



## 20. Indexing – Hashing

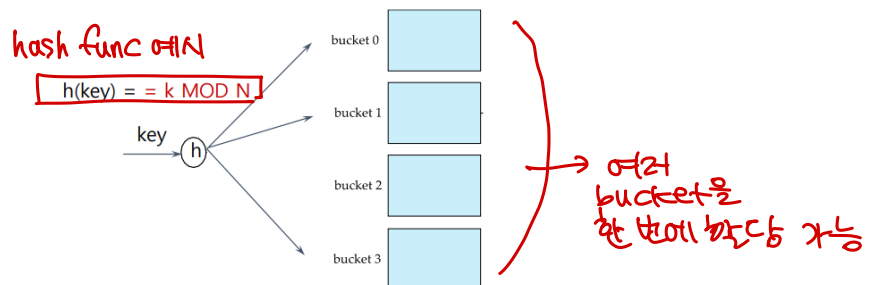
### ▼ Hashing

- bucket : data를 담는 공간 → data가 저장되는 단계 (보통 block 크기, 4KB)
- Hash function : bucket을 쭉 탐색해야 함 → data가 어디에 저장될지 짐작.
  - $h$  (hash func) : search-key 집합으로부터 bucket address 집합 매칭  
→ data 위치를 알 수 있다 (bucket 번호를 알 수 있다!)
- 다른 search key가 같은 bucket에 mapping 될 수 있음  
⇒ 모든 bucket은 sequentially하게 탐색 가능
- hash index : record pointer가 있는 entry 저장되어 있는 bucket 번호? 위치?
- hash-file-organization → bucket에 record 저장되어 있는 구조

### ▼ Static Hashing

bucket의 개수가 fixed

\*data 찾을 때 순서  
① hash function  
② bucket 위치 확인



→ bucket이 sequentially하게 할당되어 있으며 할당된 공간은 절대 free x (재할당 불가)

### ▼ Handling of Bucket overflows

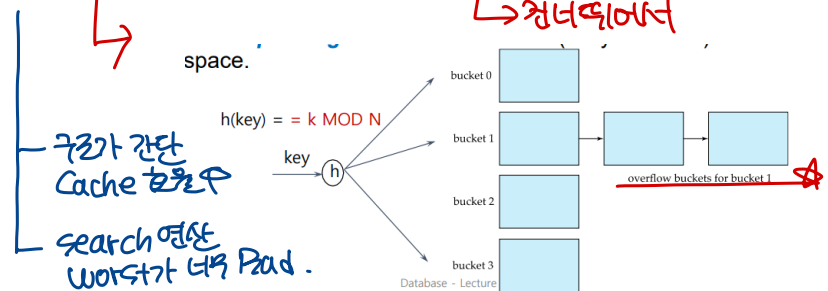
- Bucket overflow 발생 이유
  - 1. insufficient buckets : bucket이 너무 작거나 소수임
  - 2. hash function 오류(skew in distribution of records)
    - skew 발생 이유

1. multiple records : 같은 search-key를 가짐  $\Rightarrow$  미리 막을 방법이 없음
2. 선택된 hash function이 너무 단순하거나 non-uniform한 분포를 가짐

$\Rightarrow$  bucket overflow를 줄일 수 있지만 아예 일어나지 않게 할 수는 없음

$\rightarrow$  overflow buckets를 사용하여 조절함

- ① **Overflow Chaining**
  - 새로운 bucket 설정 + bucket이랑 linked list로 연결
  - Closed Addressing(Closed Hashing) : overflow bucket이 linked list 안에 chained된 형태
  - ② **Open Hashing** : overflow bucket 사용하지 않음  $\rightarrow$  다른 bucket에 data 삽입
- $\Rightarrow$  **Linear Probing** : space를 가지는 다음 bucket 사용(in cyclic order)



### ▼ Example of Hash file organization $\Rightarrow$ record 자체가 저장됨

$\rightarrow$  instructor file

key: dept\_name

bucket represent

$h \rightarrow \text{dept\_name modulo } 10$   
 $\rightarrow$  어떤 값을 return

ex)  $h(\text{Music}) = 1$

$h(\text{History}) = 2$

$h(\text{Physics}) = 3$

$h(\text{Elec. Eng.}) = 3$

$\Rightarrow h$ : uniform distribution 제공 못해!

bucket 0	bucket 4
	12121 Wu Finance 90000
	76543 Singh Finance 80000
bucket 1	bucket 5
15151 Mozart Music 40000	76766 Crick Biology 72000
bucket 2	bucket 6
32343 El Said History 80000	10101 Srinivasan Comp. Sci. 65000
58583 Califieri History 60000	45565 Katz Comp. Sci. 75000
bucket 3	bucket 7
22222 Einstein Physics 95000	
33456 Gold Physics 87000	
98345 Kim Elec. Eng. 80000	

$\rightarrow$  만약 bucket 3을 탐색  $\rightarrow$  3번 탐색  $\rightarrow$  uniform x  
반면에 bucket 7은 x

### ▼ Deficiencies of Static Hashing **중요점은?** $\Rightarrow$ 추가 연산 필요 x

- static hashing  $\rightarrow$  func(h)가 bucket address의 고정된 set에 맞게 search-key mapping

$\rightarrow$  h대로 계산해서 나옴  $\times$  찾아야 함!!

- database
  - bucket이 너무 작음 → overflow때매 성능 저하
  - bucket이 너무 큼 → underfull + disk space 낭비

Solution

⇒ Rehashing

- 새로운 hash function으로 re-organization
  - but, expensive

⇒ better solution: bucket수가 동적으로 수정 가능하도록!

## ▼ Dynamic Hashing

Linear Hashing :  $h()$  수정 Put, 중간 계층 X  
 Extendible Hashing :  $h()$  수정 X Put " O

- Basic Idea ⇒ Bucket (dir) 2배로 증가!

Idea Indirection level을 높이자!

- dir pointer 사용
- bucket 수 늘릴
- overflow된 bucket만 split

⇒ Hash function이 어떻게 조정해야하냐에 따라 달림

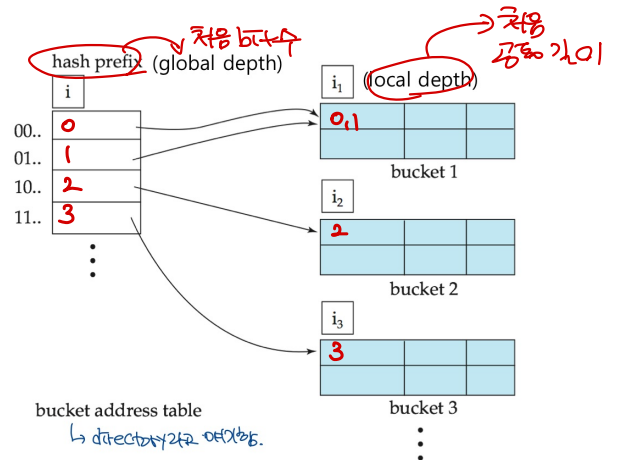
## ▼ Extendible Hashing

hash key의 postfix (prefix) 사용

- postfix한 길이 =  $i$  라고 해보자

- $Dir = 2^i$
  - $i = 0 \rightarrow$  크기감
  - db가 커질수록  $i$ 도 커짐

- General한 구조



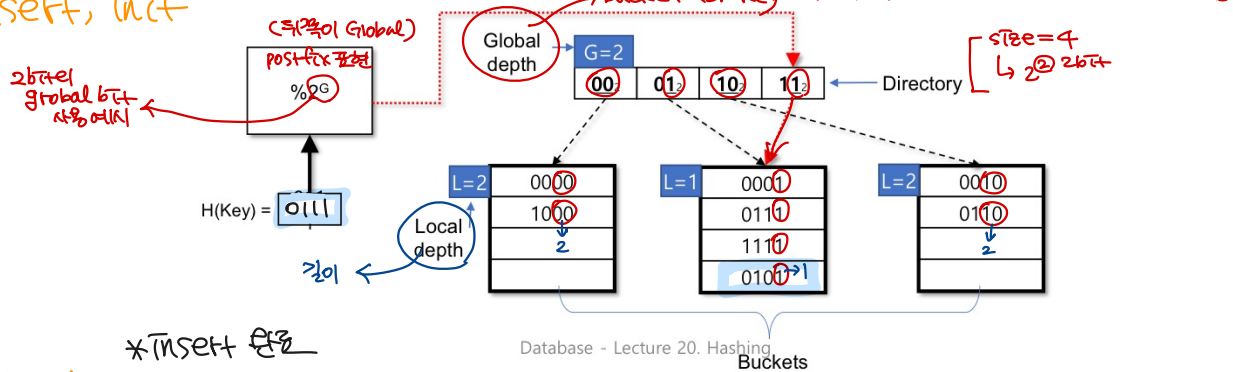
In this structure,  $i_2 = i_3 = i$  whereas  $i_1 = i - 1$

이전 가지고 있음

- Multiple Dir Entries : single bucket point

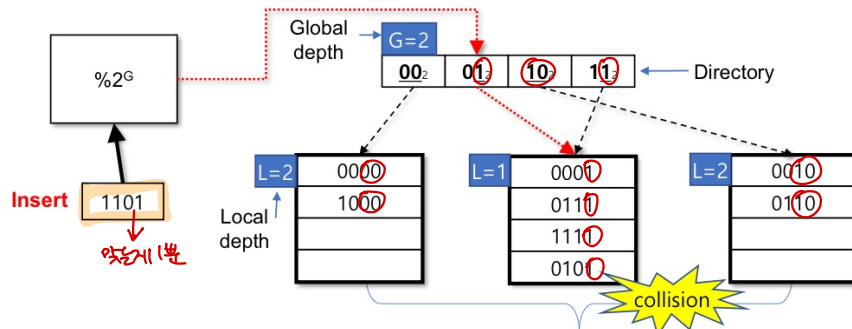
⇒ 실제 bucket 수  $< 2^i$   
 (동적으로 변경될 수 있음)

## ① Insert, $\overline{M}$



\* Insert 순서

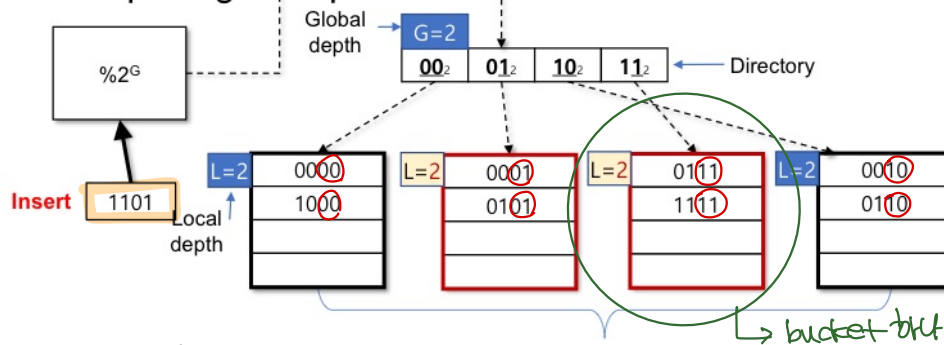
①  $G > L$   
 $\Rightarrow$  bucket 사용 할당



\* Split 시 확인 필요한 것  
 : Local depth, global depth  
 같은지 확인한 것!  
 $\Rightarrow$  다르다면 bucket만  
 같다면 d가 증가 2배

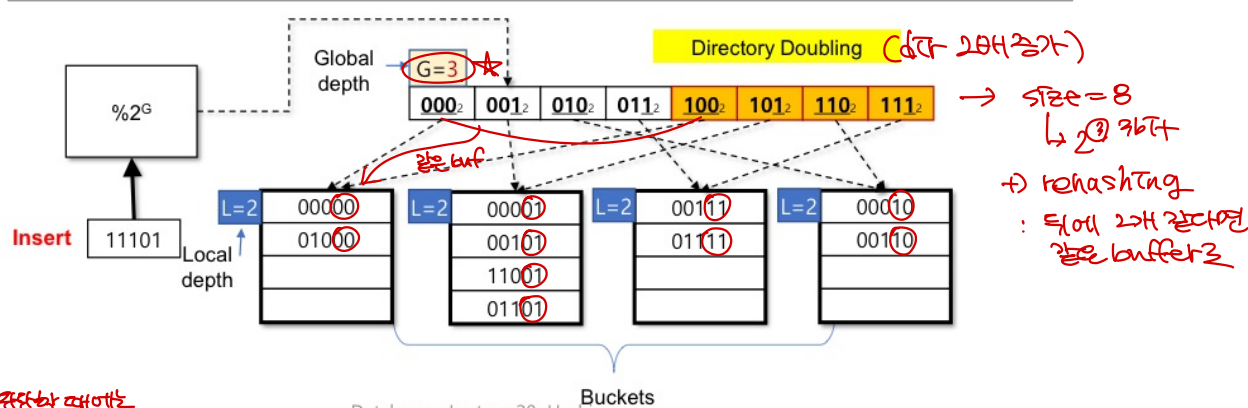
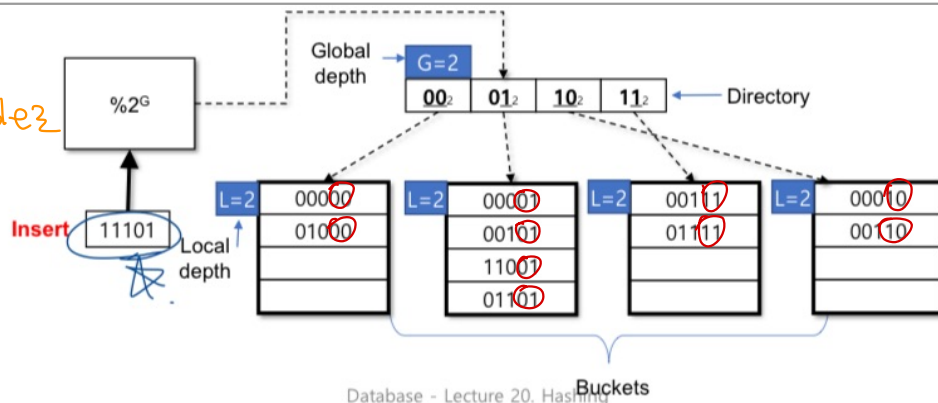
\* Bucket Split 필요

urther splitting is required if the bucket is still full



\* Bucket Split 필요

②  $G = L$   
 $\Rightarrow$  directory double 2배



\* Local depth를 판단할 때에는  
 global key를 잘 파악하고!!

• Extendible vs. Other

- ① 장점
  - [ Hash 함수이 file 크기에 비해 나빠지지 않음
  - space overhead minimum
- ② 단점
  - [ Level이 Indirection이 추가적일 필요 (record 찾기 위함)
  - [ Bucket Address Table → 매우 커질 것임 (memory가 더 크지!)
  - ↳ Disk에 큰 공간 할당함
  - ⇒ B+ tree structure 사용 → bucket address table에!
  - ⇒ But, 정적 사이즈가 영장 커졌을 때만 효과가 있다..

▼ Linear Hashing

directory 사용하지 않고 long overflow chain 문제 해결!

⇒ 일회적인 overflow page 사용됨

split 하기 위한 bucket ⇒ RR로 선택됨

⇒ "next"가 현재 가지고 있는 bucket을 split한 뒤 next pointer ++

• hash function family 사용 :  $h_0, h_1, h_2 \dots$

- [  $h_i(\text{key}) = \text{key} \bmod N$
- [  $N$  :  $i$ th bucket 개수 (2의 배수)
- [  $h_{i+1} \leftarrow h_i \rightarrow$  (다트 doubling을 함)

• Algorithm → round 기반

- [  $i$ th Level :  $N \times 2^{\text{Level}}$  →  $h_{\text{level}}$  구함!
- [ Bucket 0 ~ (Next-1)까지 split
- [ Next ~  $N_{\text{level}}$ 은 아직 split X
- [ Next =  $N_{\text{level}}$
- [ 다음 round : Level ++ , Next = 0

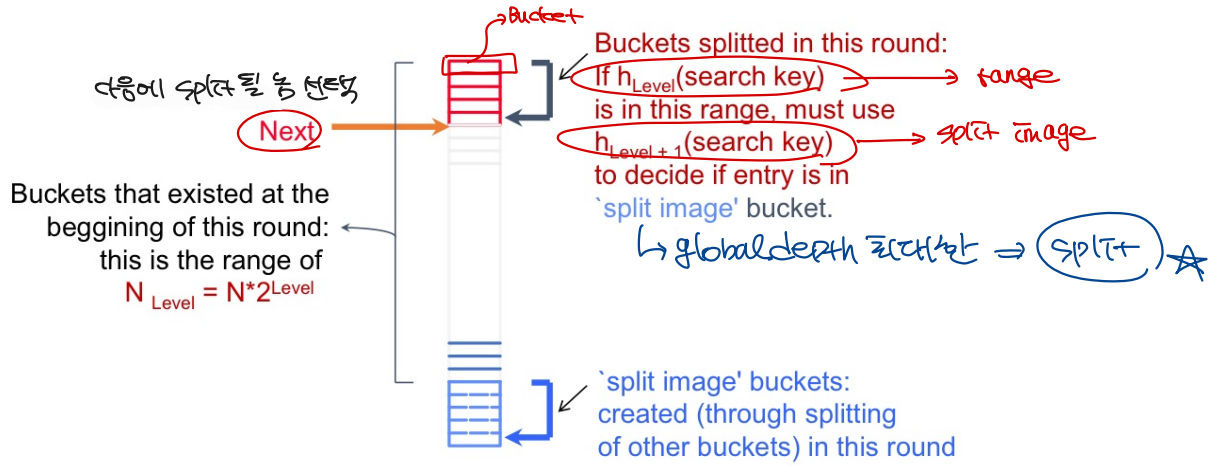
• Insert → long overflow 발생 가능

Bucket 찾을

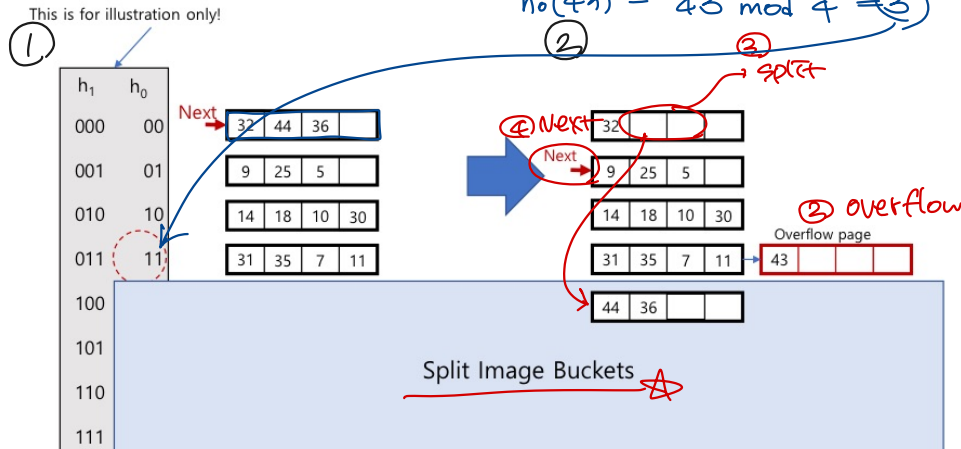
↳ if fit → DONE

↳ else

- ↳ overflow된 것이 아님!
  - ↳ overflow page 찾고 insert data ★
  - ↳ next bucket split → Next ++
  - ↳ re-distribute entries :  $h(\text{level}+1)$  사용!! ★

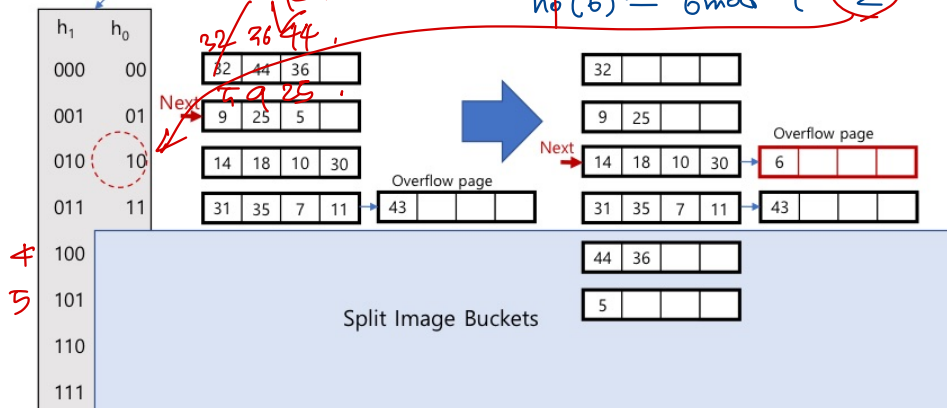


Insert 43 (101011<sub>(2)</sub>) (Level 0, N=4) ①  $2^{Level} = N \cdot 2^{Level} = 4$   
 $h_0(43) = 43 \bmod 4 = 3$

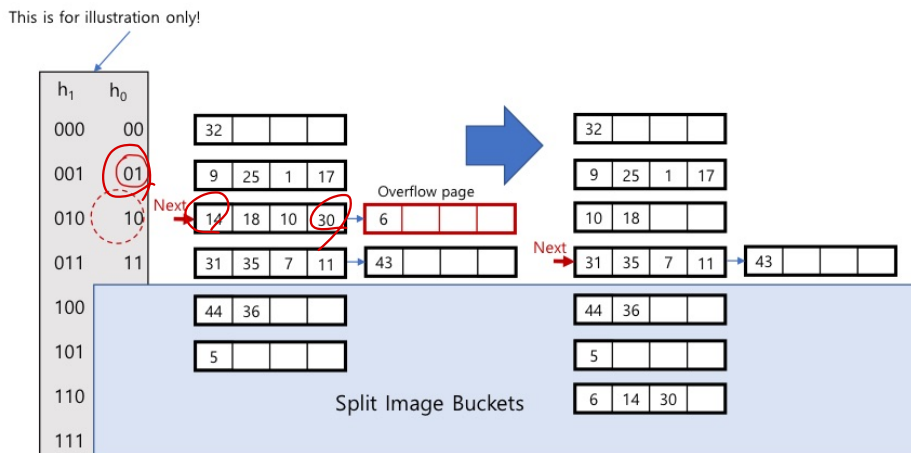


This is for illustration only!

Insert 6 (110<sub>(2)</sub>) (Level 0, N=4) ①  $2^{Level} = 4$   
 $h_0(6) = 6 \bmod 4 = 2$



Insert 65 (100001<sub>(2)</sub>) (Level 0, N=4)  $h_0(65) = 65 \div 4 = 1$



## • searching

i) entry  $r \rightarrow h_{level}(r)$

ii) 아직 01 2원소이므로 split X

① if  $h_{level}(r) \geq N_{level} \rightarrow r$ 은 0 (bucket에 있음)

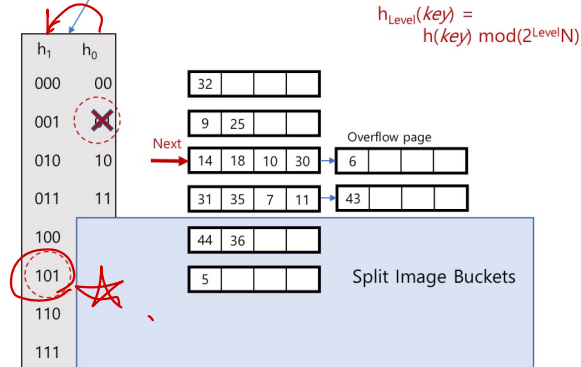
iii) " " " "

①  $h_{level}(r)$  혹은  $h_{level}(r) + N_{level}$

②  $h_{level+1}(r)$ 로 찾은 것을 재사용해야 함.

Search 5 (00101<sub>(2)</sub>) (Level 0, N=4)

This is for illustration only!



\* re-organization  $\rightarrow$  cost  $\uparrow$

Insertion, deletion, 상대적인 Frequent