

+ "equal" $v = 25$

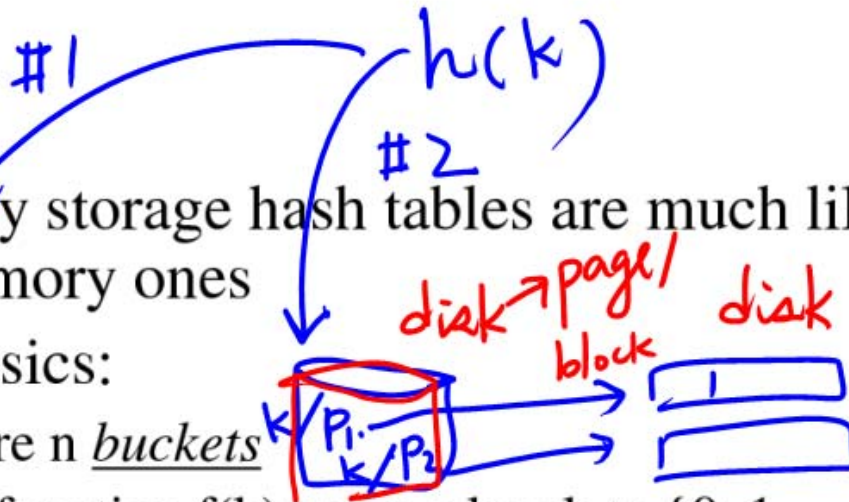
+ range $v \geq 25$

+ $v = 25$ ✓

X $v \geq 25$?

Hash Tables

- Secondary storage hash tables are much like main memory ones
- Recall basics:
 - There are n buckets
 - A hash function $f(k)$ maps a key k to $\{0, 1, \dots, n-1\}$
 - Store in bucket $f(k)$ a pointer to record with key k
- Secondary storage: bucket = block, use overflow blocks when needed



Hash Table Example

- Assume 1 bucket (block) stores 2 keys + pointers

- $h(\underline{e})=0$
- $h(b)=h(f)=1$
- $h(g)=2$
- $h(a)=h(c)=3$

$$h("x") = 2$$

nothing.

select
from
where
name='eric'

4 buckets.

name attr

0	e	→ [Eric Wolff]
1	b	
	f	
2	g	
	x?	
3	a	
	c	

Insertion in Hash Table

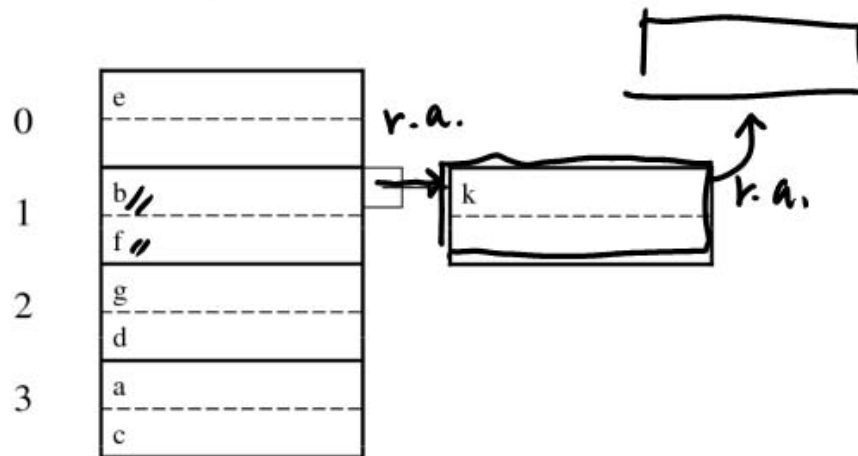
- Place in right bucket, if space
- E.g. $h(d)=2$
 $\swarrow \quad \searrow$

0	e	
1	b	f
# 2	g	d
3	a	c

} block

Insertion in Hash Table

- Create overflow block, if no space
- E.g. $h(\underline{k})=1$



- More over-flow blocks may be needed

Hash Table Performance

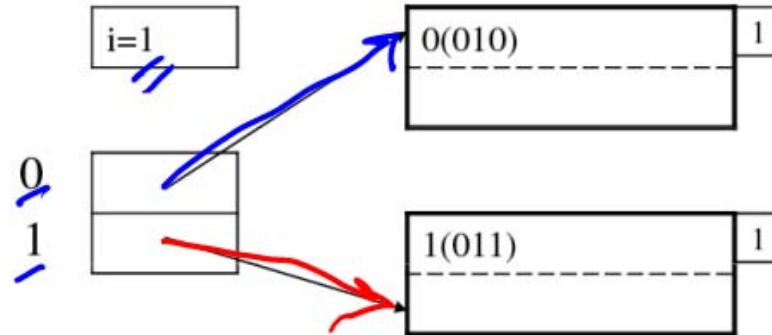
- Excellent, if no overflow blocks
- Degrades considerably when number of keys exceeds the number of buckets (I.e. many overflow blocks).

① Extensible Hash Table ② Linear

- Allows hash table to grow, to avoid performance degradation
 - Assume a hash function h that returns numbers in $\{0, \dots, 2^k - 1\}$
 - Start with $n = 2^i \ll 2^k$, only look at first i most significant bits
- Handwritten notes and diagrams:
- k bits $k=4$
 - 0000 ... 1111 (with a red bracket indicating 16 buckets)
 - i bits (circled in red)
 - MSB (Most Significant Bit) and LSB (Least Significant Bit) labels with an arrow pointing to the right.
 - Linear (underlined)
 - Diagram showing hash values h_1 and h_2 with their first i bits highlighted in blue.

Extensible Hash Table

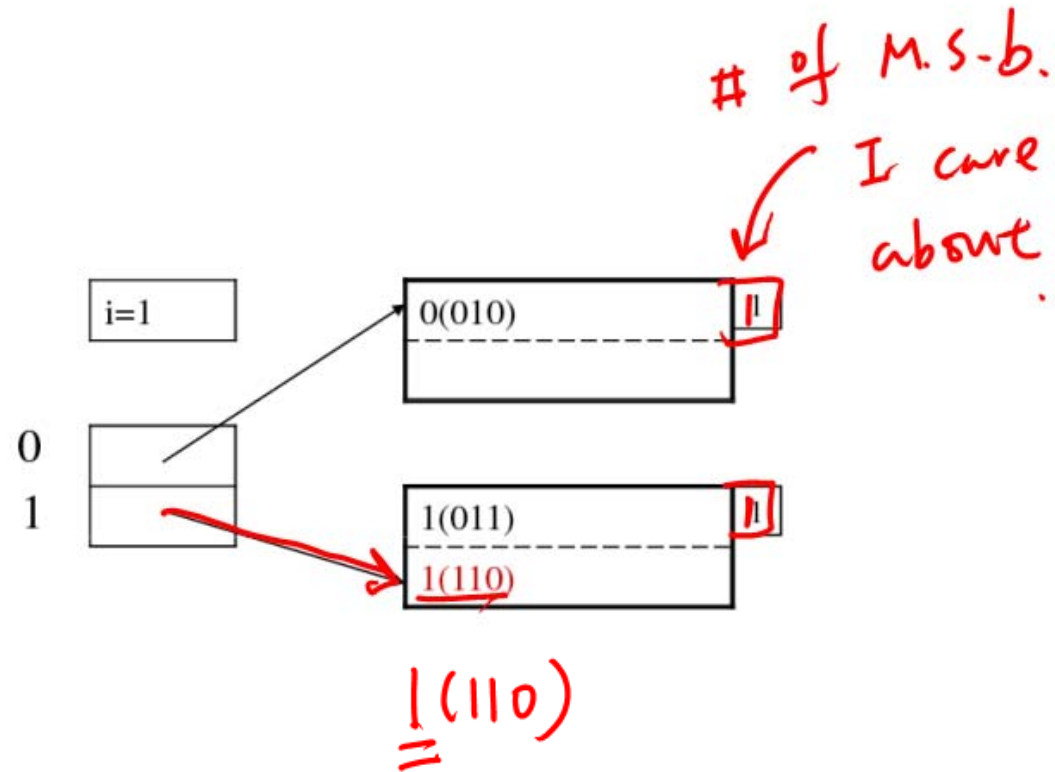
- E.g. $i=1$, $n=2$, $k=4$ $\begin{array}{c} 0010 \\ 1011 \end{array}$
- $n = 2^i$ $\begin{array}{c} 1 \\ 2 \end{array}$ max



- Note: we only look at the first bit (0 or 1)

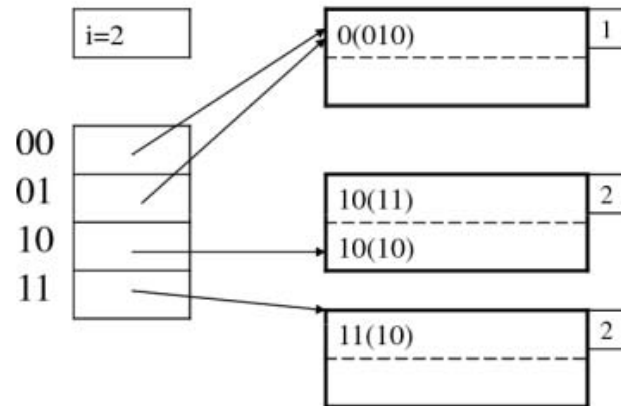
Insertion in Extensible Hash Table

- Insert 1110



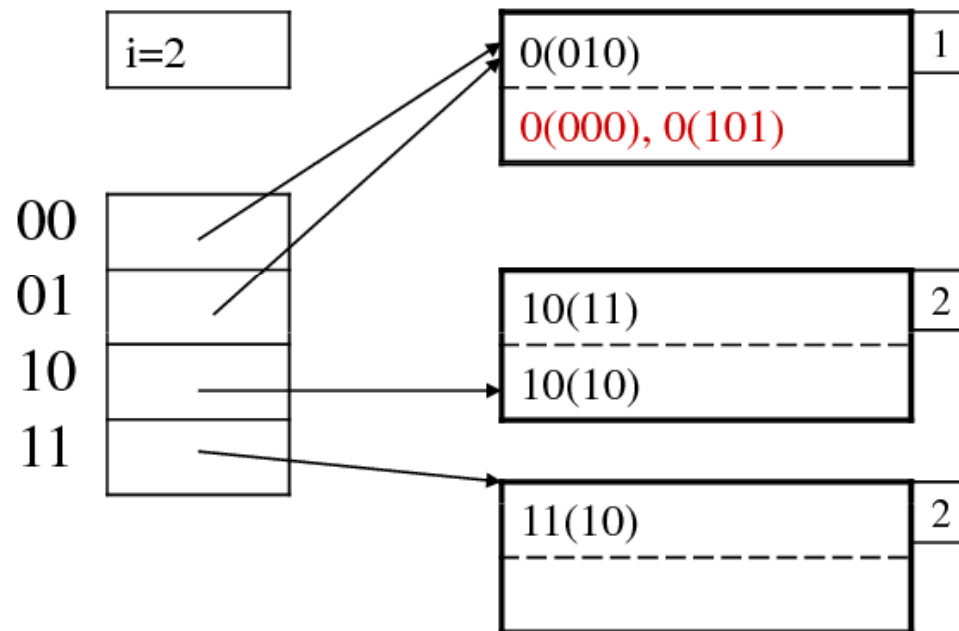
Insertion in Extensible Hash Table

- Now insert 1010 (cont.)



Insertion in Extensible Hash Table

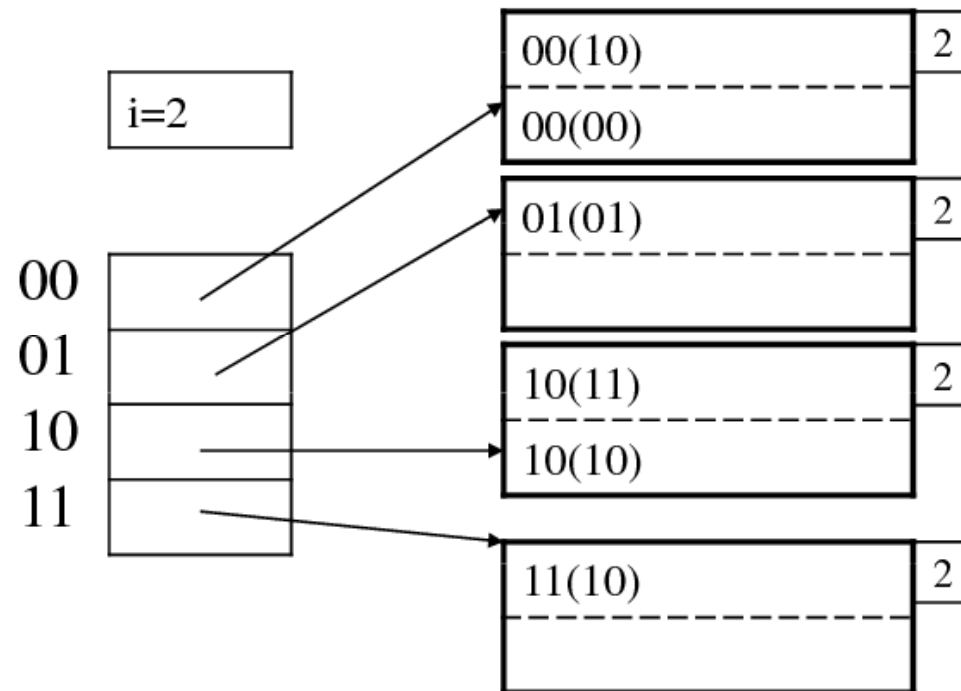
- Now insert 0000, then 0101



- Need to split block

Insertion in Extensible Hash Table

- After splitting the block



Performance Extensible Hash Table

- No overflow blocks: access always one read
- BUT:
 - Extensions can be ^{2x} costly and ^{no longer fit,} disruptive
 - After an extension table may no longer fit in memory

CS411

Database Systems

8: Query Processing

Why Do We Learn This?

Engine

top-3

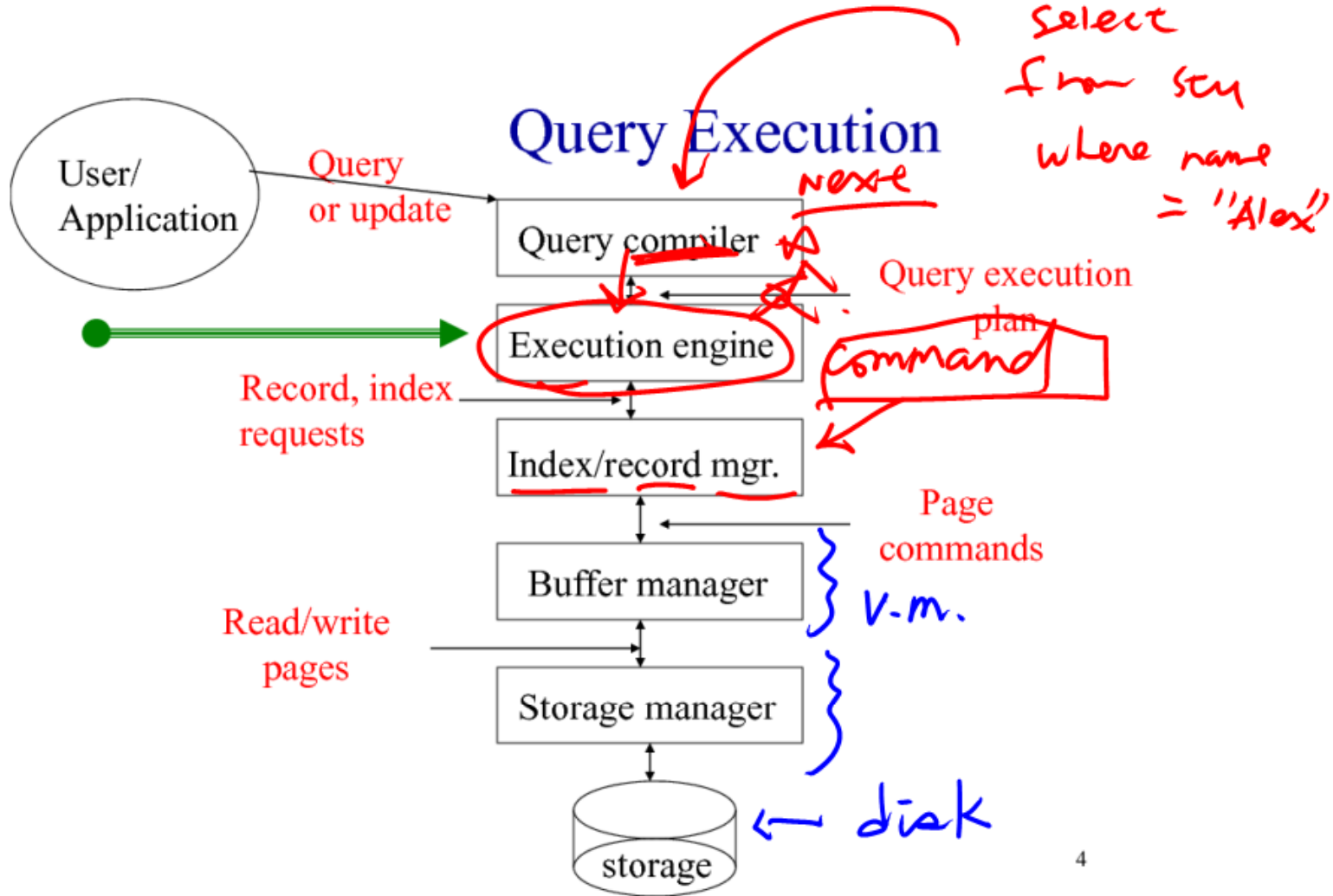
① Work at G/
M\$/Oracle

"performance",
"performance",
"performance".

② insight.
of SQL tune

Outline

- Logical/physical operators
- Cost parameters and sorting
- One-pass algorithms
- Nested-loop joins
- Two-pass algorithms



Logical v.s. Physical Operators

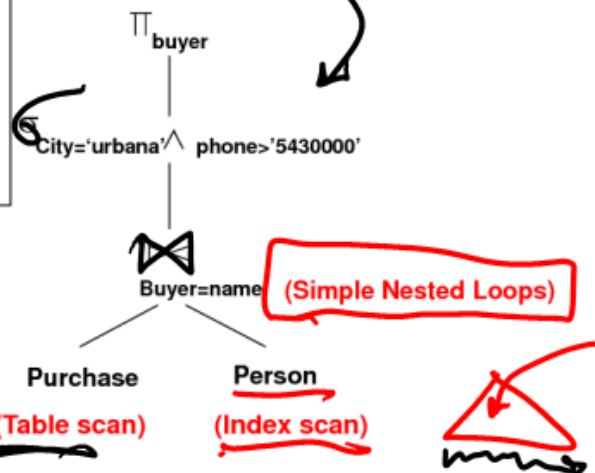
- Logical operators \bowtie, σ, π, R
 - what they do
 - e.g., union, selection, project, join, grouping
- Physical operators
 - how they do it \bowtie sort-merge
 - e.g., nested loop join, sort-merge join, hash join, index join



Query Execution Plans

compiler

```
SELECT S.sname
FROM Purchase P, Person Q
WHERE P.buyer=Q.name AND
      Q.city='urbana' AND
      Q.phone > '5430000'
```



Query Plan:

- logical tree
- implementation choice at every node
- scheduling of operations.

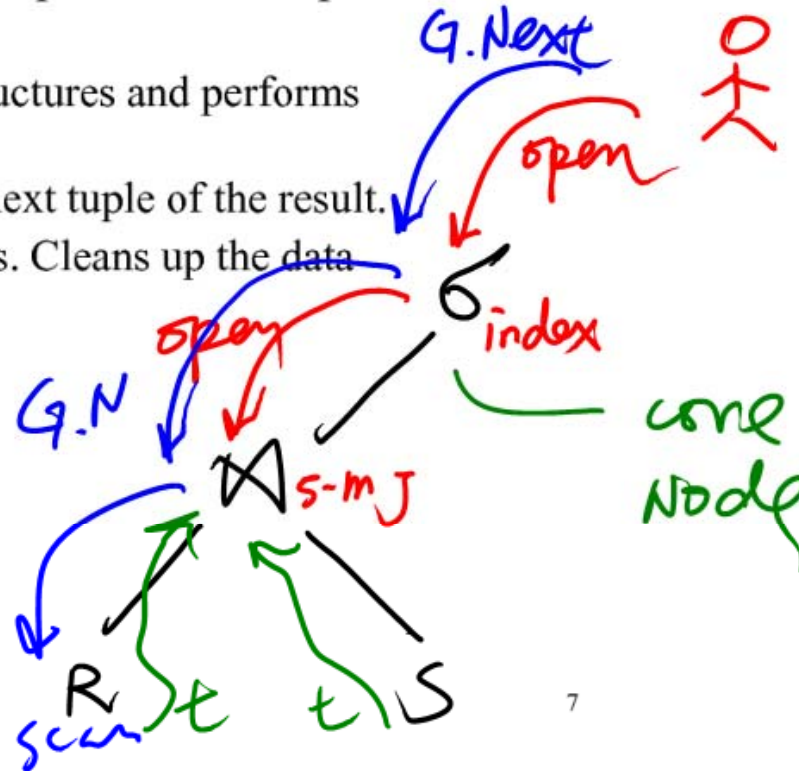
Some operators are from relational algebra, and others (e.g., scan, group) are not.

How do We Combine Operations?

- The **iterator model**. Each operation is implemented by 3 functions:

- **Open**: sets up the data structures and performs initializations
- **GetNext**: returns the the next tuple of the result.
- **Close**: ends the operations. Cleans up the data structures.

- Enables pipelining!



Cost Parameters



R



8M

- Cost parameters
 - M = number of blocks that fit in main memory
 - $B(R)$ = number of blocks holding R
 - $T(R)$ = number of tuples in R
 - $V(R, a)$ = number of distinct values of the attribute a
- Estimating the cost:
 - Important in optimization (next lecture)
 - Compute I/O cost only
 - We compute the cost to *read* the tables
 - We don't compute the cost to *write* the result (because pipelining)

all alg $[B(R) < M]$
lg Sorting

- Two pass multi-way merge sort
- Step 1:
 - Read M blocks at a time, sort, write
 - Result: have runs of length M on disk
- Step 2:
 - Merge $M-1$ at a time, write to disk
 - Result: have runs of length $M(M-1) \approx M^2$
- Cost: $3B(R)$, Assumption: $B(R) \leq M^2$

lg²

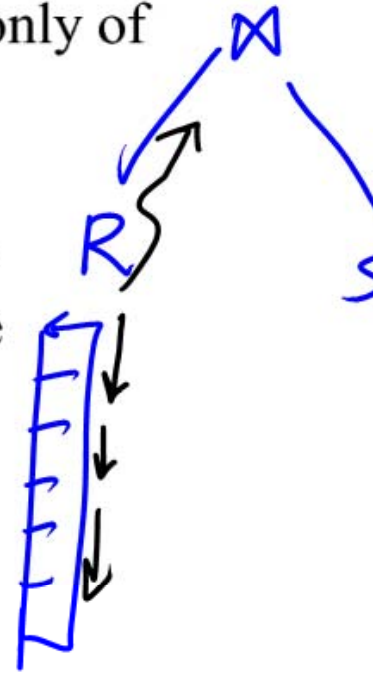
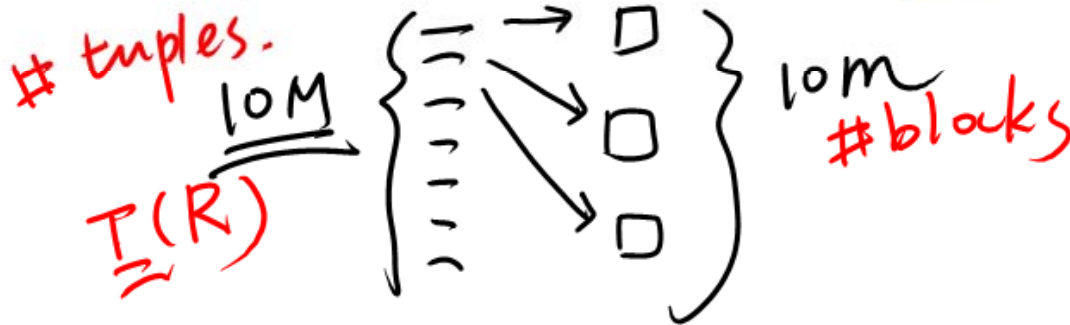
[illegible]

③ Requirement $20 \leq 50$

$$\frac{B(R)}{M} \leq M \Rightarrow B(R) \leq M^2 \quad \underline{\underline{100 \text{ blocks}}}$$

Scanning Tables

- The table is clustered (I.e. blocks consists only of records from this table): $B(R)$
 - Table-scan: if we know where the blocks are
 - Index scan: if we have index to find the blocks
- The table is unclustered (e.g. its records are placed on blocks with other tables) $T(R)$
 - May need one read for each record $T(R)$



Cost of the Scan Operator

- Clustered relation:
 - Table scan: $B(R)$; to sort: $3B(R)$
 - Index scan: $B(R)$; to sort: $B(R)$ or $3B(R)$
- Unclustered relation
 - $T(R)$; to sort: $T(R) + 2B(R)$

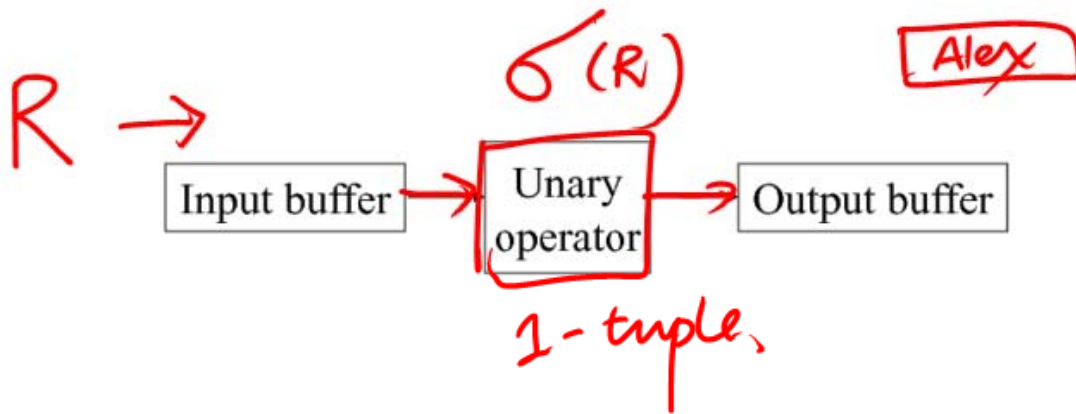
One pass algorithms

One-pass Algorithms

Selection $\sigma(R)$, projection $\Pi(R)$

- Both are tuple-at-a-Time algorithms
- Cost: $B(R)$

$\sigma_{name = 'alex'}$ Student



One-pass Algorithms

Duplicate elimination $\delta(R)$

- Need to keep a dictionary in memory:

- balanced search tree
- hash table
- etc

- Cost: $B(R)$

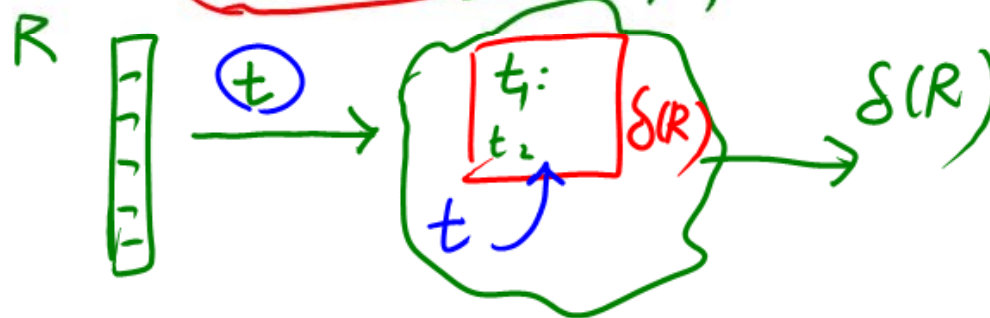
- Assumption: $B(\delta(R)) \ll M$

Select distinct ...
from ...
where

1-pass
exp



\Rightarrow 2-pass



One-pass Algorithms

Grouping: $\gamma_{\text{city, sum(price)}}(R)$

- Need to keep a dictionary in memory
- Also store the sum(price) for each city
- Cost: $B(R)$
- Assumption: number of cities fits in memory

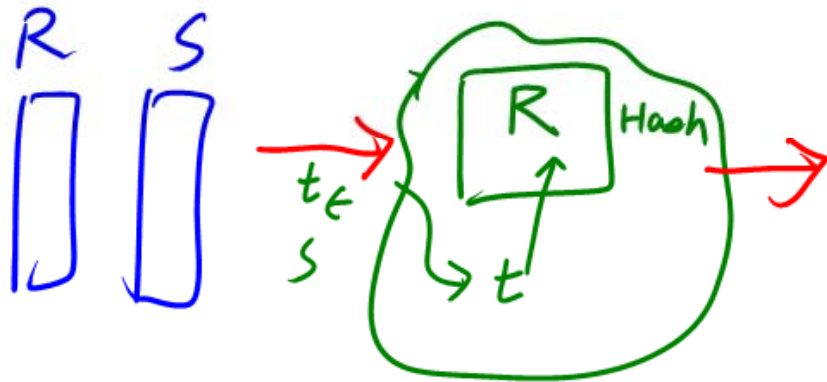
One-pass Algorithms

Binary operations: $R \cap S$, $R \cup S$, $R - S$

- Assumption: $\min(B(R), B(S)) \leq M$
- Scan one table first, then the next, eliminate duplicates

- Cost: $B(R) + B(S)$

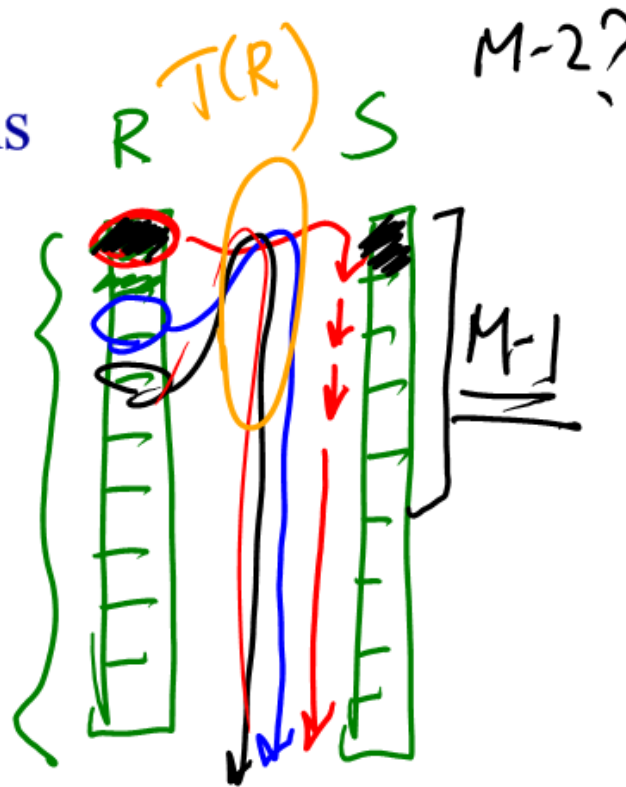
$B(R) \leq M$ or $B(S) \leq M$



Nested Loop Joins

- Tuple-based nested loop $R \bowtie S$
- R =outer relation, S =inner relation

for each tuple r in R do
 for each tuple s in S do
 if r and s join then output (r,s)



- Cost: $T(R) T(S)$, sometimes $T(R) B(S)$

$$B(R) + \cancel{T(R)} \times \cancel{T(S)} \quad (B(S))$$

R

- $B(s)$

shared by all R

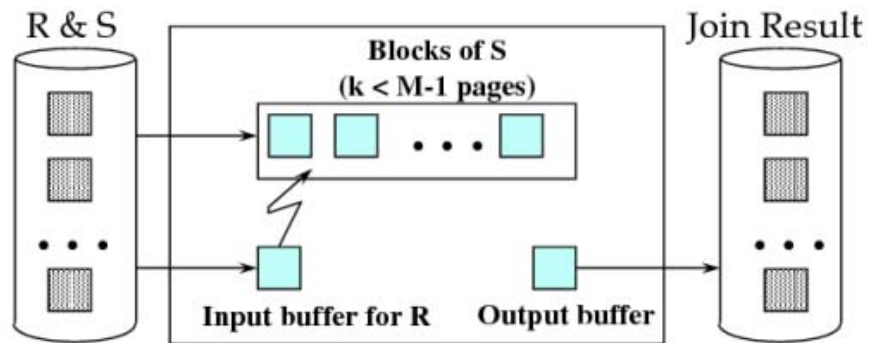
M-1

M-1


$$\frac{B(s)}{M}$$

18

Nested Loop Joins



Nested Loop Joins

- Block-based Nested Loop Join
- Cost:
 - Read S once: cost $B(S)$
 - Outer loop runs $B(S)/(M-1)$ times, and each time need to read R: costs $B(S)B(R)/(M-1)$
 - Total cost: $B(S) + B(S)B(R)/(M-1)$
- Notice: it is better to iterate over the smaller relation first— i.e., S smaller