



Organização de Computadores Digitais

Capítulo 4 - Memória Cache



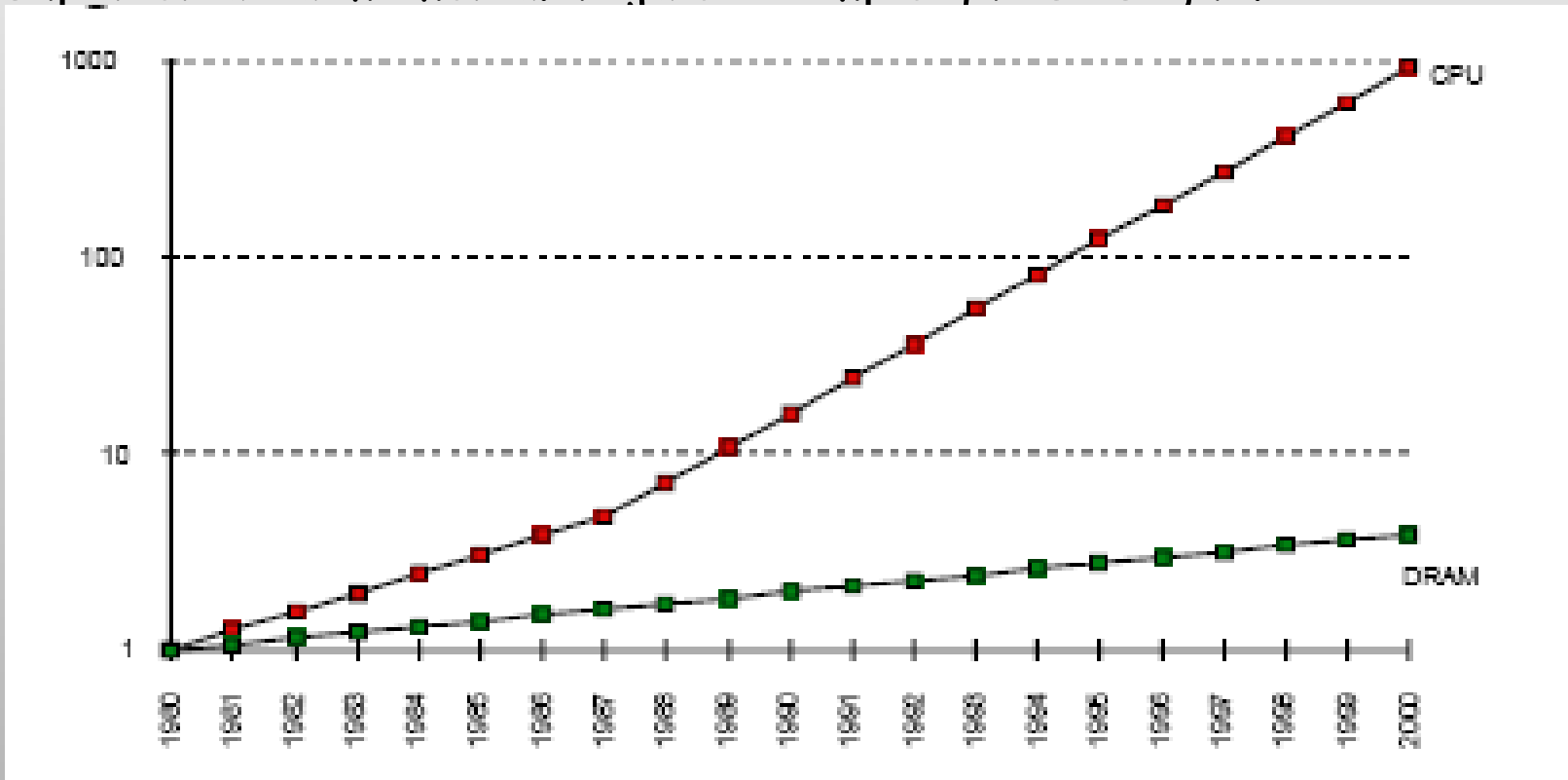
Porque a memória é importante?

- Os dispositivos de memória servem para armazenar programas (instruções) e dados (sequências de palavras em memória que não incluem qualquer código de operação)

0 0 0

Porque a memória é importante?

- O desempenho dos processadores aumentou a uma taxa muito superior a da memória. Veja o exemplo que se segue.





Características

- o Localização
- o Capacidade
- o Unidade de transferência
- o Método de acesso
- o Desempenho
- o Características físicas
- o Tipo físico
- o Organização



Localização

- A organização da memória rege-se por uma hierarquia de sistemas de memória. As principais são as seguintes:
 - A memória interna refere-se à memória principal, diretamente acessada pela CPU
 - A CPU necessita da sua própria memória local, na forma de registradores
 - A memória externa (memória secundária) é referida a todos os dispositivos usados para armazenar dados e que acessem a CPU através de controladores/módulos de Entrada/Saída



Capacidade

- o Expressa em bytes ou words para a memória interna
- o Os tamanhos mais usuais de palavras são de 8,16, 32 bits
- o A capacidade de memória externa é expressa somente em bytes



Unidade de Transferência

o Interna

- o No nível da memória interna, refere-se ao número de linhas de dados que entram e saem do módulo de memória. Geralmente é igual à largura da palavra, mas nem sempre
- o É geralmente determinada pelo largura do barramento de dados

o Externa

- o É geralmente um bloco que é maior do que uma palavra

o Unidade Endereçável

- o O menor local de memória que pode ser endereçado

0 0 0

Métodos de Acesso(1)

- o Sequencial
 - o Começa no início e lê os dados sequencialmente
 - o O tempo de acesso depende da localização dos dados e da posição anterior
 - o Por exemplo, uma fita cassete
- o Direto
 - o Blocos individuais têm endereços únicos
 - o O acesso é feito por salto e por procura sequencial
 - o O tempo de acesso depende na localização dos dados e da localização anterior
 - o Por exemplo, um disco

0 0 0

Métodos de Acesso(2)

o Aleatório

- o Endereços individuais identificam locais exatos
- o O tempo de acesso é independente do local ou de acessos anteriores
- o Por exemplo, a RAM

o Associativa

- o Os dados são localizados por comparação de parte dos dados armazenados
- o O tempo de acesso é independente do local ou de acesso anteriores
- o Por exemplo, a cache



Desempenho (1)

- o Tempo de acesso (Access time)
 - o Para uma memória de acesso aleatório: é o tempo que demora a executar uma operação de leitura ou escrita, ou seja, é o tempo a partir do instante que o endereço é apresentado à memória até ao instante que os dados foram arquivados ou tenham ficado disponíveis para uso
 - o Para os outros tipos de memória, é o tempo que demora o mecanismo de leitura-escrita a posicionar-se na localização desejada



Desempenho (2)

- Tempo de ciclo de memória (Memory cycle time)
 - Compreende o tempo de acesso e o tempo necessário a memória para recuperar antes que outro acesso possa ser realizado
 - Tempo de ciclo de memória = acesso + recuperar
- Taxa de transferência (Transfer Rate)
 - Taxa a qual os dados podem ser transferidos
 - Para memória aleatórias é $1/(\text{Tempo de Ciclo})$
 - Para as outras memórias: $TN = TA + N/R$
 - TN - Tempo médio de leitura ou escrita de N bits
 - TA - Tempo médio de acesso
 - N - Número de bits
 - R - Taxa de transferência, em bits por segundo (bps)



Tipos de memória física

- o Semicondutora
 - o RAM
- o Magnética
 - o Disco & Fita
- o Ótica
 - o CD & DVD
- o Outros
 - o Magneto -optica



Memória

Questões importantes

- As questões no desenho de uma memória de um computador podem ser resumidas a três questões
 - Dimensão
 - Velocidade
 - Custo
- A relação entre estas características terá de ser considerada na escolha da memória
 - Quanto menor o tempo de acesso, maior o custo por bit
 - Quanto maior a capacidade, menor o custo por bit
 - Quanto maior a capacidade, maior o tempo de acesso



Hierarquia de memória

- o A forma escolhida para ter em consideração todas estas relações foi não ter apenas uma única memória mas sim ter uma hierarquia de memórias
- o À medida que descemos na hierarquia da memória temos:
 - o custo por bit decrescente
 - o capacidade crescente
 - o tempo de acesso crescente
 - o frequência de acesso à memória pelo CPU decrescente

0 0 0

Lista da Hierarquia

Maior
capacidade



Menor
custo
por bit

- o Registradores
- o L1 Cache
- o L2 Cache
- o Memória principal
- o Cache de disco
- o Disco
- o Óptico
- o Fita

Mais
Rápida

Maior
Frequência
de Uso

o o o

Hierarquia de memória

Tempo de acesso

Mais
Rápida



- o Registradores 5 nanoseg
- o L1 Cache
- o L2 Cache
- o Memória principal 70 nanoseg
- o Cache de disco
- o Disco 10.000.000 nanoseg (10 miliseg)
- o Óptico
- o Fita 10.000.000.000 nanoseg (10 seg)



Memória Cache



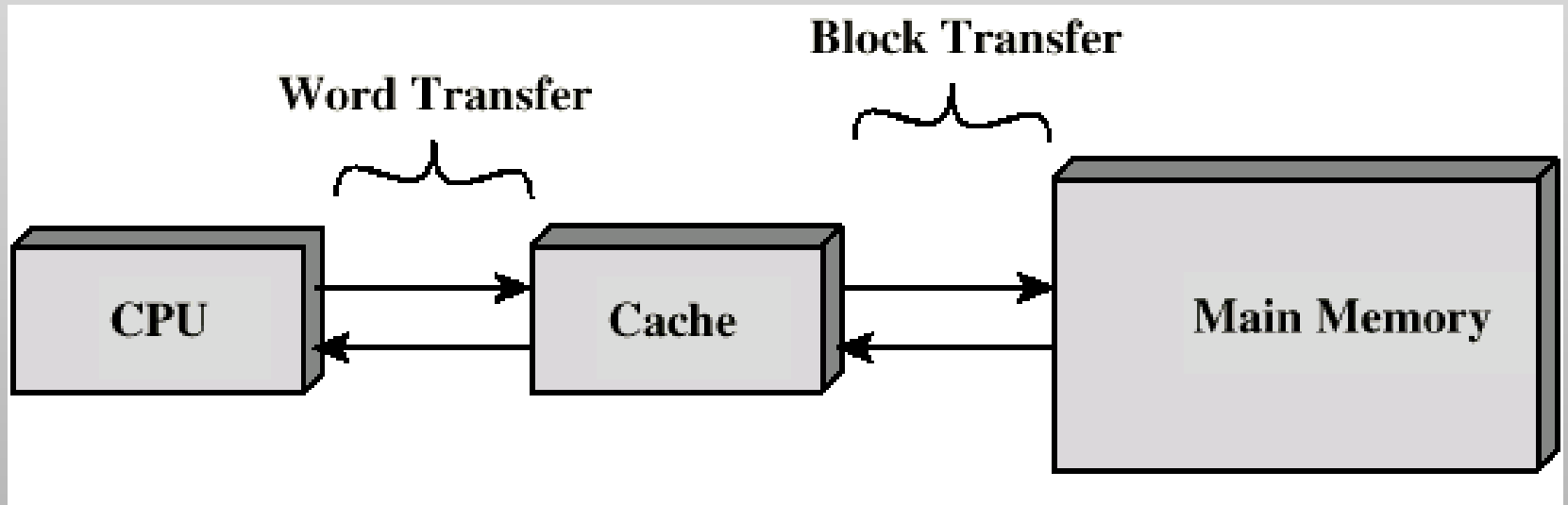
Memória Cache

- O tempo de execução de uma instrução é inferior ao tempo de acesso à memória
- A taxa que o CPU pode executar as instruções é claramente limitada pelo ciclo de acesso à memória
- Como o tempo de acesso à memória não acompanha o tempo de ciclo do processador, a solução consiste em colocar uma memória pequena e rápida entre o processador e a memória principal
- Esta memória é denominada memória cache.

0 0 0

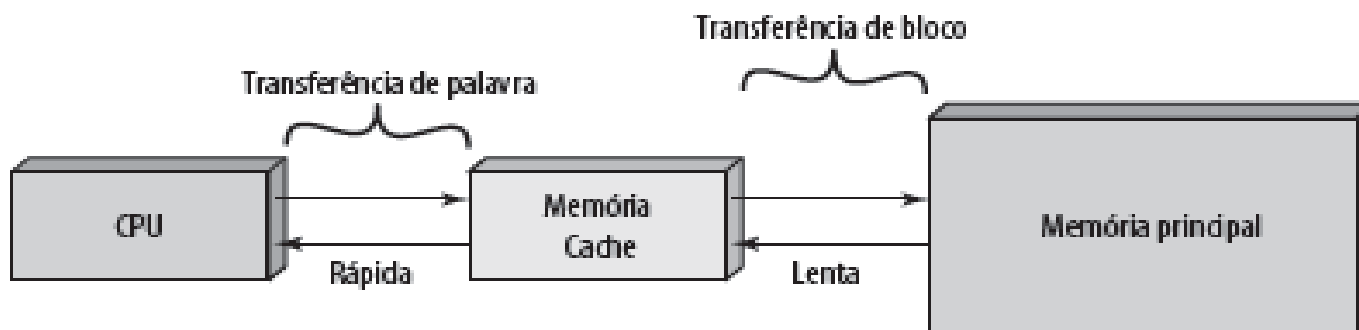
Cache

- o Pequena memória rápida
- o Posiciona-se entre a memória principal e a CPU
- o Pode estar localizada na CPU ou num módulo

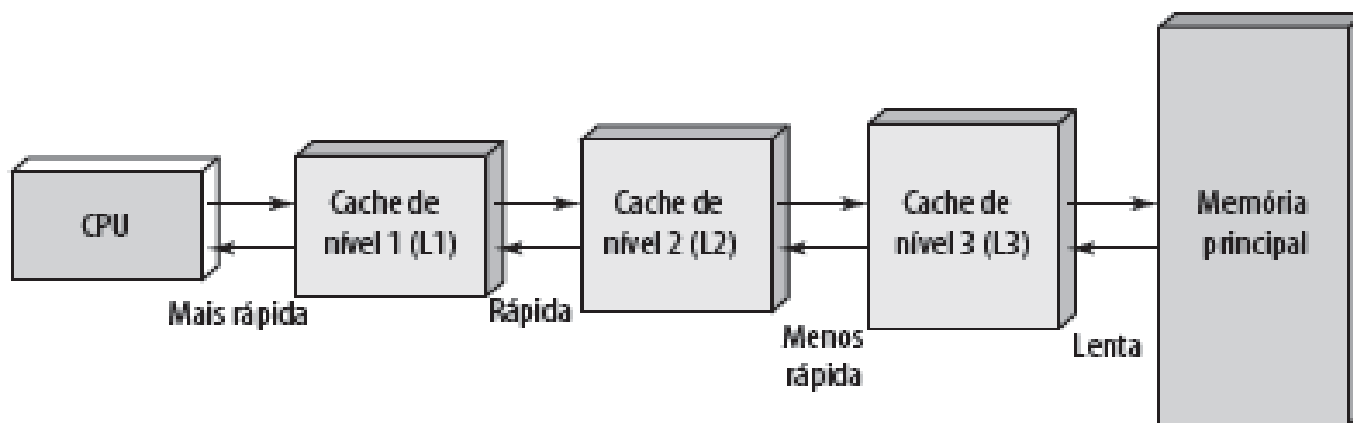


0 0 0

Cache e memória principal



(a) Cache única



(b) Organização de cache em três níveis



Operações da Cache

- A CPU pede dados a memória principal
- A cache é pesquisada a procura dos dados
- Se estes estão presentes, os dados são obtidos rapidamente
- Se não estão presentes, um bloco de dados é lido da memória principal para a cache
- Os dados são então entregues da cache a CPU
- A cache inclui identificadores (tags) para identificar os blocos de memória principal presentes na cache



Estrutura da Cache

- Para uma memória principal de 2^n palavras endereçadas, dividida em blocos de K palavras ($M=2^n/K$ blocos), a memória cache é organizada em C slots de tamanho fixo que contêm K palavras cada
- Como é óbvio $C \ll M$
- Por razões de eficiência, as leituras para (e da) cache são feitas em bloco e não palavra a palavra



Estrutura da Cache

- o Como existem mais blocos de memória que slots ($C \ll M$), um slot individual não pode estar permanentemente dedicado a um bloco
- o Então cada slot inclui uma tag que identifica o bloco que está ser guardado
- o Esta tag é usualmente parte do endereço da memória principal



Estrutura da Cache

- o A CPU gera um endereço de uma palavra para ser lida
- o Se a palavra estiver em cache, é entregue a CPU, caso contrário, o bloco que contém a palavra é carregado da cache e entregue a CPU



Cache Design

- o Tamanho
- o Função de Mapeamento
 - o Direto, associativo, associativo por conjuntos
- o Algoritmo de Substituição
 - o LRU, FIFO, LFU
- o Regras de Escrita
 - o Write-through, write-back, write-once
- o Dimensão dos Blocos
- o Número de Caches
 - o Um ou dois níveis
 - o Unificada ou separada

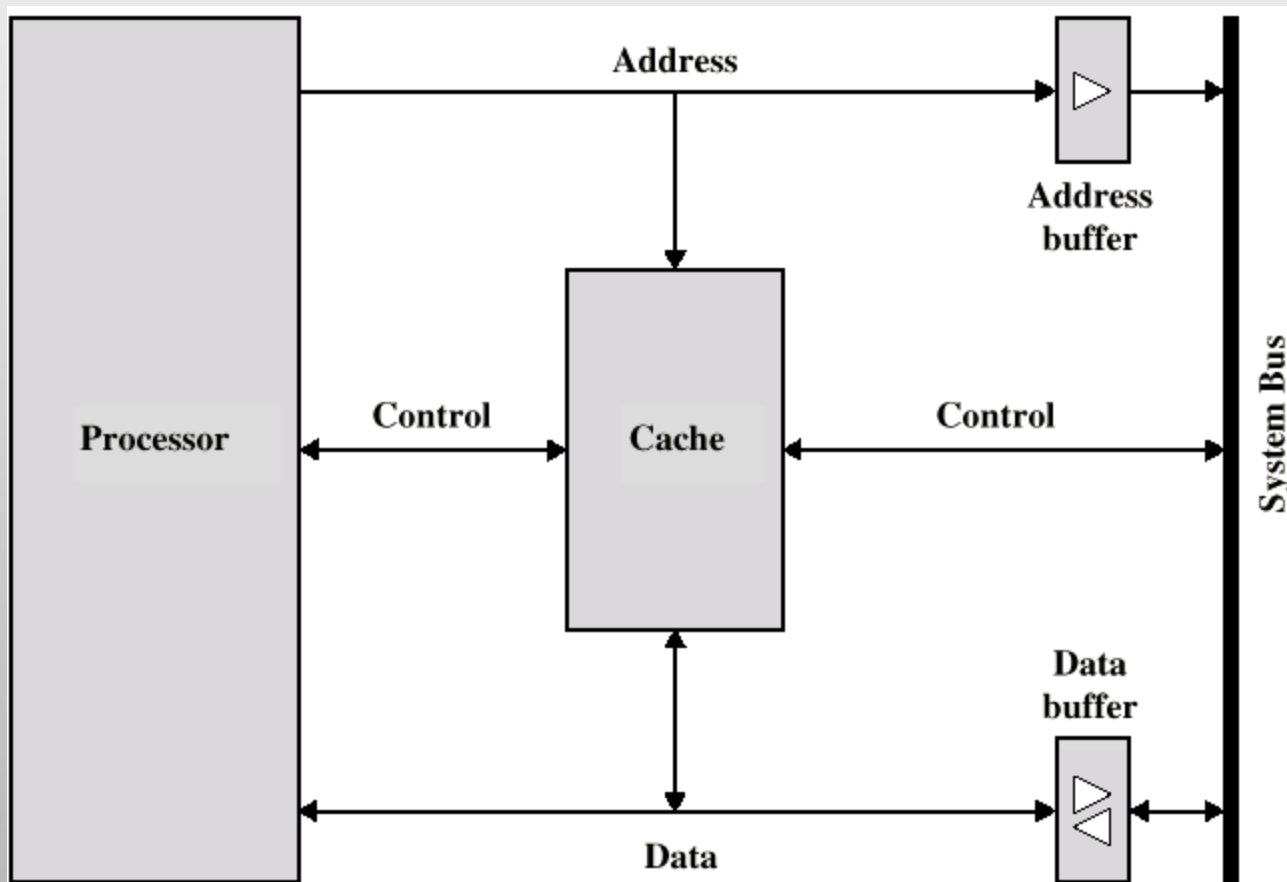


Tamanho é importante

- o Custo
 - o Ter uma maior cache é mais dispendioso
- o Velocidade
 - o Uma maior cache aumenta a velocidade (até um limite)
 - o Pesquisar a cache a procura de dados demora tempo

0 0 0

Organização Típica da Cache





Função de Mapeamento

Exemplo

- o 16MBytes memória principal
- o Endereços de 24 bit
 - o $(2^{24}=16M)$
- o Cache de 64kByte
- o Cache com blocos de 4 bytes
 - o i.e. cache tem 16k (2^{14}) linhas de 4 bytes

0 0 0

Mapeamento Direto



- Tamanho da cache = 32KB
- Tamanho do bloco = 16 bytes
 - i.e. cache tem 2K (2^{11}) linhas de 16 bytes
- Memória principal = 16 MB
 - Endereços de 24 bit ($2^{24}=16M$)
 - 2^{20} (1M) blocos de 2^4 bytes



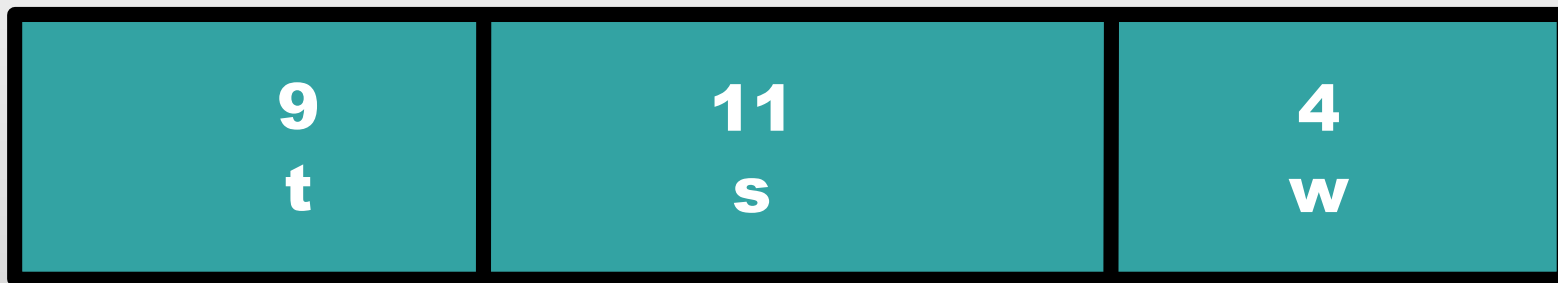
Mapeamento Direto

- Cada bloco da memória principal é mapeado a uma só linha da cache
 - i.e. se um bloco está na cache, este tem de estar numa posição específica
- Fórmula para mapear um bloco de memória a uma linha de cache :
$$i = j \bmod c$$
 - i = Número da Linha da Cache
 - j = Número do Bloco da Memória Principal
 - c = Número de Linhas na Cache

o o o

Estrutura dos Endereços

Mapeamento Direto



- Os w bits menos significativos identificam uma palavra numa linha da cache
- Os próximos s bits significativos identificam em que slot um endereço é mapeado
- Os últimos t bits significativos identificam um bloco de memória

0 0 0

Mapeamento Direto com $C=4$

○ Linha da Cache	Bloco de Memória Principal
○ 0	0, 4, 8, ...
○ 1	1, 5, 9, ...
○ 2	2, 6, 10, ...
○ 3	3, 7, 11, ...
○ Regra Geral:	
○ 0	0, C , $2C$, $3C$, ...
○ 1	1, $C+1$, $2C+1$, $3C+1$, ...
○ 2	2, $C+2$, $2C+2$, $3C+2$, ...
○ 3	3, $C+3$, $2C+3$, $3C+3$, ...

o o o

Estrutura dos Endereços

Mapeamento Direto

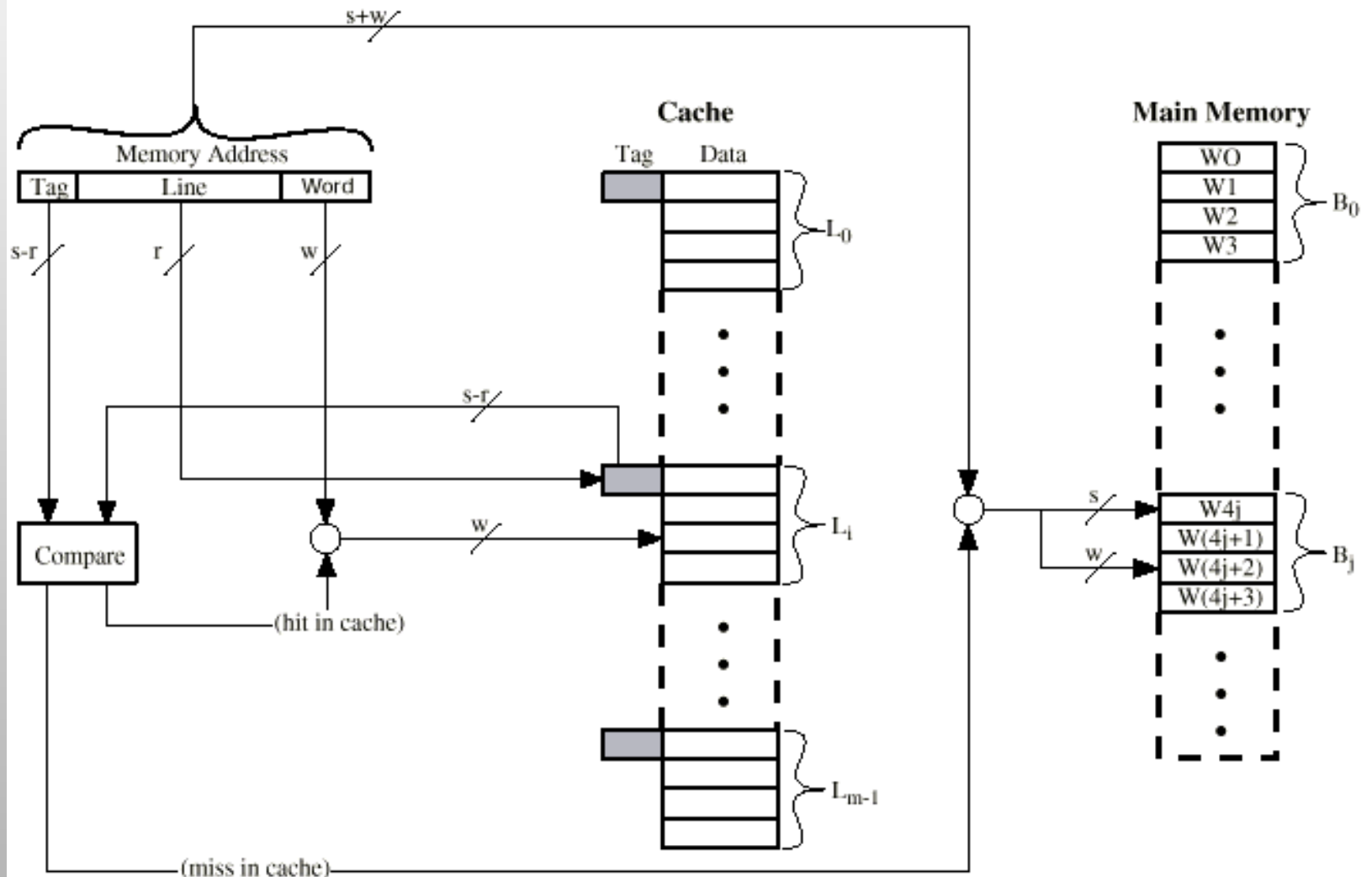


- o Endereços de 24 bit
- o Identificador de palavra: 2 bits (blocos de 4 bytes)
- o Identificador de bloco: 22 bit
 - o Tag de 8 bit (=22-14)
 - o Linhas de 14 bit
- o Dois blocos na mesma linha nunca têm o mesmo Tag
- o Encontram-se dados na cache procurando a linha (r) e comparando a Tag (s-r)

Mapeamento Direto

Mapeamento Direto

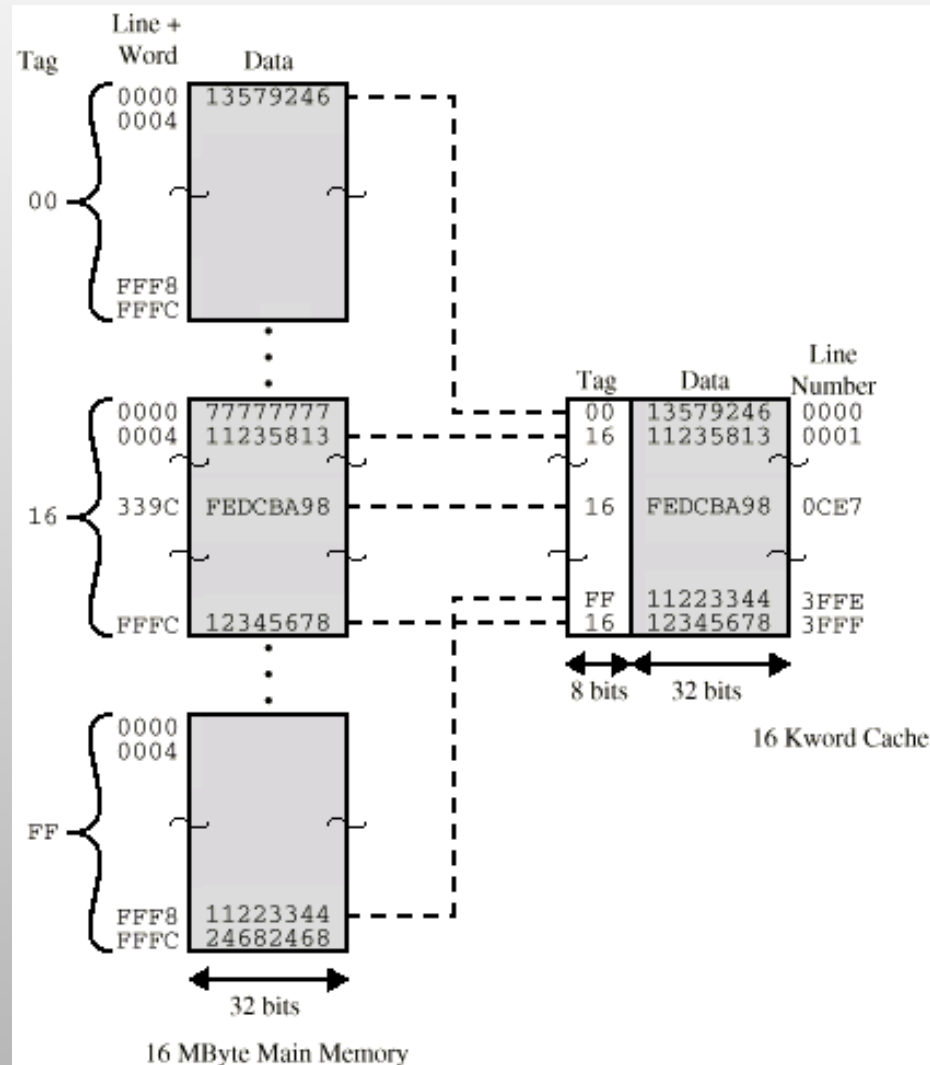
Mapeamento Direto



0 0 0

Direct Mapping

Exemplo





Mapeamento Direto

prós & contras

- o Simples
- o Barato
- o Localização fixa para os blocos
 - o Se um programa acede a dois blocos que mapeiam na mesma linha repetidamente, os cache misses tornam-se muito elevados

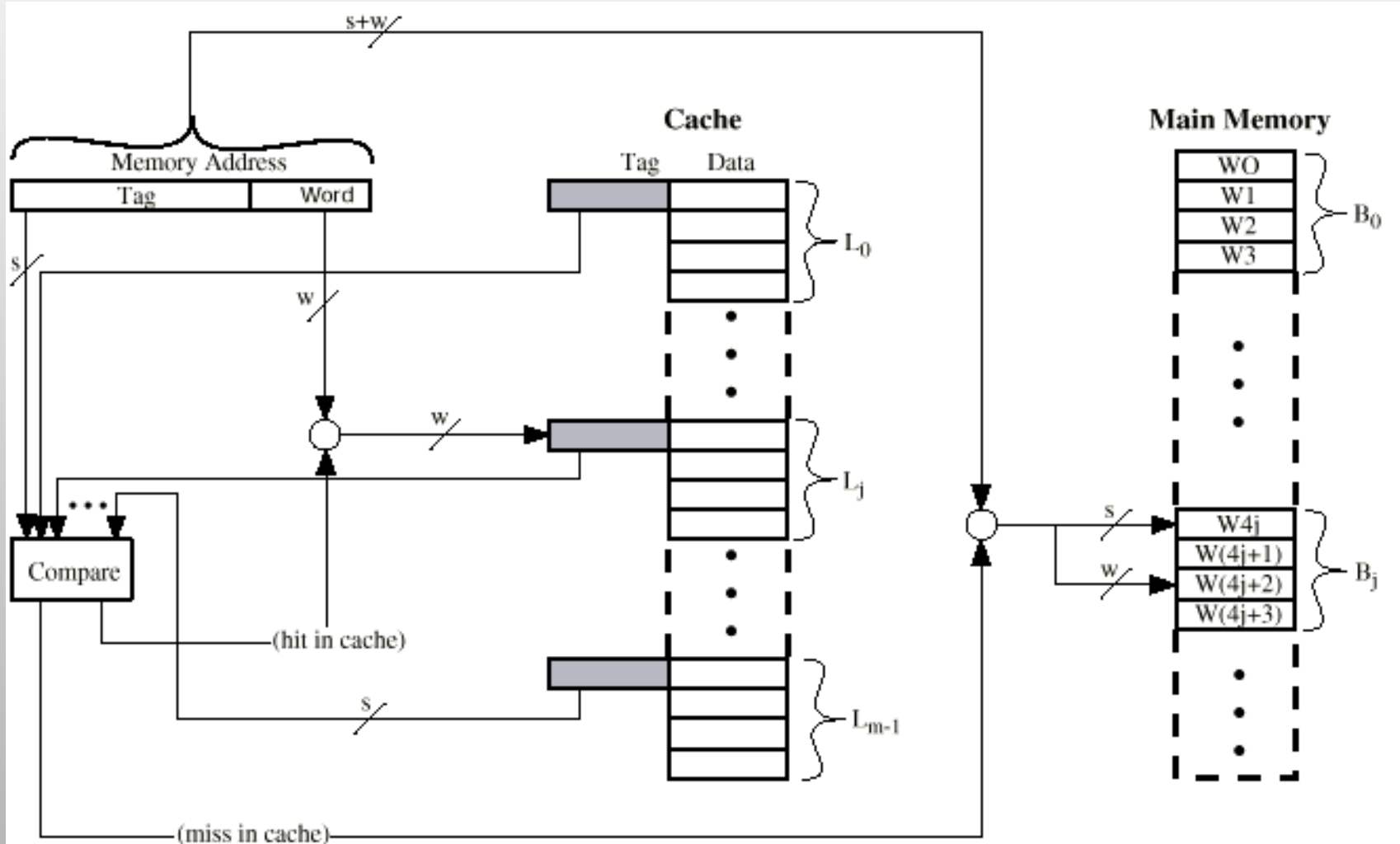


Mapeamento Associativo

- o Um bloco de memória pode ser carregado em qualquer linha da cache
- o Um endereço de memória é composto por tag e palavra
- o O tag identifica o bloco de memória
- o Cada linha de uma tag é examinada a procura de um match

o o o

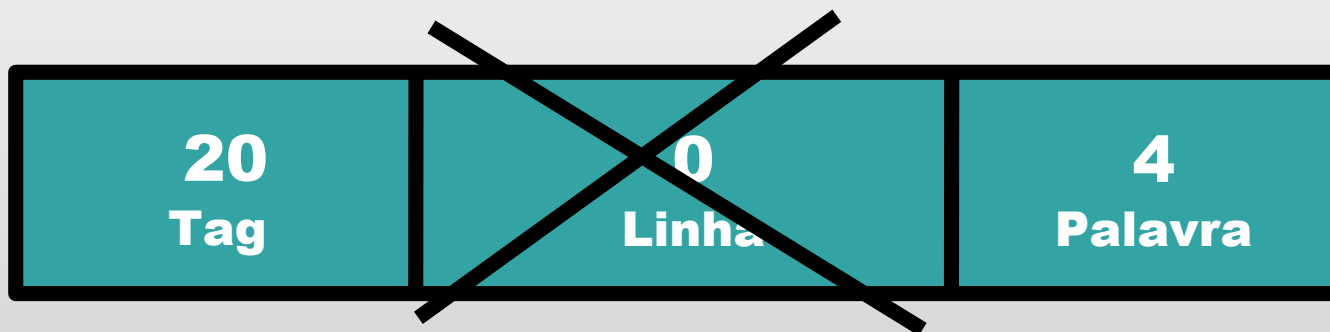
Organização de Cache Associativa



o o o

Estrutura dos Endereços

Mapeamento Associativo

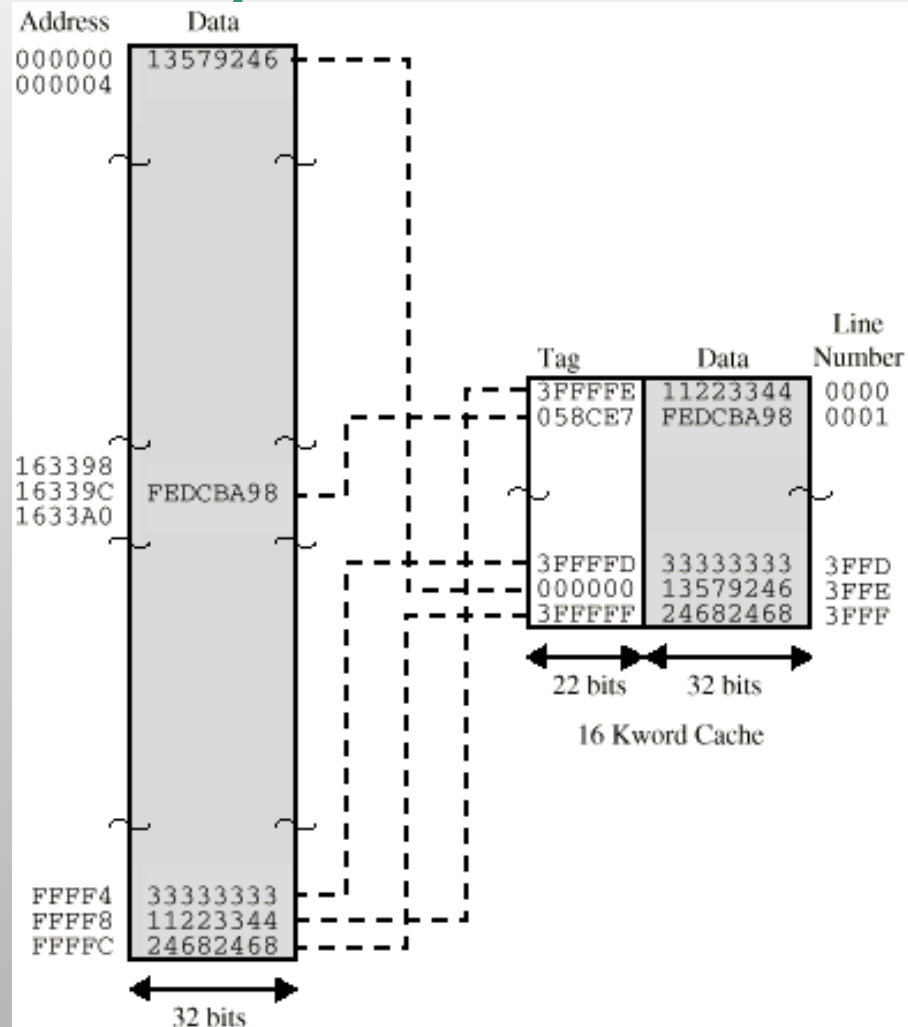


- o Tamanho da cache = 32KB
- o Tamanho do bloco = 16 bytes
- o Memória principal = 16 MB
- o 2^{24} bytes

o o o

Mapeamento Associativo

Exemplo





Mapeamento Associativo

- A pesquisa da cache é cara
 - Idealmente seria necessário um circuito que pode simultaneamente examinar todas as tags ao mesmo tempo
 - Implicaria a utilização de muitos circuitos
 - custo elevado
- Como qualquer dado pode ser expulso da cache é necessário desenvolver regras de substituição



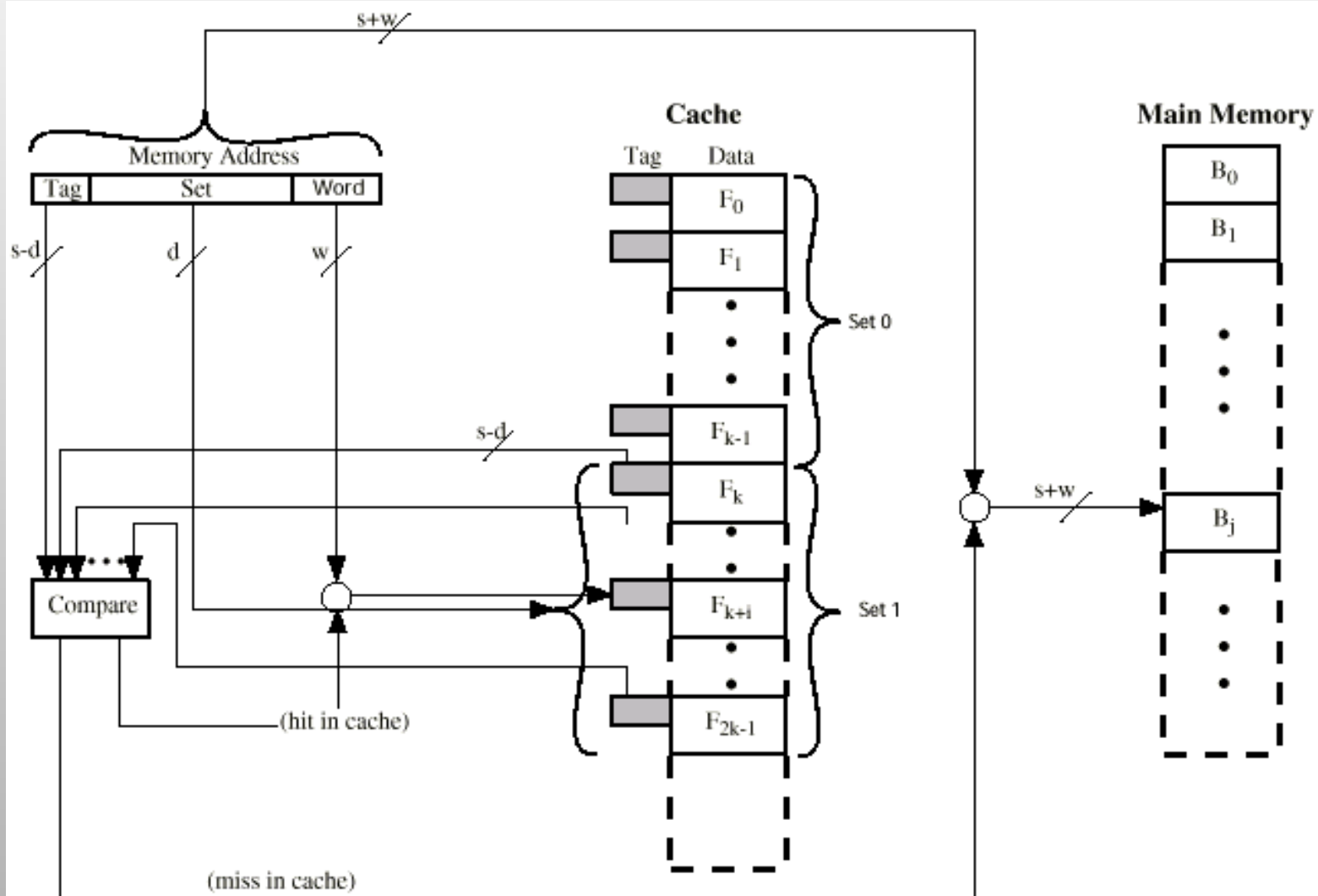
Mapeamento Associativo por Conjuntos

- Divide a cache em v conjuntos de k linhas
 - $m = v \times k$
 - $i = j \text{ modulo } v$, onde
 - i = número do conjunto na cache
 - j = número do bloco na memória principal
 - m = número de linhas da Cache
- Com este tipo de mapeamento, o bloco B pode ser mapeado em qualquer linha do conjunto i
- Permite capturar as vantagens dos dois métodos anteriores
 - Nos casos extremos se $v = m$, $k = 1$ (direta); se $v = 1$, $k =$

m (associativa)

0 0 0

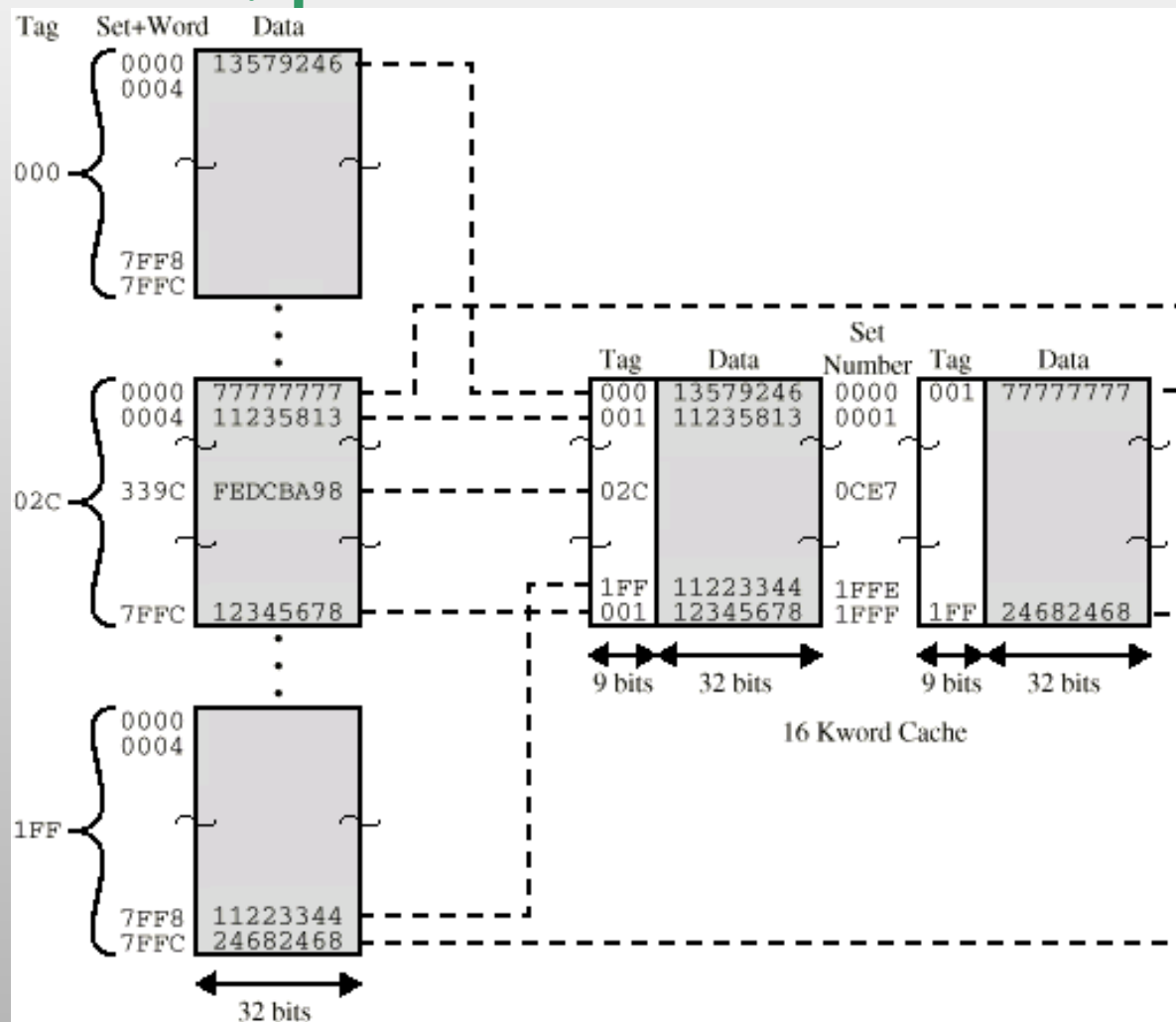
Organização de Cache com Mapeamento Associativo por Conjunto



o o o

Mapeamento Associativo

Exemplo





Algoritmos de Substituição

- Qual o bloco a substituir? (se o bloco não está na cache)
- Mapeamento Direto
 - Só existe uma opção
- Mapeamento Associativo
 - Cada bloco é um candidato
- Associativa de nível k
 - Um de k

0 0 0

Algoritmos de Substituição

- o A substituição de blocos na cache pode ser
- o implementada por diversos algoritmos
 - o LRU (least recently used)
 - o FIFO (First-in First Out)
 - o LFU (least frequently used)
 - o Random (aleatoriamente)



Política de Escrita

- Existem duas técnicas de escrita de blocos modificados na memória principal
 - Write-Through - as operações são executadas em simultâneo na cache e na memória
 - Write-back - as operações são executadas apenas na cache e atualizando a memória principal apenas quando o bloco é substituído



Tamanho de linha

- Recupere não apenas a palavra desejada, mas também uma série de palavras adjacentes.
- Tamanho de bloco aumentado aumentará razão de acerto a princípio.
 - O princípio da localidade.
- Razão de acerto diminuirá à medida que o bloco se tornar ainda maior.
 - Probabilidade de uso de informações recém- buscadas torna-se menor que probabilidade de reutilizar informações substituídas.



- o Blocos maiores:
 - o Reduzem número de blocos que cabem na cache.
 - o Dados sobrescritos pouco depois de serem buscados.
 - o Cada palavra adicional é menos local, de modo que é menos provável de ser necessária.
- o Nenhum valor ideal definitivo foi descoberto.
- o 8 a 64 bytes parece ser razoável.
- o Para sistemas HP, 64 e 128 bytes mais comum.



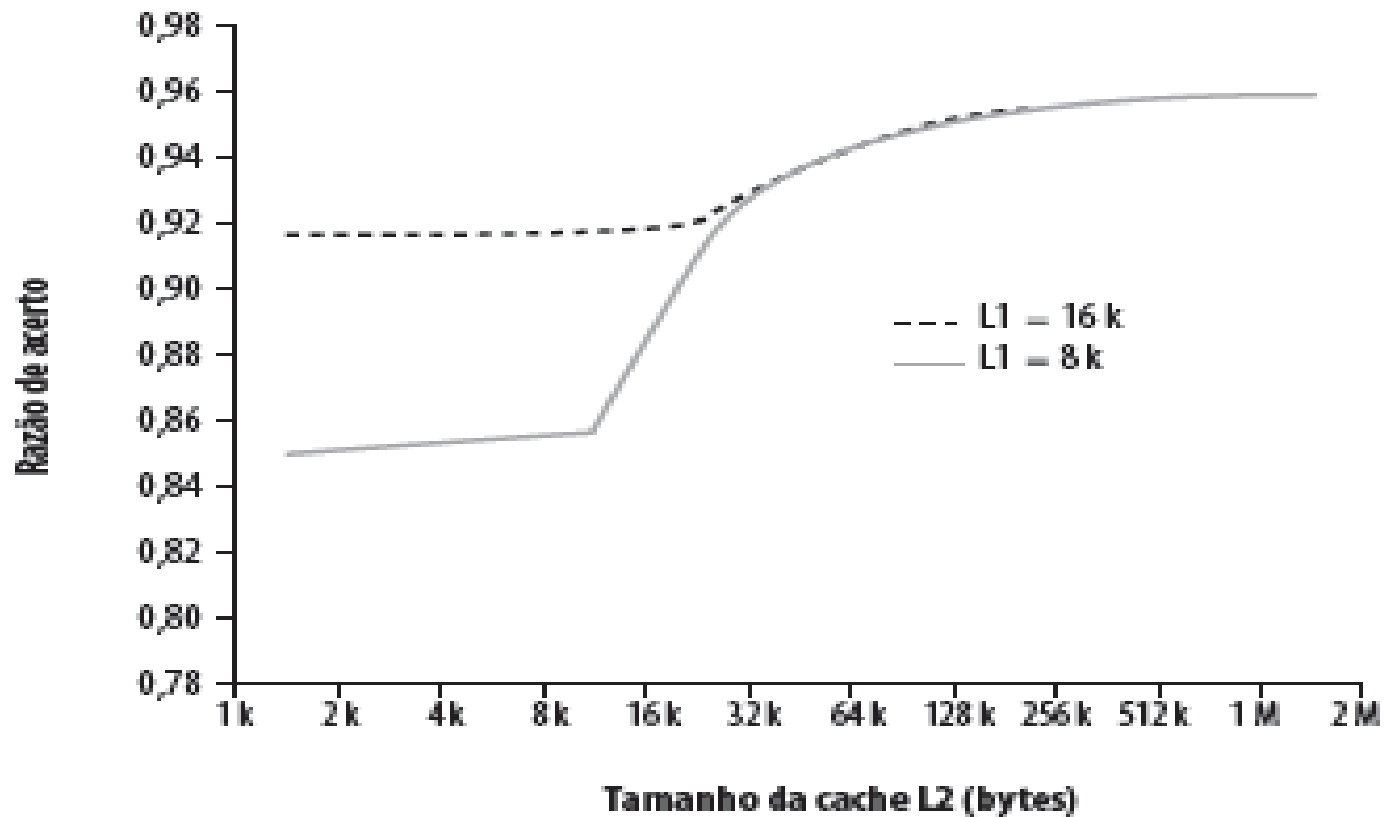
Caches multinível

- Alta densidade lógica permite caches no chip.
 - Mais rápido que acesso ao barramento.
 - Libera barramento para outras transferências.
- Comum usar cache dentro e fora do chip.
 - L1 no chip, L2 fora do chip na RAM estática.
 - Acesso L2 muito mais rápido que DRAM ou ROM.
 - L2 normalmente usa caminho de dados separado.
 - L2 pode agora estar no chip.
 - Resultando em cache L3.
 - Acesso ao barramento agora no chip.

0 0 0

Razão de acerto total (L1 & L2)

Para L1 de 8 KB e 16 KB





Caches unificadas *versus* separadas

- Uma cache para dados e instruções ou duas, uma para dados e uma para instruções.
- Vantagens da cache unificada:
 - Maior taxa de acerto.
 - Equilibra carga entre buscas de instrução e dados.
 - Apenas uma cache para projetar e implementar.
- Vantagens da cache separada:
 - Elimina disputa pela cache entre a unidade de busca/decodificação de instrução e a unidade de execução.
 - Importante no pipeline de instruções.