# Write-Optimized Dynamic Hashing for Persistent Memory

*Moohyeon Nam*[+], *Hokeun Cha*[δ], *Young - ri Choi*[+], *Sam H. Noh*[+], *Beomseok Nam*[δ]

UNIST (*Ulsan National Institute of Science and Technology*)[+]
Sungkyunkwan University (*SKKU*)[δ]

# Outline

## Background

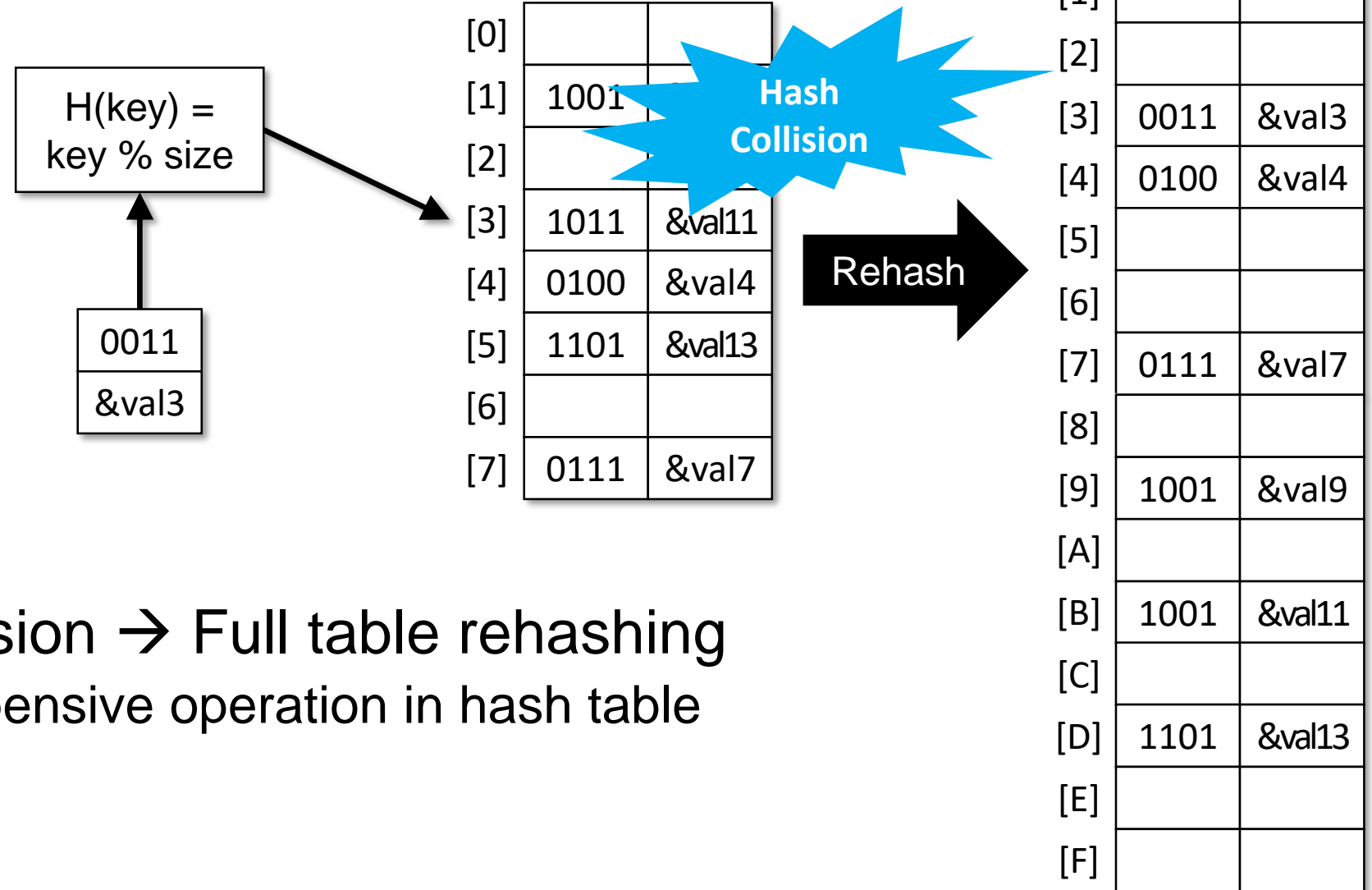- Static Hashing
- Extendible Hashing
- Persistent Memory

## Cacheline-Conscious Extendible Hashing

- Challenges and Contributions
- 3-Level Structure of CCEH
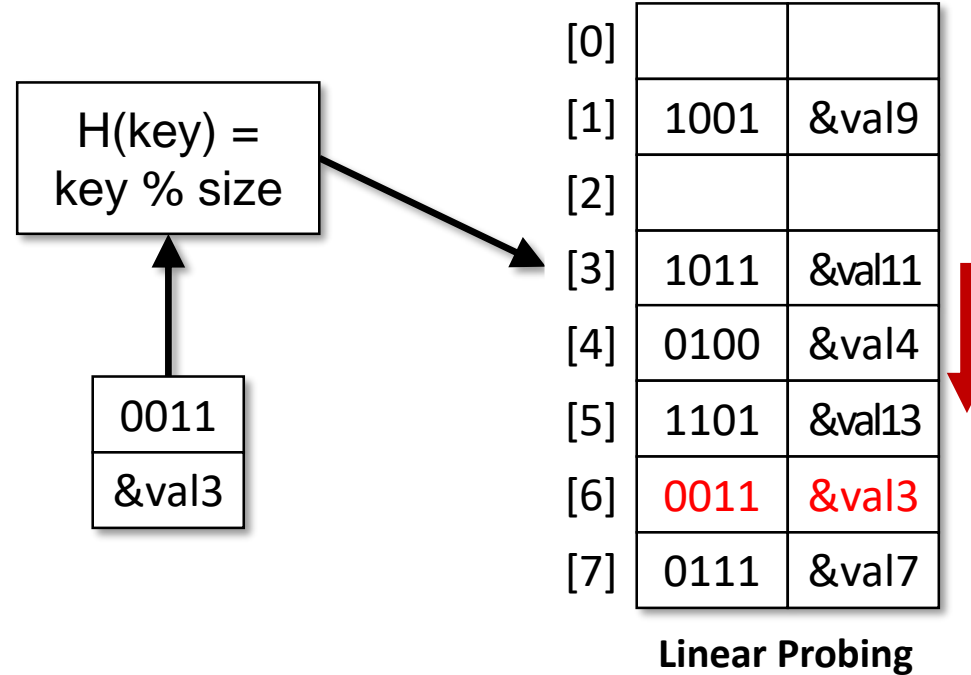- Failure-atomic Directory Update

## Evaluation

## Conclusion

| | | |
|---|---|---|
| [0] | | |
| [1] | 1001 | |
| [2] | | |
| [3] | 1011 | &val11 |
| [4] | 0100 | &val4 |
| [5] | 1101 | &val13 |
| [6] | | |
| [7] | 0111 | &val7 |

H(key) = key % size

0011
&val3

**Hash Collision**

Rehash

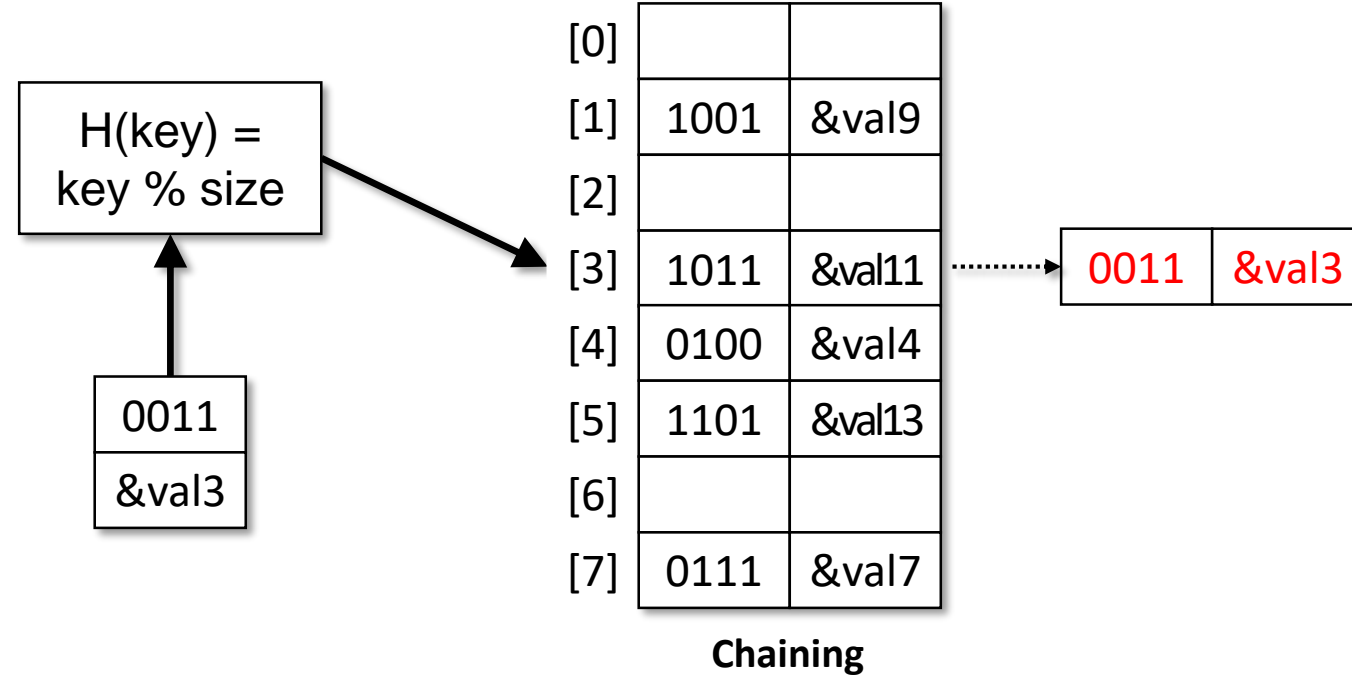| | | |
|---|---|---|
| [0] | | |
| [1] | | |
| [2] | | |
| [3] | 0011 | &val3 |
| [4] | 0100 | &val4 |
| [5] | | |
| [6] | | |
| [7] | 0111 | &val7 |
| [8] | | |
| [9] | 1001 | &val9 |
| [A] | | |
| [B] | 1001 | &val11 |
| [C] | | |
| [D] | 1101 | &val13 |
| [E] | | |
| [F] | | |

- Hash key collision → Full table rehashing
  - The most expensive operation in hash table

3

# Background:
# Static Hashing

|  |  |
|---|---|
| [0] |  |  |
| [1] | 1001 | &val9 |
| [2] |  |  |
| [3] | 1011 | &val11 |
| [4] | 0100 | &val4 |
| [5] | 1101 | &val13 |
| [6] | 0011 | &val3 |
| [7] | 0111 | &val7 |

H(key) =
key % size

0011
&val3

**Linear Probing**

- To avoid full table rehashing:
  - Linear probing
  - Chaining
  - Double hashing such as Cuckoo hashing

# Background:
# Static Hashing

H(key) =
key % size

| | |
|---|---|
| [0] | |
| [1] 1001 | &val9 |
| [2] | |
| [3] 1011 | &val11 |
| [4] 0100 | &val4 |
| [5] 1101 | &val13 |
| [6] | |
| [7] 0111 | &val7 |

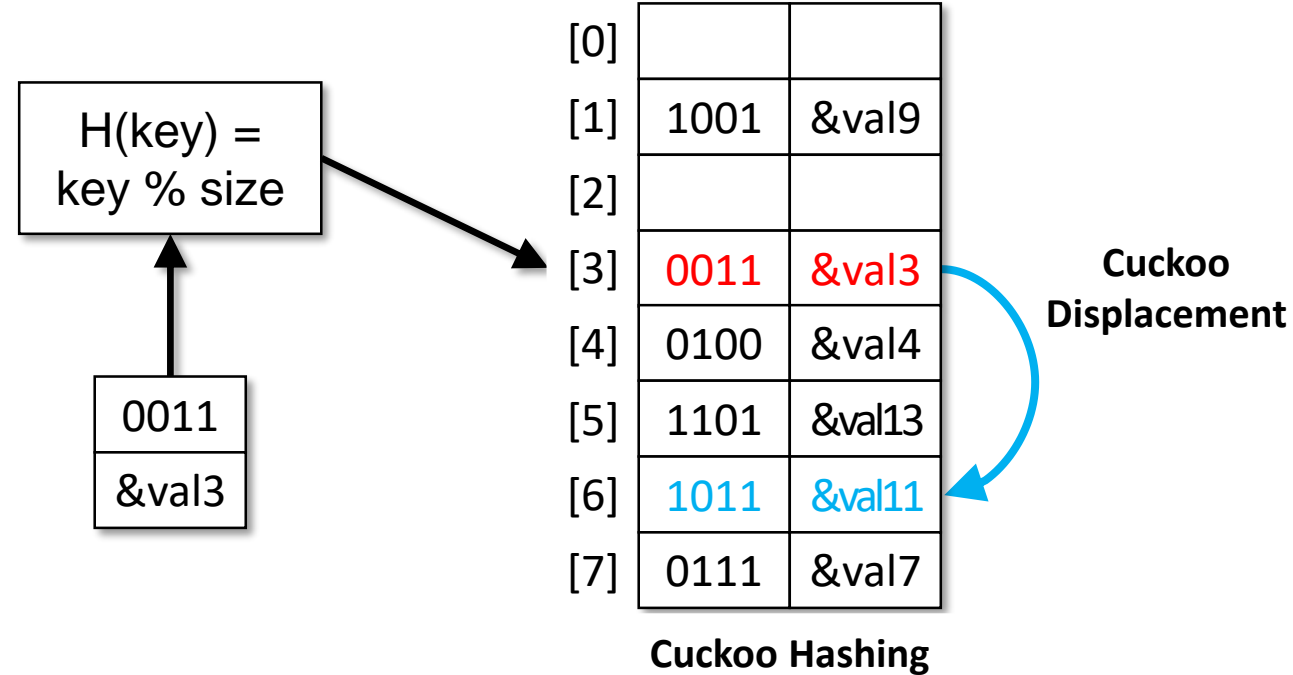| 0011 | &val3 |
|---|---|

| 0011 | &val3 |
|---|---|

**Chaining**

- To avoid full table rehashing:
  - Linear probing
  - Chaining
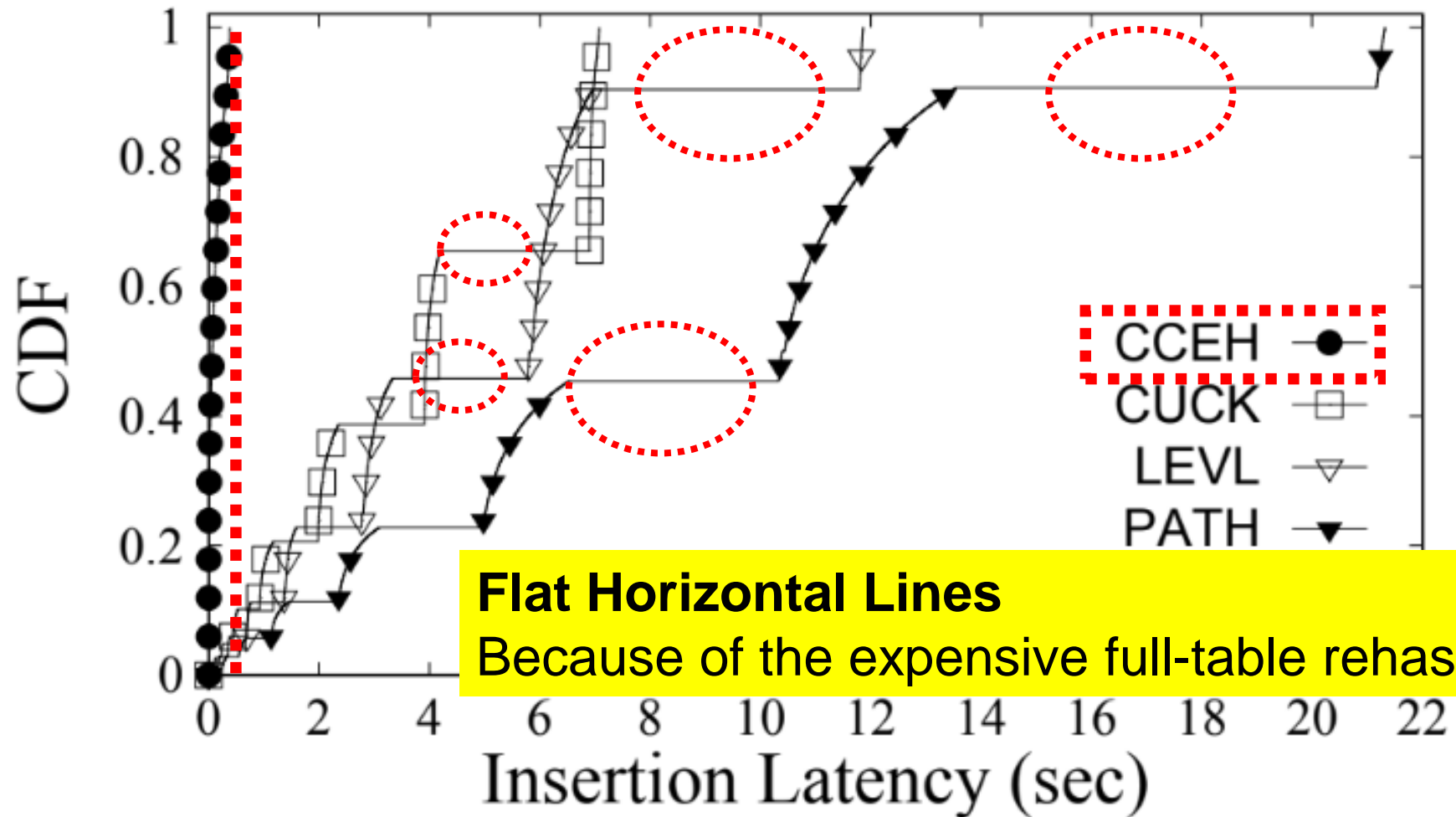  - Double hashing such as Cuckoo hashing

# Background:
# Static Hashing



$H(key) =$ key % size

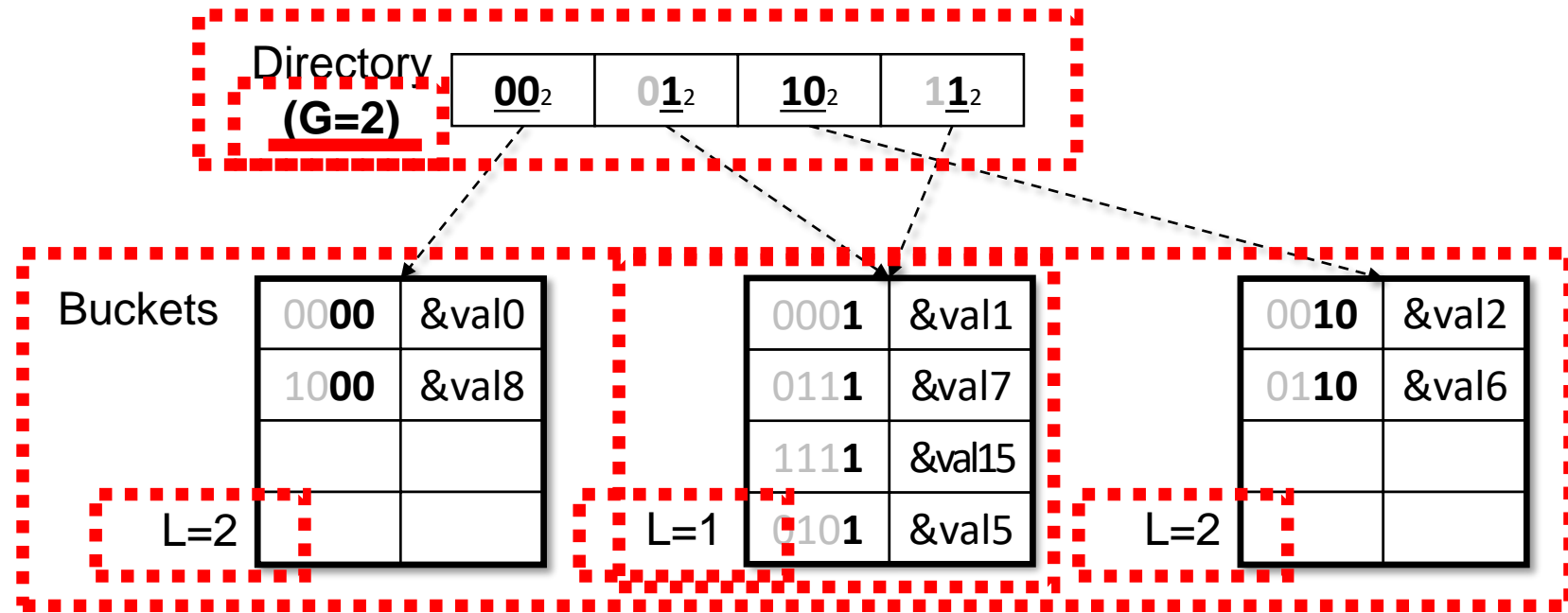| | | |
|---|---|---|
| [0] | | |
| [1] | 1001 | &val9 |
| [2] | | |
| [3] | 0011 | &val3 |
| [4] | 0100 | &val4 |
| [5] | 1101 | &val13 |
| [6] | 1011 | &val11 |
| [7] | 0111 | &val7 |

0011
&val3

**Cuckoo Displacement**

**Cuckoo Hashing**

- To avoid full table rehashing:
  - Linear probing
  - Chaining
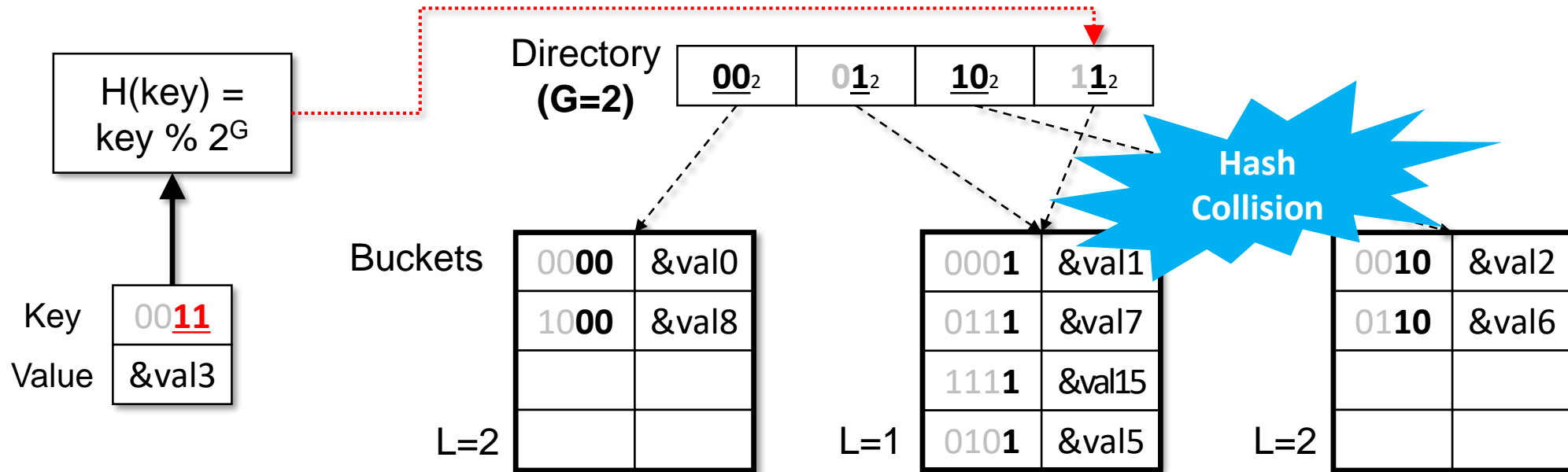  - Double hashing such as Cuckoo hashing

6

# Insertion Latency CDF



**Flat Horizontal Lines**
Because of the expensive full-table rehashing

# Background:
# Disk-based Extendible Hashing

Hash Function:
$H(key) = key \% 2^G$

Directory
**(G=2)**

| $\underline{00}_2$ | $0\underline{1}_2$ | $\underline{10}_2$ | $1\underline{1}_2$ |

Buckets

| 00**00** | &val0 |
|----------|-------|
| 10**00** | &val8 |
|          |       |
| L=2      |       |

| 000**1** | &val1 |
|----------|-------|
| 011**1** | &val7 |
| 111**1** | &val15 |
| L=1  010**1** | &val5 |

| 00**10** | &val2 |
|----------|-------|
| 01**10** | &val6 |
|          |       |
| L=2      |       |

- Dynamically splits one bucket or merges two buckets at a time

# Extendible Hashing – Insertion

$$H(key) = key \% 2^G$$

Key | $00\underline{\textbf{11}}$
Value | &val3

Directory **(G=2)**

| $\underline{\textbf{00}}_2$ | $0\underline{\textbf{1}}_2$ | $\underline{\textbf{10}}_2$ | $1\underline{\textbf{1}}_2$ |

**Hash Collision**

Buckets

| 00**00** | &val0 |
| 10**00** | &val8 |
| | |
| | |

L=2

| 000**1** | &val1 |
| 011**1** | &val7 |
| 111**1** | &val15 |
| 010**1** | &val5 |

L=1

| 00**10** | &val2 |
| 01**10** | &val6 |
| | |
| | |

L=2

# Background:
# Extendible Hashing – Bucket Split

$$H(key) = key \% 2^G$$

Key    0011

Value   &val3

Directory **(G=2)**

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |

Buckets

| 00**00** | &val0 |
|---|---|
| 10**00** | &val8 |
| | |
| | |

L=2

| 00**01** | &val1 |
|---|---|
| 01**11** | &val7 |
| 11**11** | &val15 |
| 01**01** | **Bucket Split** |

L=1

| 00**10** | &val2 |
|---|---|
| 01**10** | &val6 |
| | |
| | |

L=2

- Only overflown bucket is modified

# Background:
# Extendible Hashing – Bucket Split

$$H(key) = key \% 2^G$$

| Key | 0011 |
|---|---|
| Value | &val3 |

Directory (G=2)

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |
|---|---|---|---|

**Update Directory**

Buckets

| 00**00** | &val0 |
|---|---|
| 10**00** | &val8 |
| | |
| | |

L=2

| 00**10** | &val2 |
|---|---|
| 01**10** | &val6 |
| | |
| | |

L=2

| 00**01** | &val1 |
|---|---|
| 01**01** | &val5 |
| | |
| | |

L=2

| 01**11** | &val7 |
|---|---|
| 11**11** | &val15 |
| 00**11** | &val3 |
| | |

L=2

- Update Directory
  - At least two pointers need to be updated

# Background:
# Extendible Hashing – Directory doubling

H(key) = key % $2^G$
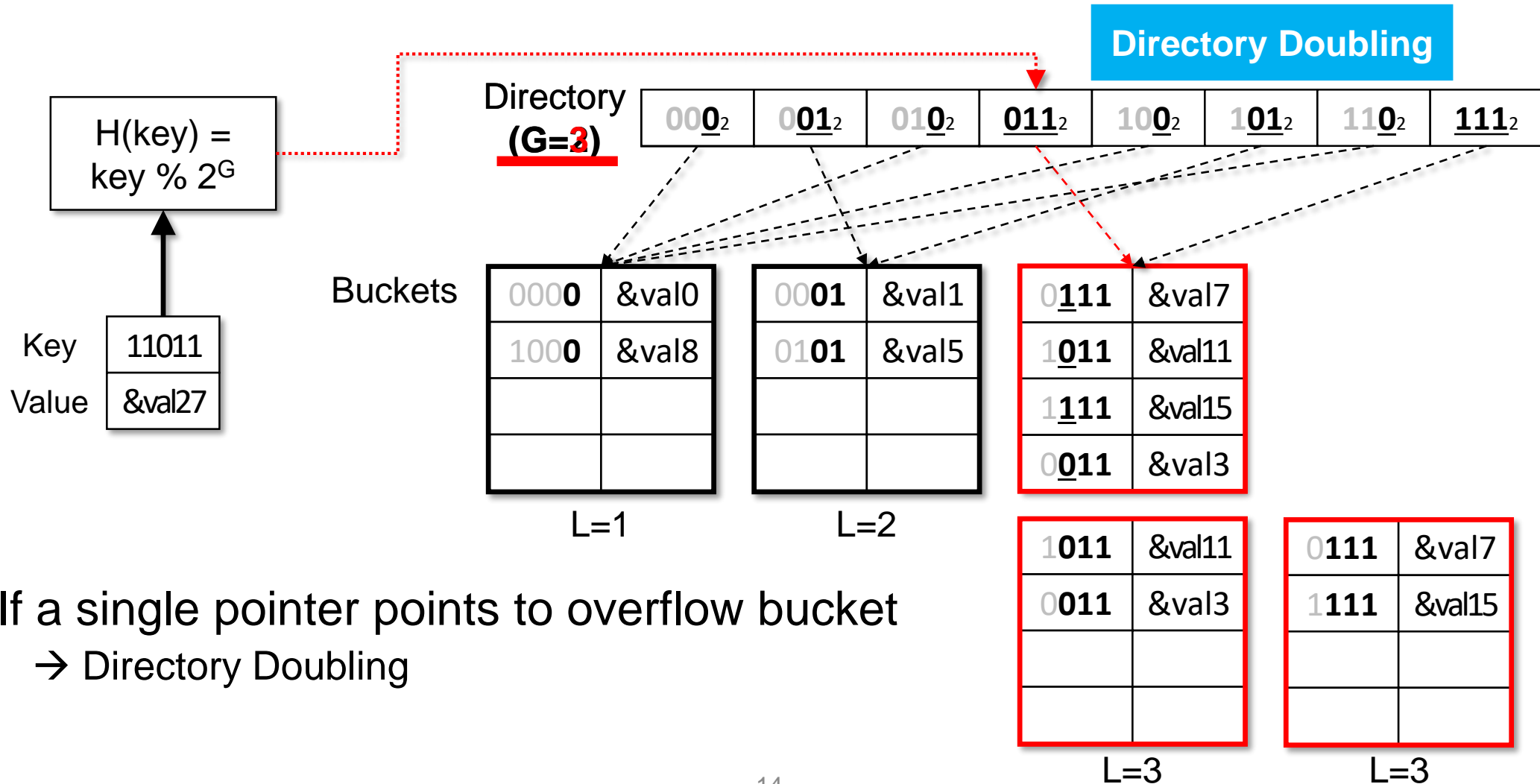
Key | 110**11**
Value | &val27

Directory **(G=2)**

| $0\underline{0}_2$ | $0\underline{1}_2$ | $1\underline{0}_2$ | $1\underline{1}_2$ |

Buckets

| 000**0** | &val0 |
|---|---|
| 100**0** | &val8 |
| | |
| | |

L=1

| 00**01** | &val1 |
|---|---|
| 01**01** | &val5 |
| | |
| | |

L=2

| 01**11** | &val7 |
|---|---|
| 10**11** | &val11 |
| 11**11** | &val15 |
| 00**11** | &val3 |

L=2

**Hash Collision**

- If a single pointer points to overflow bucket
  - → Directory Doubling

# Extendible Hashing – Directory doubling



| H(key) = key % $2^G$ |
|---|

| Key | 11011 |
|---|---|
| Value | &val27 |

Directory **(G=2)**

| $0\underline{0}_2$ | $\underline{01}_2$ | $1\underline{0}_2$ | $\underline{11}_2$ |
|---|---|---|---|

Buckets

| 000**0** | &val0 |
|---|---|
| 100**0** | &val8 |
| | |
| | |

L=1

| 00**01** | &val1 |
|---|---|
| 01**01** | &val5 |
| | |
| | |

L=2

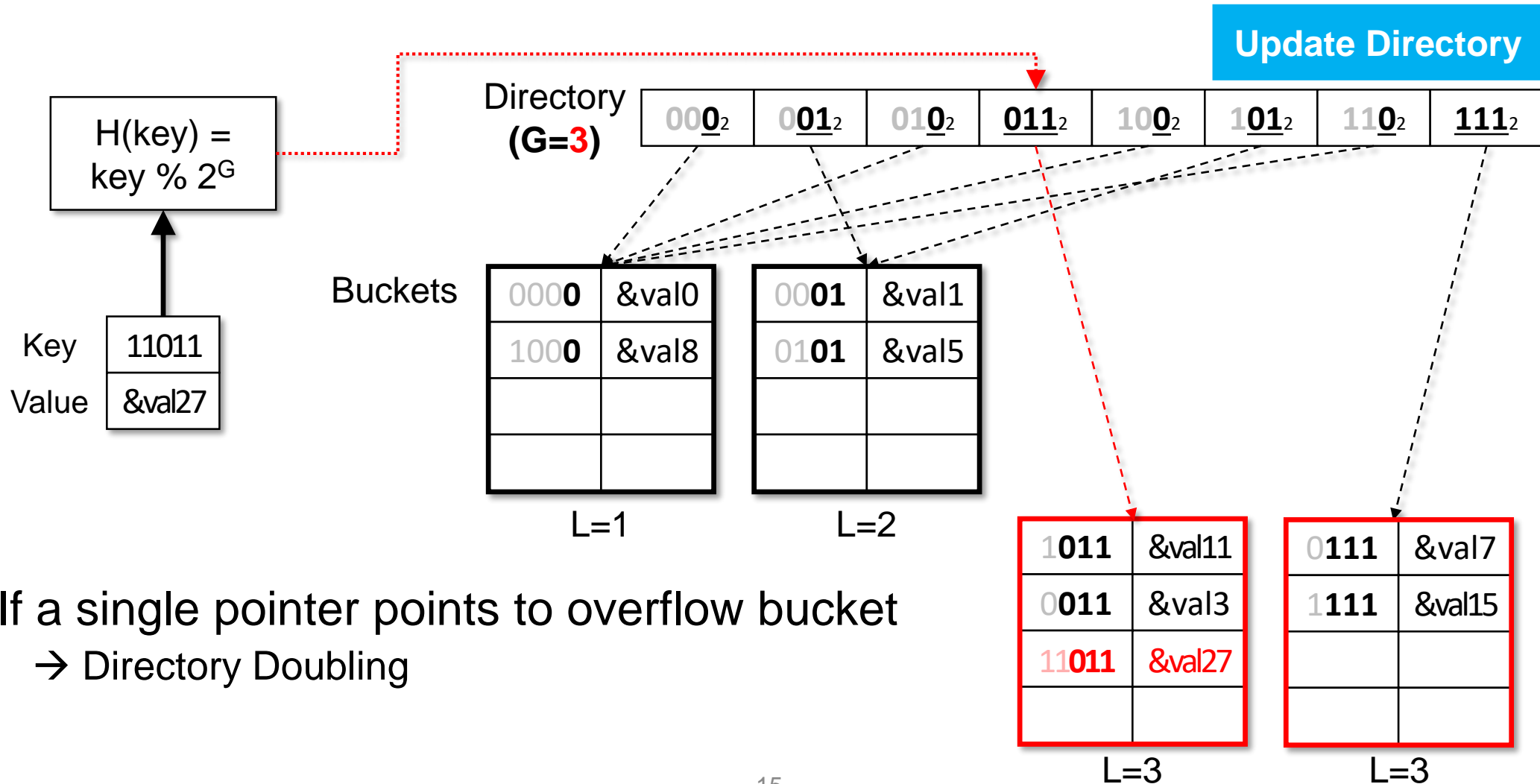| 01**11** | &val7 |
|---|---|
| 10**11** | &val11 |
| 11**11** | &val15 |
| 00**11** | &val3 |

L=2

**Bucket Split**

- If a single pointer points to overflow bucket
  - → Directory Doubling
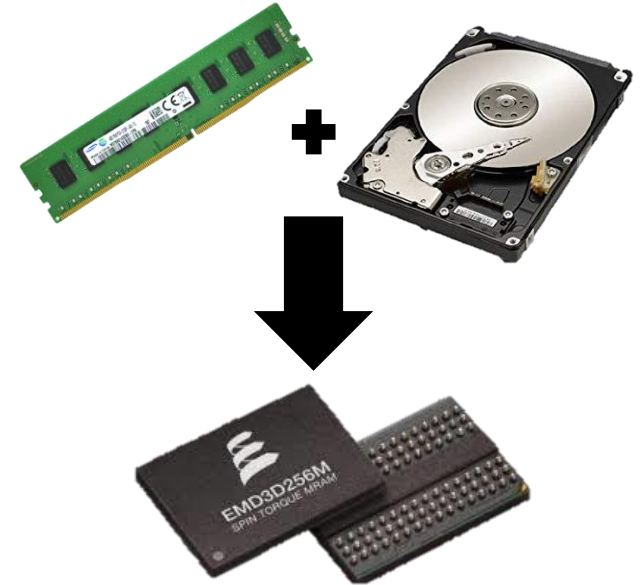
# Background:
# Extendible Hashing – Directory doubling



Directory Doubling

$H(key) = key \% 2^G$

Key | 11011

Value | &val27

Directory **(G=3)**

| $000\underline{0}_2$ | $00\underline{1}_2$ | $01\underline{0}_2$ | $\underline{011}_2$ | $10\underline{0}_2$ | $1\underline{01}_2$ | $11\underline{0}_2$ | $\underline{111}_2$ |

Buckets

| $000\mathbf{0}$ | &val0 |
|---|---|
| $100\mathbf{0}$ | &val8 |
| | |
| | |

L=1

| $00\mathbf{01}$ | &val1 |
|---|---|
| $01\mathbf{01}$ | &val5 |
| | |
| | |

L=2

| $0\mathbf{111}$ | &val7 |
|---|---|
| $1\underline{011}$ | &val11 |
| $1\mathbf{111}$ | &val15 |
| $0\underline{011}$ | &val3 |

| $1\mathbf{011}$ | &val11 |
|---|---|
| $0\underline{011}$ | &val3 |
| | |
| | |

L=3

| $0\mathbf{111}$ | &val7 |
|---|---|
| $1\mathbf{111}$ | &val15 |
| | |
| | |

L=3

- If a single pointer points to overflow bucket
  → Directory Doubling

# Background:
# Extendible Hashing – Directory doubling

**Update Directory**

H(key) = key % 2$^G$

Directory (G=**3**)

| 000$_2$ | 001$_2$ | 010$_2$ | 011$_2$ | 100$_2$ | 101$_2$ | 110$_2$ | 111$_2$ |
|---|---|---|---|---|---|---|---|

Key 11011

Value &val27

Buckets

| 000**0** | &val0 |
|---|---|
| 100**0** | &val8 |
| | |
| | |

L=1

| 00**01** | &val1 |
|---|---|
| 01**01** | &val5 |
| | |
| | |

L=2

| 1**011** | &val11 |
|---|---|
| 0**011** | &val3 |
| 11**011** | &val27 |
| | |

L=3

| 0**111** | &val7 |
|---|---|
| 1**111** | &val15 |
| | |
| | |

L=3

- If a single pointer points to overflow bucket
  → Directory Doubling

# Persistent Memory

## Characteristics
- High performance – Comparable to DRAM
- Byte-addressability – As DRAM
- Persistence – As storage devices (HDD/SSD)

## Challenges
- Atomic unit of writes → 8-bytes
- Data transfer unit between CPU cache and PM → 64 byte cacheline
- Order of memory writes is not guaranteed

# Outline

**~~Background~~**

- ~~Static Hashing~~
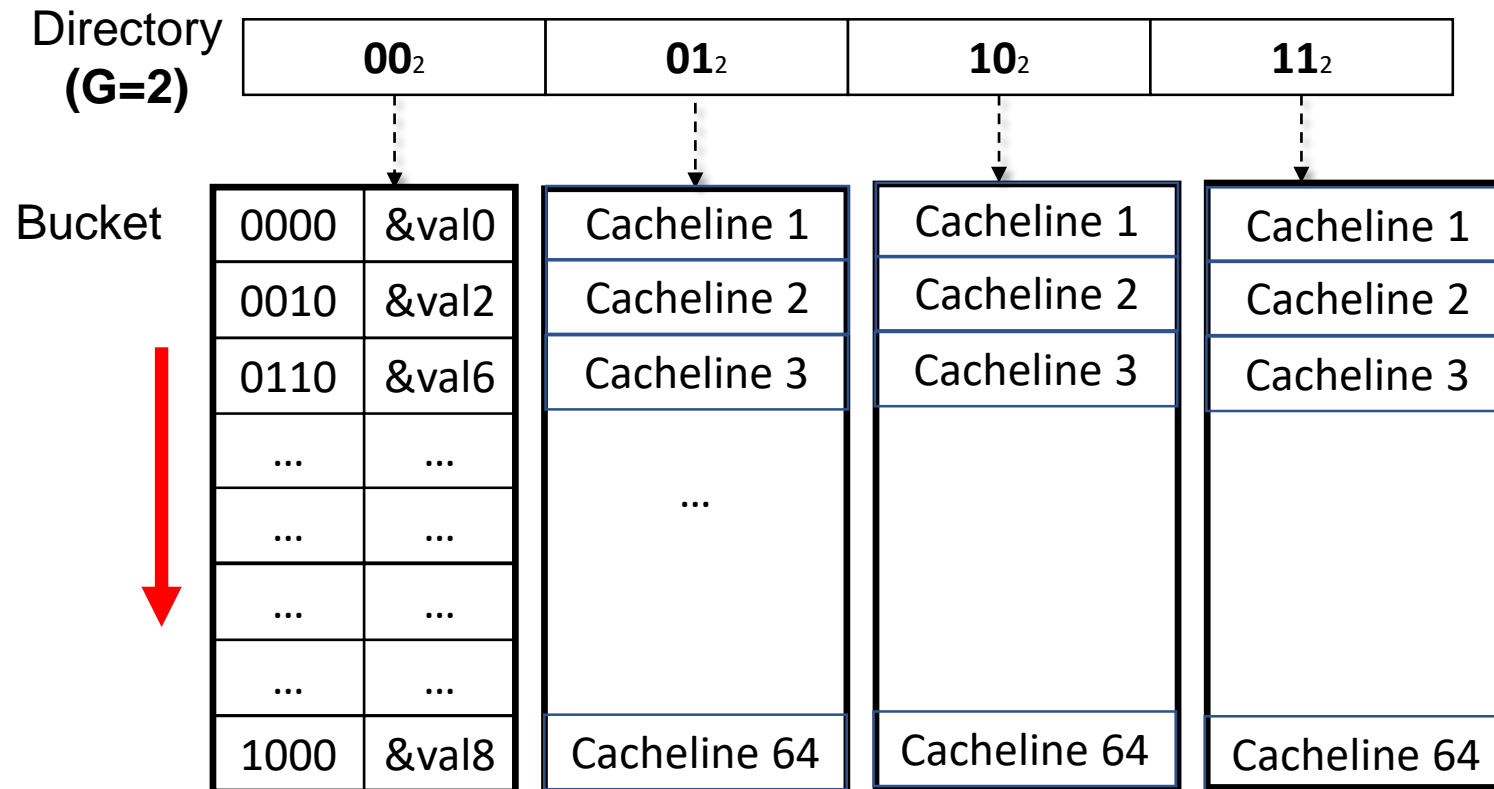- ~~Extendible Hashing~~
- ~~Persistent Memory~~

**Cacheline-Conscious Extendible Hashing**

- Challenges and Contributions
- 3-Level Structure of CCEH
- Failure-atomic Directory Update

**Evaluation**

**Conclusion**

# Challenge in In-Memory Extendible Hashing

**Directory (G=2)**

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |
|---|---|---|---|

**Bucket**

| 0000 | &val0 |
|---|---|
| 0010 | &val2 |
| 0110 | &val6 |
| ... | ... |
| ... | ... |
| ... | ... |
| ... | ... |
| 1000 | &val8 |

| Cacheline 1 |
|---|
| Cacheline 2 |
| Cacheline 3 |
| ... |
| Cacheline 64 |

| Cacheline 1 |
|---|
| Cacheline 2 |
| Cacheline 3 |
| |
| Cacheline 64 |

| Cacheline 1 |
|---|
| Cacheline 2 |
| Cacheline 3 |
| |
| Cacheline 64 |

**Problems**
- Page-sized bucket → 64 cacheline accesses per bucket

# Challenge in In-Memory Extendible Hashing



Directory (G=20)

Buckets

00..00

Cacheline-Size

**Problems**

Cacheline-sized small bucket → a large directory
(8 byte pointer per cacheline)

# Challenge in Extendible Hashing on PM

Directory
**(G=2)**

| **$000_2$** | **$001_2$** | **$010_2$** | **$011_2$** |
|---|---|---|---|

Buckets

| 0000 | &val0 |
|---|---|
| 1000 | &val8 |
| | |
| | |

| 0010 | &val1 |
|---|---|
| 0110 | &val5 |
| | |
| | |

| 1001 | &val11 |
|---|---|
| 0001 | &val3 |
| 11001 | &val27 |
| | |

| 0111 | &val7 |
|---|---|
| 1111 | &val15 |
| | |
| | |

**Problems**

Split operation updates multiple pointers → Not Failure-Atomic

# Contributions

## 3-Level Structure

→ Introduces an intermediate level, *Segment*

→ Lookup via only **two cacheline accesses**

## Failure-atomic Directory Updates

→ Introduces **the split buddy tree** to **manage split history**

## Failure-atomic Segment Split

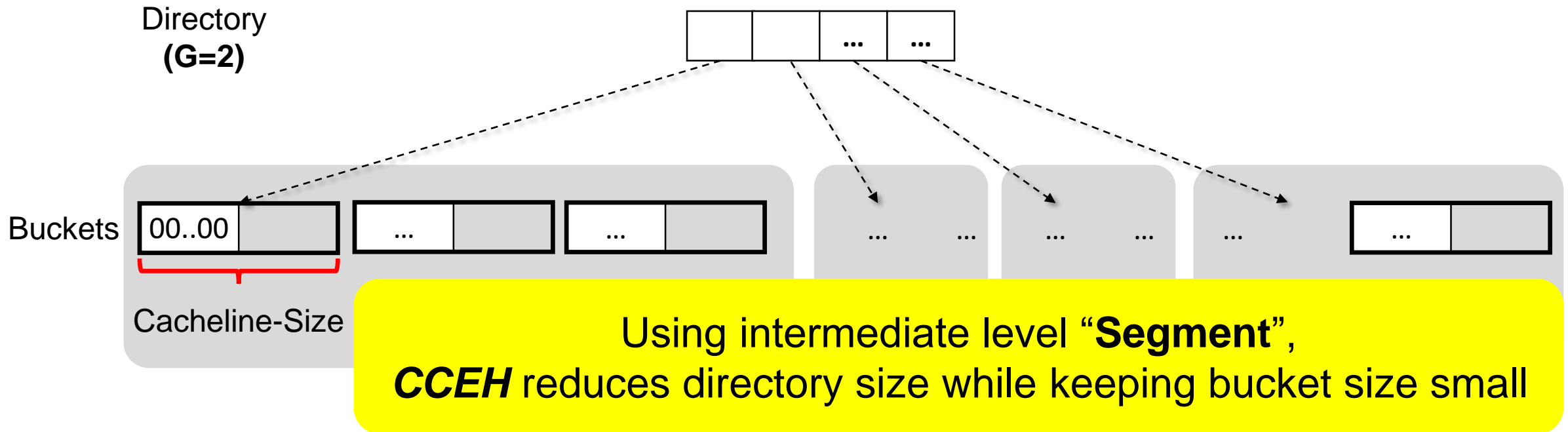→ *Lazy deletion* scheme to **minimize dirty writes**

# Segment: Intemediate Level
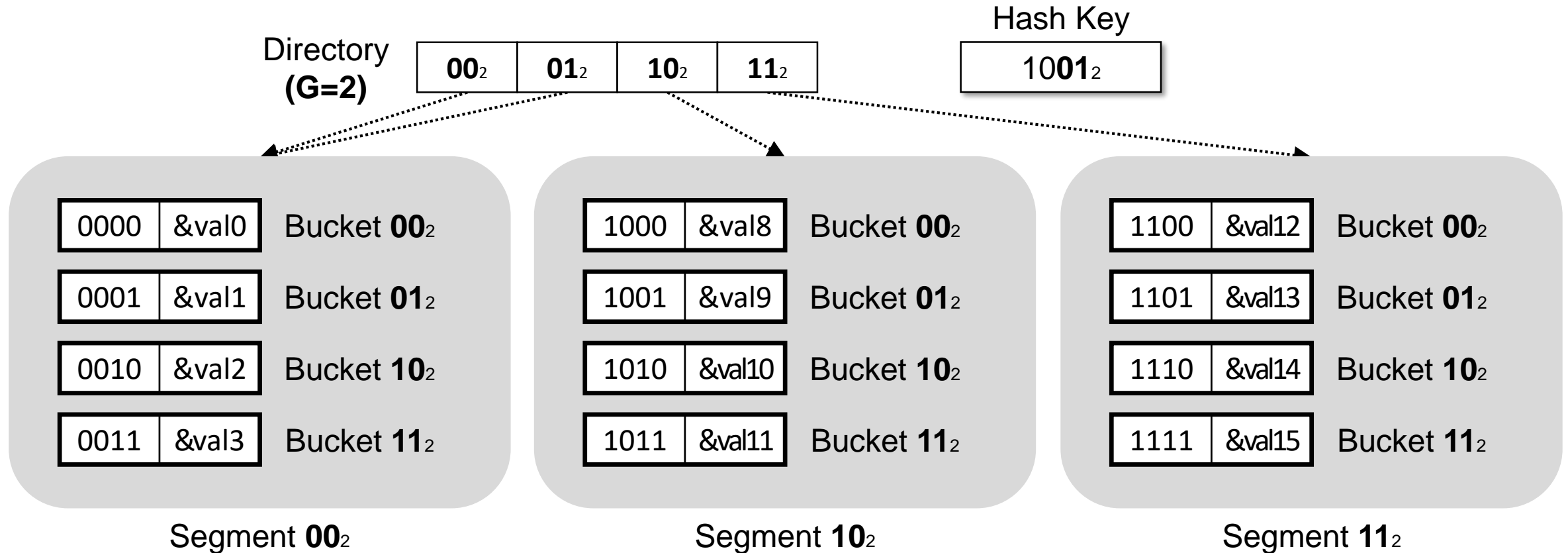


- A group of multiple cacheline-sized buckets = Segment
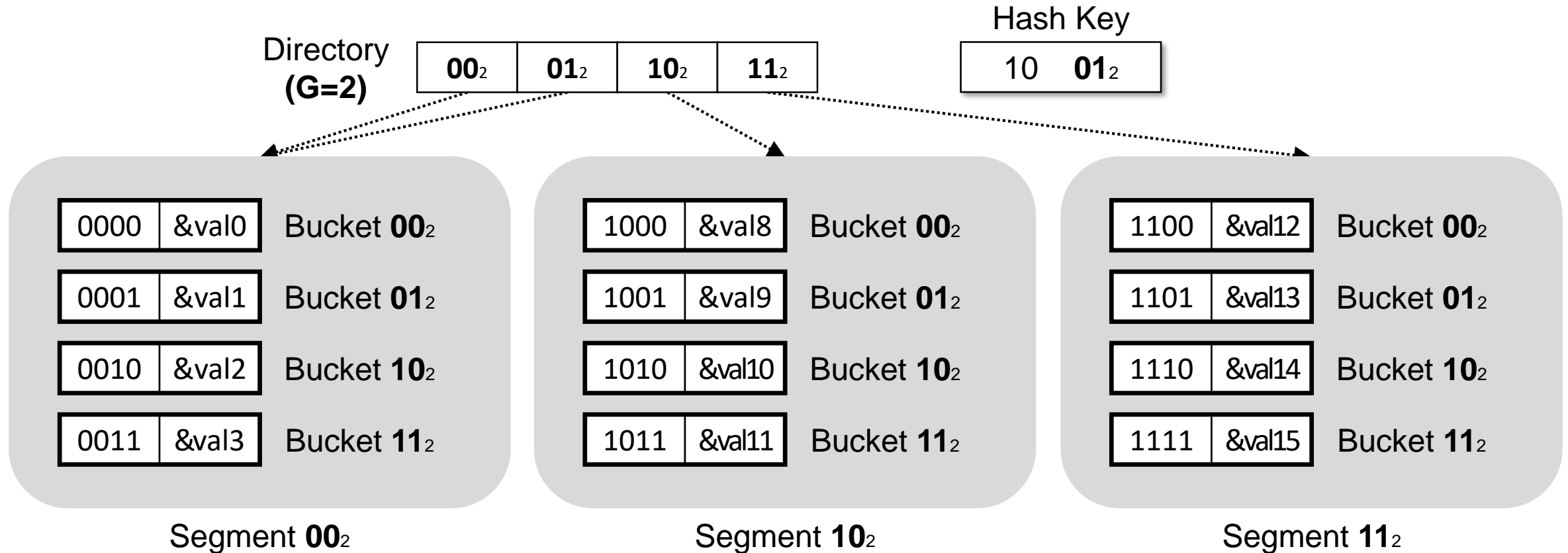
# Segment: Intemediate Level

Directory
**(G=2)**

Buckets | 00..00 | ... | ... | ... | ... | ... | ... | ... | ... |

Cacheline-Size

Using intermediate level "**Segment**",
*CCEH* reduces directory size while keeping bucket size small

## 3-Level Structure
Directory → Segment → Cacheline-sized Bucket
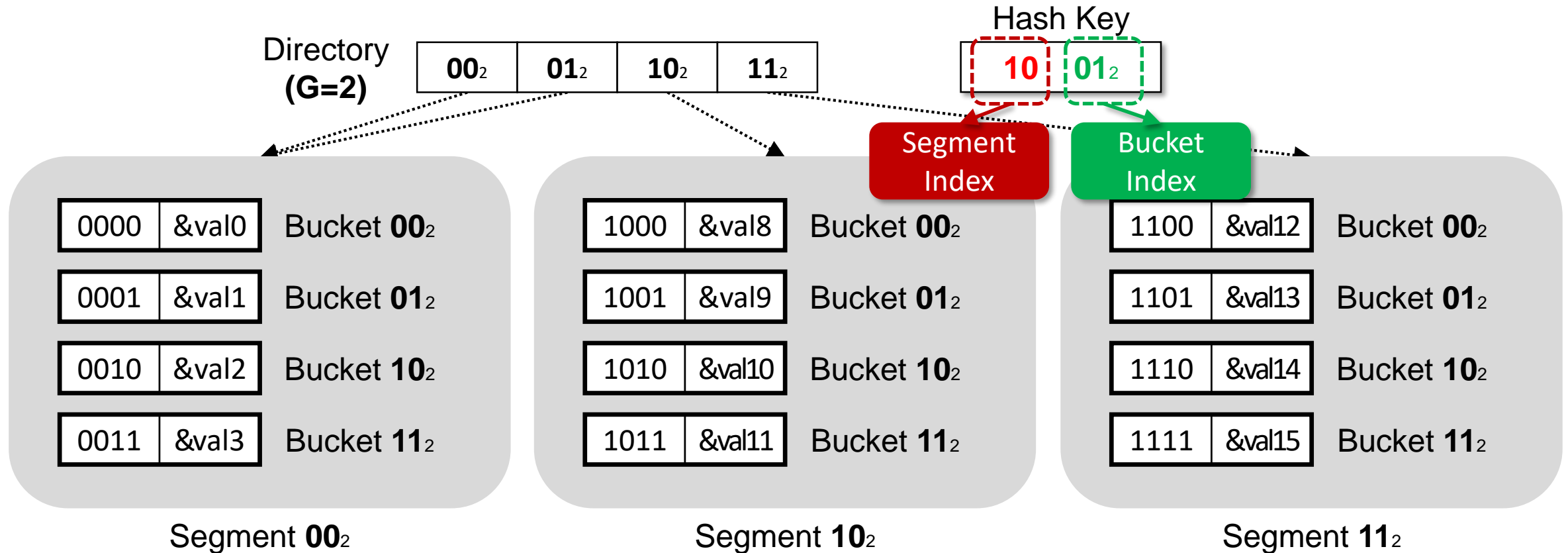
# Minimize Cacheline Accesses in Segment

Directory (G=2)

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |

Hash Key

| $10\mathbf{01}_2$ |

| 0000 | &val0 | Bucket $00_2$
| 0001 | &val1 | Bucket $01_2$
| 0010 | &val2 | Bucket $10_2$
| 0011 | &val3 | Bucket $11_2$

| 1000 | &val8 | Bucket $00_2$
| 1001 | &val9 | Bucket $01_2$
| 1010 | &val10 | Bucket $10_2$
| 1011 | &val11 | Bucket $11_2$

| 1100 | &val12 | Bucket $00_2$
| 1101 | &val13 | Bucket $01_2$
| 1110 | &val14 | Bucket $10_2$
| 1111 | &val15 | Bucket $11_2$

Segment $00_2$

Segment $10_2$

Segment $11_2$

**Q**: With large segments, how can we minimize cacheline accesses?

# Minimize Cacheline Accesses in Segment

Directory (G=2)

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |
|---|---|---|---|

Hash Key

| 10 $01_2$ |
|---|

| 0000 | &val0 | Bucket $00_2$ |
|---|---|---|
| 0001 | &val1 | Bucket $01_2$ |
| 0010 | &val2 | Bucket $10_2$ |
| 0011 | &val3 | Bucket $11_2$ |

Segment $00_2$

| 1000 | &val8 | Bucket $00_2$ |
|---|---|---|
| 1001 | &val9 | Bucket $01_2$ |
| 1010 | &val10 | Bucket $10_2$ |
| 1011 | &val11 | Bucket $11_2$ |

Segment $10_2$

| 1100 | &val12 | Bucket $00_2$ |
|---|---|---|
| 1101 | &val13 | Bucket $01_2$ |
| 1110 | &val14 | Bucket $10_2$ |
| 1111 | &val15 | Bucket $11_2$ |

Segment $11_2$

**Q**: With large segments, how can we minimize cacheline accesses?

# Minimize Cacheline Accesses in Segment

Directory (G=2)

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |

Hash Key

| $10$ | $01_2$ | |

**Segment Index**

**Bucket Index**

| 0000 | &val0 | Bucket $00_2$ |
| 0001 | &val1 | Bucket $01_2$ |
| 0010 | &val2 | Bucket $10_2$ |
| 0011 | &val3 | Bucket $11_2$ |

Segment $00_2$

| 1000 | &val8 | Bucket $00_2$ |
| 1001 | &val9 | Bucket $01_2$ |
| 1010 | &val10 | Bucket $10_2$ |
| 1011 | &val11 | Bucket $11_2$ |

Segment $10_2$

| 1100 | &val12 | Bucket $00_2$ |
| 1101 | &val13 | Bucket $01_2$ |
| 1110 | &val14 | Bucket $10_2$ |
| 1111 | &val15 | Bucket $11_2$ |

Segment $11_2$

**Q**: With large segments, how can we minimize cacheline accesses?

# Minimize Cacheline Accesses in Segment

Directory (G=2)

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |

Hash Key

| **10** **$01_2$** |

| 0000 | &val0 | Bucket $00_2$
| 0001 | &val1 | Bucket $01_2$
| 0010 | &val2 | Bucket $10_2$
| 0011 | &val3 | Bucket $11_2$

| 1000 | &val8 | Bucket $00_2$
| 1001 | &val9 | Bucket $01_2$
| 1010 | &val10 | Bucket $10_2$
| 1011 | &val11 | Bucket $11_2$

| 1100 | &val12 | Bucket $00_2$
| 1101 | &val13 | Bucket $01_2$
| 1110 | &val14 | Bucket $10_2$
| 1111 | &val15 | Bucket $11_2$

Segment $11_2$

Use hash key as index for both directory and segment
→ No need to access irrelevant buckets

# Contributions

**3-Level Structure**

    → Introduces an intermediate level, *Segment*

    → Lookup via only *two cacheline accesses*

**Failure-atomic Directory Updates**

    → Introduces *the split buddy tree* to **manage split history**

**Failure-atomic Segment Split**

    → *Lazy deletion* scheme to **minimize dirty writes**

# Recovery: Split History Buddy Tree in CCEH



Using **MSB** segment index, split segments are pointed by **adjacent directory** entries

# Recovery: Split History Buddy Tree in CCEH

Global Depth:

Depth:   0    1    2    3

**Directory (G=3)**

| | |
|---|---|
| S1 | L=1 |
| S1 | L=1 |
| S1 | L=1 |
| **S5** | **L=2** |
| S2 | L=2 |
| S2 | L=2 |
| S3 | L=3 |
| S4 | L=3 |

S1 → S1 → S5

S2 → S2 → S3 → S3 → S4

Suppose a system crashes while **S5** splits

# Recovery: Split History Buddy Tree in CCEH

Global Depth:

Depth:    0        1        2        3

Directory (G=3)

| | |
|---|---|
| S1 | L=1 |
| S1 | L=1 |
| S1 | L=1 |
| S5 | L=2 |
| S2 | L=2 |
| S2 | L=2 |
| S3 | L=3 |
| S4 | L=3 |

$$Stride = 2^{G-L}$$

Global Depth **G** = 3
Local Depth **L** = 1
*Stride* = 4

- Each **segment** must be pointed by $2^{G-L}$ directory entries
- If not, rollback the split

# Recovery: Split History Buddy Tree in CCEH

Global Depth:

Depth:    0        1        2        3



Directory (G=3)

| | |
|---|---|
| S1 | L=1 |
| S1 | L=1 |
| S1 | L=1 |
| **S1** | **L=1** |
| S2 | L=2 |
| S2 | L=2 |
| S3 | L=3 |
| S4 | L=3 |

$$Stride = 2^{G-L}$$

Global Depth **G** = 3
Local Depth **L** = 1
*Stride* = 4

- Each **segment** must be pointed by $2^{G-L}$ directory entries
- If not, rollback the split

# Contributions

## 3-Level Structure

→ Introduces an intermediate level, *Segment*

→ Lookup via only **two cacheline accesses**

## Failure-atomic Directory Updates

→ Introduces **the split buddy tree** to **manage split history**

## Failure-atomic Segment Split

→ ***Lazy deletion*** scheme to **minimize dirty writes**

33

# Segment Split: Legacy CoW

Directory

| $00_2$ | $01_2$ | $10_2$ | $11_2$ |
|---|---|---|---|

**Copy-on-Write Split**
→ Lock-Free Search is enabled

| Local Depth = 1 | |
|---|---|
| $00_2$ | **01**000   &val8 |
| $01_2$ | **00**001   &val1 |
| $10_2$ | **01**010   &val10 |
| $11_2$ | **00**011   &val3 |

Segment $0_2$

| Local Depth = **2** | |
|---|---|
| $00_2$ | |
| $01_2$ | |
| $10_2$ | |
| $11_2$ | |

Segment **00**$_2$

| Local Depth = **2** | |
|---|---|
| $00_2$ | |
| $01_2$ | |
| $10_2$ | |
| $11_2$ | |

Segment **01**$_2$

# Segment Split: Lazy Deletion

Directory | $00_2$ | $01_2$ | $10_2$ | $11_2$

**Lazy Deletion**
  → Minimizes dirty writes

Local Depth = **2**

**Single Cacheline-flush invalidates all the migrated data**

| $00_2$ | ~~01000~~ | ~~&val8~~ |
| $01_2$ | **00**001 | &val1 |
| $10_2$ | ~~01010~~ | ~~&val10~~ |
| $11_2$ | **00**011 | &val3 |

Segment **$00_2$**

Local Depth = **2**

| $00_2$ | | |
| $01_2$ | | |
| $10_2$ | | |
| $11_2$ | | |

Segment **$01_2$**

35

# Outline

**~~Background~~**

- ~~Static Hashing~~
- ~~Extendible Hashing~~
- ~~Persistent Memory~~

**~~Cacheline-Conscious Extendible Hashing~~**

- ~~Challenges and Contributions~~
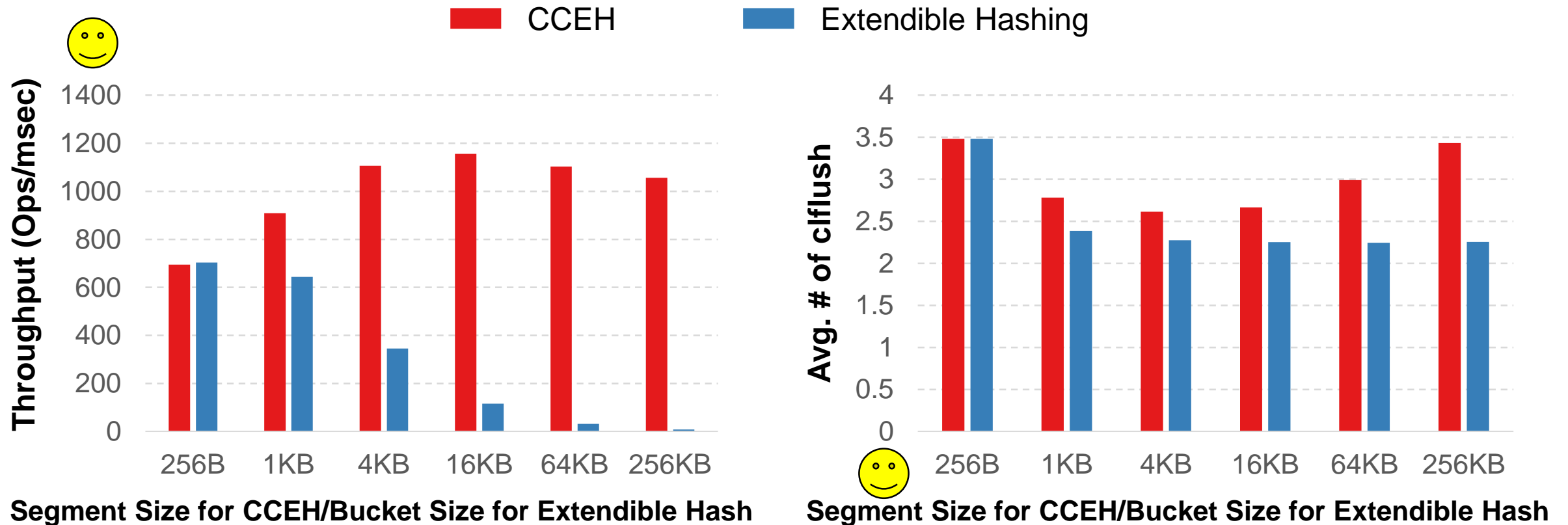- ~~3-Level Structure of CCEH~~
- ~~Failure-atomic Directory Update~~

**Evaluation**

**Conclusion**

# Experimental Setup

| | |
|---|---|
| **CPU** | 2x Intel Xeon Haswell-Ex E7-4809 v3<br>$\rightarrow$ 8 cores, 2.0 GHz<br>$\rightarrow$ 20MB L3 cache |
| **Memory** | 64GB of DDR3 DRAM |
| **PM** | **Quartz:** A DRAM-based PM latency emulator<br>* To emulate write latency, we inject stall cycle after each *clflush* instructions |
| **Workload** | 160 Million random number dataset |

# CCEH VS Legacy Extendible Hash



**CCEH**     **Extendible Hashing**

Throughput (Ops/msec) vs Segment Size for CCEH/Bucket Size for Extendible Hash

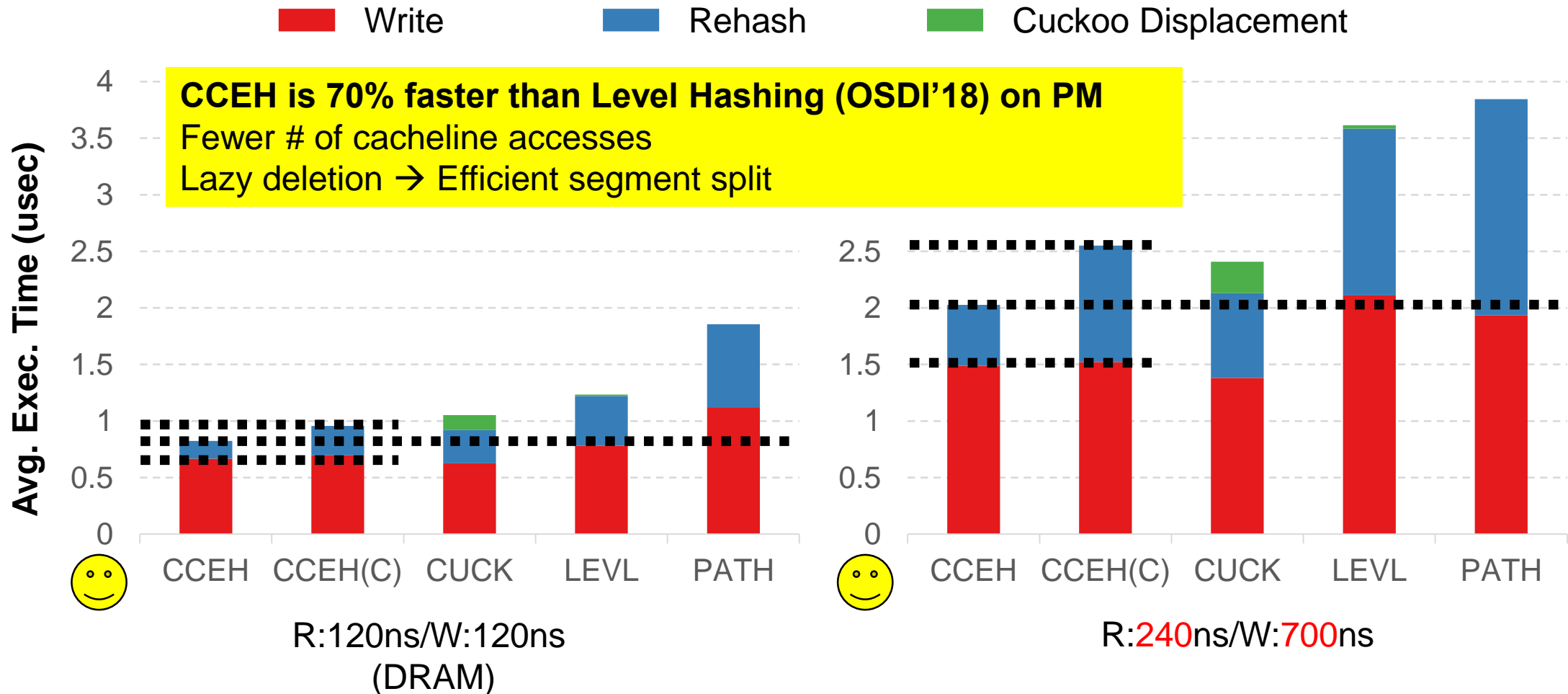Avg. # of clflush vs Segment Size for CCEH/Bucket Size for Extendible Hash

**CCEH compared to legacy Extendible Hashing**
Cons: Low utilization and more cacheline flushes due to hash collisions
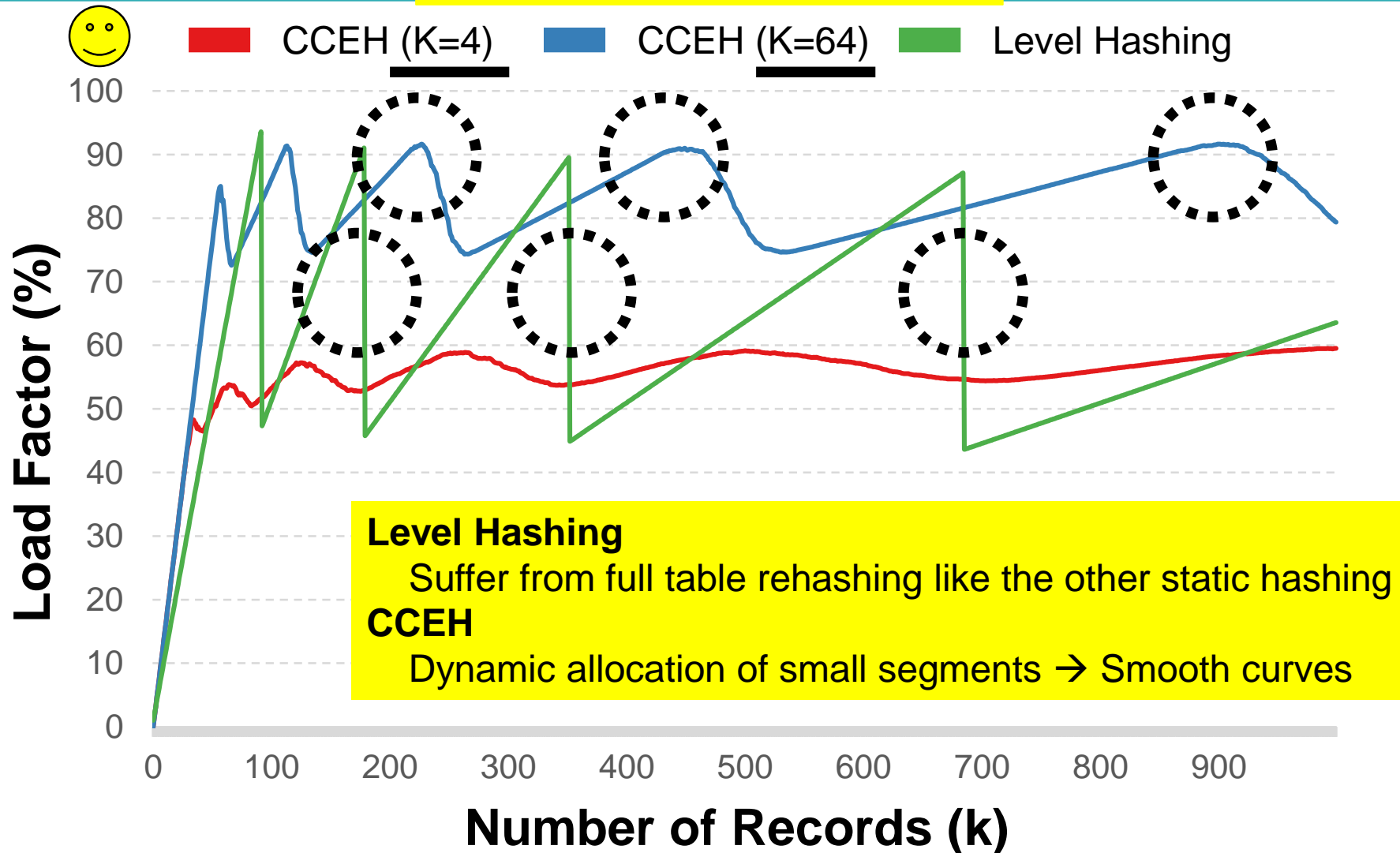Pros: Constant number of cacheline accesses with varying directory size

# Insertion Performance Breakdown

**Write** ▮    **Rehash** ▮    **Cuckoo Displacement** ▮

**CCEH is 70% faster than Level Hashing (OSDI'18) on PM**
Fewer # of cacheline accesses
Lazy deletion → Efficient segment split

Avg. Exec. Time (usec)

CCEH   CCEH(C)   CUCK   LEVL   PATH

R:120ns/W:120ns
(DRAM)

CCEH   CCEH(C)   CUCK   LEVL   PATH

R:240ns/W:700ns

# Load Factor

# Conclusion

**Cacheline-Conscious Extendible Hashing (CCEH)**
- 3-Level Structure
  - Introduced an intermediate level, Segment
  - Constant Lookup: Only two cacheline accesses → Write-Optimal

- Failure-Atomic Write-Optimal Lazy Deletion → Minimize I/O
- Failure-Atomic Directory Updates → Log-less directory update

Disk-based hashing needs to be modified for PM to make effective use of cachelines.

Source Codes:  http://github.com/DICL/CCEH

# Question?