

זיכרון מטמון - Cache

מה זו היררכיית זיכרון?

בגישה של היררכיית זיכרון, הזיכרון מחולק למספר רמות – הרמה התחתונה היא הרמה הגדולה והאיטית יחסית והרמה העליונה (cache) היא הרמה הקטנה והמהירה ביותר. כל המידע ברמת זיכרון נתונה, נמצא בכל אחת מהרמות שמתחתיה. מעבר של נתונים מרמה אחת לשנייה נעשה ביחידות של בלוקים (ע"מ לנצל את תכונת המקומיות במרחב, ובכך לחסוך בכמות הגישות לזיכרונות איטיים).

למה צריך זיכרון מטמון?

במצב האידיאלי היינו רוצים מעבד עם זכרון יחיד (לפשטות), הגדול ביותר, המהיר ביותר והזול ביותר. אלא שאילוצים טכנולוגיים גורמים לכך שככל שהזכרון גדול יותר, כך הוא איטי יותר וככל שהזכרון קטן יותר כך הוא מהיר יותר ויקר יותר.. והמטרה היא להשיג מערכת זיכרון מהירה **כמעט** כמו המהירה ביותר, וזולה כמעט כמו הזולה ביותר.

- 1) for (I=0;I<1000;I++)
- 2) X+=arr[I];

כלל אצבע – המערכת מבלה 80% מהזמן בביצוע של 20% מהקוד (במקרה זה שורות 1,2) **מקומיות בזמן** – תוכניות משתמשות פעם נוספת בנתונים ופקודות שבהם הן השתמשו לאחרונה (בתוכנית שבדוגמא המשתנה I בשימוש פעמים רבות, והפקודות עצמן מובאות שוב ושוב). זיכרון מטמון קטן ומהיר שיכיל את הפקודות המרכיבות את הקוד שבלולאה יהיה ניצול טוב של תכונת המקומיות בזמן.

מקומיות במרחב - ל-2 משתנים שכתובותיהן נמצאות אחת ליד השנייה, תהייה התייחסות קרובה בזמן (כל כתובות המערך arr סמוכות אחת לשנייה). מערכת הזיכרון עובדת ביחידות של בלוקים של Bytes, ובכך מנצלת את תכונת המקומיות במרחב.

נוסחא להערכת ביצועי זיכרון היררכי:

Average Memory Access Time =

$$\text{hit time} + \text{miss rate} \times \text{miss penalty}$$

hit time – זמן קריאת/כתיבת נתון לבלוק בזיכרון המטמון, כאשר הבלוק נמצא.

miss rate – החלק של הגישות לזיכרון שבהן הבלוק שאילו אנו רוצים לקרוא/לכתוב אינו נמצא ב-cache.

miss penalty – הזמן הלוך להביא בלוק מהזיכרון הראשי ל-cache.

כדי להבין את אופן פעולת היררכיית זיכרון, נענה לגביה על 4 שאלות:

- 1) להיכן ב- Cache יכול להתמפות בלוק מהזיכרון?
- 2) כיצד מוצאים בלוק ב-Cache?
- 3) איזה בלוק מחליפים במקרה שהבלוק לא נמצא ב-Cache (miss), וצריך להביאו מהזיכרון?.
- 4) מה קורה בכתיבה?

1) להיכן ב- Cache יכול להתמפות בלוק מהזיכרון?

קיימות 3 שיטות למיפוי בלוק מהזיכרון הראשי ל- Cache

1. מיפוי ישיר (direct mapped) – בלוק מהזיכרון ממופה לבלוק אחד ויחיד בזיכרון המטמון. פונקציית המיפוי הינה:

$$\left(\begin{array}{c} \text{מספר הבלוק} \\ \text{בזיכרון} \end{array} \right) \bmod \left(\begin{array}{c} \text{מספר הבלוקים} \\ \text{ב-cache} \end{array} \right)$$

חסרון עיקרי – אם במקרה משתמשים תכופות ב-2 כתובות שמתמפות לאתו בלוק ב-cache, יוצר מצב שבו נביא בלוק מהזיכרון, ומייד אח"כ נביא בלוק אחר שידרוס אותו וחוזר חלילה, לדוגמא בתוכנית:

```
for (I=0;I<10;I++)
{
    x1+=arr1[I];
    x2+=arr2[I];
}
```

יתרון עיקרי – פשטות במימוש.

2. Fully associative – בשיטה זו ניתן למפות כל בלוק בזיכרון לכל בלוק ב-cache. חסרון עיקרי – סיבוכיות במימוש

יתרון עיקרי – גמישות רבה בבחירת הבלוק להחלפה.

3. n way set associative – בלוק מהזיכרון ימופה לקבוצת (set) בלוקים ב-cache, בתוך הקבוצה בגודל n בלוקים, ניתן למפות לכל אחד מ-n הבלוקים.

פונקציית המיפוי של בלוק לקבוצה הינה:

$$\left(\begin{array}{c} \text{מספר הבלוק} \\ \text{בזיכרון} \end{array} \right) \bmod \left(\begin{array}{c} \text{מספר הקבוצות} \\ \text{ב-cache} \end{array} \right)$$

פיתרון זה מהווה פשרה בין סיבוכיות מימוש לגמישות בבחירת הבלוק להחלפה.

הערה:

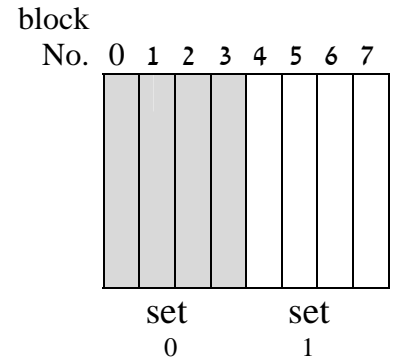
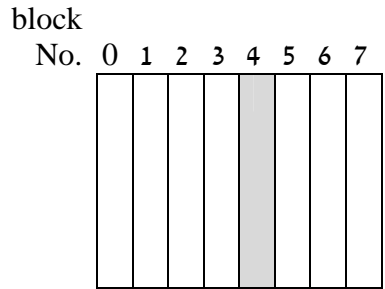
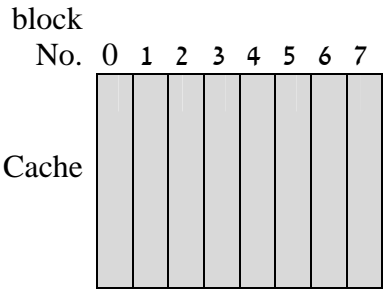
ניתן להתייחס למיפוי ישיר כ-1 way set associative, ול- Fully associative כ-M way set associative, כאשר M הינו מספר הבלוקים ב-cache.

דוגמא למיפוי בלוק מס' 12 בזיכרון הראשי ל-cache:

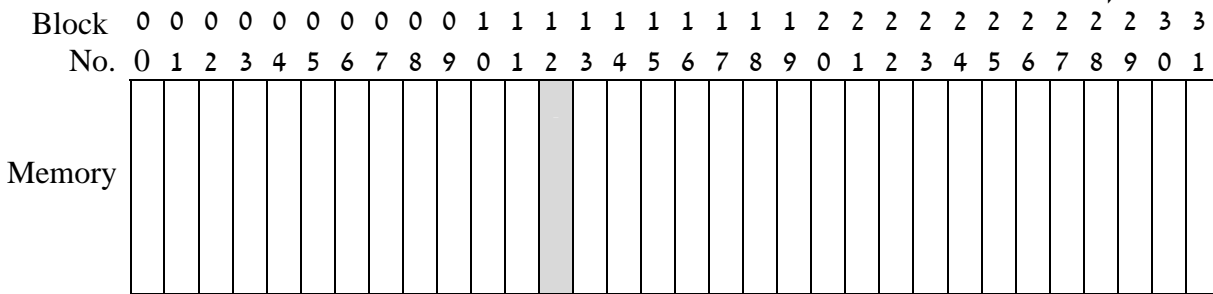
Fully associative:
בלוק 12 יכול להתמפות לכל אחד מהבלוקים ב-cache

Direct Mapped:
בלוק 12 יכול להתמפות רק לבלוק מס' 4 ב-cache
($12 \bmod 8 = 4$)

4-way set associative
בלוק 12 יכול להתמפות לכל אחד מ-4 הבלוקים שב-set 0.
($12 \bmod 2 = 0$)



הזיכרון הראשי:



2) כיצד מוצאים בלוק ב-Cache?

נניח:

גודל הזיכרון: $2^{32}=4\text{GBytes}$, לכן כתובת של Byte הינה בת 32 ביטים.

גודל בלוק: $2^3=8\text{Bytes}$

גודל ה-cache – $8\text{Kbyte} = 2^{10} = 1024$ בלוקים.

שיטת מיפוי n way set associative

מס' הקבוצות (sets) ב-cache = $n/\text{מס' הבלוקים ב-cache}$

מבנה הכתובת:

Byte Address		
Block Frame (= Block address)		Offset
Tag	Index	
מה שנשאר	$\log_2(\text{cache-ב-קבי})$	$\log_2(\text{block size}) = 3 \text{ bits}$

Offset – בוחר את ה-Byte המתאים מתוך הבלוק.

Index - בוחר את הקבוצה (set) המתאימה.

Tag – נשמר יחד עם הבלוק ב-cache.

כאשר מאחסנים בלוק ב-cache שומרים יחד איתו גם valid bit שערכו 1 כאשר יש מידע אמיתי בבלוק, ואת שדה ה-Tag.

אין צורך, כמובן, לשמור את שדה ה-offset, כי כל הבלוק נשמר כיחידה אחת. אין צורך לשמור את שדה ה-Index, כי ממילא הימצאות של בלוק בקבוצה x ב-cache, מחייבת ששדה האינדקס של אותו בלוק יהיה x.

כדי למצוא Byte ב-cache, נלך אל ה-set המתאים בעזרת ה-Index, נעבור על כל n הבלוקים שב-set (במקביל) ונשווה את ה-Tag של הכתובת הרצויה ל-Tag ששמור יחד עם הבלוק ב-cache. אם יש התאמה של ה-Tag ו- $\text{valid bit}=1$ אז נמצא ה-Byte וזהו cache hit, אחרת (אם אין התאמה של ה-Tag או ש- $\text{valid bit}=0$) אז ה-Byte הרצוי אינו ב-cache, ויש להביאו מהזיכרון וזהו cache miss.

3) איזה בלוק מחליפים במקרה שהבלוק לא נמצא ב-Cache (miss), וצריך להביאו מהזיכרון?

כאשר המעבד מבקש לקרוא נתון מבלוק שלא נמצא ב-cache, הבלוק מובא מהזיכרון ובהנחה שה-cache מלא, יש לבחור איזה בלוק ב-cache ל"זרוק" על מנת לפנות מקום לבלוק החדש. במקרה של מיפוי ישיר יש רק בלוק אחד ויחיד ואין אפשרות בחירה (מה שמפשט את המימוש), במקרה של n way set associative, יש n בלוקים שניתן להחליף.

אסטרטגיות החלפה עיקריות:

- 1) בחירה אקראית מתוך n האפשרויות (יתרון - פשוטות מימוש).
- 2) LRU (least recently used) – "זריקת" הבלוק שלא נעשה בו שימוש הכי הרבה זמן (יתרון – miss rate נמוך יותר. חיסרון – נדרש מנגנון למעקב אחרי השימוש בבלוקים) כדי להינות מ-2 העולמות מממשים בד"כ מדיניות שהיא קירוב של מדיניות ה-LRU.

4) מה קורה בכתיבה?

קיימות 2 מדיניות כתיבה עיקריות לזיכרון המטמון:

1. **Write through** – כתיבת המידע המועדכן גם ל-cache וגם לזיכרון הראשי. ב-WT המעבד חייב להמתין עד לסיום הכתיבה של הבלוק לזיכרון הראשי. ניתן להמנע מהמתנה זו ע"י הוספת Write Buffer שמאפשר למעבד להמשיך בביצוע פקודות, בזמן שהחוצץ כותב את הבלוק לזיכרון.
2. **Write Back** – כתיבת המידע המועדכן רק לבלוק ב-cache. הבלוק המתאים לזיכרון יעודכן רק כאשר הבלוק "יזרק" מה-cache. בכדי לדעת אם בלוק ב-cache השתנה, מוסיפים לכל בלוק dirty bit, שערכו 1 כאשר הבלוק עודכן – ביט זה חוסך כתיבות לזיכרון של בלוקים שלא נעשה בהם כל שינוי.

יתרונות מדיניות ה-Write Back:

זמן הכתיבה תלוי במהירות ה-cache בלבד.

מספר כתיבות לאותו בלוק מצריכות כתיבה אחת של הבלוק לזיכרון הראשי.

יתרונות מדיניות ה-Write Through:

miss בקריאה לא גורמים לכתיבה לזיכרון הראשי של הבלוק ש"נזרק" מה-cache. קל יותר למימוש.

במקרה שהבלוק שאותו רוצים לעדכן לא נמצא ב-cache (Write miss), משתמשים באחת מהמדיניות הבאות:

1. **Write Allocate** – הבלוק הדרוש מועבר מהזיכרון הראשי ל-cache, ולאחר מכן פועלים לפי מדיניות הכתיבה שנבחרה (WB או WT).
2. **No Write Allocate** – הבלוק מועדכן בזיכרון הראשי, אך אינו מועבר ל-cache.

Index = 2, לכן נבדקים כל (2) הבלוקים שב-2 set, מאחר שה-0 valid bit יש Read miss. cache-ה שולח סיגנל stall שעוצר את ה-CPU (עד שתושלם הבאת הנתון מהזיכרון הראשי) בלוק מס' 14 (בו נמצא 117 Byte) מובא מהזיכרון הראשי ומושם באחד מ-2 הבלוקים שב-2 set. ה-0 valid bit של הבלוק ב-cache מתעדכן ל-1, וה-3 Tag מתעדכן עפ"י ה-3 Tag של הבלוק (Tag = 3) מתבצעת חזרה על תהליך הקריאה שהפעם מצליחה והבלוק מועבר ל-CPU, ממנו נלקח הבית ה-5.

אחרי הקריאה:

cache-ה:

2 way set associative

block No.	0	1	2	3	4	5	6	7
valid bit	0	0	0	0	1	0	0	0
Tag 3 bits	x	x	x	x	3	x	x	x
Block data	x	x	x	x	בלוק נדרש	x	x	x
	set 0		set 1		set 2		set 3	

Block No.	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	2	2	2	2	2	2	2	2	2	2	3	3		
	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	

