

Med flerniv  caching blir ligningen for   beregne den effektive aksesstiden mer komplisert fordi vi opererer med flere treffrater. Jeg har valgt   ikke diskutere formler for effektiv aksess-tid for flerniv  caching i dette kurset.

2.1.2. Splittet cache kontra enhetlig cache

En tradisjonell cache brukes b de til instruksjoner og data. Dette kalles en *enhetlig cache* (*unified cache* p  engelsk), og har den fordelen at den automatisk avpasser antall instruksjoner kontra antall data i cachen.

Det vil si at dersom et program oftere henter instruksjoner enn data fra minnet, s  tilpasser cachen seg dette ved at den automatisk inneholder flere instruksjoner. Det motsatte gjelder selvf lgelig ogs , hvis programmet henter mer data enn instruksjoner. Vi har alts  en dynamisk fordeling av data og instruksjoner. I tillegg har en enhetlig cache den fordelen at det er enklere   designe og implementere  n cache enn to.

Motstykket til en slik enhetlig cache er det vi kaller en *splittet cache*. Splittet cache vil si at cachen best r av to selvstendige del-cacher; en for instruksjoner og en for data. Dette kan v re en stor fordel i enkelte tilfeller, nemlig dersom vi skulle f  bruk for b de en instruksjon og data samtidig. Men n r skulle et slikt behov kunne oppst ? Jo, det skal vi se mer i detalj p  senere i kurset, men saken er at moderne prosessorer jobber med flere instruksjoner samtidig. Dette gjelder spesielt s kalte *superskalare prosessorer*; det vil si prosessorer som inneholder mer enn en prosesserende enhet, og som utf rer flere instruksjoner samtidig. P  slike prosessorer hender det ofte at  n prosesserende enhet vil lese en ny instruksjon akkurat samtidig med at en annen prosesserende enhet vil lese eller skrive data. Dette er fullt mulig med en delt cache siden de to del-cachene har hver sine forbindelseslinjer til prosessoren, og dermed kan aksesseres samtidig.

2.1.3. Egen kontra felles (delt) cache

Moderne prosessorer har flere kjerner. Hva dette betyr skal vi se p  senere i kurset. Men i kortversjon betyr det at flere prosessorer er bygget sammen i en og samme krets. Med en fire-kjernes prosessor vil operativsystemet oppfatte det som en maskin med fire prosessorer. Selv om det bare er en prosessorkrets.

P  en slik fler-kjerne prosessor kan enkelte cacher v re felles for alle kjerner, mens andre cacher bare brukes av en kjerne.

P  nyere prosessorer ser vi ofte at hver kjerne har sin egen L1- og L2-cache. Og s  er det en egen L3-cache som brukes av alle kjernene; en s kalt *felles cache* eller *delt cache*.

2.1.4. Sett-st rrelse

Den som designer et cache-system m  velge hvor mange linjer det skal v re i hvert sett. Dette kalles sett-st rrelse, og som det fremg r av kapittel 1.8.2 er den ofte 2, 4, 8 eller 16. N  vet vi at det finnes flere niv  med cacher p  en prosessor, og vi kan diskutere sett-st rrelsen litt mer inng ende.

En h y verdi p  sett-st rrelsen vil selvsagt gj re cachen lite s rbar for trashing. Sannsynligheten for   trashing vil jo v re lavere dess h yere sett-st rrelsen er.

Problemet med stor sett-st rrelse er imidlertid at cachen blir mer komplisert   lage, og det er vanskelig   f  cachen s rlig hurtig fordi den krever mer kontroll-elektronikk. (De som leser Vedlegg A til denne leksjonen vil f.eks se at det kreves flere bits i tag-en til cachen n r sett-st rrelsen  ker.) Alt dette gj r at straffen for   bomme blir st rre ved store sett-st rrelser.

Avveiningen blir derfor mellom straffen for   bomme kontra gevinsten ved   treffe.

Siden straffen for   bomme er lavere for cachene som befinner seg n r prosessoren vil L1-cachen ha en lav sett-st rrelse, f.eks. to eller fire, mens sett-st rrelsen  ker for cacher lenger bort. L3-cachen er ofte en 16-veis sett-assosiativ cache

2.1.5. Utviklingstrender innen cache

Siden gapet mellom ytelsen til CPU og ytelsen til prim rminnet bare fortsetter    ke, vil cache bli viktigere og viktigere. De trendene vi har sett de senere  r er at cache-mengden  ker og at det tas i bruk flere niv  av cache.

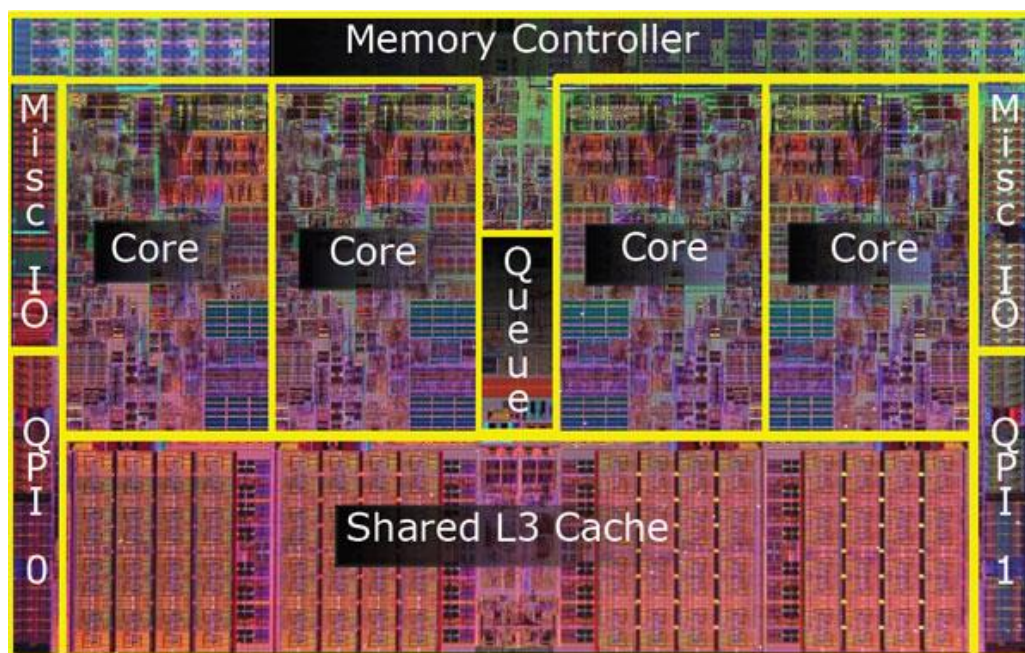
St rre og st rre totalmengde cache. Fra en total mengde p  noen KB cache p  sent 80-tall, via noen hundretalls KB p  90-tallet, ser vi n  at totalmengde cache m les i 10-talls MB p  moderne prosessorer.

Antall cache-niv   ker. P  de fleste plattformer er L1-cachen bare litt st rre enn den var p  90-tallet, 32 KB til instruksjons-cache og 32 KB til data-cache er typisk. L2-cachen er oftest fra en halv MB til noen f  MB. Mens man ved st rre cache-behov heller baserer seg p  et 3. cache-niv : en L3-cache som kan v re temmelig stor: fra 2 MB og oppover til flere 10-talls MB.

Totalmengde cache er sv rt avhengig av ytelsen til CPU. Lavpris-prosessorer har mye mindre cache enn dyre prosessorer med h y ytelse.

Eksempel p  moderne cache, Intel i7 Sandy Bridge

Sandy Bridge er en nyere prosessorarkitektur fra Intel. Denne prosessoren har fire kjerner mer hver sin L1- og L2-cacher. I tillegg har den en stor felles L3-cache. Innmaten er vist i Figur 4.



Figur 4 Innmaten i Intel i7 (sandy bridge). Bildet viser en Sandy Bridge prosessor. Det er fire kjerner (og hver kjerne har sin egen L1 og L2 cache), en delt L3 cache (8 MB) og innebygget minne-kontroller. Legg merke til hvor stor del av prosessoren som brukes til cache. Bildet er hentet fra <http://www.sharkyextreme.com/hardware/cpu/article.php/3782516/Intel-Core-i7965-XE--Core-i7920-Review.htm>

Et eksempel på cache-konfigurasjon på en i7-prosessor kan være:

Blokk-størrelse: 64 Byte

L1-cache. Hver kjerne har:

Instruksjons-cache: 32 KB, 8-veis sett-assosiativ cache.

Data-cache: 32 KB, 8-veis sett-assosiativ.

L2-cache. Hver kjerne har:

Enhetlig cache: 256 KB, 8-veis sett-assosiativ cache.

L3-cache. Kjernene deler på en:

Delt enhetlig cache: 8 MB, 16-veis sett-assosiativ cache.

2.2. Sammenkobling av delene

I de to siste leksjonene har vi sett på primærminne og på cache. Disse samarbeider om å fore CPU med instruksjoner og data. I neste leksjon skal vi se nærmere på hvordan dette foregår. Denne delen av datateknikk kalles gjerne Systemarkitektur.

Vedlegg A. Mer om mapping-funksjoner.

I leksjonen har vi sett at n r en blokk kopieres til cachen, s  brukes en mappingfunksjon til   bestemme hvor i cachen blokka skal ligge.

N  skal vi se n rmere p  de tre viktigste mappingfunksjonene, nemlig full-assosiativ cache, ikke-assosiativ cache (direkte mapping) og sett-assosiativ cache.

Det vi skal se p  er hvordan CPU greier holde oversikt over hvilke blokker som til en hver tid ligger i de ulike cache-linjer. Vi m  huske at informasjonen om hvilken blokk som ligger i en linje finnes i tagen til linjen.

Full-assosiativ cache

Her kan en blokk legges i en hvilken som helst linje. Ulempen med metoden er at n r CPU skal hente data, s  vet den ikke p  forh nd hvilken linje den kan forvente   finne blokka.. Derfor m  cache sjekke tag-en til hver eneste linje. I praksis tar dette for langt tid. Vi tar en titt p  dette tilfellet likevel.

Tag ved full-assosiativ cache: et enkelt eksempel

Vi skal se hvordan adressene til den blokken som ligger i cache kan brukes til   bygge opp tagen. Vi begynner med et enkelt eksempel for   se prinsippet - og deretter ser vi p  et mer realistisk eksempel.

I Figur 5 ser vi et minne som best r av 16 lokasjoner. Vi antar at cache bruker en blokk-st rrelse p  2 lokasjoner. Da kan minnet ogs  sees p  som 8 blokker. Cachen har to linjer som hver har plass til en blokk. Vi ser at blokk 0 er kopiert til linje 0 i cachen.

Hvordan skal tagen se ut? Jo, se p  Figur 6 og legg merke til at de tre mest signifikante¹ bitene i adressene utgj r blokknummeret. Derfor kan vi bruke disse til tagen. I dette tilfellet benyttes ogs  de tre mest signifikante bitene til tag. Den minst signifikante biten angir bare bytenummer innenfor blokken.

Tag ved full-assosiativ cache: et realistisk eksempel

La oss n  se p  et mer realistisk eksempel. Vi antar en maskin med f lgende data:

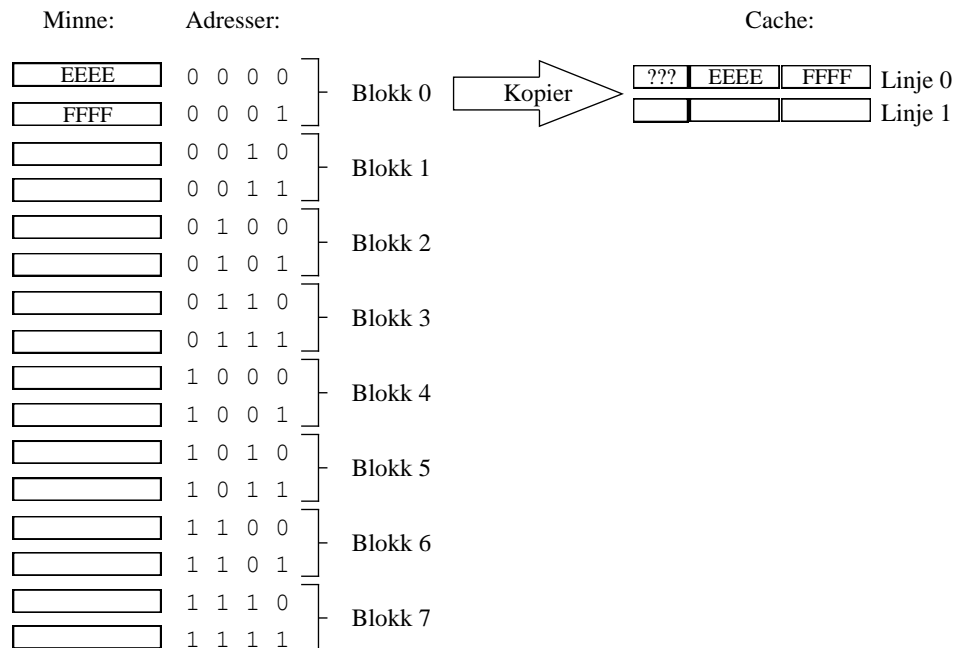
Prim rminne : 4 MB (krever 22 bits adresser)
Cache st rrelse : 4 KB
Blokkst rrelse : 4 B

Da vil antall linjer i cachen v re: $C = 4 \text{ KB} / 4 \text{ B} = 1 \text{ K} = 1024$

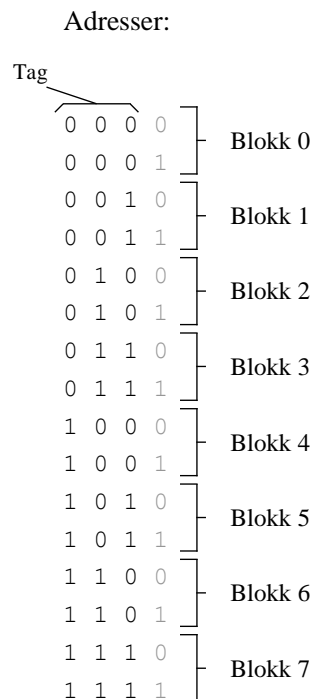
Isteden for   se p  minnet som 4M lokasjoner   en byte, kan vi se p  det som 1 M blokker som hver er 4 bytes. I en gitt linje kan det ligge et vilk rlig av disse 1M blokkene. For   identifisere hvilken blokk som er inne i en linje m  vi bruke de 20 mest signifikante bitene i adressen som tag. De to laveste bitene, ...xx00, ...xx01, ...xx10 og ...xx11, skiller bare mellom

¹ Den mest signifikante biten (MSB - Most Significant Bit) er biten helt til venstre i et bitm nster. Den minst signifikante biten (LSB - Least Significant Bit) er biten helt til h yre. Mest signifikant betyr at en endring av bit-verdien (fra 0 til 1 eller motsatt) vil gi st rst for ndring i bitm nsterets verdi. Minst signifikant betyr at en endring av bit-verdien vil gi minst endring i bitm nsterets verdi.

blokkens 4 bytes, som alle hentes inn samlet. De er derfor ikke n dvendige for   skille mellom ulike blokkene i cach .



Figur 5 En forenklet fremstilling av minne og cache. P  figuren er blokk 0 kopiert til linje 0 i cach , og sp rsm let er: hvordan skal Tag se ut for at vi skal kunne identifisere hvilken blokk som ligger der? Tagen er merket med tre sp rsm lstegn.



Figur 6 De tre mest signifikante bitene utgj r blokknummeret, og kan brukes til tag. Den minst signifikante biten angir bare bytenummer innenfor blokken.

Ikke-assosiativ cache (direkte mapping)

For   slippe   ha sammenligningslogikk for hver eneste linje (1024 sammenligningskretser i eksempelet ovenfor) kan vi isteden la hver blokk i minnet ha en *fast* plass i cache. Dette kalles *ikke-assosiativ cache*, eller *direkte-mappet cache*. Da kan en blokk bare legges p  en bestemt linje i cachen, men det er mange blokker som m  dele p  hver linje. N  blir det mye enklere for cachen. Den vet p  forh nd hvilken linje den skal sjekke for   finne ut om blokkene som ettersp rres ligger i cache.

Tag ved ikke-assosiativ cache: et enkelt eksempel

Vi ser igjen p  det forenklete tilfellet med bare 16 minnelokasjoner, og to cache-linjer. Vi bestemmer oss for at blokkene 0, 2, 4 og 6 (like nummer) skal "tilh re" linje 0 i cachen, mens blokkene 1, 3, 5 og 7 (odde nummer) skal tilh re linje 1. Som vi ser av Figur 7 vil bit nummer 1 n  angi hvilken linje blokken skal ligge i, mens bit nummer 0 fortsatt angir bytenummer innenfor blokken. De  vrige bitene brukes til tag. N r CPU skal sjekke om en blokk ligger i cachen, vet den om den skal sjekke en like eller odde blokk - og kan g  direkte til riktig linje. Derfor er det tilstrekkelig med bare to bits til tag i dette tilfellet.

Adresser:

Tag	Linjenr	
0 0	0 0	Blokk 0
0 0	0 1	
0 0	1 0	Blokk 1
0 0	1 1	
0 1	0 0	Blokk 2
0 1	0 1	
0 1	1 0	Blokk 3
0 1	1 1	
1 0	0 0	Blokk 4
1 0	0 1	
1 0	1 0	Blokk 5
1 0	1 1	
1 1	0 0	Blokk 6
1 1	0 1	
1 1	1 0	Blokk 7
1 1	1 1	

Figur 7 Tag for ikke-assosiativ cache. Vi ser at de to mest signifikante bitene brukes til tag. Bit nummer 1 er linjenummer der blokken skal legges, og bit nummer 0 (den minst signifikante) angir fortsatt bare bytenummer innenfor blokken.

Eksempel p  tag:

Hva blir tag hvis blokk 6 ligger i cachen? Jo, vi ser fra Figur 7 at da blir tag lik 11₂, som er de to mest signifikante bitene i adressene til blokk 6.

Enn hvis blokk 3 ligger i cachen? Da ser vi p  samme m te at tag blir 01₂.

Men hvordan sjekkes det om en blokk ligger i cache?

N r CPU ettersp r en minnelokasjon m  elektronikken alts  f rst sjekke om denne minnelokasjonen allerede er kopiert til cachen. La oss se p  et eksempel med v r enkle maskin med 16 ord i minnet og to-linjers cache:

Anta at blokk 6 og 3 ligger i cachen, slik som i rammen rett ovenfor. Da er tag lik henholdsvis 11₂ og 01₂. Blokk 6 ligger i linje 0 siden 6 er et partall. Blokk 3 ligger i linje 1 siden 3 er et odde tall.

Eksempel p  treff:

Anta at CPU vil lese lokasjon 7 fra minnet. Cachen vet at dersom lokasjon 7 (med adresse 0111₂) finnes i cachen, s  ligger den i linje 1 (siden bit nr 1 i adressen er lik 1). Derfor vil elektronikken sjekke linje 1. Tagen til linje 1 er lik de to mest signifikante bitene i adressen som CPU ettersp r. Vi har et treff, og cachen returnerer den rette lokasjonen fra cache-linjen.

Eksempel p  bom:

Anta at CPU ettersp r lokasjon 14 (=1110₂) fra minnet. Cachen vet at dersom lokasjon 14 finnes i cachen, s  ligger den i linje 1 (siden bit nr 1 i adressen er lik 1). Derfor vil elektronikken sjekke linje 1, den finner at tagen til linje 1 ikke er lik de to mest signifikante bitene i adressen som CPU ettersp r. Vi har en bom, og elektronikken m  helt ut i minnet for   hente en ny blokk inn til linje 1.

Tag ved ikke-assosiativ cache: et realistisk eksempel

Vi fortsetter med det mer realistiske eksempelet. Vi hadde et minne p  4 MB. Det vil si at hver adresse er 22 bits. Vi antar 4 bytes st rrelse p  cache-blokkene. De to laveste adresse-bitene identifiserer fortsatt hver enkelt byte i blokka. N  lar vi de *neste* bitene bestemme hvor en blokk skal legges n r den kopieres inn i cachen. Som vi husker hadde vi 1024 linjer ($C = 1024$):

Blokk 0 (byte 0...0000₁₆ til 0...0003₁₆ i minnet) har plass i linje 0

Blokk 1 (byte 0...0004₁₆ til 0...0008₁₆) har plass i linje 1

:

:

:

Blokk 1023 (byte 0...0FFC₁₆ til 0...0FFF₁₆) har plass i linje 1023.

Blokk 1024 (byte 0...1000₁₆ til 0...1003₁₆) har plass i linje 0, og deler alts  plass med blokk 0

Blokk 1025 deler linje 1 med blokk 1 osv.

Siden linje 0 bare brukes av blokk 0, 1024, 2048, ..., og linje 1 bare brukes av blokk 1, 1025, 2049... osv. er det ikke n dvendig   ta med adressebits 2 til 11 i tag-feltet. Tag-feltet er bare p  de 10 mest signifikante bits, det vil si bit 12 til 21.

Eksempel p  bruk av ikke-assosiativ cache

La oss se p  et eksempel hvor CPU vil lese byten p  den 22-bits adressen (bin rt): 0000101011100000110101. Da deles adressen opp slik:

0000101011	1000001101	01
⤴ tag-verdi	⤴ Linjenr	⤴ byte-nummer innen blokk (brukes ikke i tag)

Ordet h rer alts  hjemme i linje 525 (1000001101 bin rt). N r akkurat denne blokken er inne i cache, vil tag-feltet p  denne linjen ha verdien 43 (0000101011 bin rt).

Fordelen med en direkte-mapping-cache er foruten plassbesparelsen ved kortere tag-felt at det er tilstrekkelig med en eneste sammenlignings-krets for   finne ut om det riktige ordet ligger i cache. Det er heller ikke n dvendig   velge noen utskiftingsalgoritme, siden det er absolutt gitt hvilken linje som skal brukes.

Sett-assosiativ cache

For   unng  trashing har vi sett at det er lurt   bruke en sett-assosiativ cache; blokkene som tilh rer samme linje f r flere linjer   dele p , s  det blir plass til begge.

Som i en direkte-mappet cache brukes en del av minne-adressen for   bestemme hvor i cache data skal lagres, men n  bestemmer adressen bare hvilket *sett* data skal legges i og ikke hvilken linje i settet. Som alltid gj r vi det slik at de laveste adresse-bitene bare bestemmer byte-nummer innen blokken og kan ignoreres, men n  lar vi de neste bitene velge sett: Hvis vi har 256 sett vil blokk 0 caches til sett 0, blokk 1 til sett 1 ... blokk 255 til sett 255, blokk 256 til sett 0, blokk 257 til sett 1...

Vi har f rre sett enn vi hadde linjer i en direkte-mappet cache. Derfor m  vi bruke flere bits i tag-feltet - en bit mer ved en 2-veis cache, to bits mer ved 4-veis.

Realistisk eksempel p  bruk av sett-assosiativ cache

Vi antar som f r et minne med 1M blokker, og en cache med 1024 linjer. Men n  er cachen organisert som 256 sett   4 linjer. Vi bruker de 8 bitene fra 2 til 9 i adressen til   velge cache-sett. La oss se p  et eksempel hvor CPU vil lese byten p  22 bits adressen (bin rt): 0000101011100000110101. Da deles adressen opp slik:

000010101110	00001101	01
⤴ tag-verdi	⤴ sett-nummer	⤴ byte-nummer innen blokk (brukes ikke i tag)

Ordet h rer alts  hjemme i sett 13 (1101 bin rt), dvs. linje 52, 53, 54 og 55. Hvis ordet er inne i cache vil tag-feltet p  en av disse fire ha verdien 175 (000010101110 bin rt).

N r CPU adresserer data unders ker den tag-feltet til alle linjer i settet for   se hvor ordet ligger (eller om det m  hentes inn). Da er det nok med s  mange sammenligningskretser som antall blokker i et sett. N r en blokk hentes inn m  det velges en av de 2, 4,... linjene i settet, og det er n dvendig   implementere en utskiftings-algoritme som velger et «offer».

En sett-assosiativ cache gir en rimelig god avveining mellom kompleksitet og ytelse. Kravene til elektronikk er moderate, og faren for kollisjoner relativt liten. De aller fleste nyere cache-design er 2-veis, 4-veis eller iblant 8-veis sett-assosiative.

Som en kontroll p  at du forst r dette kan du jo pr ve   l se f lgende oppgave:

Ta utgangspunkt i f lgende adresse:

0100110010010010011011

Anta en sett-assosiativ cache slik den ble beskrevet i forrige ramme. Vis at:

1. Denne adressen tilh rer sett nr $26_{16} = 38_{10}$.
2. At sett 38_{10} består av cache-linjene 152, 153, 154 og 155.
3. At tag til en av disse linjene m  v re $4C9_{16}$ dersom adressen ovenfor ligger i cache.