

# מערכות אחסון משני

## אמצעי אחסון משניים בעבר

Floppy disks (דיסקטים)



כרטיסים מנוקבים



CD & DVD



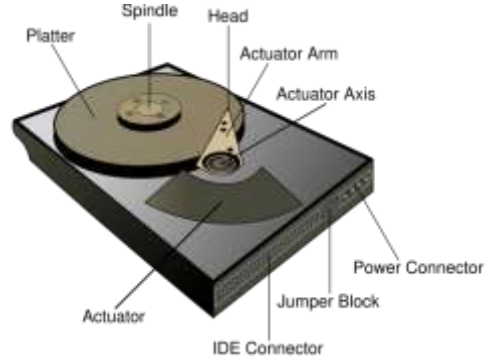
סרטים מגנטיים



סרטי נייר מנוקבים



## דיסק קשיח (Hard disk)



## Solid State Drive



## מחירי זיכרון (2012)

מהירות העברת מידע	מחיר ל-1GB	מחיר אופייני	נפח אופייני	סוג זיכרון
10 GB/Sec	30\$	120\$	4 GB	ראשי 800MHz (SDRAM)
20 MB/Sec	6\$	100\$	16 GB	Disk on key
Read 200, Write 100 MB/Sec	3\$	300\$	100 GB	Solid State Drive
100-120 MB/Sec	20 cent	100\$	500 GB	משני 7200 (hard disk) סל"ד
100 MB/Sec	10 cent	100\$	1 TB	דיסק קשיח חיצוני לגיבוי
7.8 MB/Sec	30 cent	0.2\$	700 MB	CD-ROM (52x)
10 MB/Sec	5 cent	0.25\$	4.7 GB	DVD (8x)

### 6 משטחים



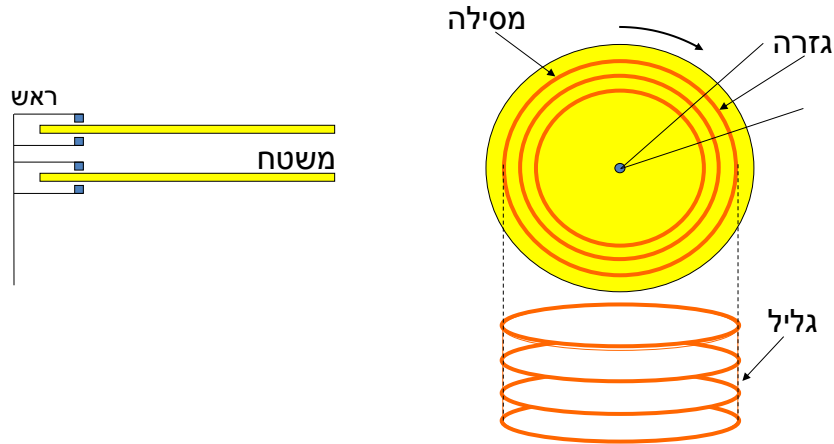
### הראש הקורא



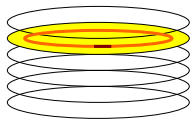
### הדיסק



## תיאור סכמתי של מבנה הדיסק

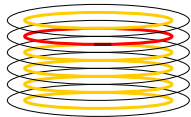


## כתובת בדיסק



**ארגון פיזי:** משטח (Plate), מסילה (Track), גזרה (Sector).

החיסרון: תמיד ישנה הזזת זרוע במעבר ממסילה למסילה.



**ארגון לוגי:** גליל (Cylinder), מסילה, גזרה.

גליל: אוסף כל המסילות בעלות רדיוס זהה. היתרון: במעבר על כל הדיסק מספר הזזות הזרוע קטן (פי מספר המשטחים).

## נתונים לדיסק אופייני (אבל ישן)

IBM Ultrastar 36z15 scsi <http://www.hgst.com/hdd/ultra/ul36z15.htm>

גודל: 4 פלטות, 8 משטחים, 10500 מסילות למשטח, 450 גזרות למסילה, 512 בתים לגזרה

קיבולת:

	512 B	גזרה:
512 B x 450 =	230,400 B = 225 KB	מסילה:
225 KB x 8 =	1800 KB = 1.8 MB	גליל:
1.8MB x 10,500 =	18900 MB = 18.9 GB	מארז:

מהירות:

סיבוב: 15,000 סל"ד = 250 סיבובים לשניה  $\Leftarrow$  4 ms לסיבוב  
 זמן תזוזה: מינימלי 0.65 ms  
 מקסימלי 6.7 ms  
 ממוצע משוקלל (יצרן) 3.4 ms

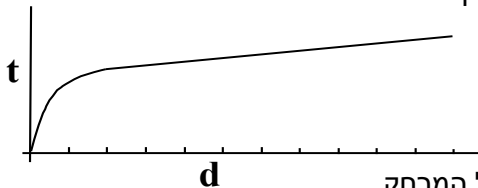
## זמן קריאת רשומה

- סכום של שלושה גורמים:
  - זמן תזוזת זרוע (seek time)
  - השהיית סיבוב (rotational delay, latency)
  - זמן קריאת גזרה - העברה (transfer time)



## I חישוב זמן התזוזה

- אם  $d$  המרחק (במסילות) בין המקום הנוכחי של הזרוע למסילת היעד.
- את רוב המרחק הזרוע עוברת במהירות קבועה ( $v$  מהירות שיוט)
- אך כדי להגיע למהירות השיוט, יש להאיץ תחילה
- גם כדי לקרוא, על הזרוע להאט לעצירה



לכן זמן התזוזה אינו פונקציה לינארית של המרחק  
נניח שזמן ההאצה למהירות שיוט הוא 1 (וכך גם זמן ההאטה)

$$t(0) = 0$$

$$t(1) = T_{\min}$$

$$t(2) = T_{\text{base}}$$

$$t(d) = (d-2) / v + T_{\text{base}} \quad (d > 2)$$

## II חישוב זמן התזוזה

דוגמה: בדיסק Ultrastar 36z15

$$V = 1800 \text{ tracks / msec}$$

$$t(1) = 0.65 \text{ msec}$$

$$t(2) = 0.87 \text{ msec}$$

$$t(10) = 8 / 1800 + 0.87 \approx 0.874 \text{ msec}$$

$$t(100) = 98 / 1800 + 0.87 \approx 0.924 \text{ msec}$$

$$t(1000) \approx 1000 / 1800 + 0.87 \approx 1.43 \text{ msec}$$

$$t(3500) \approx 3500 / 1800 + 0.87 \approx 2.8 \text{ msec}$$

זמן תזוזה מקסימלי

$$t(10500) \approx 10,500 / 1800 + 0.87 \approx 6.7 \text{ msec}$$

## מרחק תזוזה צפוי

**משפט:** נתון דיסק עם  $n$  גלילים.

אם מיקום הראש הקורא ויעד התזוזה נבחרים באופן בלתי-תלוי ובהתפלגות אחידה על פני כל הגלילים,

אז בממוצע בכל תזוזה, עוברים  $E(d) = \frac{n}{3} - O\left(\frac{1}{n}\right)$  מסילות

לדיסק בדוגמה יש 10,500 מסילות

↔ מרחק תזוזה ממוצע = 3,500 מסילות

↔ זמן תזוזה ממוצע = 2.8 msec

♣ היצרן חישב זמן תזוזה ממוצע אחרת (כנראה בעזרת ניסויים) וקיבל **3.4 msec**

## מרחק תזוזה צפוי: הוכחה

ישנם  $n$  גלילים

$P(\text{src}=i) = 1/n$  • גליל המוצא נבחר בהתפלגות אחידה:

$P(\text{dest}=i) = 1/n$  • גליל היעד נבחר בהתפלגות אחידה:

• גליל היעד אינו תלוי בגליל המוצא

$$P(\text{src}=i \wedge \text{dest}=j) = P(\text{src}=i) \cdot P(\text{dest}=i) = 1/n^2$$

$$\text{Expected distance} = \sum_d d \cdot P(d) = \sum_{1 \leq i, j \leq n} |i-j| \cdot \frac{1}{n^2} = \frac{1}{n^2} \cdot 2 \cdot \sum_{i=1}^n \sum_{j<i} j =$$

$$\frac{1}{n^2} \cdot 2 \cdot \sum_{i=1}^n \sum_{j=1}^{i-1} j = \frac{1}{n^2} \cdot 2 \cdot \sum_{i=1}^n \frac{i \cdot (i-1)}{2} = \frac{1}{n^2} \cdot \left( \sum_{i=1}^n i^2 - \sum_{i=1}^n i \right) =$$

$$\frac{1}{n^2} \cdot \left( \frac{n(n+1)(2n+1)}{6} - \frac{n(n+1)}{2} \right) = \frac{1}{n^2} \cdot \left( \frac{2n^3 - 2n}{6} \right) = \frac{n}{3} - O\left(\frac{1}{n}\right)$$

## השהיית סיבוב

15,000 סיבובים לדקה (סל"ד) = 250 סיבובים לשניה  
 $\Leftarrow$  זמן סיבוב = 4 msec  
 ממוצע השהיית הסיבוב = זמן חצי סיבוב = 2 msec



## זמן העברה

• חישוב בבתיים:  $transfer = \frac{\# \text{ bytes transferred}}{\# \text{ bytes per track}} \times (rotation \text{ time})$

• חישוב לפי גזרות (סקטורים):  
 אם יש 450 גזרות למסילה,  
 זמן קריאת גזרה =  $1 / 450$  זמן סיבוב.

עבור 15,000 סל"ד (סיבוב = 4 ms):  
 זמן העברת גזרה =  $4 / 450 = 0.0089 \text{ msec}$ .



## תרגיל כיתה

- קיבולת גליל:**  
 $512 \text{ B} \times 450 \times 8 = 1800 \text{ KB}$   
**מספר הגלילים לאחסון הקובץ:**  
 $900 \text{ MB} / 1800 \text{ KB} = 512$   
**זמן קריאת מסילה = זמן סיבוב:**  
 $4 \text{ msec}$   
**זמן קריאת גליל:**  
 $8 \times 4 \text{ msec} = 32 \text{ msec}$   
**זמן קריאת הקובץ:**  
 $512 \times 32 \text{ msec} = \text{sec } 16.4$
- קוראים באופן סידרתי קובץ של 900MB אשר מאוחסן ברציפות על גלילים סמוכים בדיסק עם:
- 10,500 גלילים
  - 8 מסילות בגליל
  - 450 גזרות למסילה
  - 512 B לגזרה
  - זמן תזוזת הזרוע:
  - מינימלית 0.65 msec
  - ממוצעת 3.4 msec
  - מהירות הסיבוב: 15,000 סל"ד

## תרגיל כיתה (המשך)

- בתחילת הקריאה, הגעה לגליל הראשון של הקובץ
  - תזוזת זרוע ממוצעת 3.4 msec
  - הגעה לתחילת המסילה (חצי סיבוב): 2 msec בממוצע
- בסיום קריאת גליל, במעבר לגליל הבא
  - תזוזת זרוע מינימלית 0.65 msec
  - הגעה לתחילת המסילה (חצי סיבוב): 2 msec בממוצע
- עבור 512 גלילים יש 511 מעברים שידרשו:  $511 \times (2 + 0.65) = 1354 \text{ msec}$
- סה"כ ל- 16.4 שחושבו בשקף הקודם יש להוסיף 1.354 שניות עבור מעברים בין גלילים.

## ואם קוראים מסילות אקראיות?

לכל מסילה יש לבצע:

$$(512/3-2)/1800+0.87=0.96 \text{ msec}$$

- תזוזה אקראית (בתוך 512/10500 מהדיסק)
- 2 msec השהיית סיבוב אקראית
- 4 msec קריאה
- 6.96 msec סה"כ

קריאת כל הקובץ:

- בקובץ  $512 \times 8 = 4096$  מסילות.
- סה"כ:  $4096 \times 6.96 \text{ msec} = 28.5 \text{ sec}$

## הכי גרוע: קריאה של גזרות אקראיות

בכל מסילה יש 450 גזרות, ועבור כל-אחת צריך לבצע:

- 0.96 msec תזוזה אקראית
- בתוך 512/10500 מהדיסק
- 2 msec השהיית סיבוב אקראית
- 0.0089 msec קריאה
- 2.9689 msec סה"כ



קריאת כל הקובץ:

- בקובץ  $450 \times 4,096 = 1,843,200$  גזרות.
- סה"כ:

$$1,843,200 \times 2.9689 \text{ msec} = 5,472 \text{ sec} = 91 \text{ min} = 1:31 \text{ hours}$$

## כתיבת רשומה בגזרה לא ריקה

- אפשר לכתוב רק גזרה שלמה.
- הפעולות שיש לבצע כאשר תוכן הגזרה לא נמצא בזיכרון:
  - חישוב כתובת היעד
  - קריאת גזרה לחוצץ
  - עדכון החוצץ
  - המתנה עד שהגזרה שוב מול הראש הקורא/כותב
  - כתיבה (העברה)
- סה"כ: כמו קריאה + סיבוב שלם נוסף.

## נתונים לדיסק אופייני (2012)

Seagate Savvio 10K.5 SAS <http://www.seagate.com/internal-hard-drives/enterprise-hard-drives/>

גודל: ~~184,799~~ ~~6~~ ~~3~~  
~~10500~~ פלטות, ~~8~~ משטחים, ~~450~~ מסילות למשטח,  
~~450~~ גזרות למסילה, 512 בתים לגזרה  
**1643 בממוצע**

קיבולת:  
 גזרה: 512 B  
 מסילה: ~~512 B x 450 = 230,400 B = 225 KB~~  
 גליל: ~~225 KB x 8 = 1800 KB = 1.8 MB~~  
 מארז: ~~1.8MB x 10,500 = 18900 MB = 18.9 GB~~  
**900 מפורמט**

מהירות:  
 סיבוב: ~~15,000~~ סל"ד = 250 סיבובים לשניה ~~6~~ ms לסיבוב  
 זמן תזוזה: מינימלי ~~0.2~~ ms  
 מקסימלי ~~6.7~~ ms  
 ממוצע משוקלל (יצרן) ~~3.4~~ ms **3.7**



## ? Old Wisdom

- זיכרון ראשי קטן ומהיר ⇔ זיכרון משני ענק ואיטי
- קריאה של בלוקים באופן רציף מהירה בהרבה
- זמן כתיבה ≈ זמן קריאה
- שיפור מועט בביצועים עם השנים



## (Solid State Drives) SSD

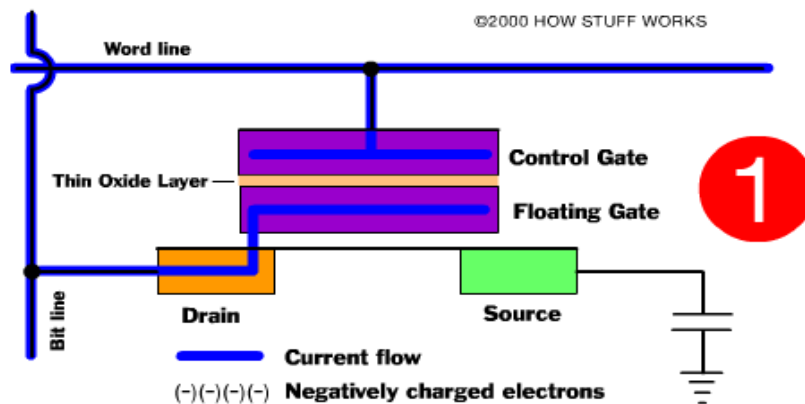
- אמצעי זיכרון משני בלא חלקים נעים
- עדיין לא "בשלים" להחליף את הדיסק הקשיח
- DRAM דורשים מתח מתמיד ולכן לא מתאימים לגיבוי או הפצה
- flash drive, כמו disk on key, יקרים ואיטיים.

## יתרונות וחסרונות של SSD

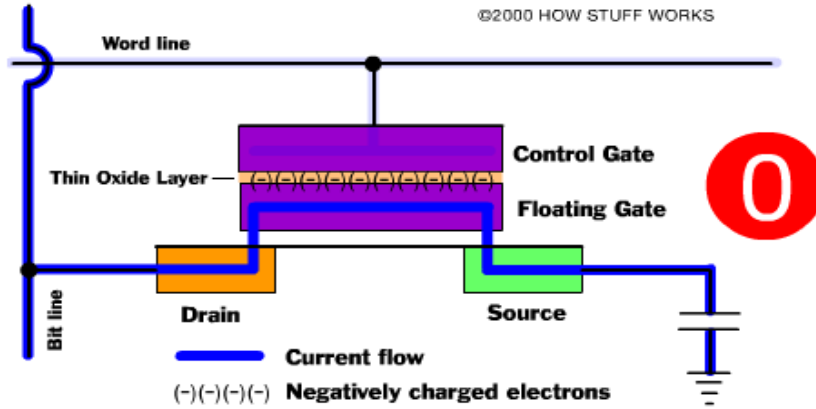
- ✗ עדכון דורש "מחיקה" של ערכים ולכן כתיבה איטית מהקריאה
- ✗ מספר כתיבות מוגבל לכל תא זיכרון
  - צריך לפזר את השימוש על פני הדיסק
- ✗ העברת נתונים איטית
- ✗ מחיר יקר ליחידת אחסון
- ✗ עדיין לא מותאם לנפחי מידע גדולים
- ✗ קריסת דיסק היא לרוב טוטאלית (יותר קשה לשחזר מידע)

- ✓ אתחול מהיר
- ✓ גישה מהירה לרשומה
- ✓ לא מתחמם
- ✓ צריכת אנרגיה נמוכה (ב-flash)
- ✓ עמידות בפני זעזועים פיסיים
- ✓ זמן הקריאה של קובץ השמור באופן רציף (סדרתי) כמעט זהה לקריאת הקובץ כאשר רשומותיו מפוזרות על פני הדיסק
- ✓ פחות שגיאות בעת קריאת הנתונים

## תא זיכרון Flash מורכב משני טרנזיסטורים

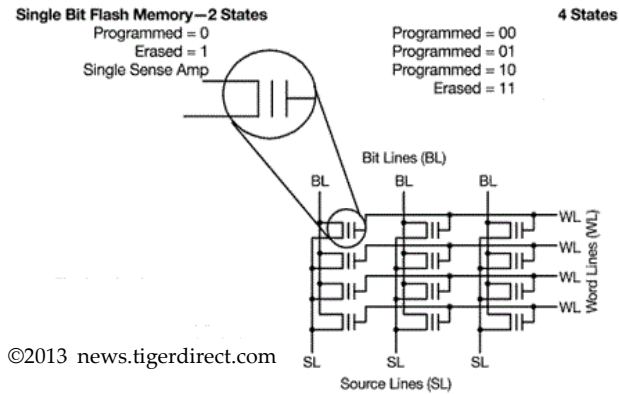


## כתיבה היא שינוי מ-1 ל-0, ושינוי הפוך דורש "מחיקה"

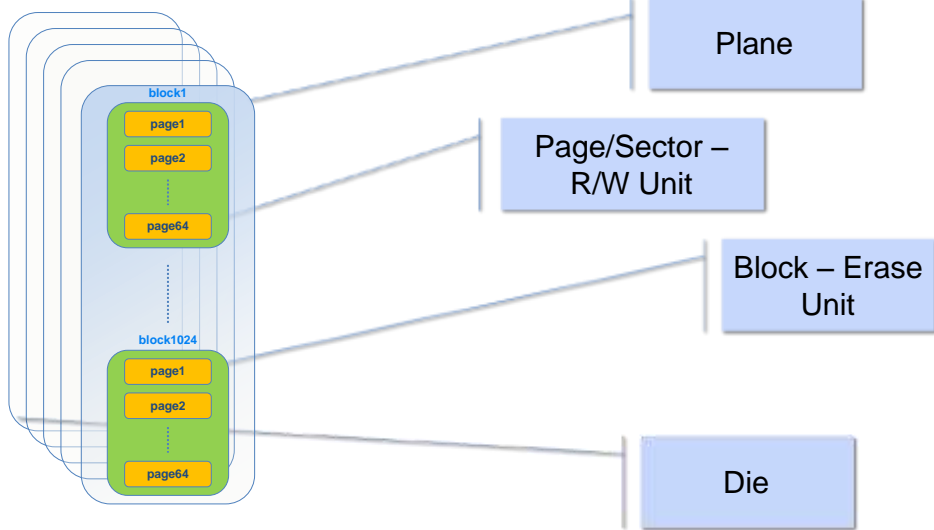


## גישה לתא

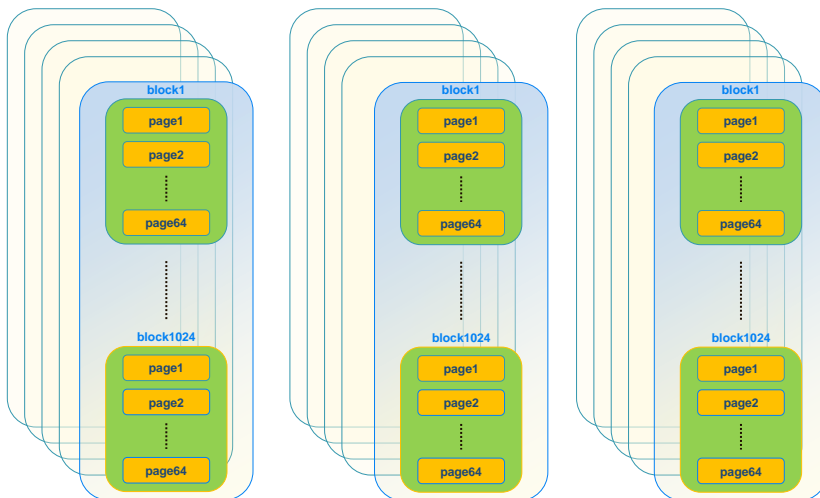
- קומבינציית המתח בין ה-word line ל-bit line מגדירה את הפעולה על התא (כתיבה, קריאה, מחיקה)



# מבנה הדיסק



# מבנה הדיסק





יחידת  
מחיקה

## פרמטרים לדוגמה

ניתן לקרוא ממשטחים שונים במקביל.

- איך זה משפיע על הביצועים?
- איך זה משפיע על ארגון קבצים?

- דף (page): 4KB
- בלוק (block): 128 דפים
- משטח (plane): 2048 בלוקים
- תבנית (die): 4 משטחים
- תבנית מכילה: 4 GB
- מארז (package): 2 תבניות
- דיסק: 10 מארזים

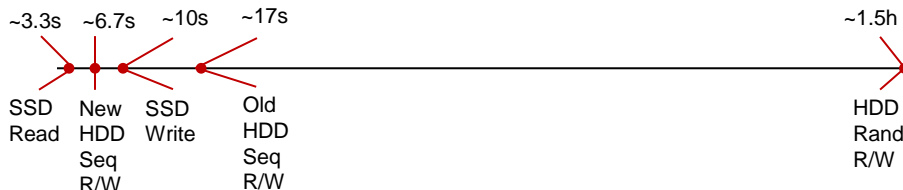
## פרמטרים ל- Intel SSD 320

- השהיית קריאה: 75 מיקרו-שניות (0.075 מילי-שניה)
- קריאה רציפה: 270 MB/Sec
- השהיית כתיבה: 90 מיקרו-שניות
- כתיבה רציפה: 90 MB/Sec
- (זמנים טיפוסיים, על פי מדידות יצרן)



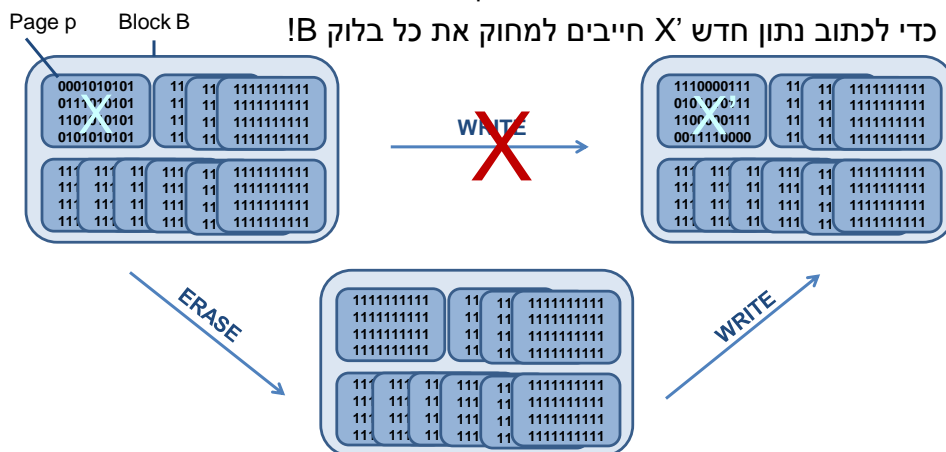
# דוגמה של קריאת/כתיבת קובץ מ/ל SSD

- קריאת 900 MB:
- זמן קריאה:  $900/270 = 3.33 \text{ Sec}$
- אין הבדל בין קובץ רציף לקובץ לא רציף
- כתיבת 900 MB:
- זמן כתיבה:  $900/90 = 10 \text{ Sec}$



# מימוש כתיבות: הבעיה עם עדכון

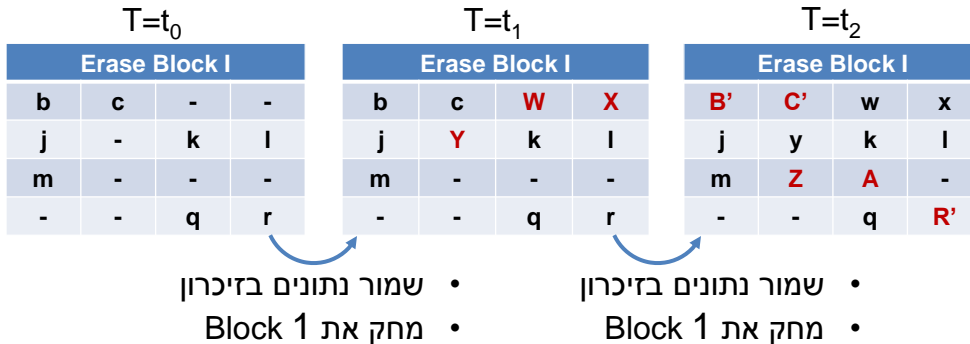
- בתוך בלוק נתונים B נמצא דף p ועליו כתוב נתון X – נתעלם כרגע משאר הדפים, ונניח שהם ריקים
- כדי לכתוב נתון חדש X' חייבים למחוק את כל בלוק B!



## מימוש כתיבות: נסיון 1 (לא יעיל)

Read/Erase/Modify/Write

- $T=t_0$ : מצב התחלתי
- $T=t_1$ : כתיבה של W, X, Y
- $T=t_2$ : כתיבה של Z, A, B', C', R'



## Read/Erase/Modify/Write

- לא יעיל
- כתיבה אטית
- יש לחכות למחיקת בלוק שלם בכל שינוי/כתיבה (לדף משומש)
- "מקצר" את חיי הדיסק
- כתיבות חוזרות לאותו בלוק

הדוגמא מבוססת על

NAND Flash Solid State Storage: Performance and Capability – an In-depth Look  
Jonathan Thatcher © 2009 SNIA

## מימוש כתיבות: נסיון 2 (לא יעיל)

Read/Modify/Write

- $T=t_0$ : מצב התחלתי
- $T=t_1$ : כתיבה של W, X, Y
- $T=t_2$ : כתיבה של Z, A, B', C', R'

$T=t_0$				$T=t_1$				$T=t_2$			
Erase Block 1				Erase Block 2				Erase Block 3			
b	c	-	-	b	c	W	X	B'	C'	w	x
j	-	k	l	j	Y	k	l	j	y	k	l
m	-	-	-	m	-	-	-	m	Z	A	-
-	-	q	r	-	-	q	r	-	-	q	R'

- שמור נתונים בזיכרון
- מחק את Block 1
- שמור נתונים בזיכרון
- מחק את Block 2

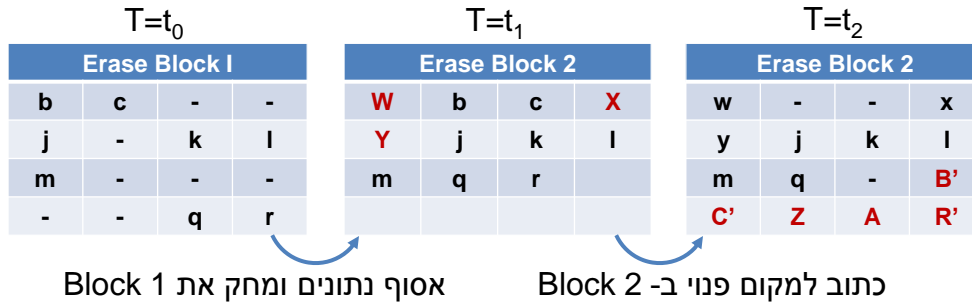
## Read/Modify/Write

- הנחה: Block 2 ו-Block 3 נמחקו מבעוד מועד
- יעיל יותר מהדוגמא הקודמת
  - אין צורך לחכות למחיקה לשם השלמת הכתיבה
- פיזור אחיד (wear leveling)
- נתונים שמעודכנים בתדירות גבוהה מפוזרים על פני הדיסק
- כתיבות מוגברות (write amplification)
  - 8 כתיבות ביוזמת המשתמש (25%)
  - 24 כתיבות תקורה בשל העברת הנתונים (75%)

## מימוש כתיבות: נסיון 3

### Garbage Collection

- $T=t_0$ : מצב התחלתי
- $T=t_1$ : כתיבה של W,X,Y
- $T=t_2$ : כתיבה של Z,A,B',C',R'



## Garbage Collection

- מתבצע ברקע
  - התדירות נקבעת על פי האלגוריתם
  - תדירות גבוהה ← כתיבות מהירות, חיי הדיסק מתקצרים
  - תדירות נמוכה ← כתיבות עלולות להתעכב, חסכון בכתיבות
  - ניתן לכוון כתיבות לבלוקים שהיו פחות בשימוש
- שיפור בתקורה
  - 8 כתיבות ביוזמת המשתמש (50%)
  - 8 כתיבות תקורה בשל העברת הנתונים (50%)