



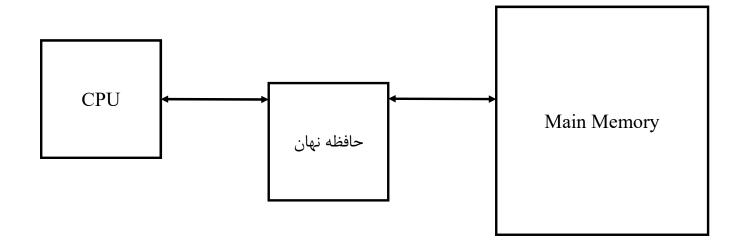


معماری کامپیوتر

جلسه هفتم: حافظه نهان-۱



- درنظر گرفتن فضایی بین پردازنده و حافظه اصلی (مشابه جعبه کوچک کتابدار)
 - باهدف ذخیرهسازی بخشی از حافظه اصلی جهت دسترسی سریعتر





- سوال اساسى: در حافظه نهان چه اطلاعاتى ذخيره كنيم؟
 - هدف اصلى: افزايش سرعت دسترسى ____ نرخ hit بالا
- طراحی بد: متوسط زمان دسترسی بیشتر نسبت به حالت دسترسی مستقیم به حافظه اصلی
 - انتخاب دادههای حافظه نهان، براساس پیشینه درخواستها از حافظه اصلی
 - استخراج ویژگیهای درخواستها با نظارت بر انتقالات داده بین پردازنده و حافظه اصلی
 - خاصیت محلیت (locality) در دادهها

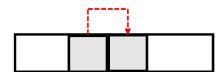


- خاصیت هم جواری (Locality):
- هم جواری مکانی (Spatial Locality)
- هم جواری زمانی (Temporal Locality)
- لحاظ کردن این خاصیت در طراحی حافظه نهان



• هم جواری مکانی (Spatial Locality)

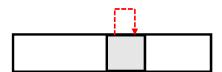
- درخواستهای متوالی از حافظه در آدرسهای نزدیک بههم قرار دارند و لذا وابستگی مکانی بین آنها قابل استخراج میباشد.
 - در حلقه for تعدادی دستورالعمل متوالی داریم
- زمانی که به یک کتاب در حیطه معماری کامپیوتر مراجعه کنیم به احتمال زیاد به کتاب دیگر در این حوزه هم نیاز داریم





• هم جواری زمانی (Temporal Locality)

- یک داده در زمانهای نزدیک بههم چندین بار مورد استفاده قرار می گیرد
 - متغیری که چندین بار در طی اجرای برنامه استفاده می شود (مثلا در حلقه ها)
- زمانی که به یک کتاب در حیطه معماری کامپیوتر مراجعه کنیم به احتمال زیاد دوباره هم به آن کتاب مراجعه می کنیم و با یکبار مسئله حل نمی شود.





همجواری مکانی:

داده: ارجاع متوالی به خانههای آرایه دستورالعمل: ارجاع به دستورات بهصورت ترتیبی

همجواری زمانی:

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;</pre>
```



• سوالهای مهم:

- دادهها را چگونه از حافظه اصلی به حافظه نهان ببریم؟
- دادههای جدید را چگونه در حافظه نهان جایگزین کنیم؟
- برای پاسخ به این دو سوال، دو بحث اساسی مطرح می شود:
 - سیاست جای دهی (Placement Policy)
 - سیاست جای گزینی (Replacement Policy)



• فرض كنيم حافظه نهان C خط (Line) به طول يك كلمه داشته باشد و حافظه اصلى C

0 Main Memory
C-1 M-1

 $C < M \bullet$

سیاست جای دهی و انواع حافظه نهان



۱- حافظه نهان نگاشت مستقیم (Direct mapping):

- اولین گام: سیاست نگاشت خانههای حافظه اصلی به حافظه نهان (Address Mapping)
 - مادامی که حافظه نهان خالی است نگاشت بدین صورت انجام می گیرد:

Cache Address = (Main Memory Address) mod (C)

- سیاست نگاشت آدرس باید پوشا و یکتا باشد
- روش بیان شده **پوشا** هست ولی **یکتا** نیست



• علت یکتا نبودن نگاشت آدرس پیشنهادی:

• ذخیره دادهای با آدرس X از حافظه اصلی در حافظه نهان

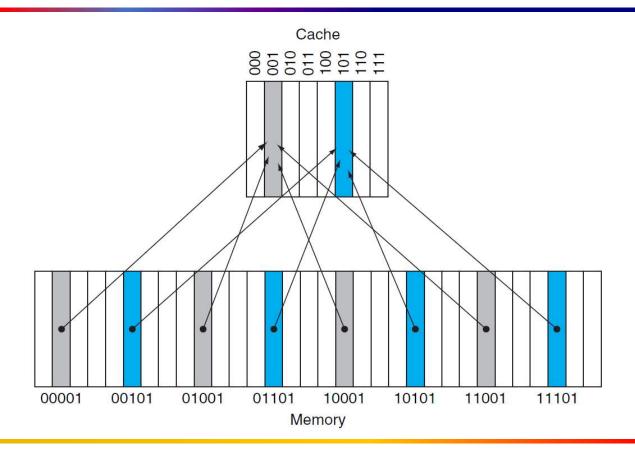
Cache address(X): X mod C

• ذخیره دادهای با آدرس X+C از حافظه اصلی در حافظه نهان

Cache address(X+C): (X+C) mod C

• هر دو به یک مکان از حافظه نهان اشاره دارند!



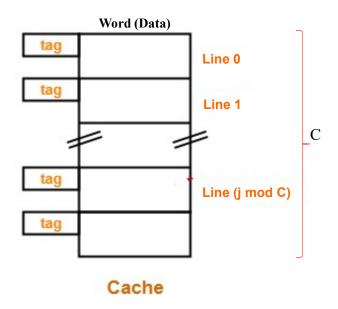




• راهکار یکتا کردن نگاشت آدرس پیشنهادی:

- ذخیره اطلاعاتی از آدرس داده در حافظه اصلی در کنار داده در حافظه نهان
- $^{f C}$ اندیس آدرس در حافظه نهان: باقی مانده تقسیم آدرس در حافظه اصلی بر
 - استفاده از خارج قسمت این تقسیم به عنوان شناسه ای از آدرس داده
 - ذخیره این شناسه به عنوان برچسب (Tag) داده در حافظه نهان







	Index	
Ī	000	1
ľ	001	1
Γ	010	1
Γ	011	1
Г	100	1
Г	101	1
Γ	110	
Г	111	1

V Tag		Data	
N	1		
N			
N			
N			
N			
N			
N			
N			

1	nd	ex
	00	0
3	00	1
	01	0
	01	1
	10	0
	10	1
	11	0
	11	1

V	Tag	Data
N		
N		
N		
N		
N		
N		
Υ	10 _{two}	Memory (10110 _{two})
N		

a. The initial state of the cache after power-on

b. After handling a miss of address (10110 $_{\mbox{two}}$)



Word address	Binary address	Hit/miss	Cache block
22	10110	Miss	110



- گام دوم: شیوه ذخیرهسازی دادهها در حافظه نهان
- در طراحی فعلی: هربار یک کلمه از حافظه اصلی به حافظه نهان منتقل می شود
 - کارایی cache وابسته به prediction مناسب درخواستها
 - همجواری زمانی و همجواری مکانی
 - اصلاح طراحی برای لحاظ کردن همجواری مکانی
 - انتقال یک بلوک B کلمهای از حافظه اصلی به حافظه نهان در هر بار دسترسی



$$M = 2^m$$

$$C = 2^c$$

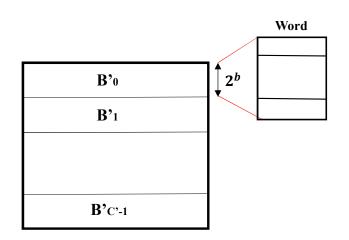
$$B = 2^b$$

• نکته مهم: M و G همگی توانهایی از ۲ هستند

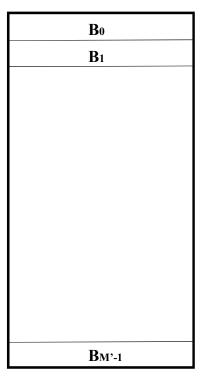
- تقسیم حافظه نهان و حافظه اصلی به بلوکهای B کلمهای
 - همگی توانهایی از ۲ هستند پس تقسیم شیفت
- در هربار انتقال کلمه از حافظه اصلی به حافظه نهان، کل بلوکی که داده در آن است منتقل میشود
- فرض: حافظه اصلی دارای M بلوک B کلمهای و حافظه نهان دارای C بلوک B کلمهای M



طول هر بلاک: 2^b کلمه



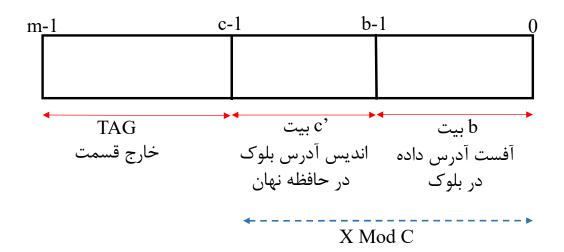
Cache



Main Memory



• قالب آدرس حافظه اصلی ورودی به حافظه نهان:





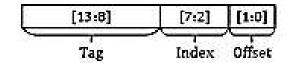
Memory Size (M) = 16 KB, m=14

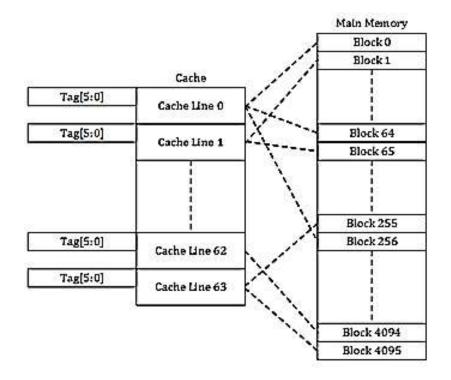
Block Size (B) = 4 B, b = 2

Cache Size (C) = 256 B, c=8

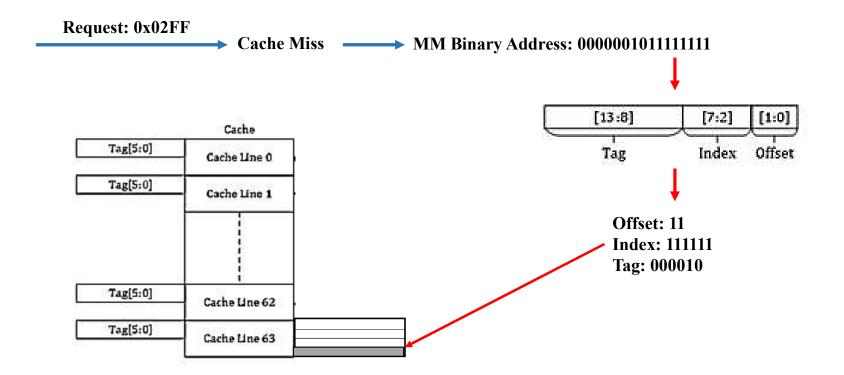
blocks in MM: M/B= 4096

blocks in cache: C / B = 64







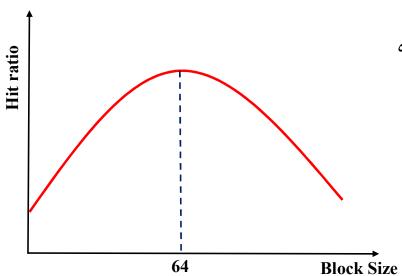




- تمامی اعضای یک بلوک حافظه نهان، Tag مشترک دارند
 - وبرای یک بلوک نگهداری یک برچسب کافی است
- روش نگاشت مستقیم در طراحی حافظه نهان، سادهترین و کمهزینهترین روش است
 - طبق آزمایشات و ارزیابیهای، hit ratio حافظه نهان طراحی شده این روش حدود ۹۲٪



- در این روش، سایز بلوک را می توان افزایش داد
 - نمودار شدن بیشتر همجواری مکانی
- کم شدن تعداد بلوکهای ذخیره شده و تنوع داده





• مدیریت درخواست نوشتن در حافظه:

(چەزمانى؟)

- مشابه درخواست خواندن ابتدا درخواست به حافظه نهان می رود:
- Miss: داده در حافظه اصلی پیدا شده و مقدار آن بهروز می گردد و به حافظه نهان منتقل می شود.
- Hit: داده در حافظه نهان موجود است و باید بهروز شود. این بهروزرسانی باید در حافظه اصلی هم انجام شود.



• مدیریت درخواست نوشتن در حافظه:

- چگونگی بهروزرسانی داده در حافظه اصلی در صورت hit شدن درخواست در حافظه نهان:
- Write through: هرگاه دادهای در حافظه نهان بهروز شد، محتویات آدرس متناظر در حافظه اصلی هم بهروز می شود
 - Write back: بهروز رسانی داده در حافظه اصلی تا زمان پاک شدن بلوک از حافظه نهان به تاخیر میافتد.
 - یک بیت کنترلی U برای هر بلوک نگه می داریم که نشانگر به روز شدن آن بلوک است
 - کارایی بهتری دارد اما ریسک از دست دادن داده جدید وجود دارد