

Information in Storage Devices  
049063 – EE Department, Technion

# LECTURE 3: SSD ACCESS

# HDD and SSD



- Has been around forever
- Improves, but looks the same
- Predictable performance
- Fast to respond
- Heavily hyped
  - High media exposure
- You know can do wonders
  - But most encounters less exciting

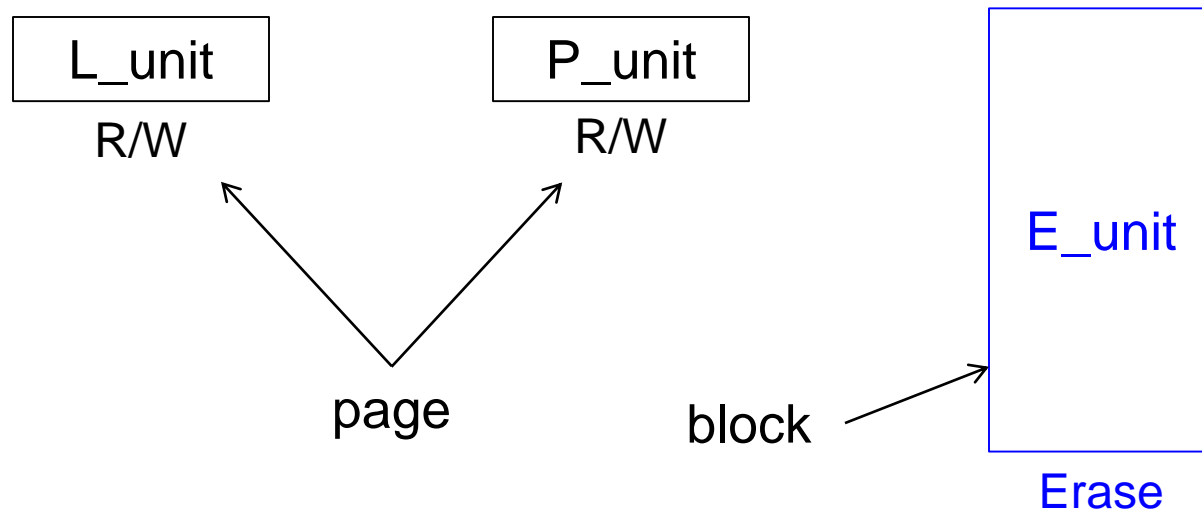
# Solid-State Drive (SSD)

- Silicon-based array of memory cells
  - with standard interface (HDD replacement)
- Invented late '90s (M-Systems, Israel)
  - Uses NAND flash for maximal density
- More expensive, but much faster
- Capacity scales by “Moore’s law”



# Flash: No Random-Access Erase

- New: Erase unit



- Physical erase of cells: only full blocks
- Write → program/erase

# No In-Place Updates

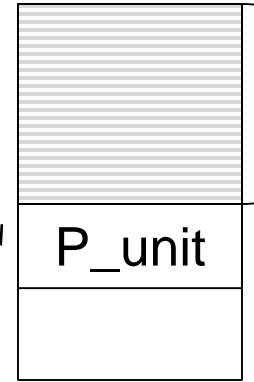
W

L\_write:

L\_unit

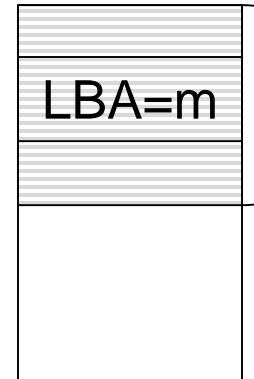
P\_write

E\_unit



used

E\_unit



used

W

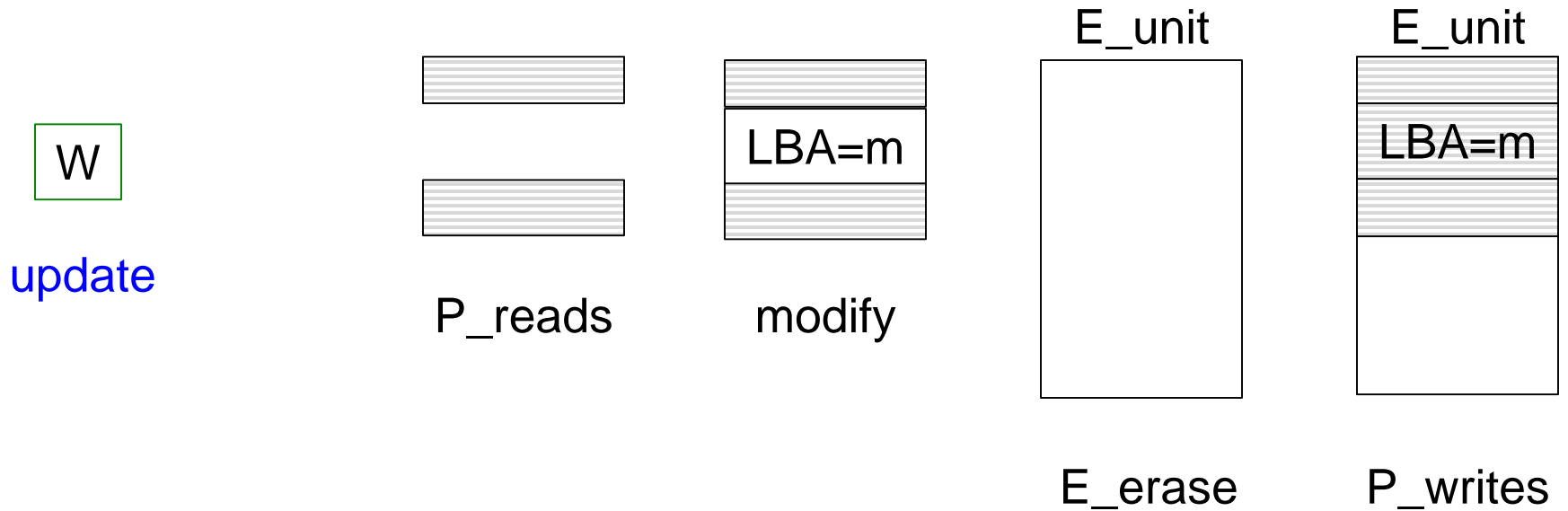
L\_write:

LBA=m

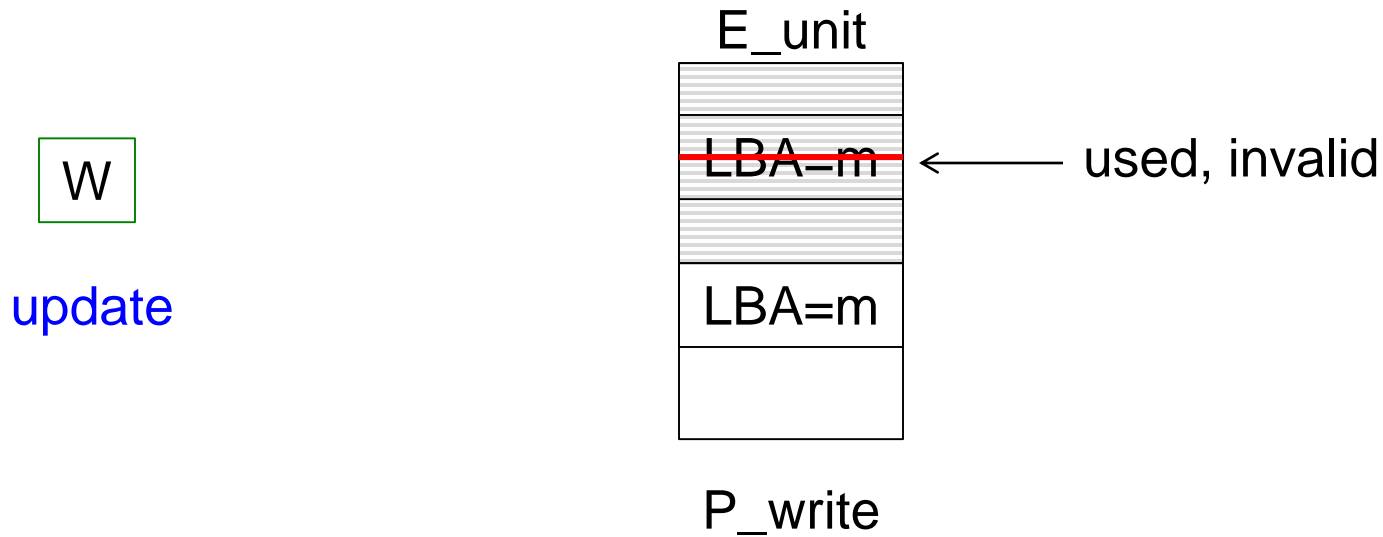
P\_write?

update

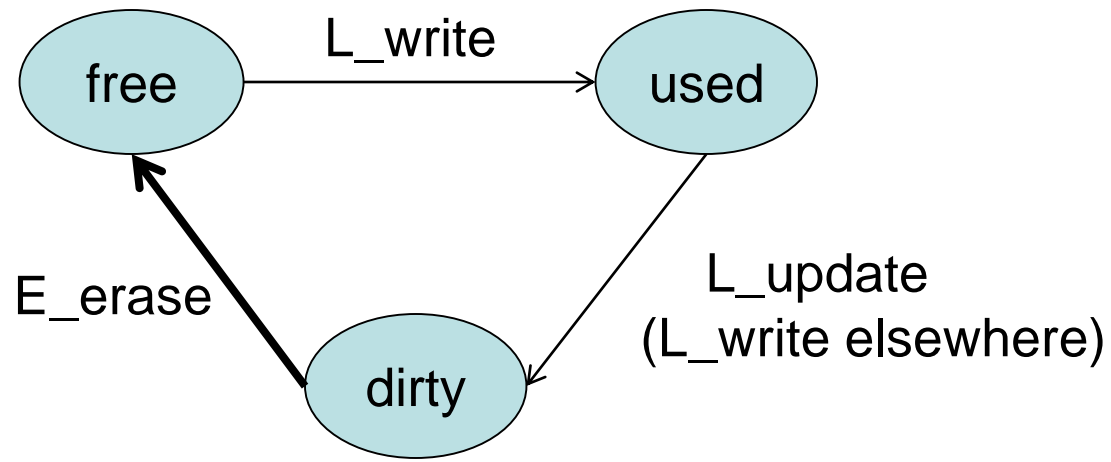
# Option 1: RMEW



# Option 2: Invalidation



# Flash State Diagram





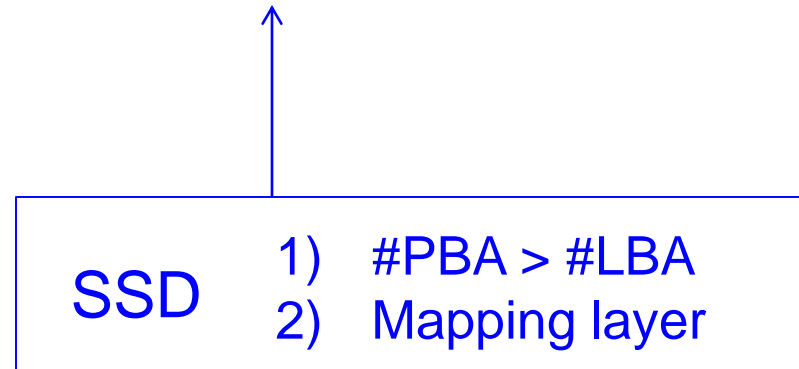
# Issues

## Option 1: RMEW

- Time
- Wear

## Option 2: Invalidation

- Over-provisioning
- Indirection



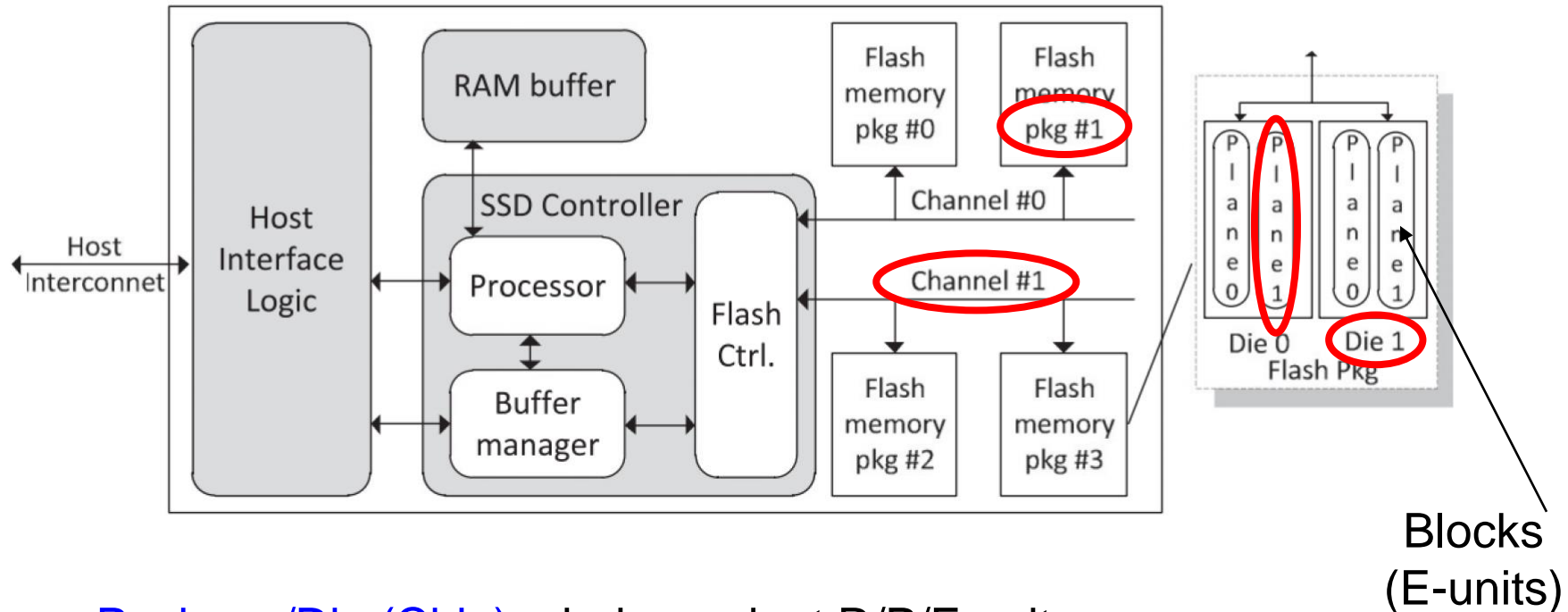
# SSD Performance vs. Technology

- **SLC**: 1 bit/cell
  - “single-level cell”
- **MLC**: 2 bits/cell
  - “multi-level cell”
- **TLC**: 3 bits/cell
  - “triple-level cell”

TABLE I  
READ/PROGRAM/ERASE TIMES FOR SLC/MLC/TLC FLASH CHIPS [10].

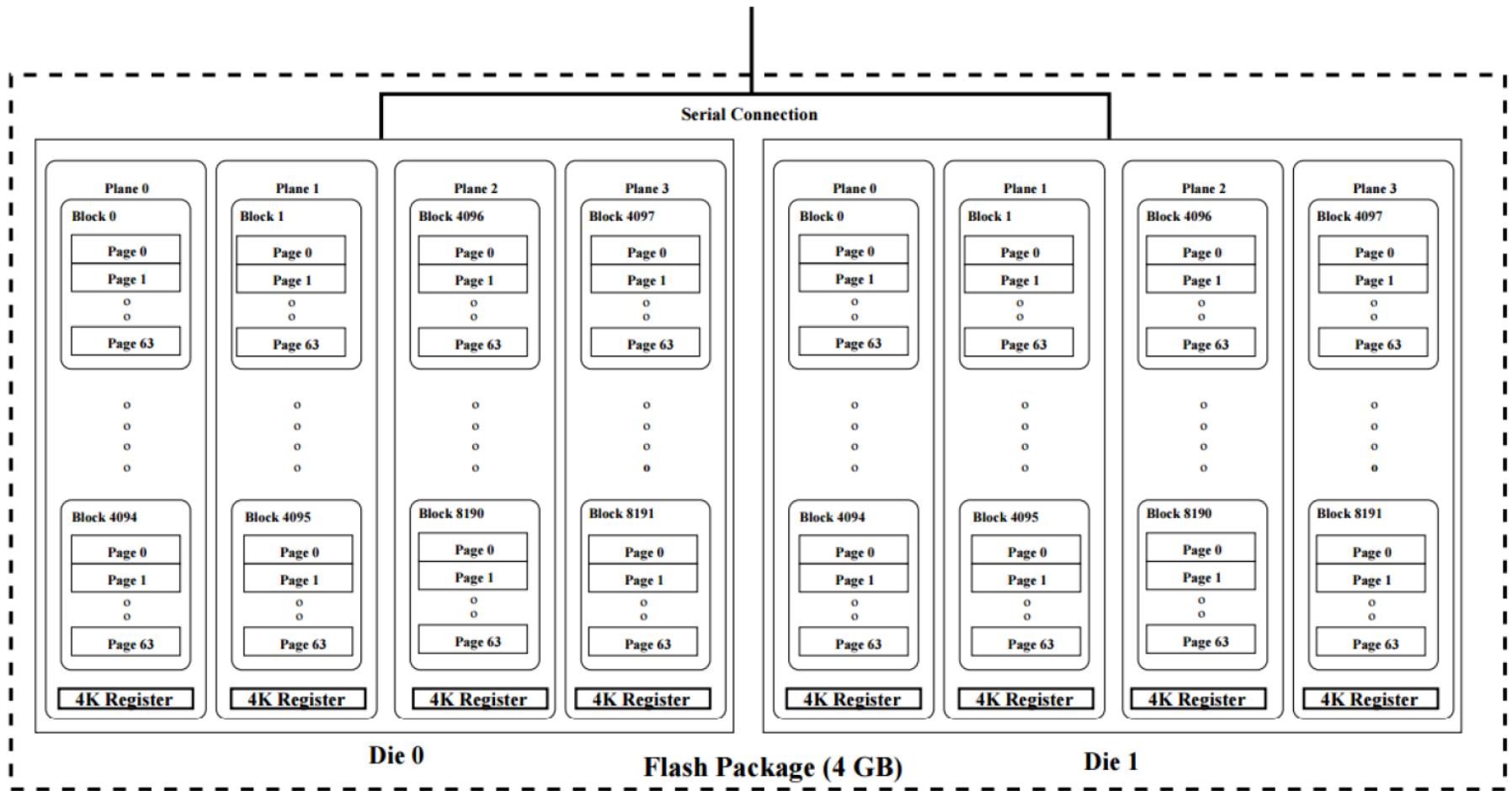
Operations	SLC chips	MLC chips	TLC chips
Random Read	20us	40us	80us
Program	260us	900us	2.3ms
Erase	2ms per block	2ms per block	10ms per block

# SSD Architecture and parallelization



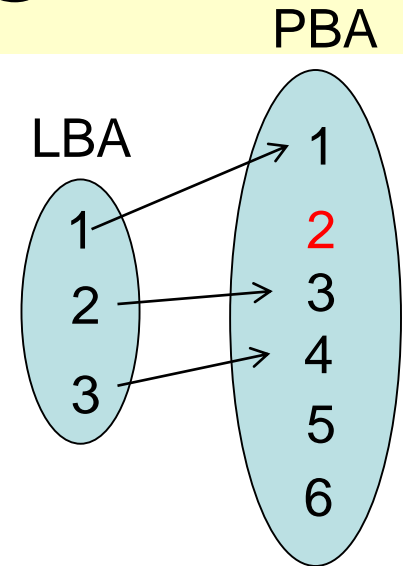
- **Package/Die (Chip)** – independent R/P/E units
- **Plane** – multiple planes of the same die can perform the same command (R/P/E) in parallel
- **Channel** – independent transfers between controller and packages
- **Parallelization:** perform commands in parallel
- **Interleaving:** Use channel for one package while the other is busy

# Flash Package Example

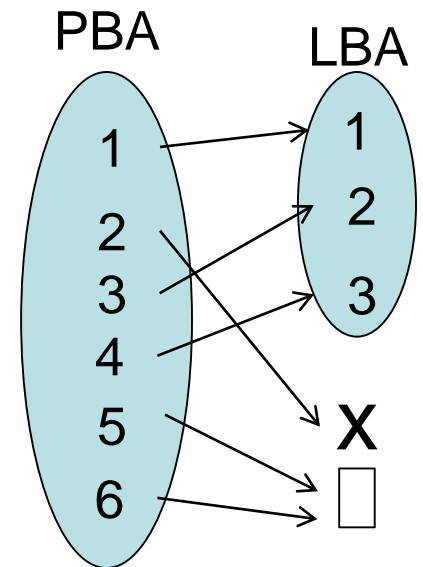


# Address Mapping

- Direct map (LBA→PBA)
  - Injective (1-1)

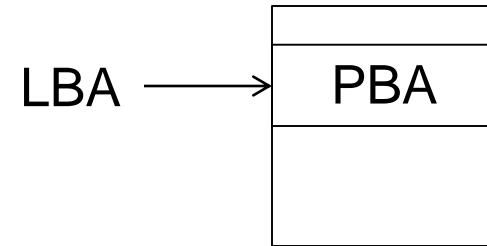


- Inverse map (PBA→LBA)
  - Range includes LBAs and
  - x: dirty
  - : free

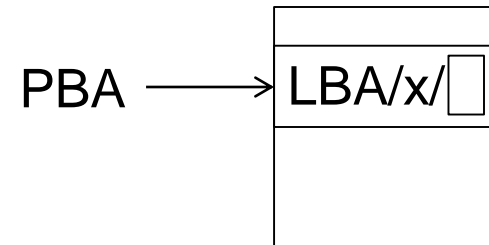


# Flash Mapping Layer

- Direct map
  - Inquired on read
  - Updated on write



- Inverse map
  - Find free PBA
  - Updated on write
  - Inquired for “clean” operations



# Clean Operations

- Reclaim dirty pages
- Also called **garbage collection**

## Procedure:

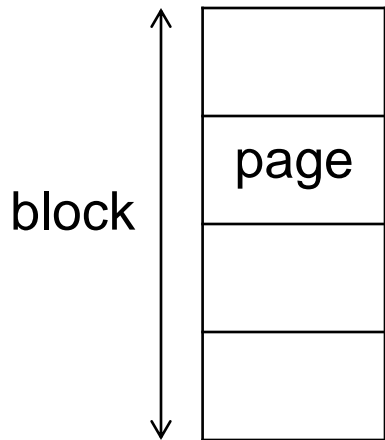
- 1) Choose an E\_unit (how?)
- 2) **Copy** all used L\_units to other E\_unit(s)
- 3) Erase E\_unit

## Objectives:

- 1) Minimize copy operations
- 2) Level the wear of E\_units

# Write Amplification (WA)

- SSD with 1 spare block



<b>0</b>	<b>4</b>	<b>8</b>	<b>12</b>	
<b>1</b>	<b>5</b>	<b>9</b>	<b>13</b>	
<b>2</b>	<b>6</b>	<b>10</b>	<b>14</b>	
<b>3</b>	<b>7</b>	<b>11</b>	<b>15</b>	



# Write Amplification (WA)

- First 4 writes

- Example:

- Writes: 0,4,8,12,

<del>0</del>	<del>4</del>	<del>8</del>	<del>12</del>	0
1	5	9	13	4
2	6	10	14	8
3	7	11	15	12

# Write Amplification (WA)

- Incoming write: LBA 1

- Example:

- Writes: 0,4,8,12,1

<del>0</del>	<del>4</del>	<del>8</del>	<del>12</del>	0
1	5	9	13	4
2	6	10	14	8
3	7	11	15	12

Choose for cleanup

Full!

# Write Amplification (WA)

- Copied 2 valid pages in the cleaned block

- Example:

- Writes: 0,4,8,12,1

<b>2</b>	<del>4</del>	<del>8</del>	<del>12</del>	0
<b>3</b>	5	9	13	4
<b>1</b>	6	10	14	8
	7	11	15	12

# WA Analysis

- Definition – Write Amplification:

The ratio of the total number of internal writes to the number of externally-requested writes.

- Notation:

Write Amplification  $WA$ ,  $WA \geq 1$

- Notes

1. Large is bad (more time, wear)
2. Depends on mapping, workload

# WA Analysis

- More notation:

T: # physical E\_units

$N_p$ : # P\_units per E\_unit

U: #L\_units/ $N_p$

$UN_p$  L\_units stored on  $TN_p$  P\_units

$S_f$ : spare factor      $S_f = \frac{T-U}{T}$       $0 \leq S_f < 1$

$\rho$ : over-provisioning      $\rho = \frac{T}{U} - 1$       $\rho \geq 0$

# Effect of Spare

<del>0</del>	<del>4</del>	<del>8</del>	<del>12</del>	0
1	5	9	13	4
2	6	10	14	8
3	7	11	15	12

Little spare

<del>0</del>	<del>4</del>	<del>8</del>	<del>12</del>	0	5
1	<del>5</del>	<del>9</del>	13	4	14
2	6	<del>10</del>	<del>14</del>	8	10
3	7	11	15	12	9

More spare

# Clean (Garbage Collection) Policies

1. LRU – least-recently written

- Pick oldest in the “log” of E\_units



2. Greedy

- Pick the E\_unit with max # dirty P\_units
- Also called min-valids

# Evaluation with Traces

BLOCK I/O TRACES [GO](#)

NFS TRACES [GO](#)

SSSI WIOCP METRICS [GO](#)

SYSTEM CALL TRACES [GO](#)

HISTORICAL SECTION [GO](#)

TOOLS [GO](#)

## MSR Cambridge Traces

The following traces are free to download under the terms of the [SNIA Trace Data Files Download License](#). Please note that cookies must be enabled within your browser in order to download traces.

For questions about downloading using a shell script, see [Using Shell Scripts](#), and for more information about downloading using a Windows batch script, see [Using Batch Scripts](#).

[View Additional Columns](#)

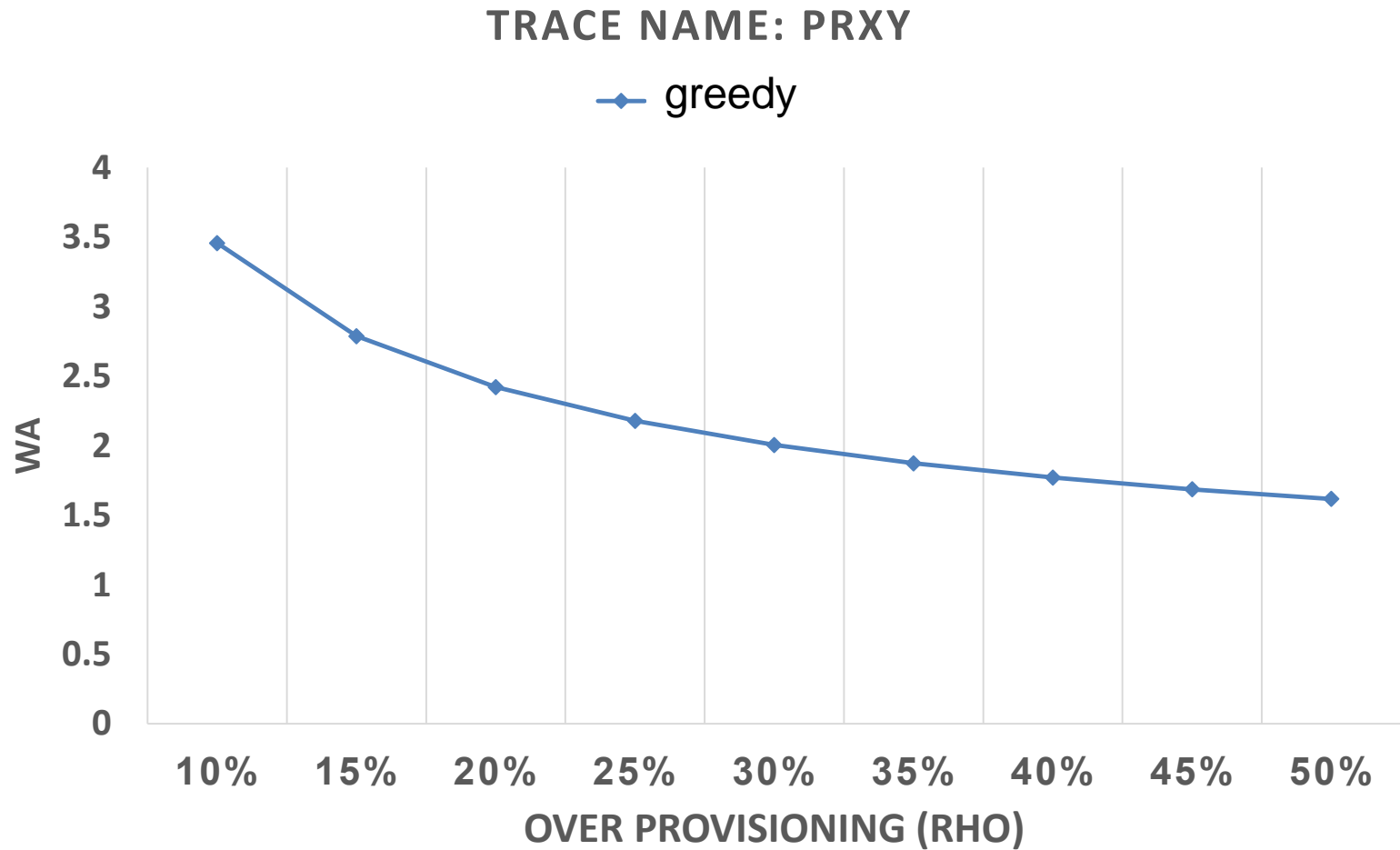
Trace Name	Details	Actions
MSR Cambridge Traces 1	1-week block I/O traces of enterprise servers at Microsoft <a href="#">Long Description</a>	<input type="text" value="Download Sample"/>
MSR Cambridge Traces 2	1-week block I/O traces of enterprise servers at Microsoft <a href="#">Long Description</a>	<input type="text" value="Download Readme"/>

I WANT TO

Choose One

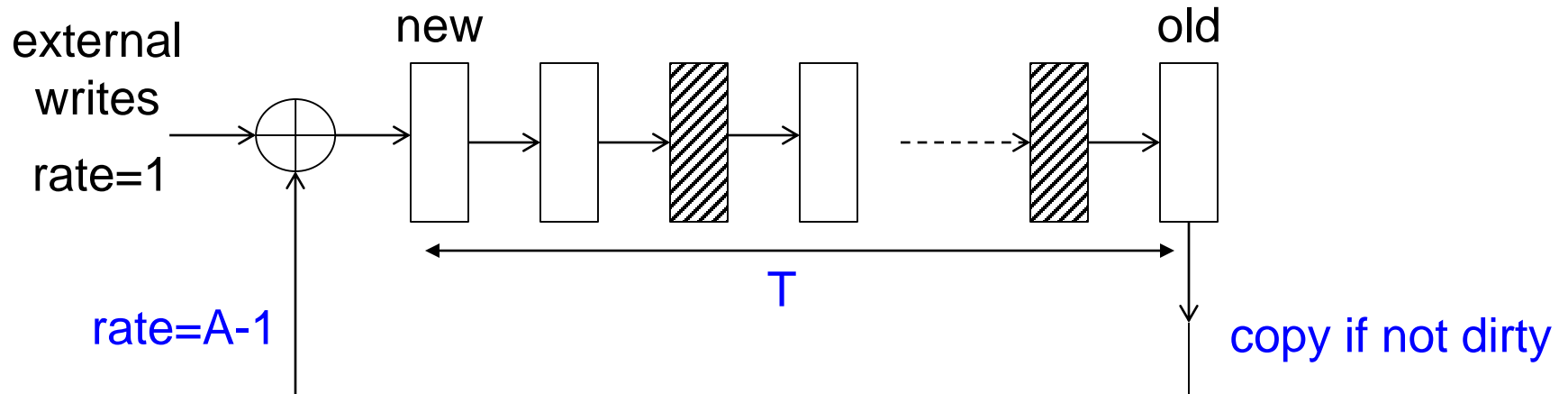


# Trace Results - Greedy



# LRU Cleaning - Model

- Clean the oldest in the E\_unit chain
- Special case:  $N_p=1$

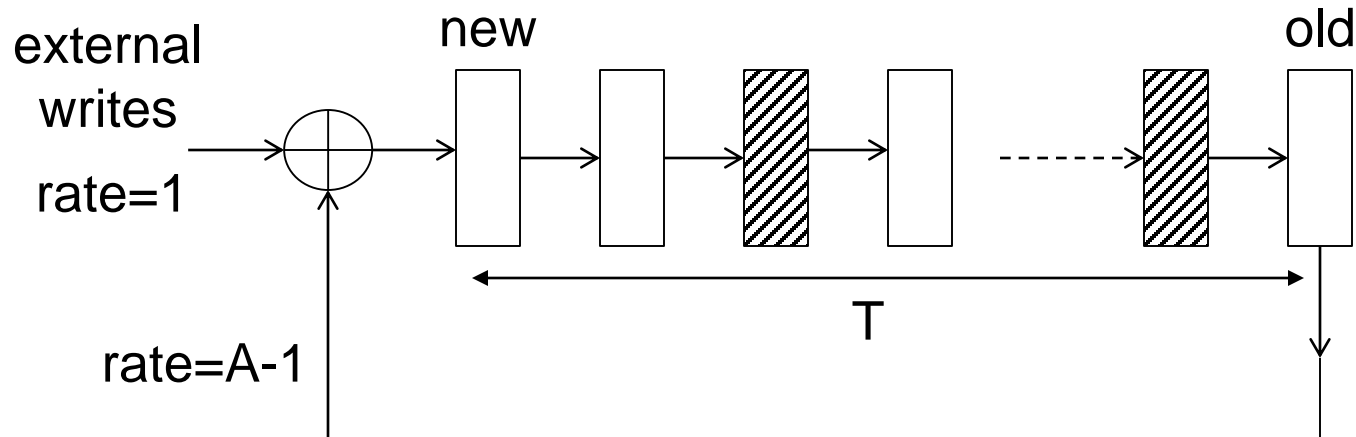


## Facts:

- T units in chain, U of them are valid
- A copy of a unit happens if **still** valid when oldest

# LRU Cleaning - Analysis

- Special case:  $N_p=1$
- Uniform workload (random write)



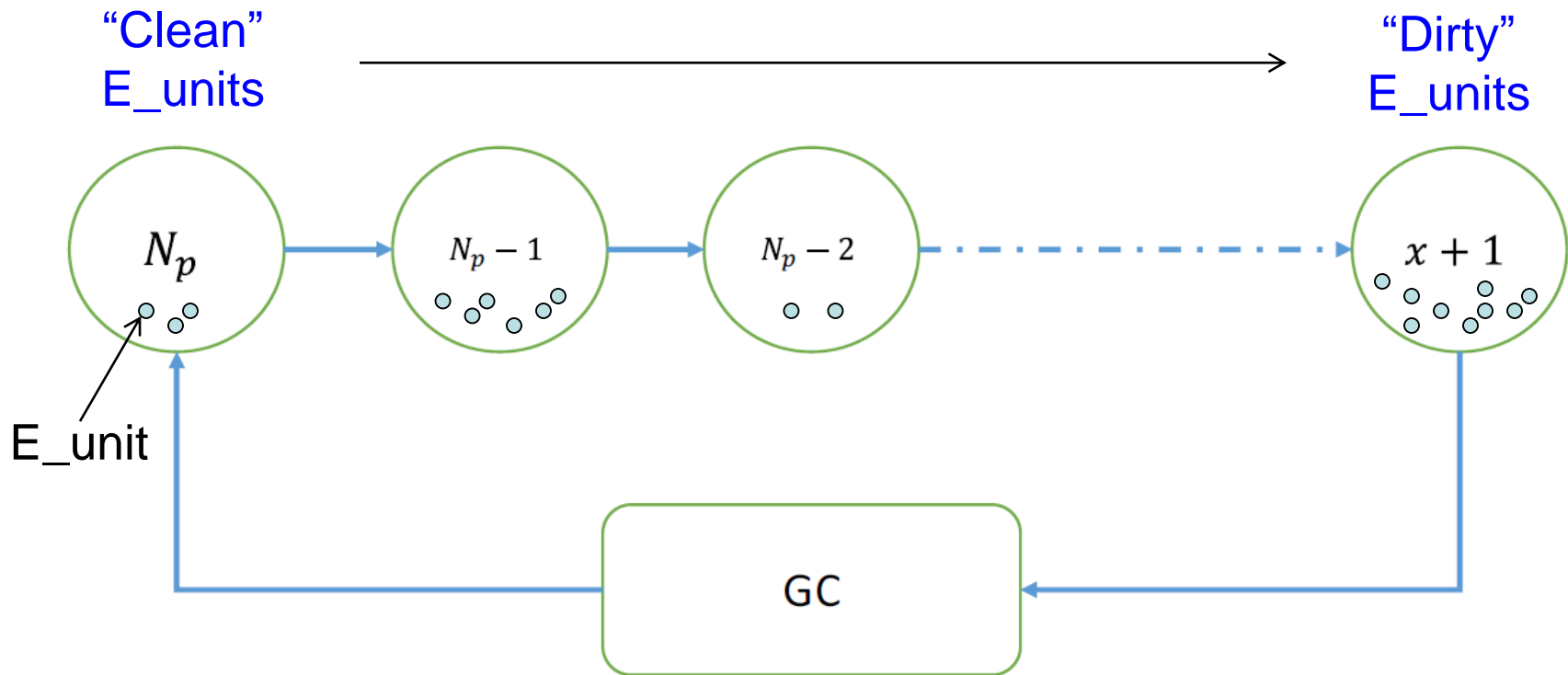
## Facts:

- $\Pr\{\text{host write is addressed to valid } E\_unit\ i\} = 1/U$ , for all  $i$
- # host writes in a full cycle  **$M=T/A$**
- The same  $L\_unit$  can cycle many times without host update
- Find  $A$  (in expectation)

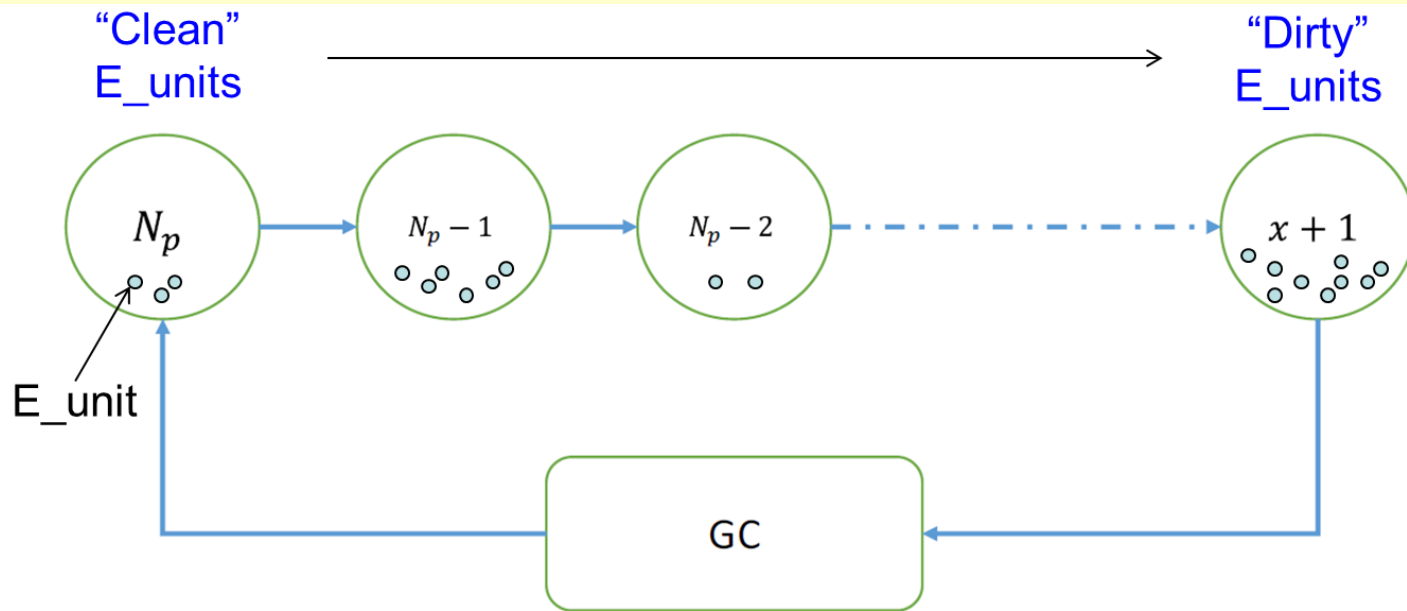


# Greedy Cleaning - Model

- Clean the E\_unit with **min # valid** (equiv. max # dirty)
- General  $N_p$ , uniform workload
- Need to calculate the expected min # valid **x**
- Markov model: states represent **# valid in E\_unit**



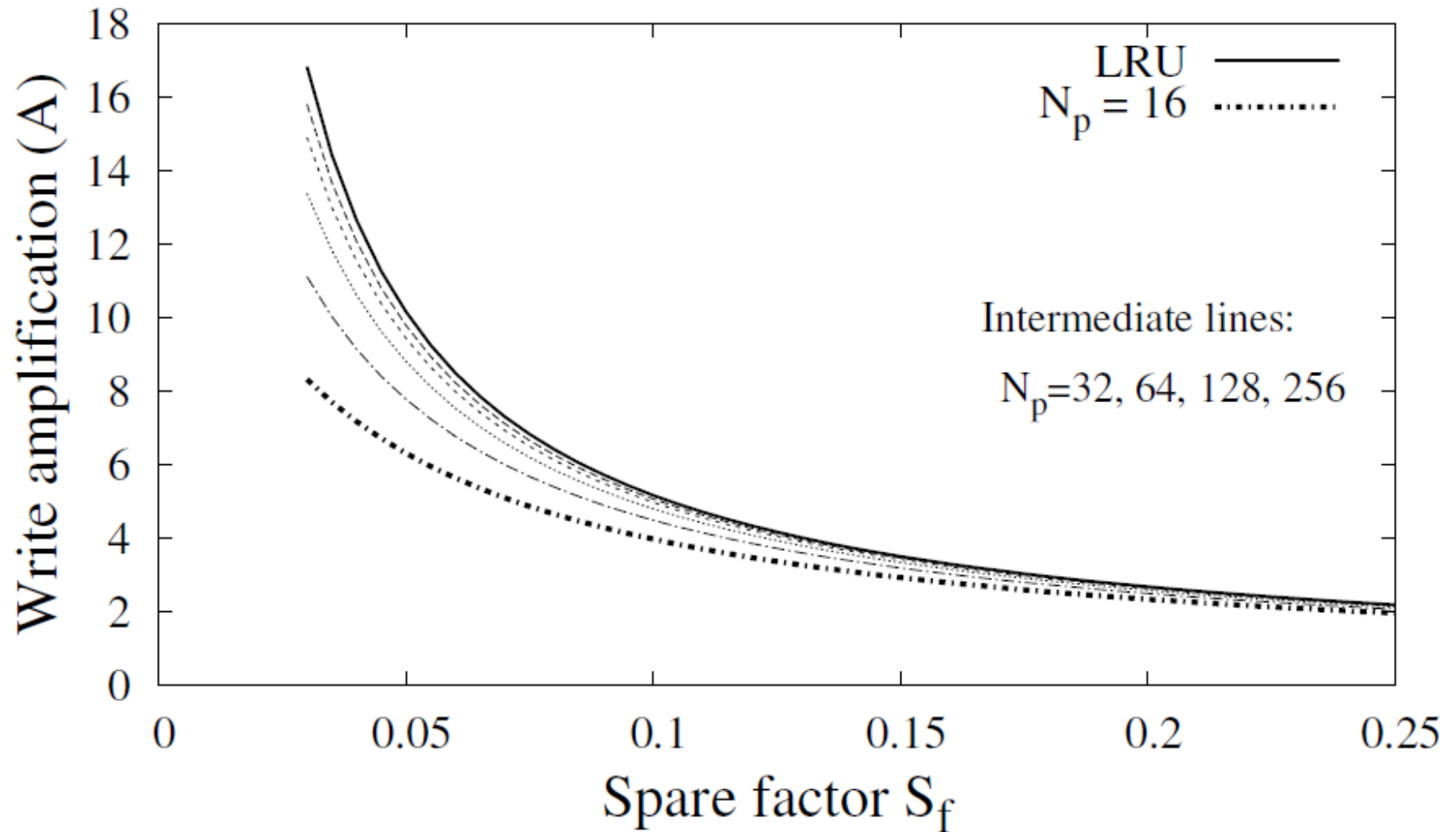
# The Greedy Markov Model



- Define fraction of E\_units in state  $i = f_i$
- Total number of valid L\_units in E\_units at state  $i = iT f_i$
- Random writes  $\rightarrow$  Transition rate from state  $i$  to  $i-1 = \frac{iT f_i}{UN_p}$
- GC transition rate  $= \frac{1}{N_p - x}$
- Find A from steady-state considerations



# Analytic Results



# Summary

- WA reduces with amount of spare
- LRU GC: simple
- Greedy GC: optimal
- Greedy approaches LRU with large  $N_p$
- Another advantage of LRU:

Wear leveling