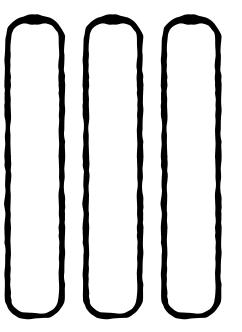
# Column-Stores



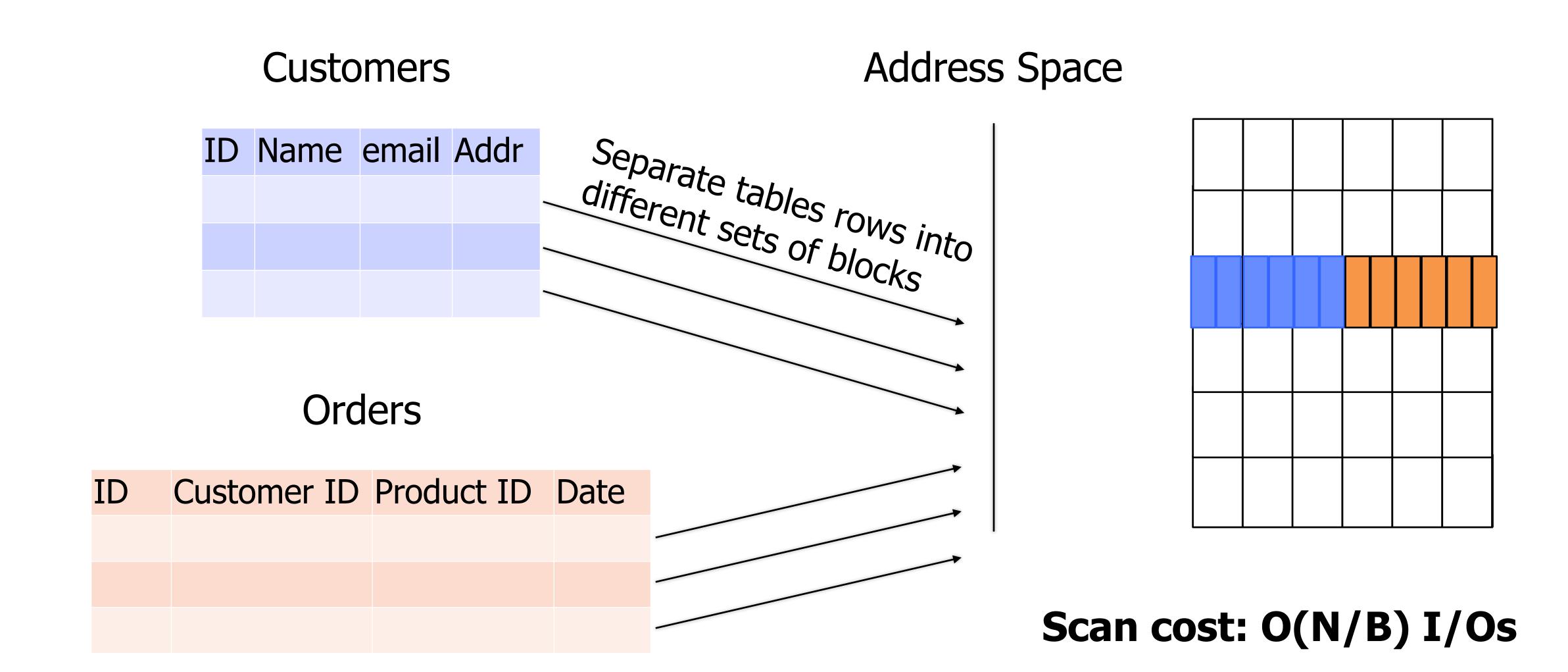
# CSC443/CSC2234 Database System Technology Niv Dayan

# **TA Office Hours**

After class today

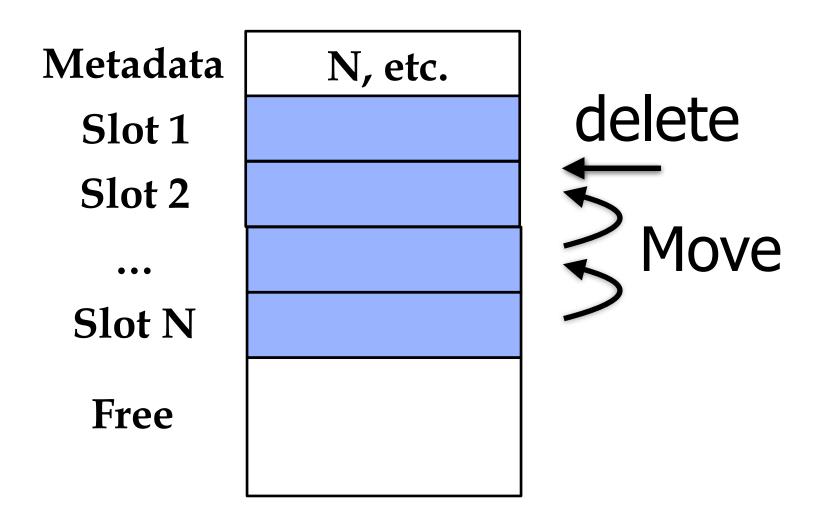
Thursday, 4-6 pm

# **Efficient Scans**



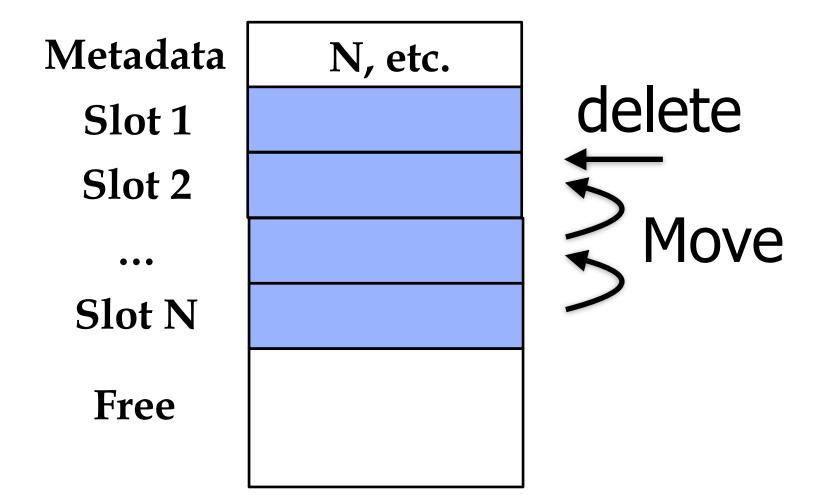
# Internal Page Organization for Fixed-Sized Rows

Entries compactly packed at front of page

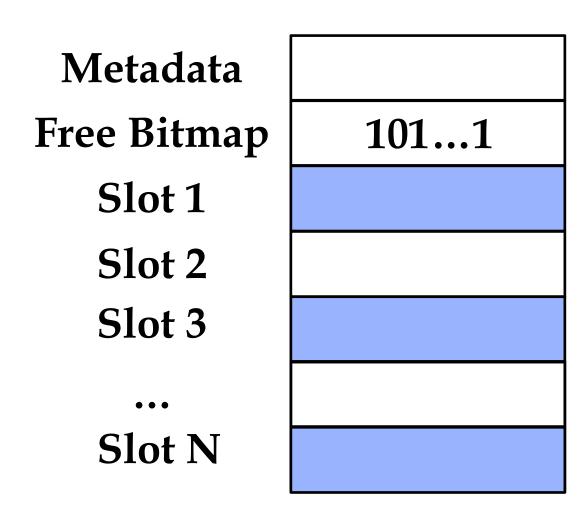


# Internal Page Organization for Fixed-Sized Rows

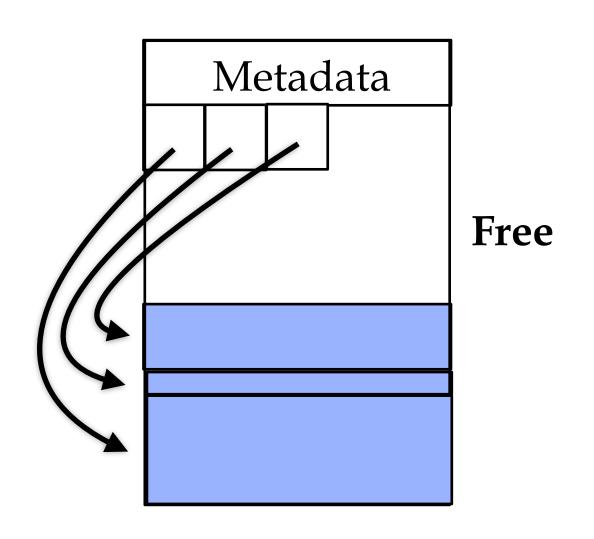
Entries compactly packed at front of page



Can also use a bitmap to mark occupied slots



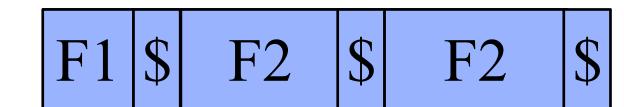
# Internal Page Organization for Variable-Sized Rows



# Variable-Sized Row Organization

Delimiters

Pointers



F1 F2 F2

Smaller
No random access

More space Random access (faster)

ID	Name	email	Addr	Salary
ID1	name1	email1	addr1	salary1
ID2	name2	email2	addr2	salary2
ID3	name3	email3	addr3	salary3

ID	Name	email	Addr	Salary

ID1 name1 email1 addr1 salary1 :ID2 name2 email2 addr2 salary2: ID3 name3 email3 addr3 salary3

Address Space



ID	Name	email	Addr	Salary

ID1 name1 email1 addr1 salary1 :ID2 name2 email2 addr2 salary2 :ID3 name3 email3 addr3 salary3

Great when queries access most columns

ID	Name	email	Addr	Salary

ID1 name1 email1 addr1 salary1 :ID2 name2 email2 addr2 salary2 :ID3 name3 email3 addr3 salary3

## Great when queries access most columns

- e.g., Select \* from Customers
- e.g., Select **Id**, **Name**, **Email**, **Addr** from Customers where salary > 10000

ID	Name	email	Addr	Salary

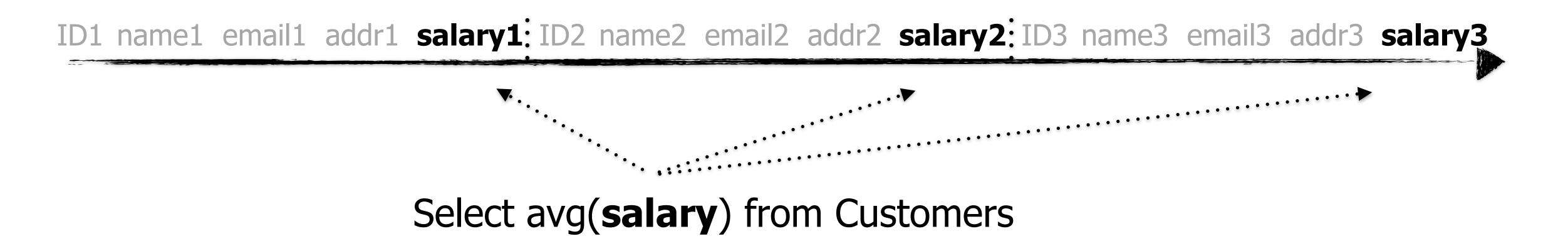
ID1 name1 email1 addr1 salary1 : ID2 name2 email2 addr2 salary2 : ID3 name3 email3 addr3 salary3

## But how about queries that examine few columns?

- e.g., Select avg(salary) from Customers
- e.g., Select sum(email) from Customers where email like '%@gmail.com';

#### **Problem?**

ID	Name	email	Addr	Salary



Problem? We only need a little of each row

ID1 name1 email1 addr1 salary1 ID2 name2 email2 addr2 salary2 ID3 name3 email3 addr3 salary3

#### 4 KB access

Select avg(salary) from Customers

Problem? We only need a little of each row

But storage access granularity is coarse

ID1 name1 email1 addr1 salary1 ID2 name2 email2 addr2 salary2 ID3 name3 email3 addr3 salary3

128B access

128B access

128B access

Select avg(salary) from Customers

Problem? We only need a little of each row

But storage access granularity is coarse

Memory can also only be accessed in cache lines

ID1 name1 email1 addr1 salary1 ID2 name2 email2 addr2 salary2 ID3 name3 email3 addr3 salary3

128B access 128B access 128B access

Select avg(salary) from Customers

Problem? We only need a little of each row

But storage access granularity is coarse

Memory can also only be accessed in cache lines

Reading more than we are interested in wastes bandwidth

ID1 name1 email1 addr1 salary1 ID2 name2 email2 addr2 salary2 ID3 name3 email3 addr3 salary3

128B access

128B access

128B access

Select avg(salary) from Customers

Problem? We only need a little of each row

But storage access granularity is coarse

Memory can also only be accessed in cache lines

Reading more than we are interested in wastes bandwidth

Solution?

ID	Name	email	Addr	Salary
ID1	name1	email1	addr1	salary1
ID2	name2	email2	addr2	salary2
ID3	name3	email3	addr3	salary3

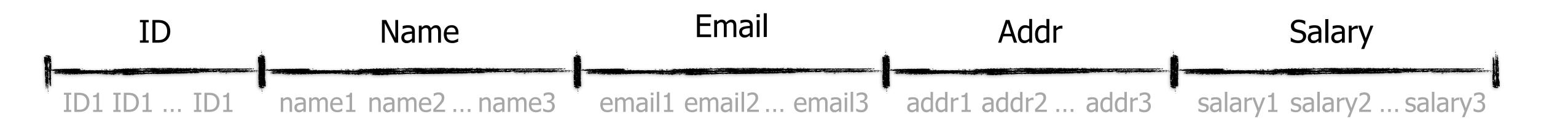
ID	Name	email	Addr	Salary

ID1 ID1 ... ID1 name1 name2 ... name3 email1 email2 ... email3 addr1 addr2 ... addr3 salary1 salary2 ... salary3

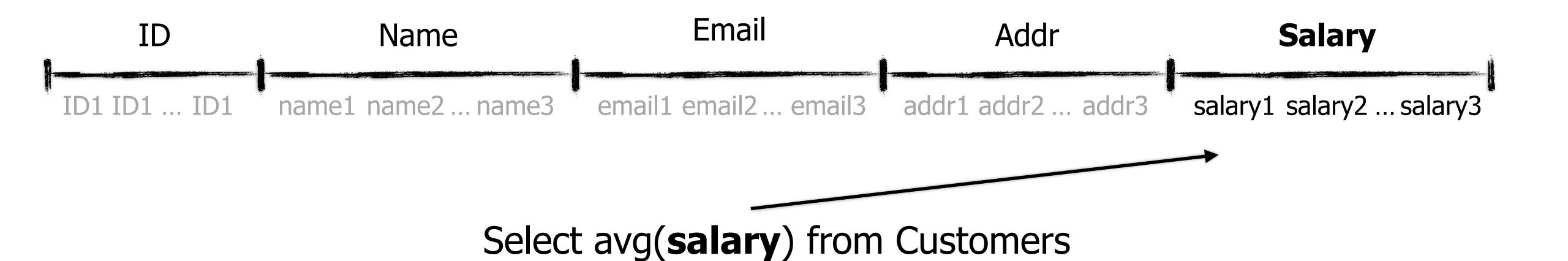
Address Space



ID	Name	email	Addr	Salary



ID	Name	email	Addr	Salary



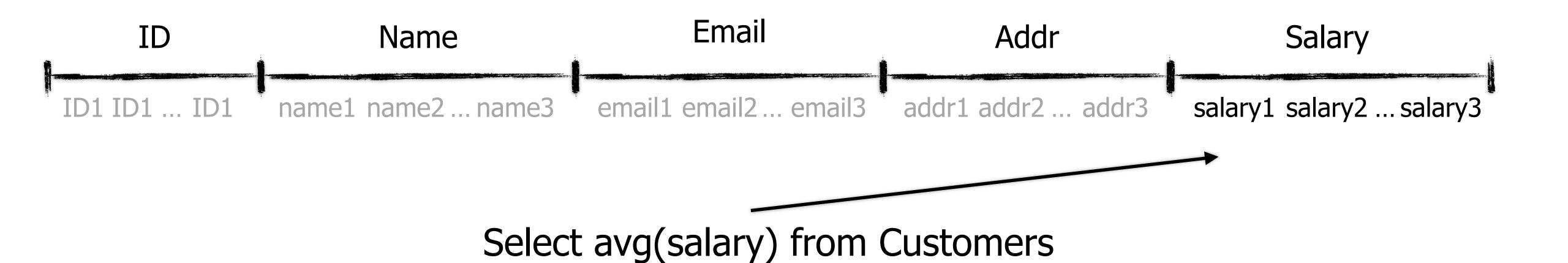
ID	Name	email	Addr	Salary



Select avg(salary) from Customers

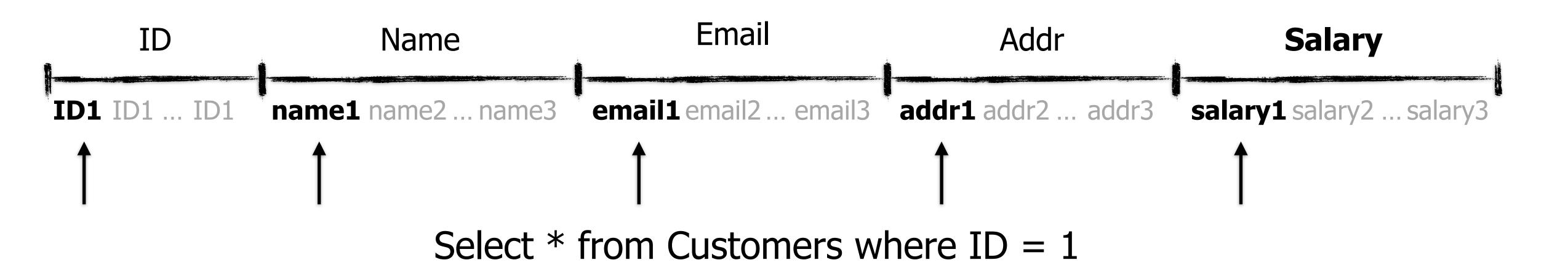
We can now skip irrelevant data!

ID	Name	email	Addr	Salary



**Any concerns?** 

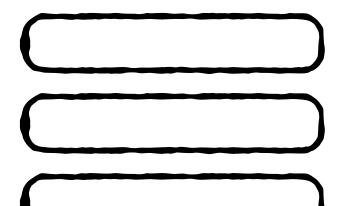
ID	Name	email	Addr	Salary



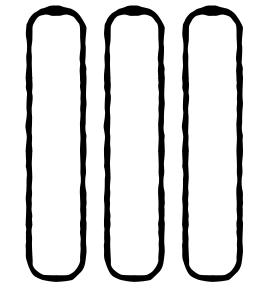
Trade-off: random I/Os for selective queries across many columns

#### The two relational database families

#### **Row-Stores**



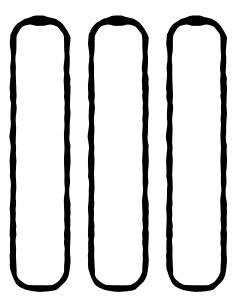
#### **Column-Stores**



Row-Stores

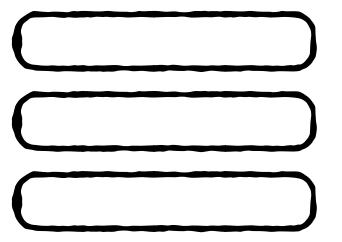
Selective queries/updates

Column-Stores



Large statistical calculations & batch updates

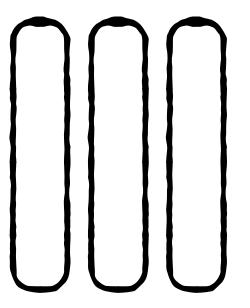
**Row-Stores** 



Selective queries/updates

Online Transaction Processing (OLTP)

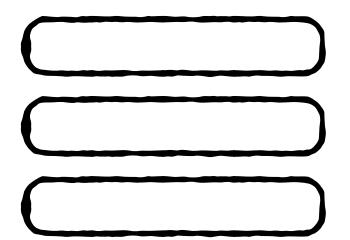
Column-Stores



Large statistical calculations & batch updates

Online Analytical Processing (OLAP)





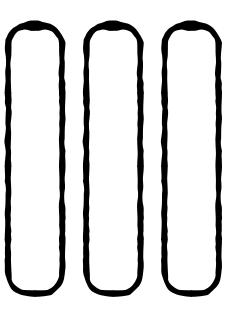
# **OLTP** examples

Banking, retail, social media





#### Column-Stores



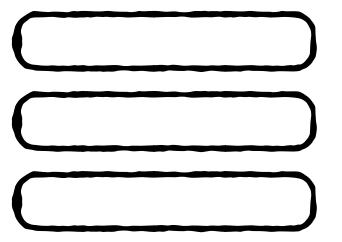
# **OLAP** examples

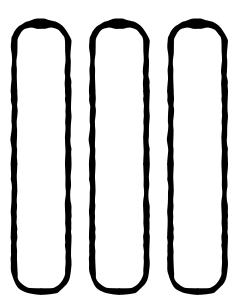
Analytics, machine learning



Row-Stores

Column-Stores

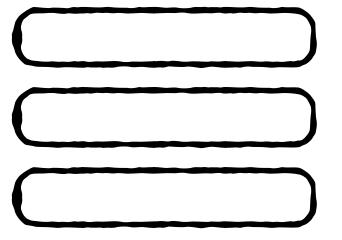




Postgres, MariaDB, etc.

MonetDB, Vectorwise, C-Store, Vertica (Today all major DBs offer this, e.g., )

**Row-Stores** 

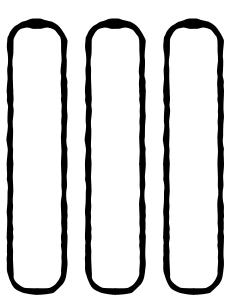


Postgres, MariaDB, etc.



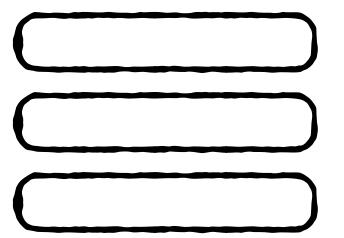
Oracle, IBM, Microsoft were here

Column-Stores



MonetDB, Vectorwise, C-Store, Vertica (Today all major DBs offer this, e.g., )



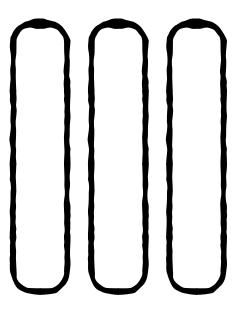


Postgres, MariaDB, etc.



Oracle, IBM, Microsoft were here But now they offer both

#### Column-Stores



MonetDB, Vectorwise, C-Store, Vertica (Today all major DBs offer this, e.g., )

First row-store
Developed and
productized

First open-source Column-store developed Initial industry adaptation of column-stores

Column-stores and row-stores are Both a norm

**1975** 

2000

2005

now

Timeline

First row-store
Developed and
productized

First open-source Column-store developed Initial industry adaptation of column-stores

Column-stores and row-stores are Both a norm

**1975** 

2000

2005

now

Why the 30 year gap?

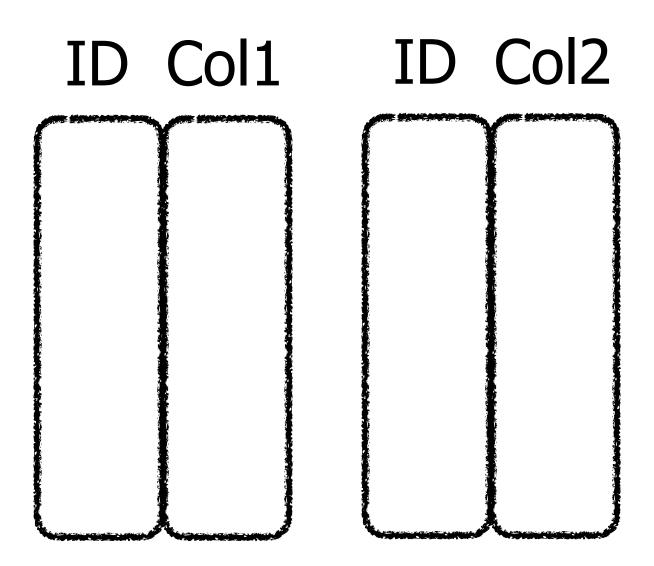
First row-store Developed and productized First open-source Column-store developed Initial industry adaptation of column-stores

Column-stores and row-stores are Both a norm

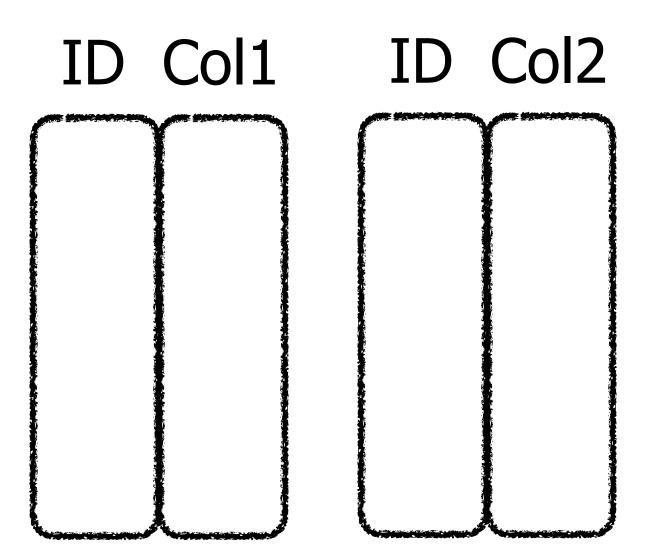
1975 2000 2005 now

Why the 30 year gap? Storage density grew, so we store far more data now. Processing it efficiently is more critical. This drives a need for specialization.

#### Should we store each column with a materialized ID?

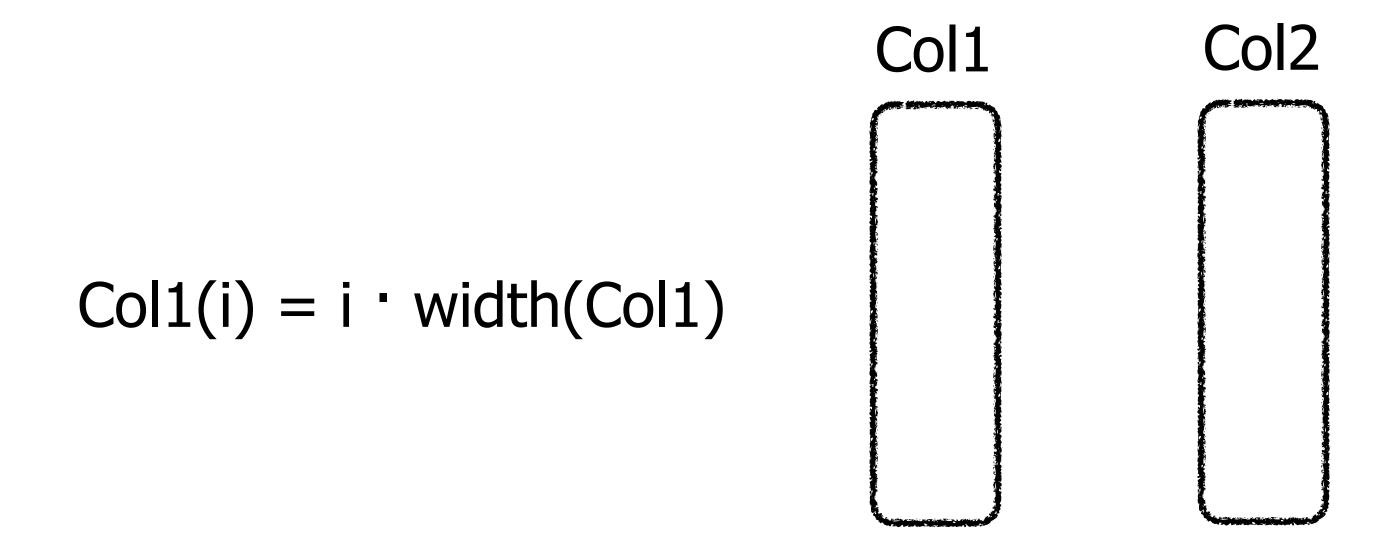


#### Should we store each column with a materialized ID?



No, this would slow down queries by reading more data.

# Instead use positional alignment

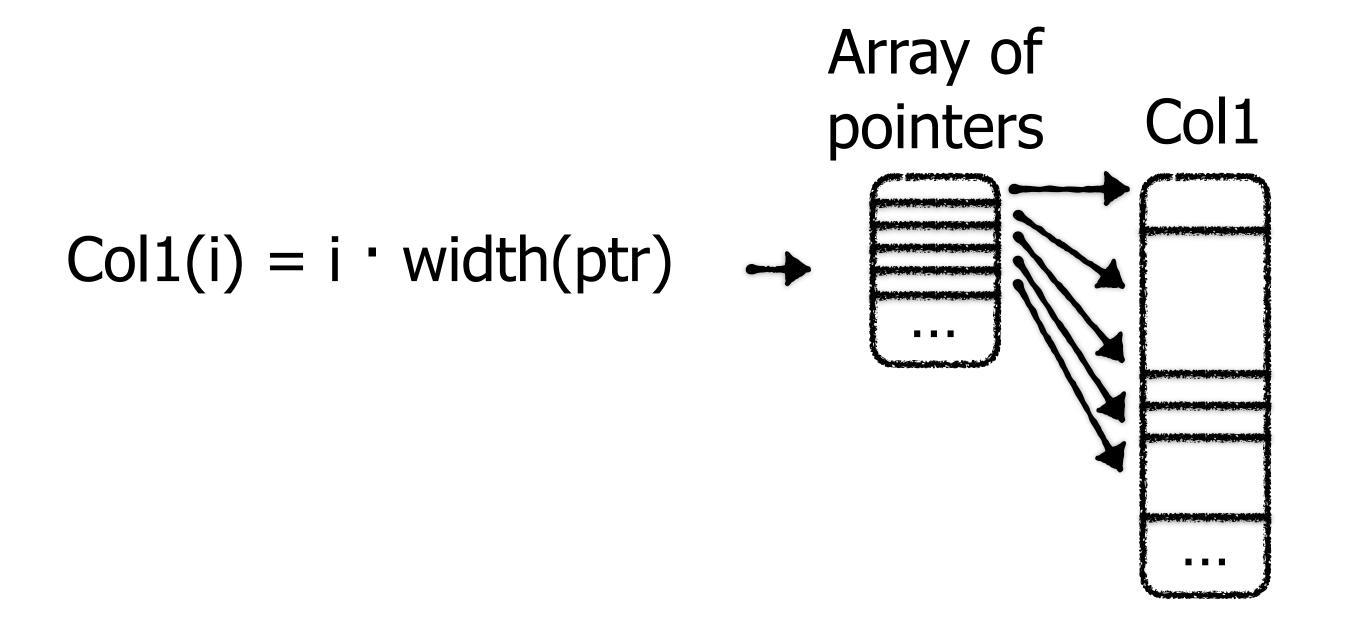


# Instead use positional alignment

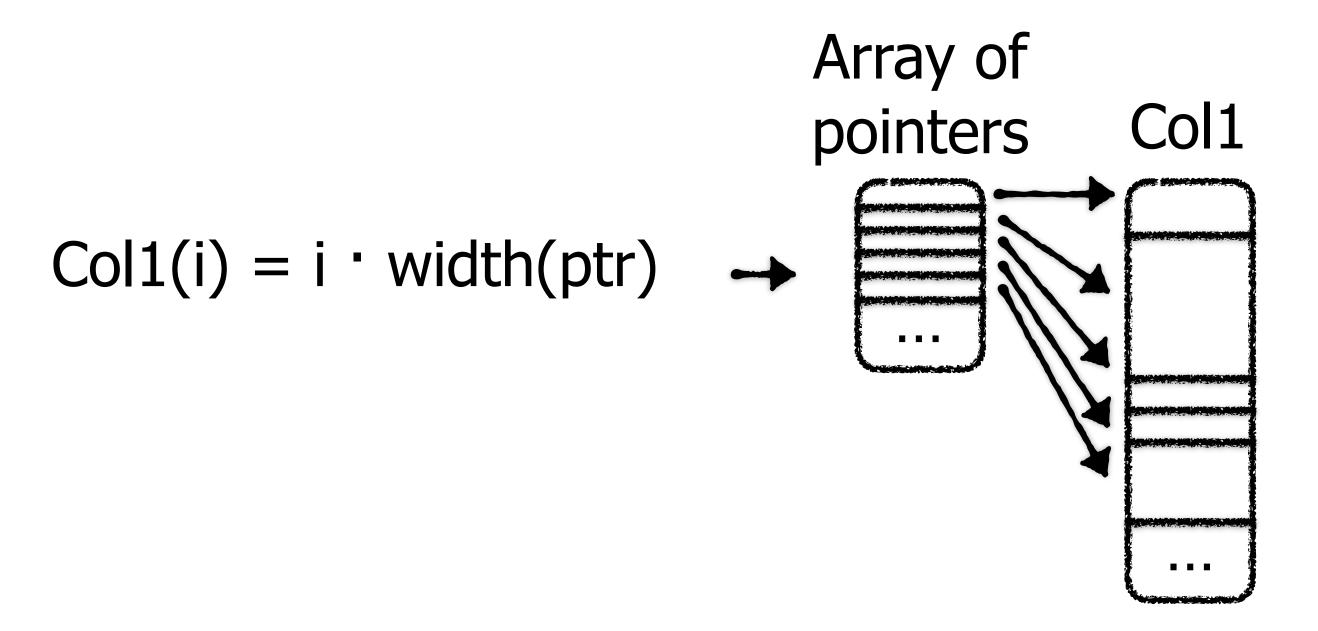
What if a column is var-length?

### In case of variable-length data, we can employ indirection

(unrecommended)

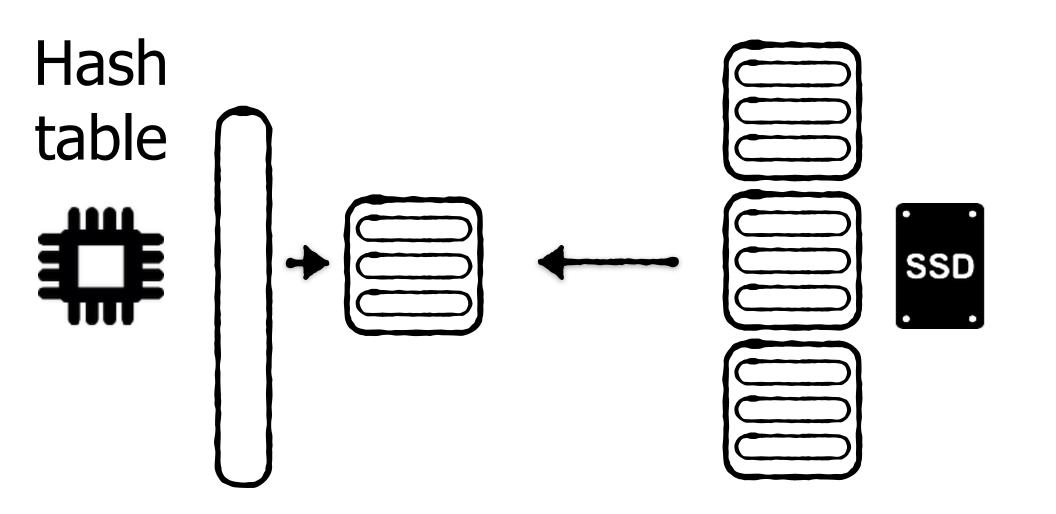


In case of variable-length data, we can employ indirection (unrecommended)

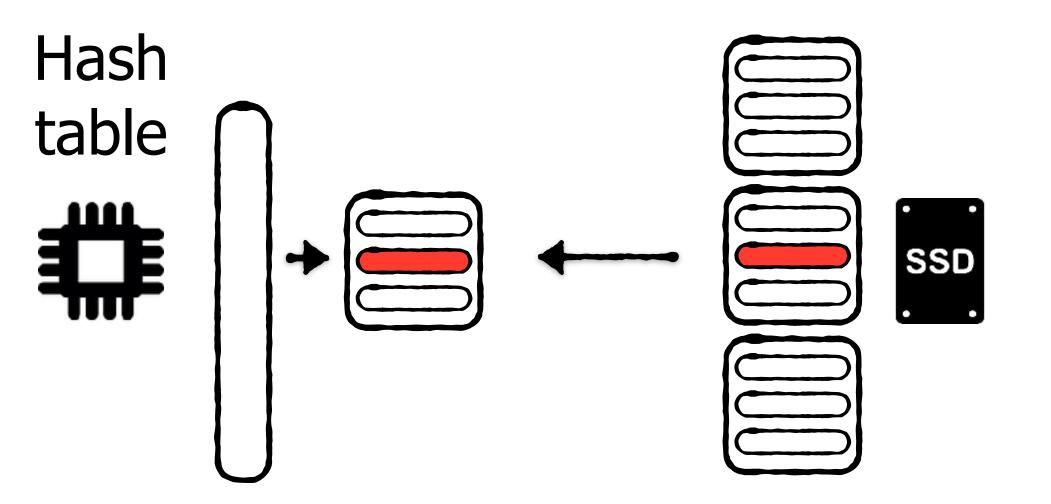


Best to define a schema with fixed-length fields if possible prevent indirection overheads

A row-store BP maps 4KB pages.



A row-store BP maps 4KB pages.

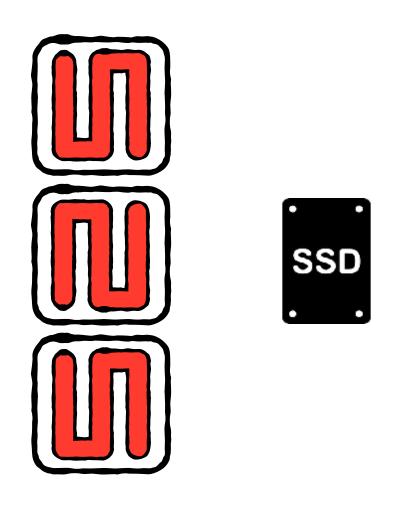


This makes sense each query is "selective" (accesses few rows)

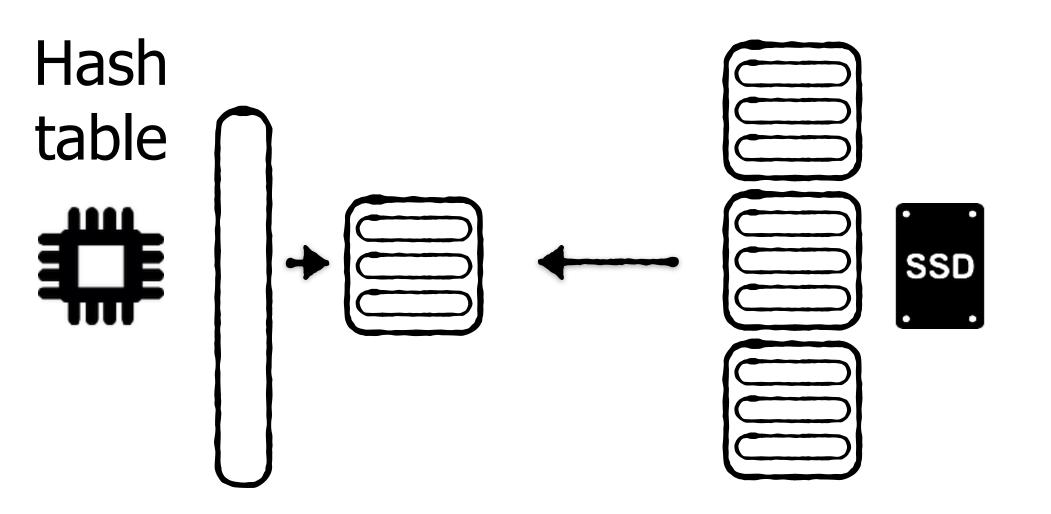
A row-store BP maps 4KB pages.

This makes sense each query is "selective" (accesses few rows)

In column-stores, each column spans multiple pages & queries are unselective

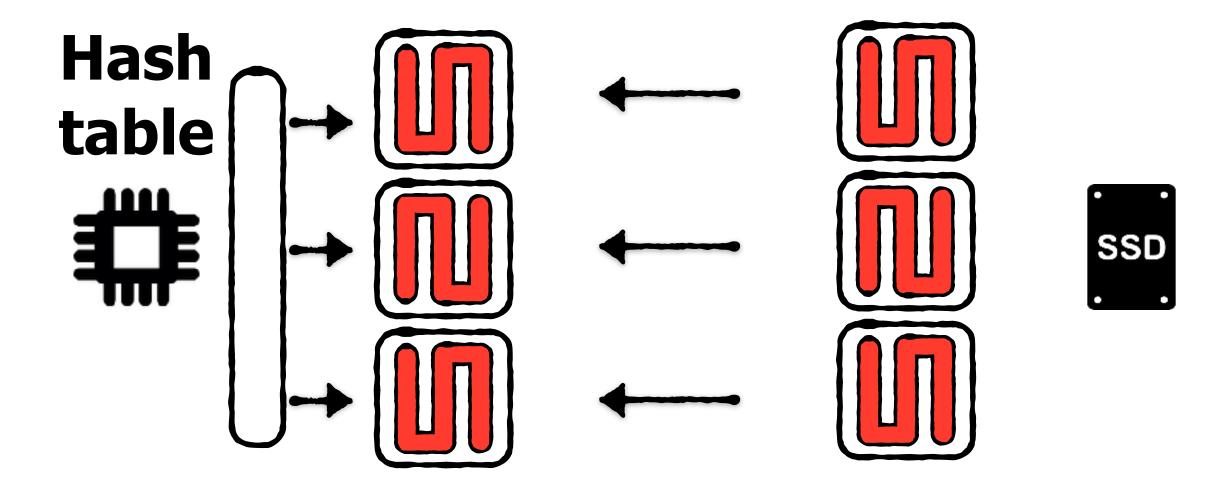


A row-store BP maps 4KB pages.



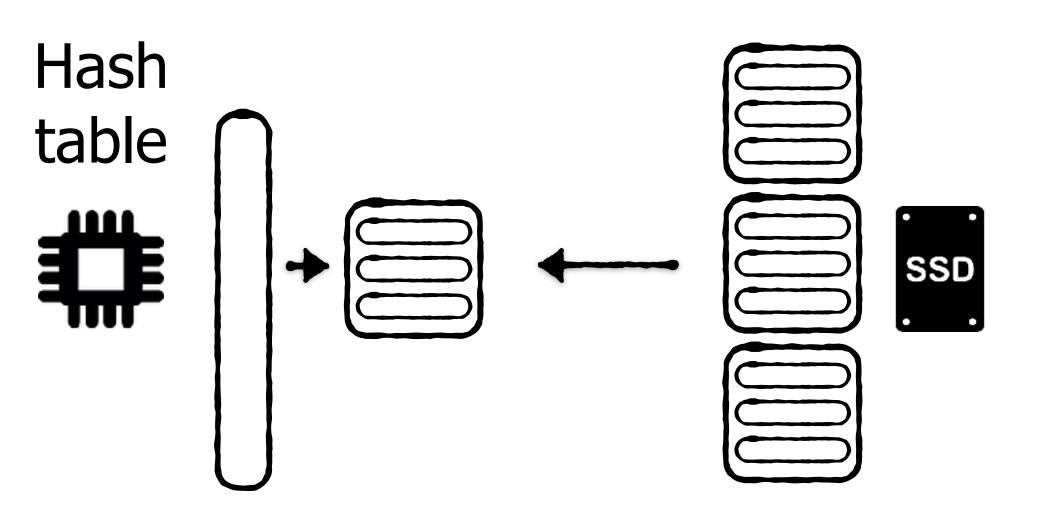
This makes sense each query is "selective" (accesses few rows)

In column-stores, each column spans multiple pages & queries are unselective



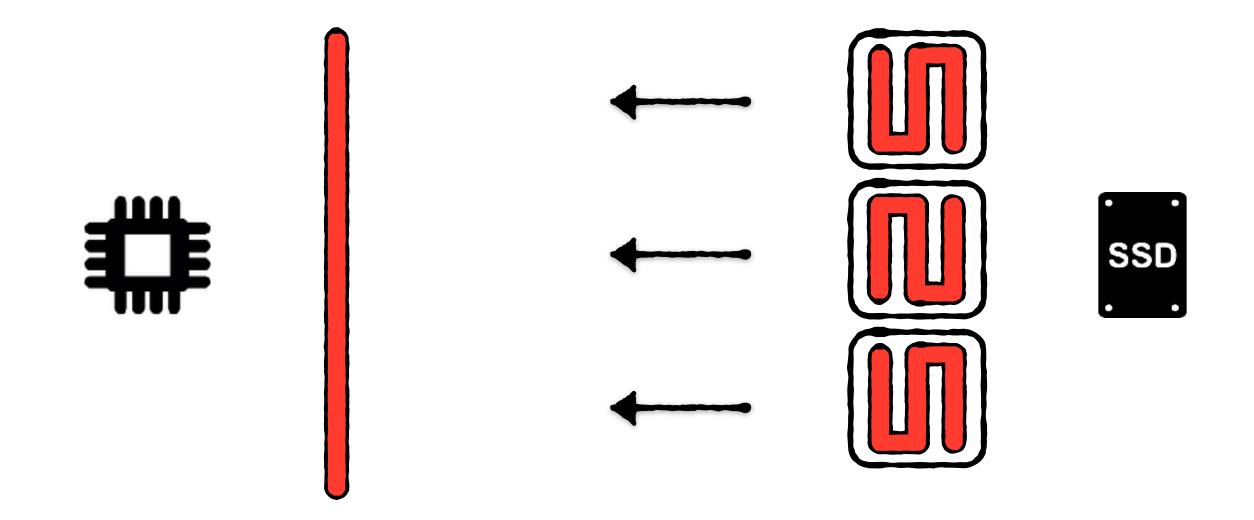
Scanning a column in memory would require hashing overheads for each page.

A row-store BP maps 4KB pages.



This makes sense each query is "selective" (accesses few rows)

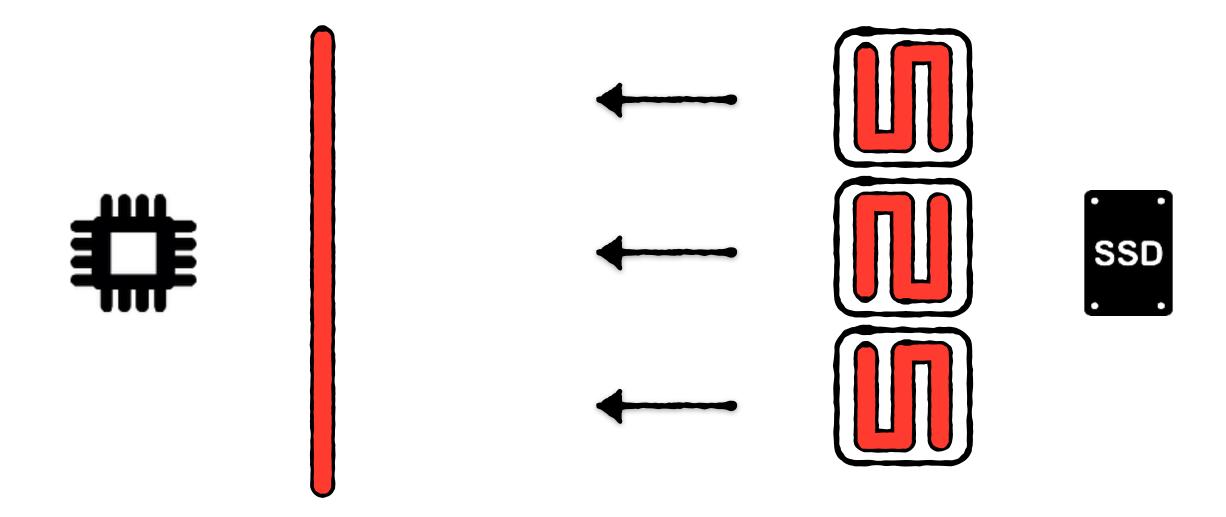
Better to read and map at columngranularity (or multiple pages thereof)



A row-store BP maps 4KB pages.

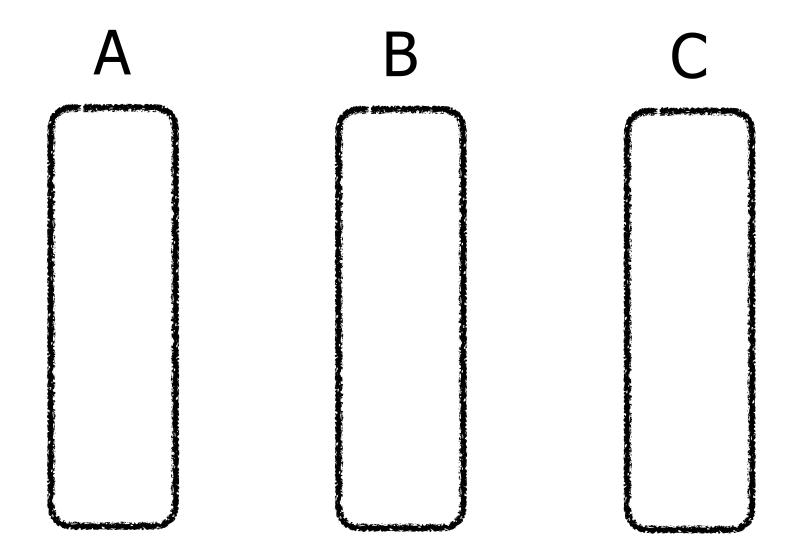
This makes sense each query is "selective" (accesses few rows)

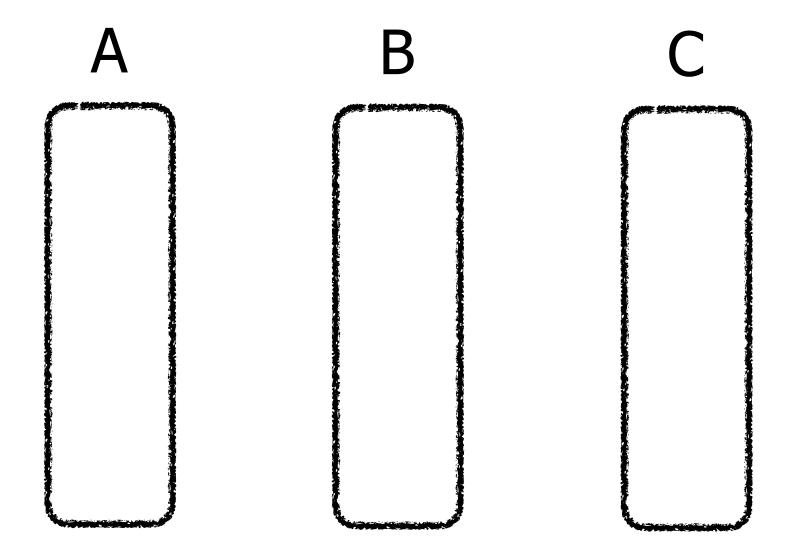
Better to read and map at columngranularity (or multiple pages thereof)



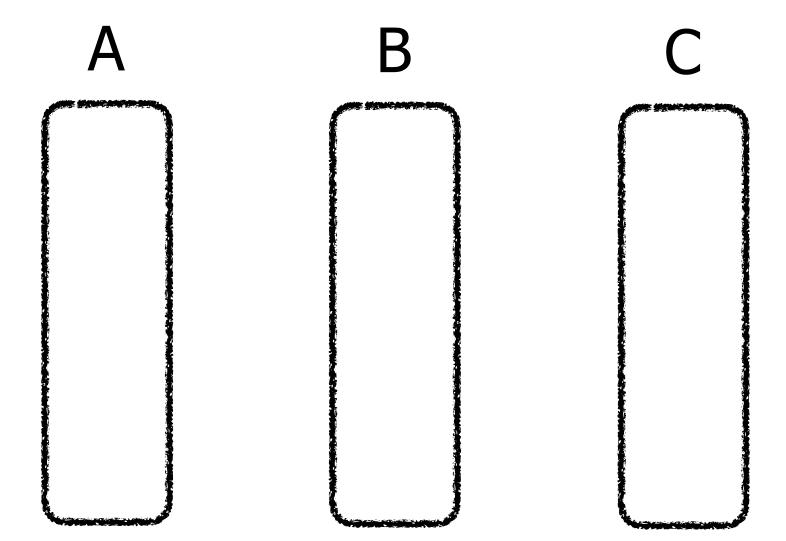
Can do this by mapping columns in virtual memory or in the BP

### How to process queries over a single table?

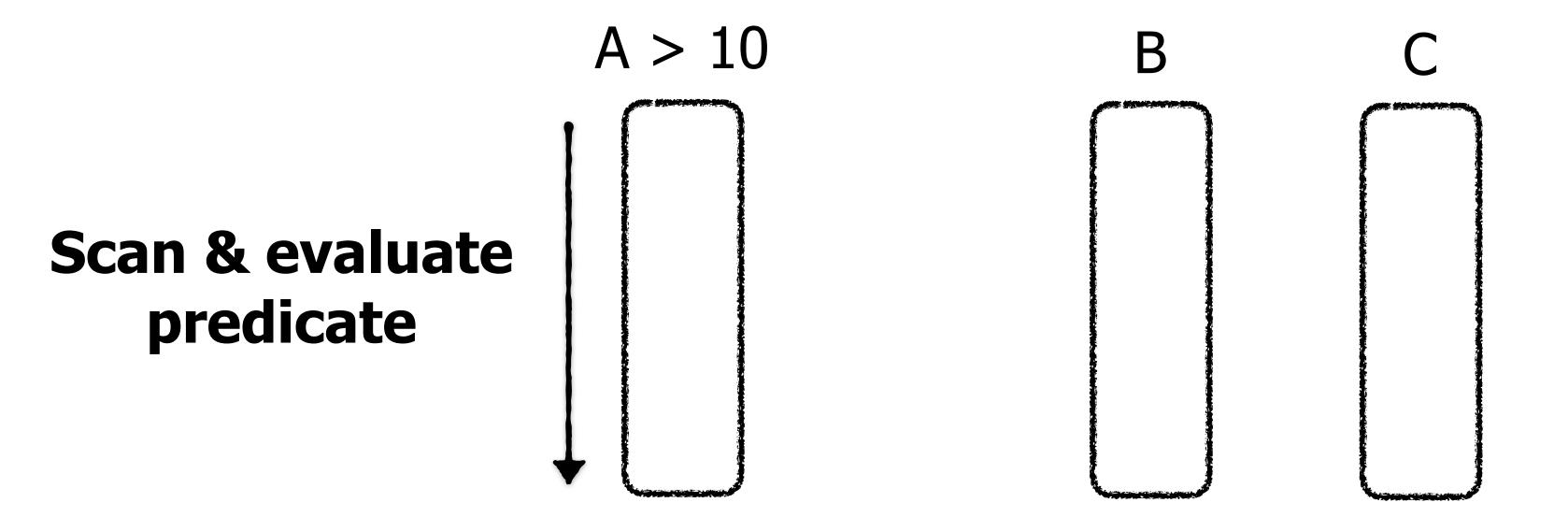




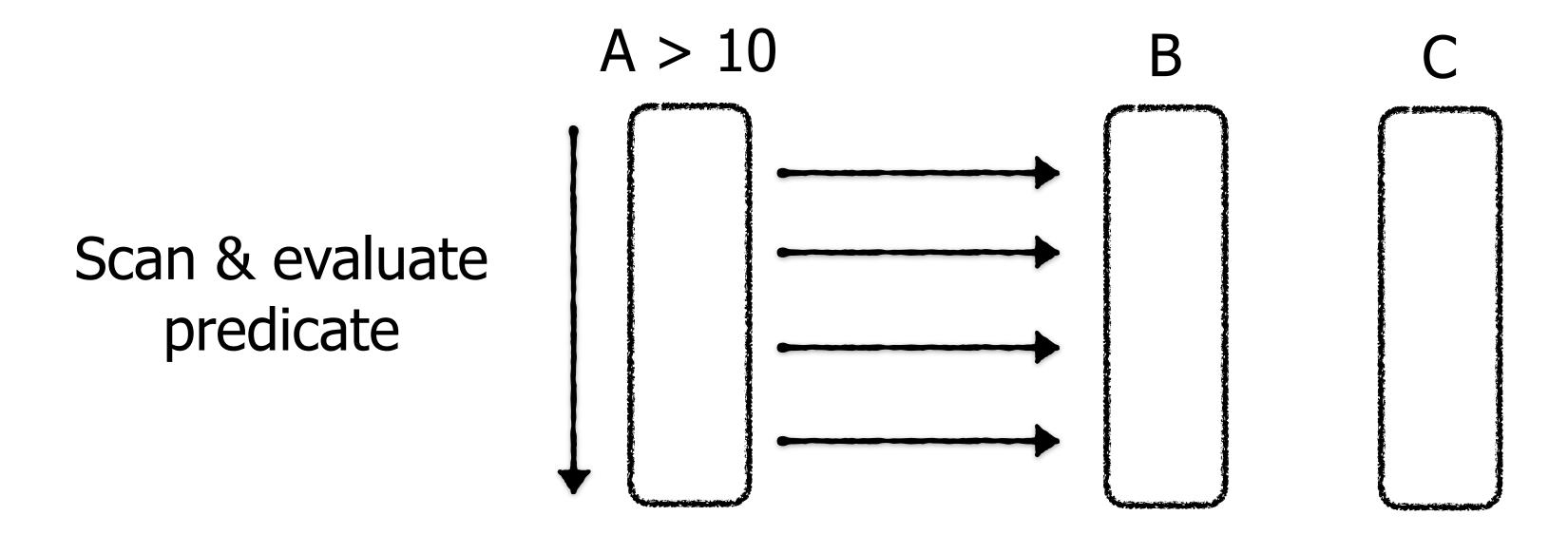
How shall we process this?



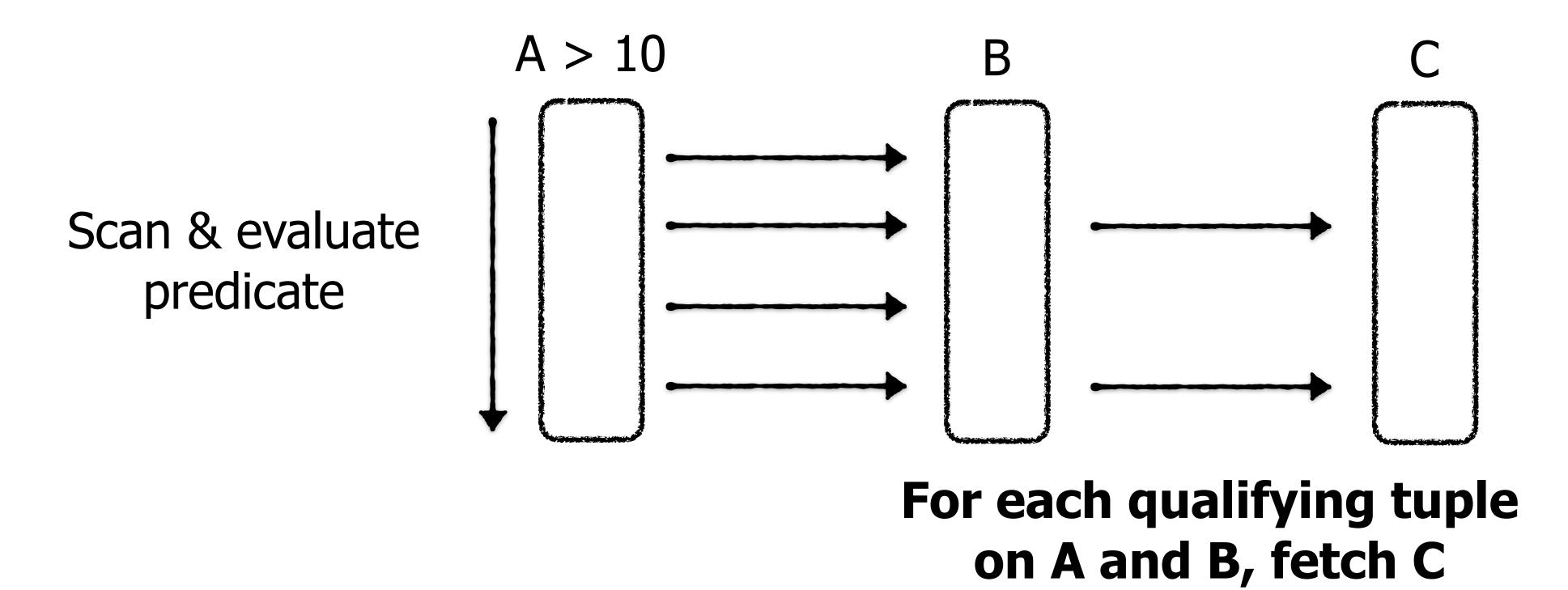
Early Materialization: return each qualifying row immediately

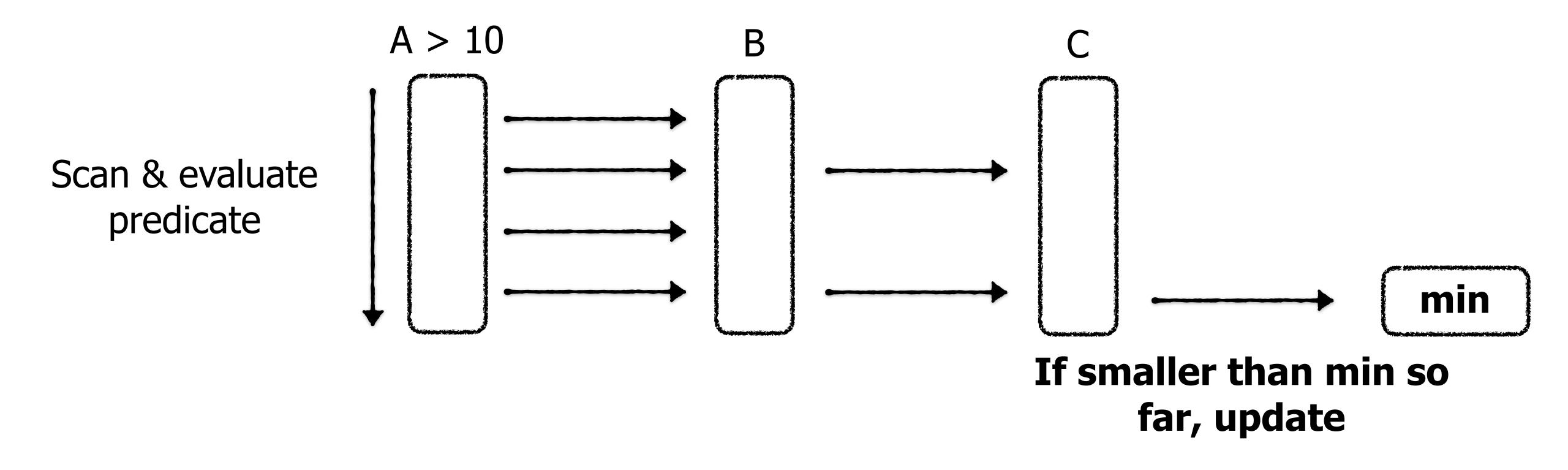


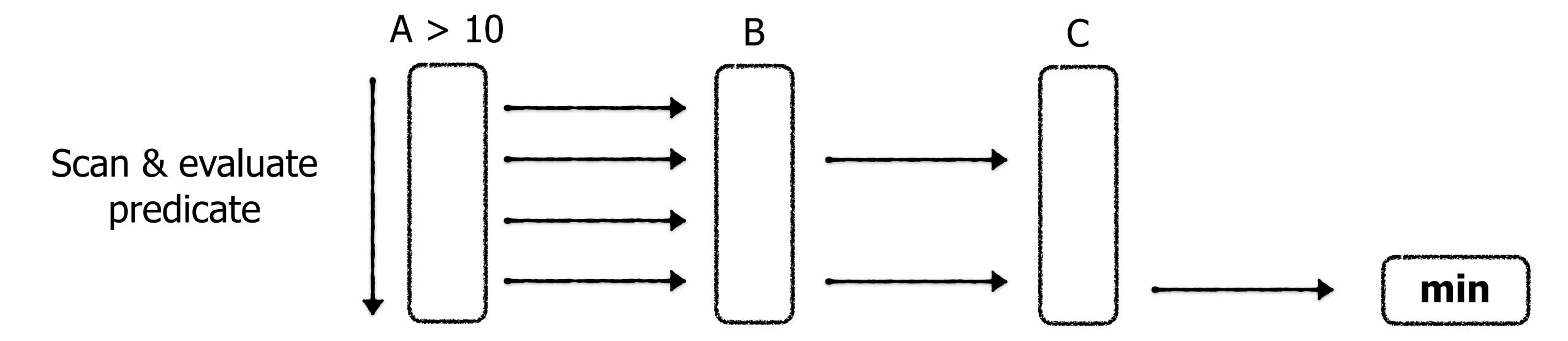
Early Materialization: return each qualifying row immediately



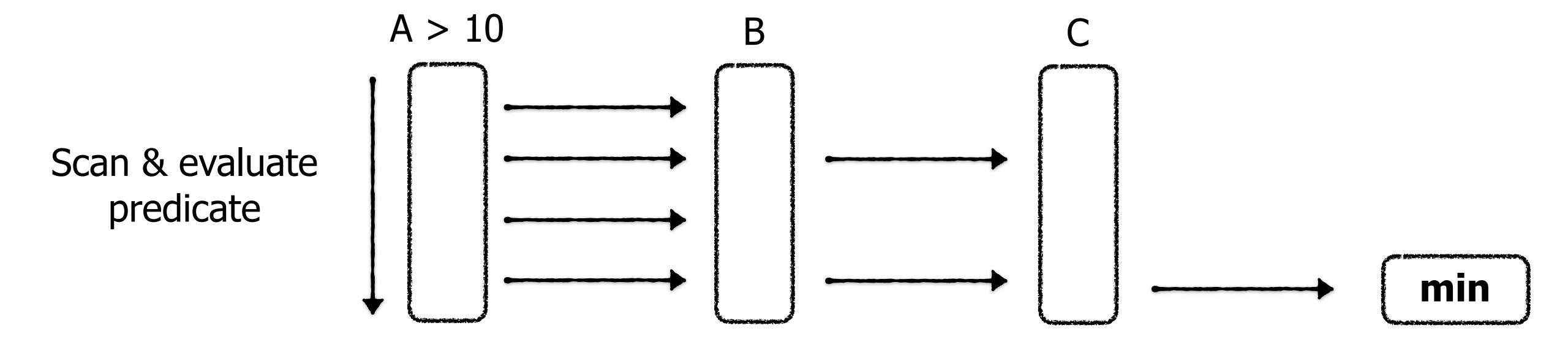
For each qualifying tuple on A, check if B also qualifies





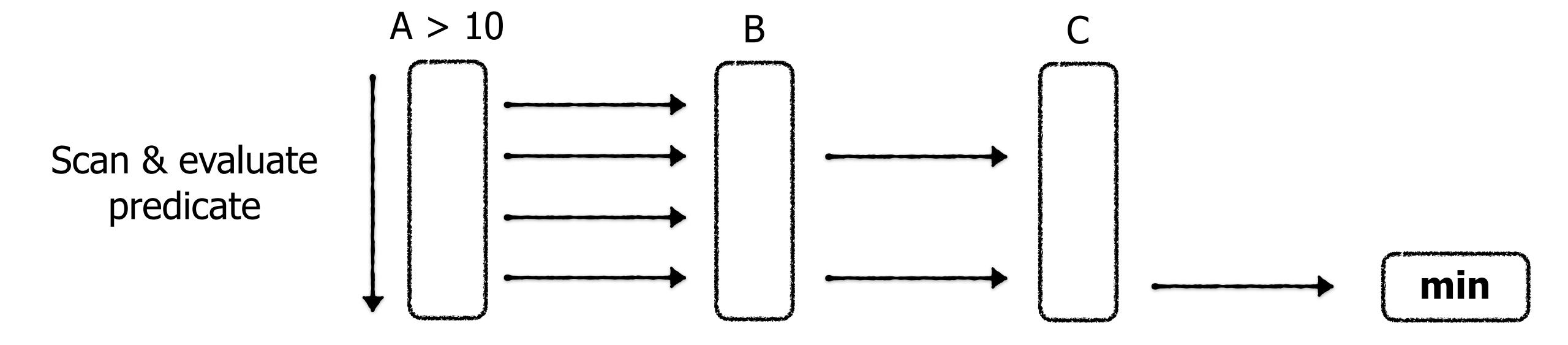


#### **Problem?**



Problem? Potentially more random I/Os

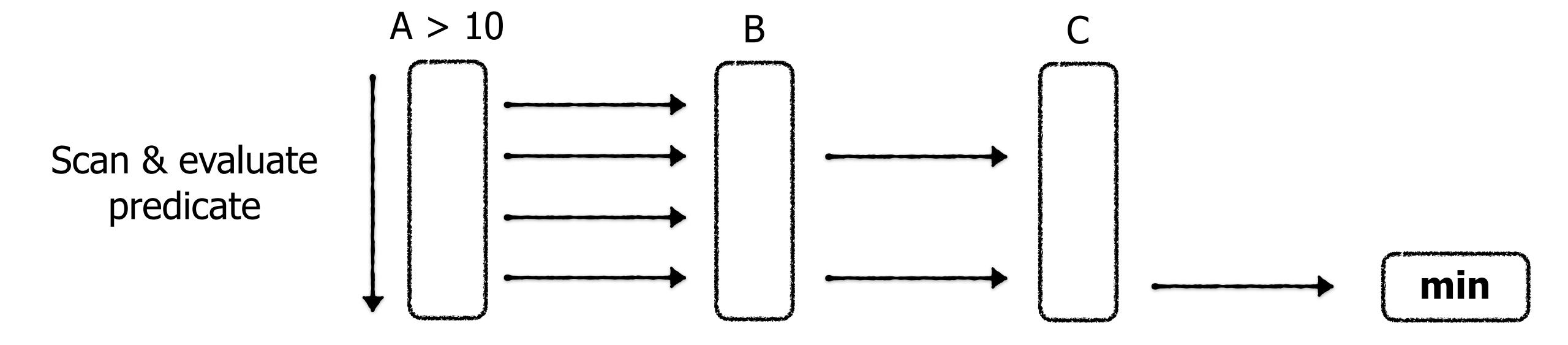
More CPU branch mispredictions



Problem? Potentially more random I/Os

More CPU branch mispredictions

Incompatible with SIMD (more later)



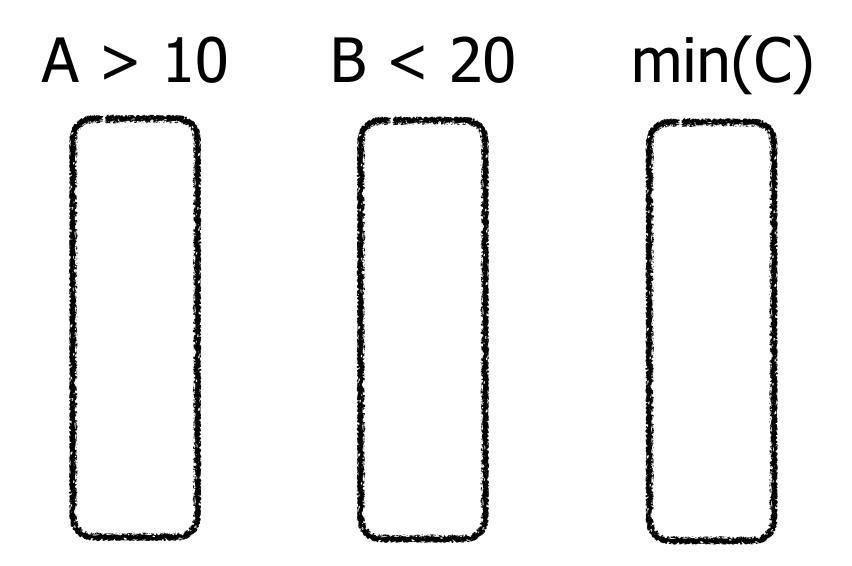
Problem? Potentially more random I/Os

More CPU branch mispredictions

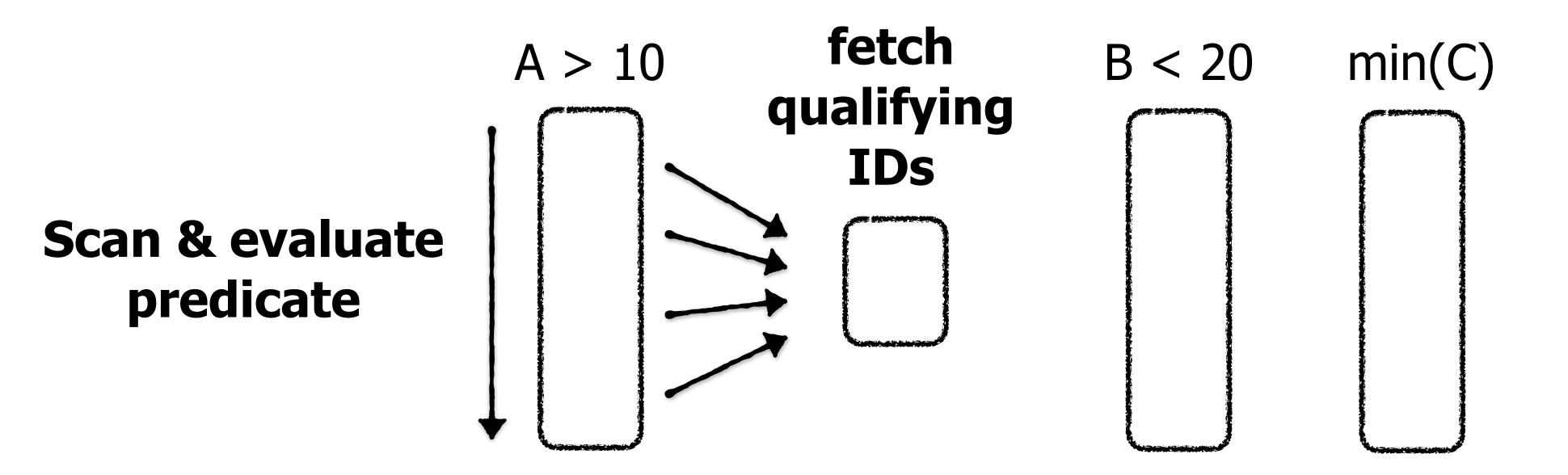
Incompatible with SIMD (more later)

**Solutions?** 

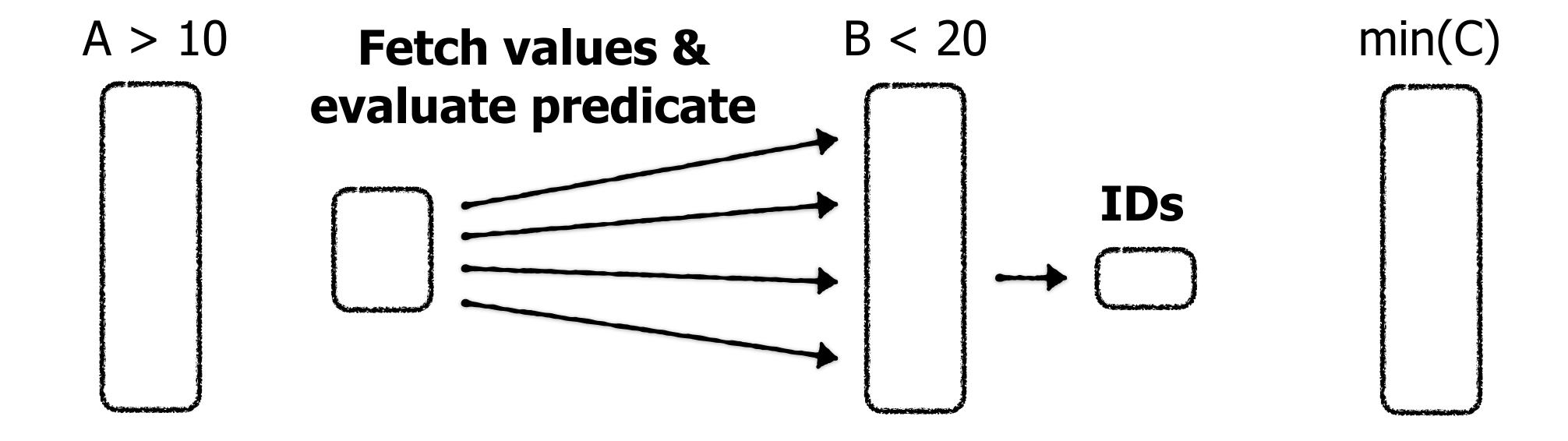
Late Materialization: process one column at a time



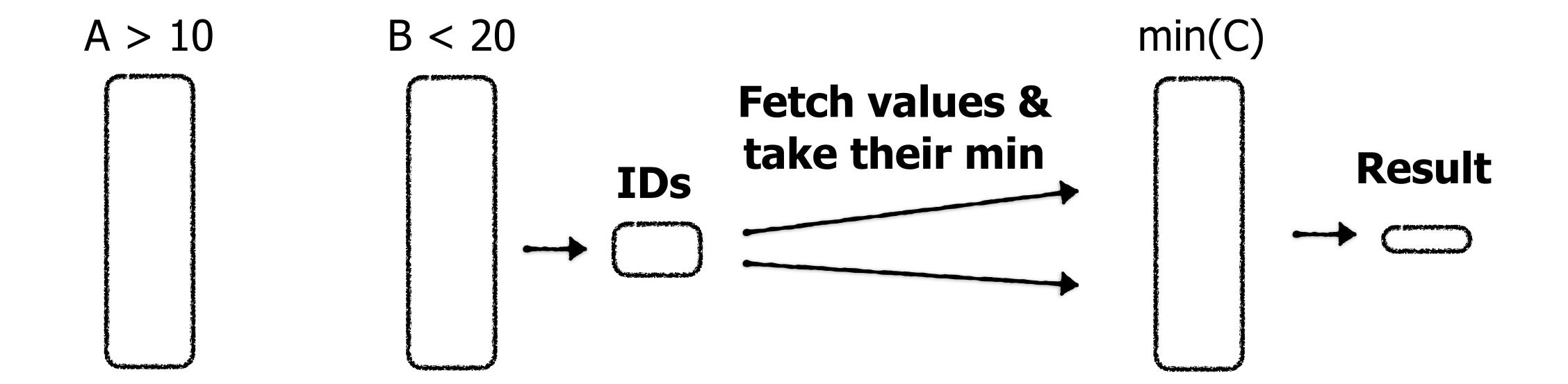
#### Late Materialization



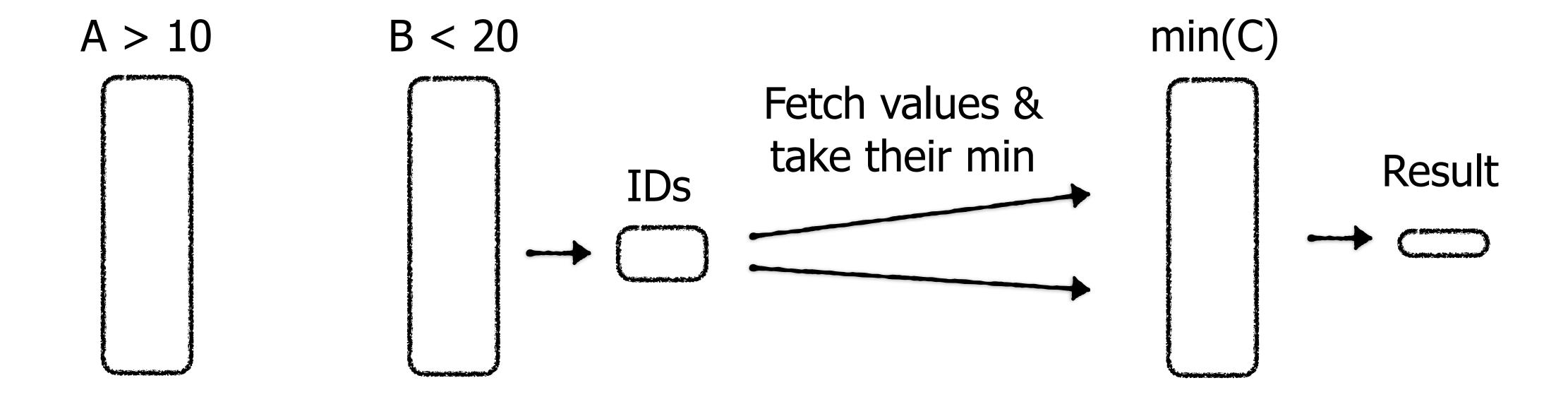
#### Late Materialization

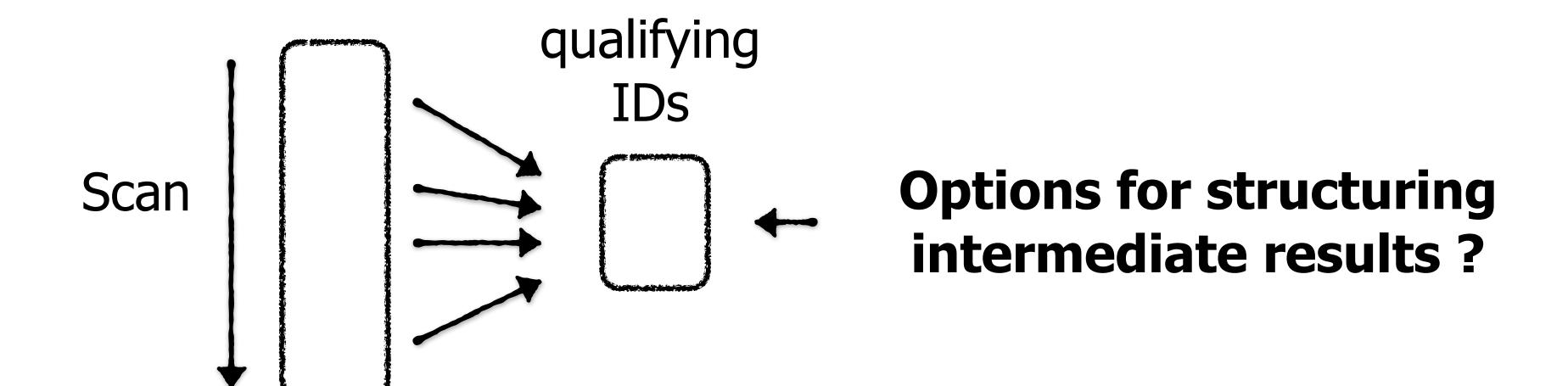


#### Late Materialization

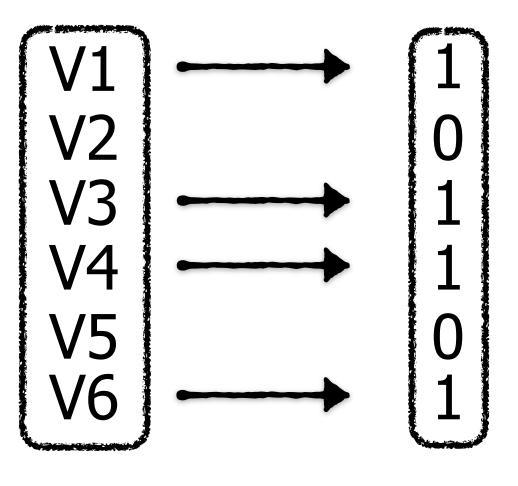


Late Materialization requires more memory to withhold intermediate results, but entails better cache behavior & CPU efficiency



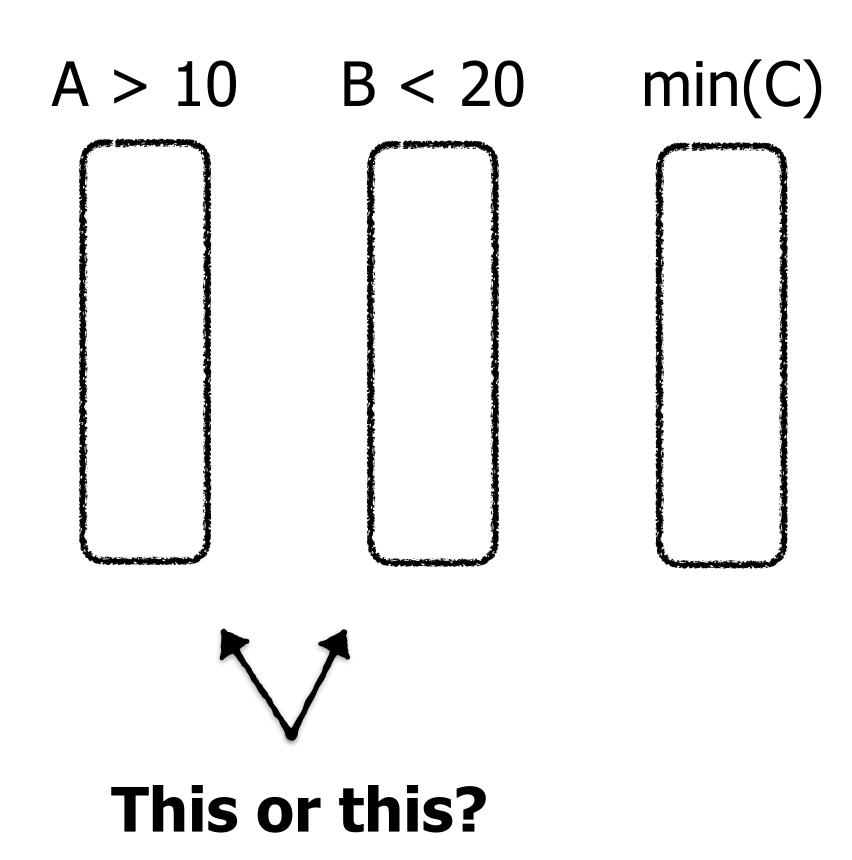




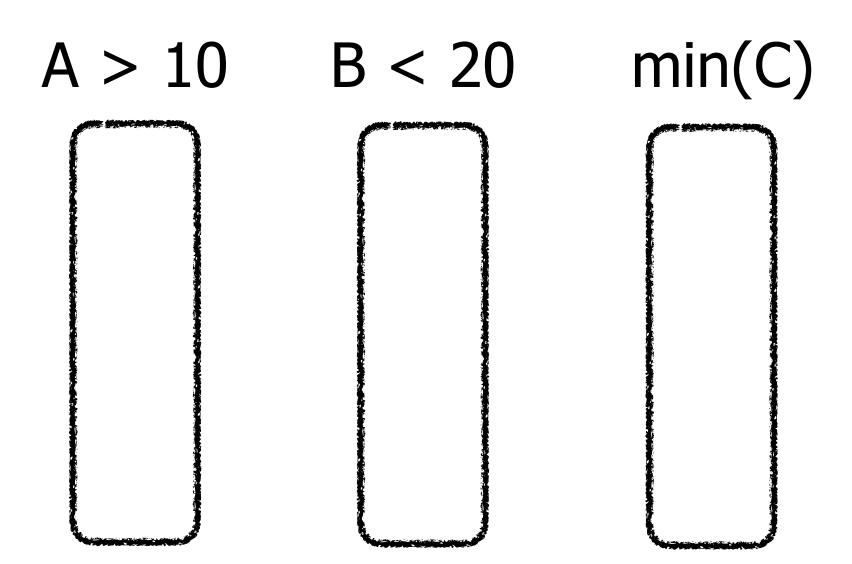


(2) Bitmap good when many results qualify

#### Which column should we filter on first?

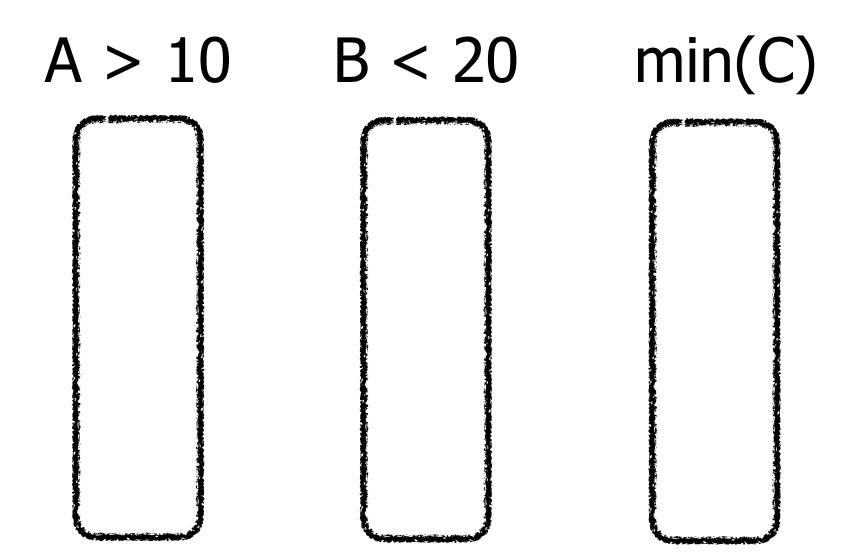


#### Which column should we filter on first?



Filter on more selective columns first to reduce size of intermediate results and thus I/O

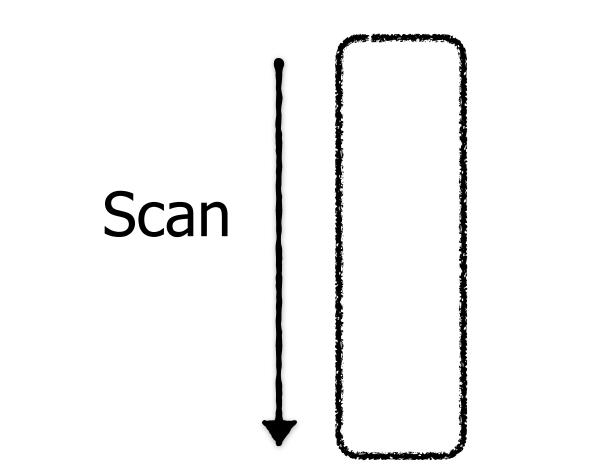
#### Which column should we filter on first?



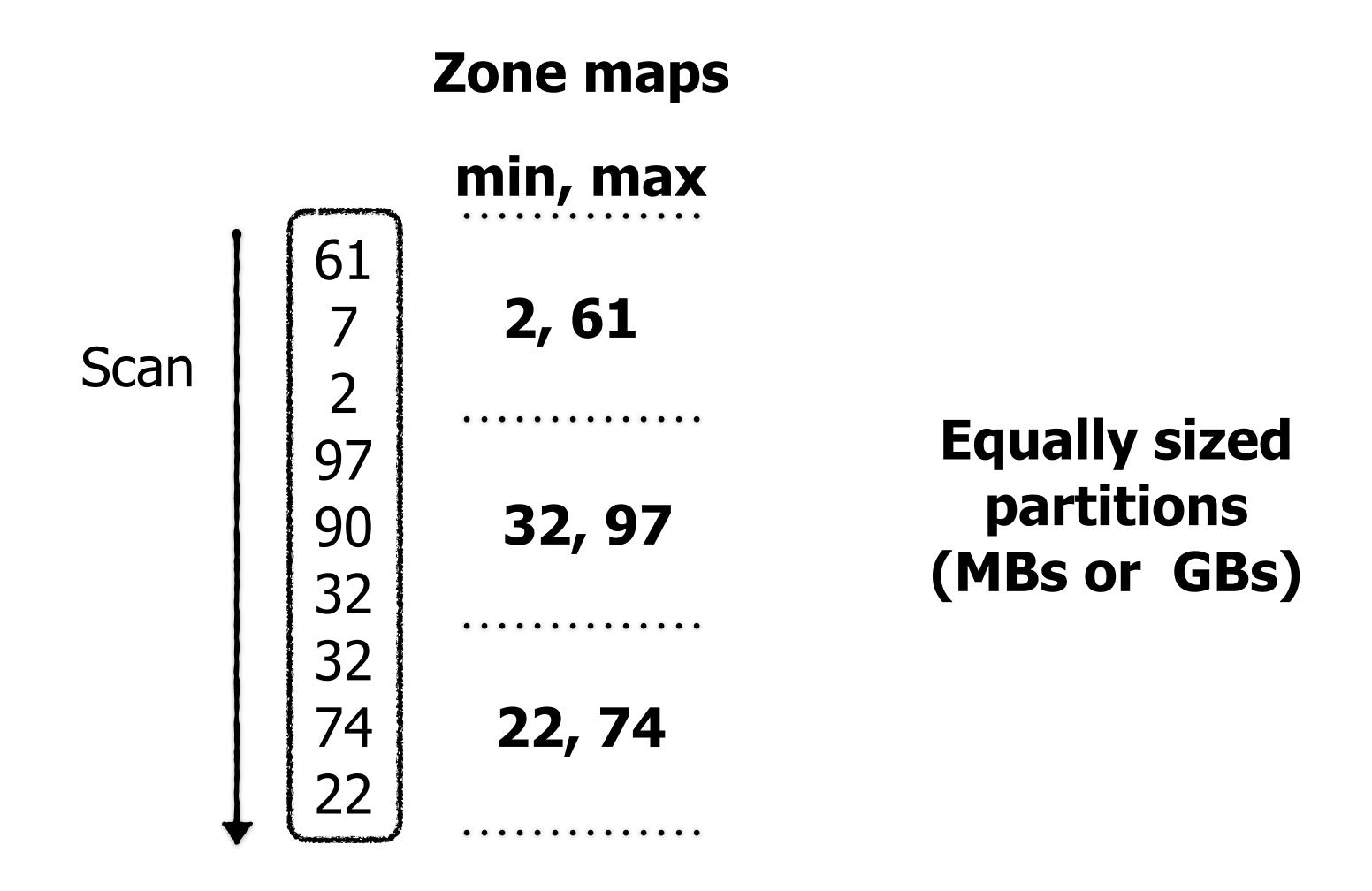
Filter on more selective columns first to reduce size of intermediate results and thus I/O

Can do this via cardinality estimation (e.g., histograms, count-min, etc)

# How to speed up scans?

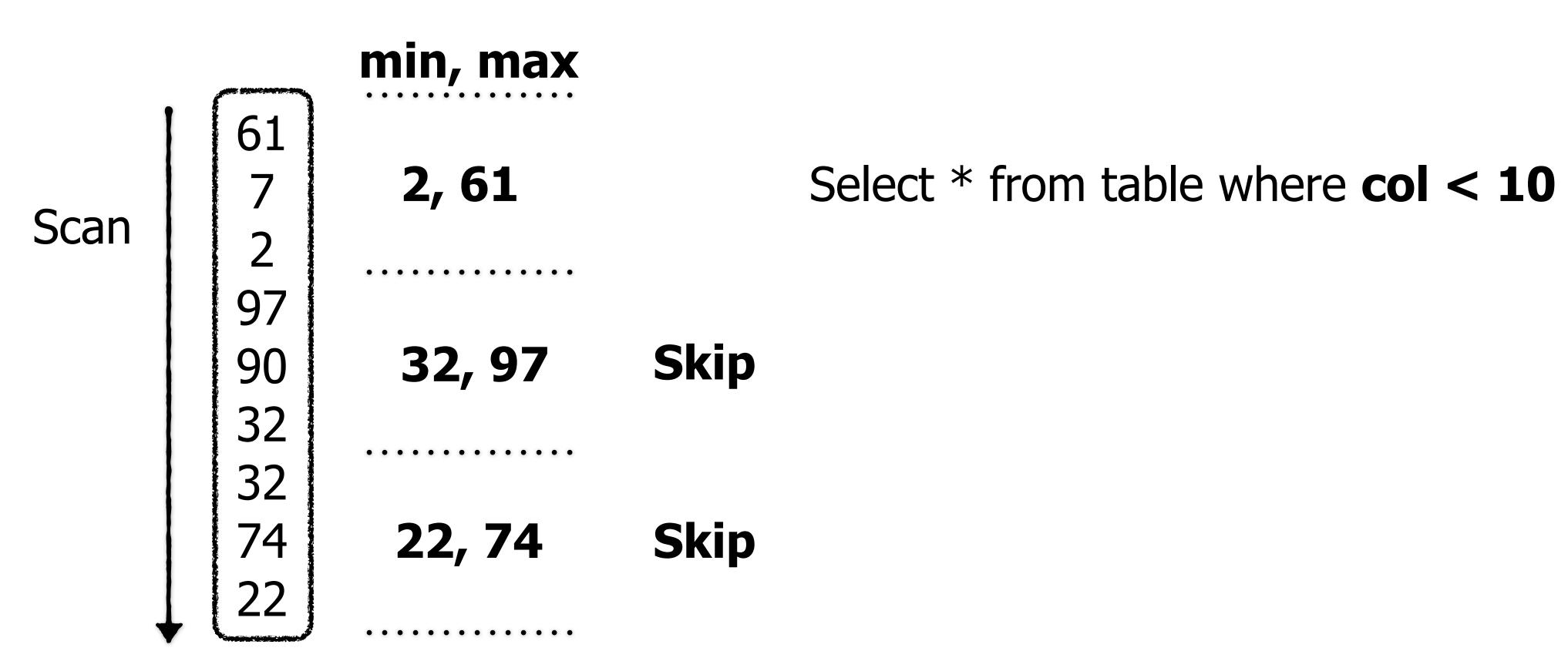


### How to speed up scans?



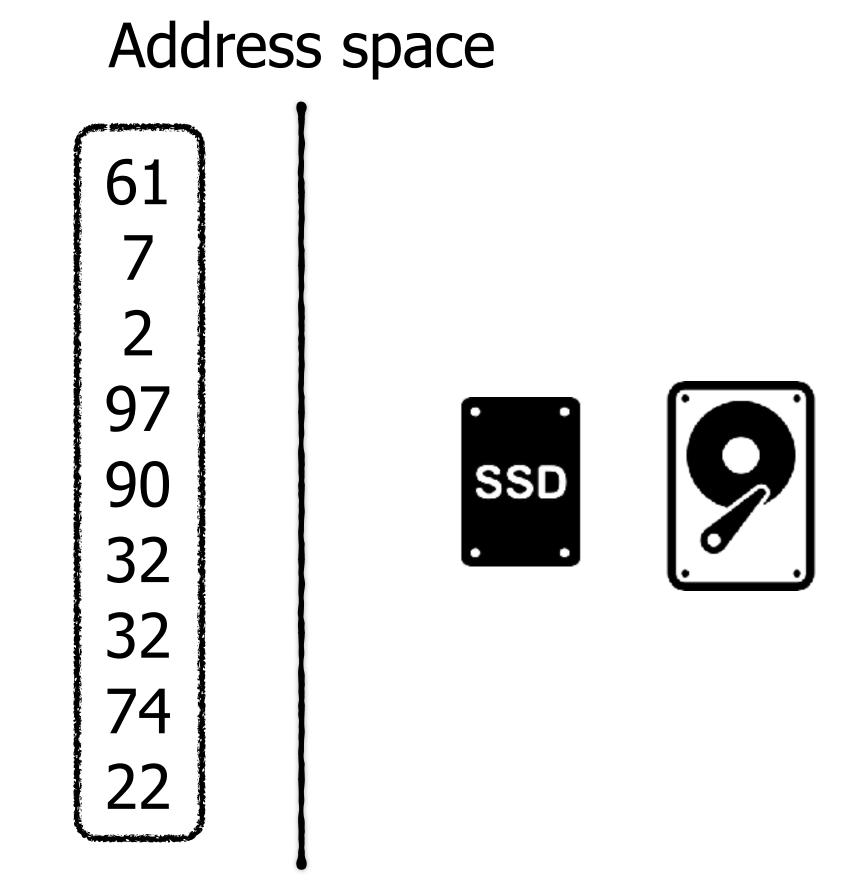
### How to speed up scans?





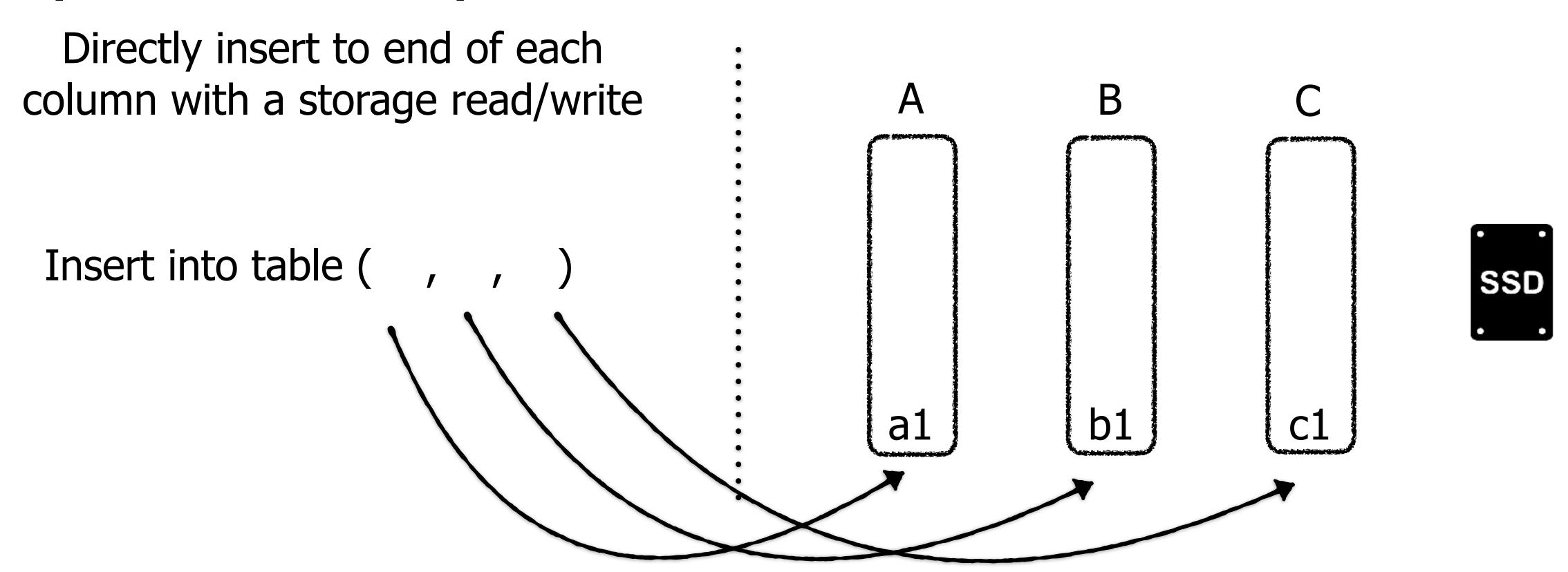
#### How to speed up scans?

#### Store columns sequentially to allow pre-fetching and I/O parallelism



Insert into table (a1,b1,c1)

#### **Option 1: In-Place Updates**

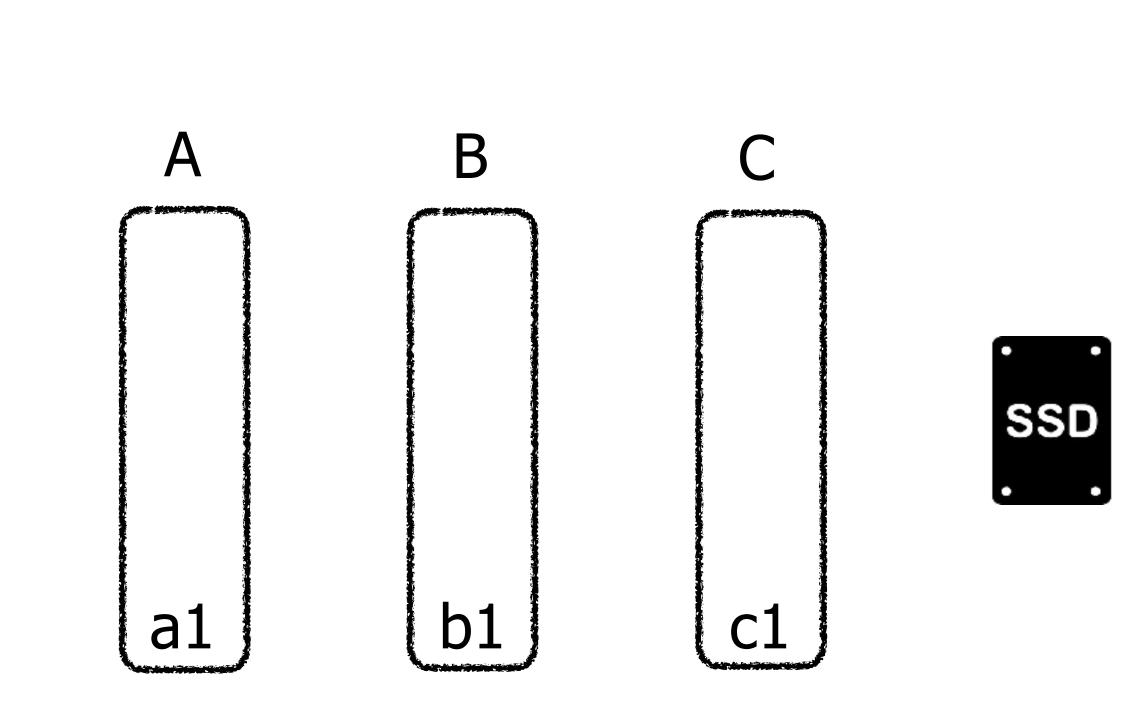


#### **Option 1: In-Place Updates**

Directly insert to end of each column with a storage read/write

Insert into table ( , , )

**Cost model?** 

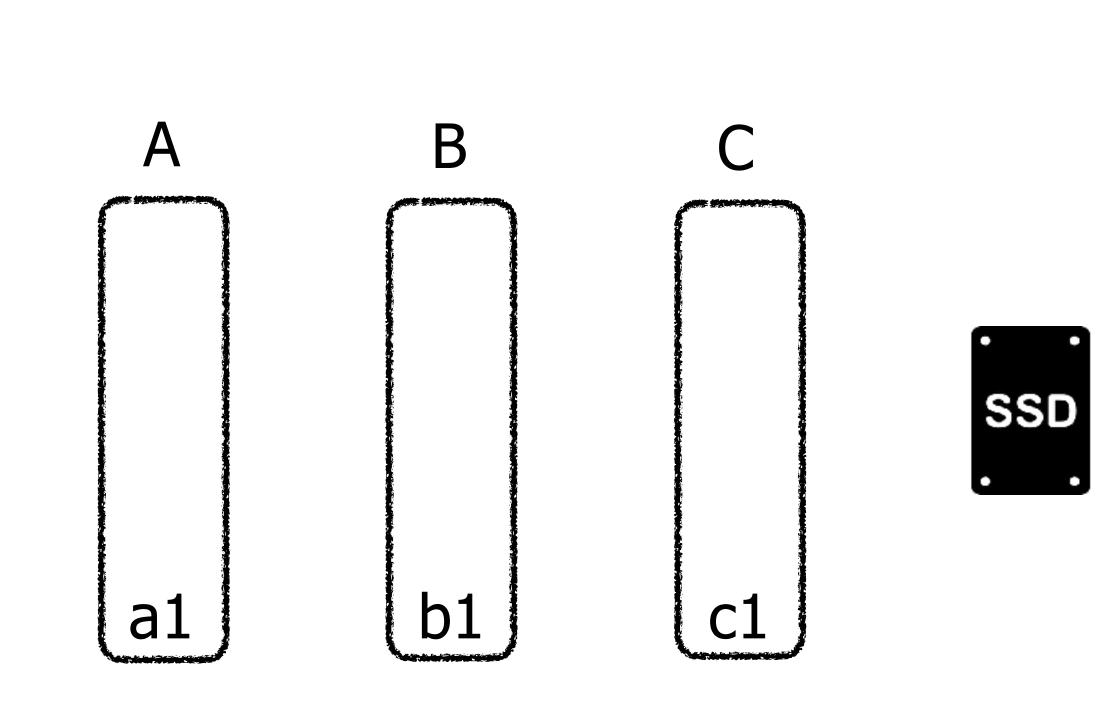


#### **Option 1: In-Place Updates**

Directly insert to end of each column with a storage read/write

Insert into table ( , , )

Cost model: O( #cols )



#### **Option 2: In-memory Buffering**

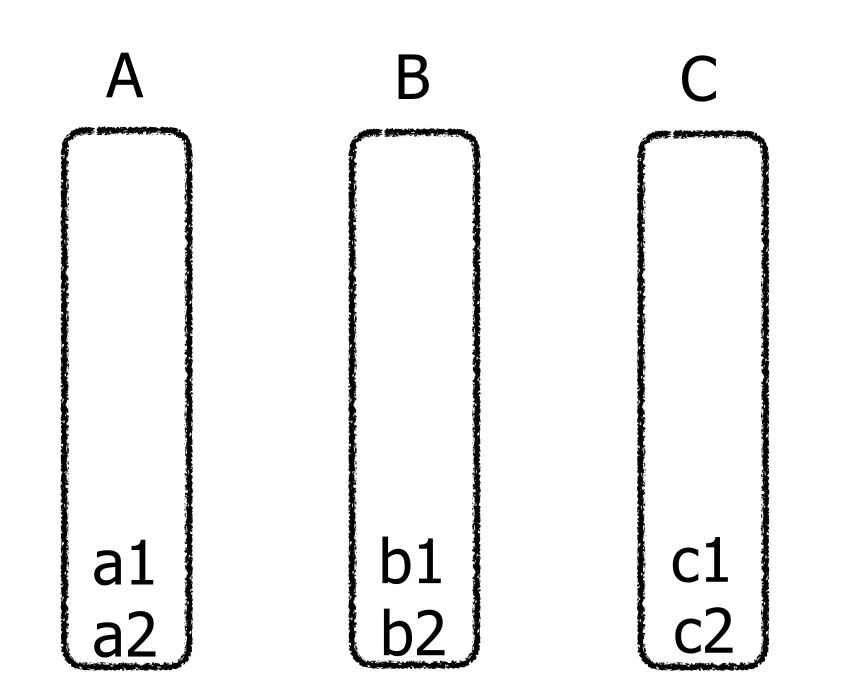
Insert into table ( , , )

Buf1 Buf2 Buf3

a1 b1 b2 c1 c2

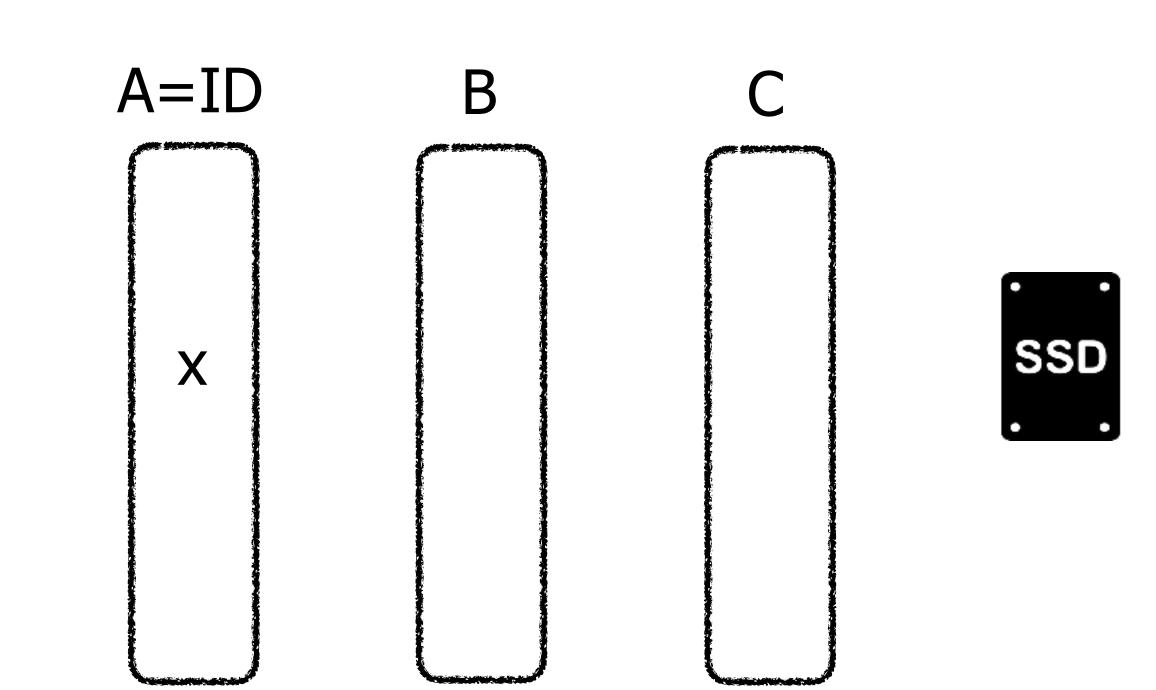
## **Option 2: In-memory Buffering**

**Cheap insertions are possible:)** 



**delete** from table where **A** = "x"

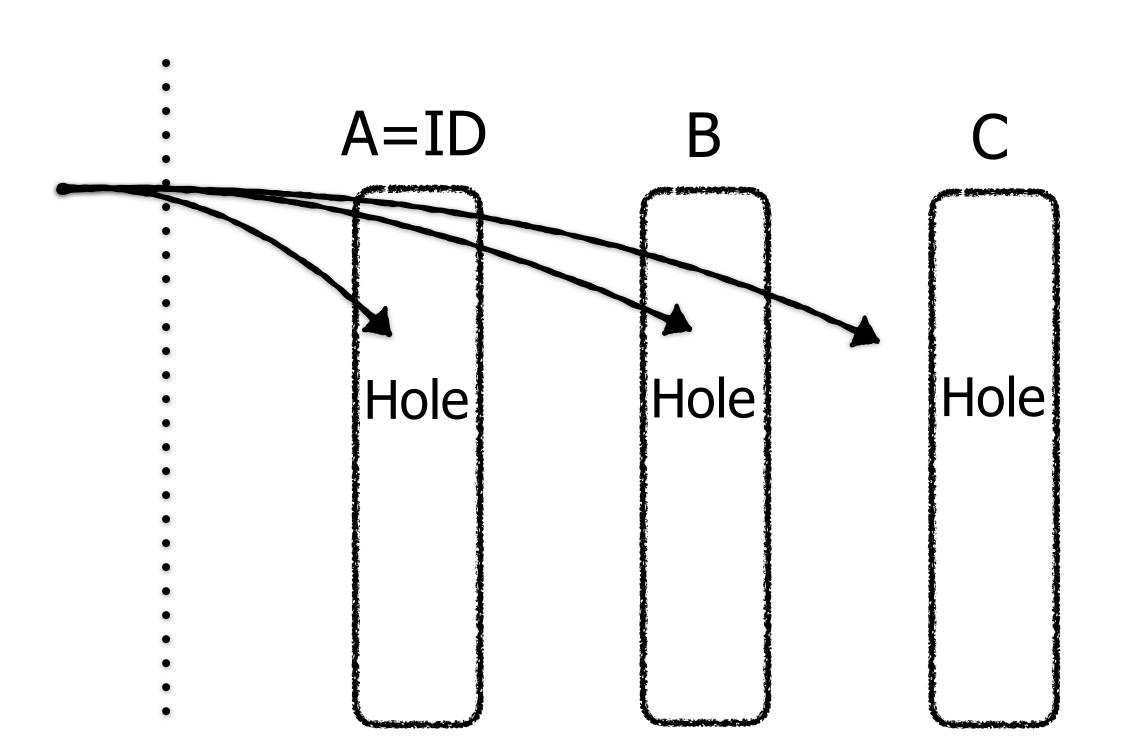
How would you do this?



## **Option 1: In-place Deletes**

delete from table where A = "x"

#### **Problems:**





Option 1: In-place Deletes

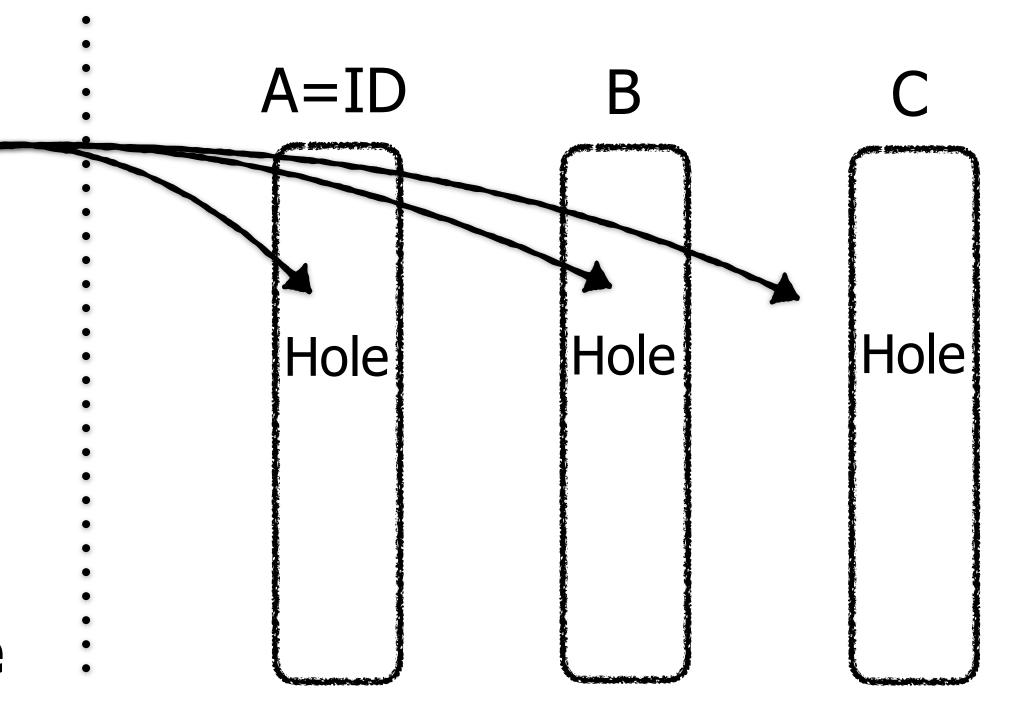
delete from table where A = "x"

Problems:

Cost: O(#cols) write I/Os

Holes slows down queries

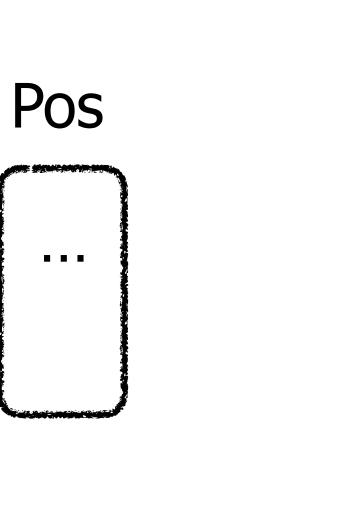
Requires 1 extra bit to note a hole

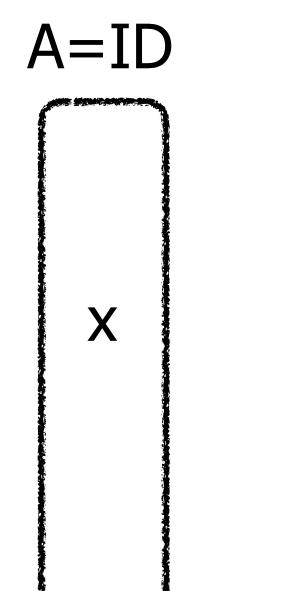


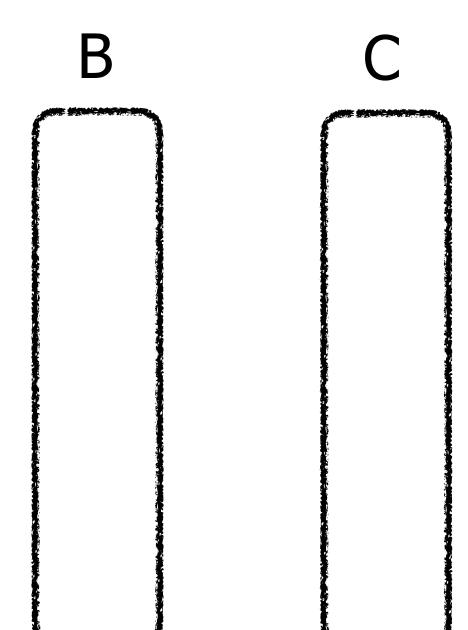


Option 2: employ delete column

delete from table where A = "x"









Option 2: employ delete column

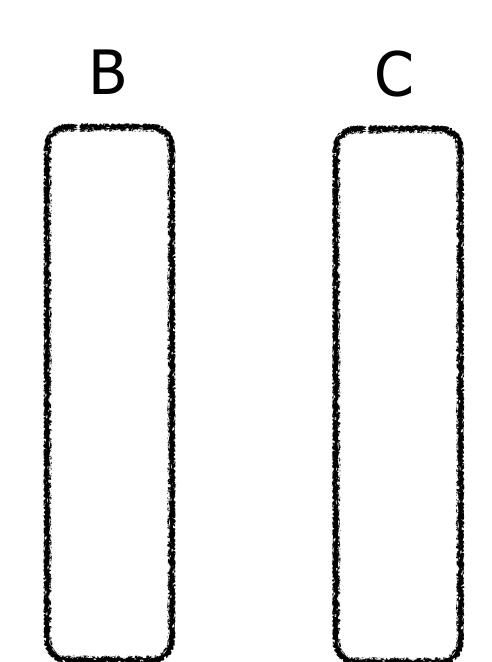
delete from table where A = "" A = ID X X X X SSD

Option 2: employ delete column

delete from table where A = " "

Pos

A=ID
X

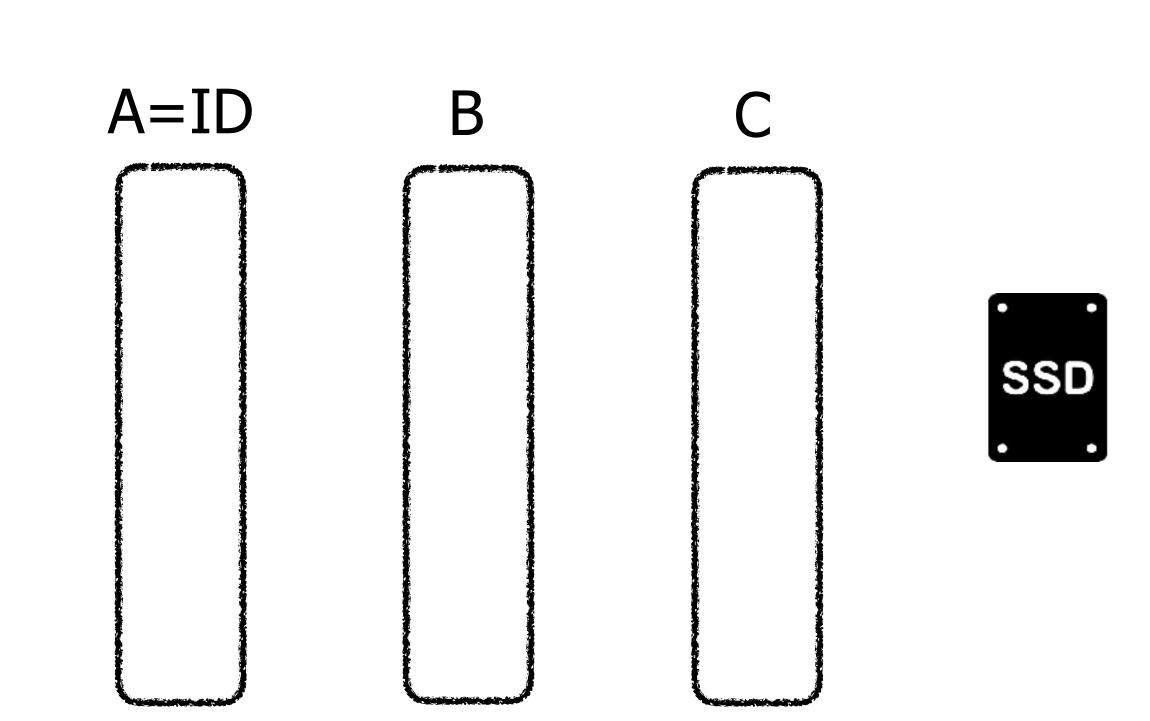




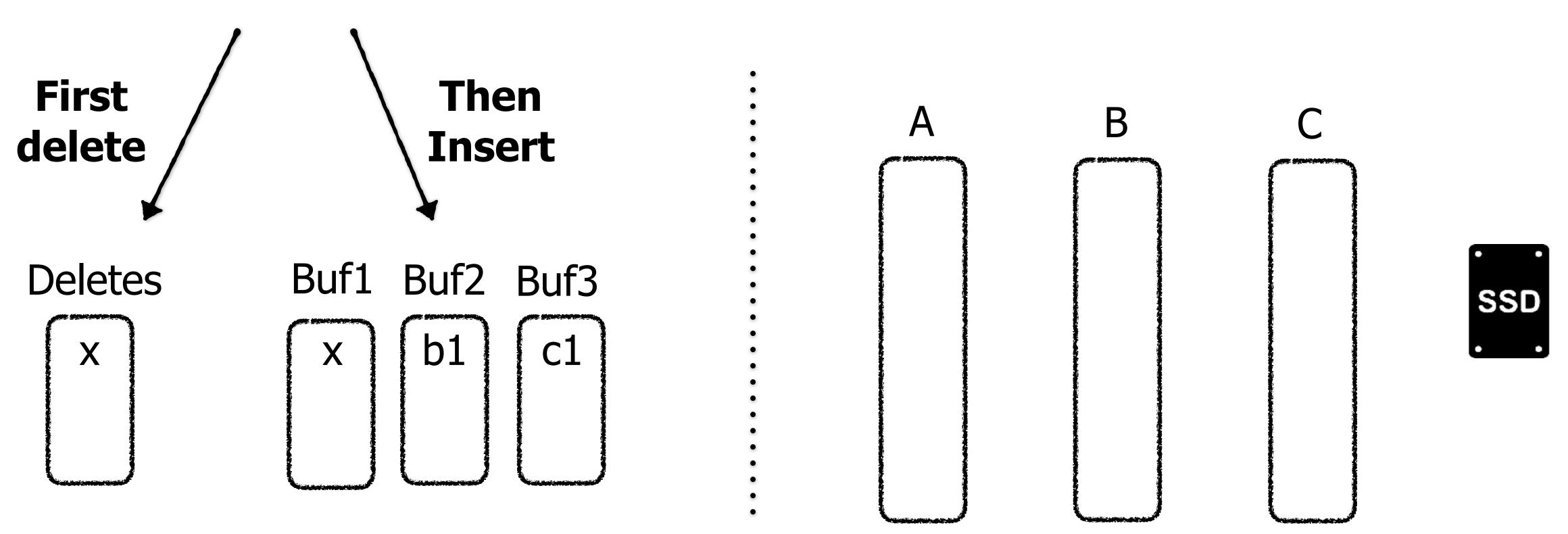
**Eventually merge with column** 

update table set B="b1", C="c1" where A="x"

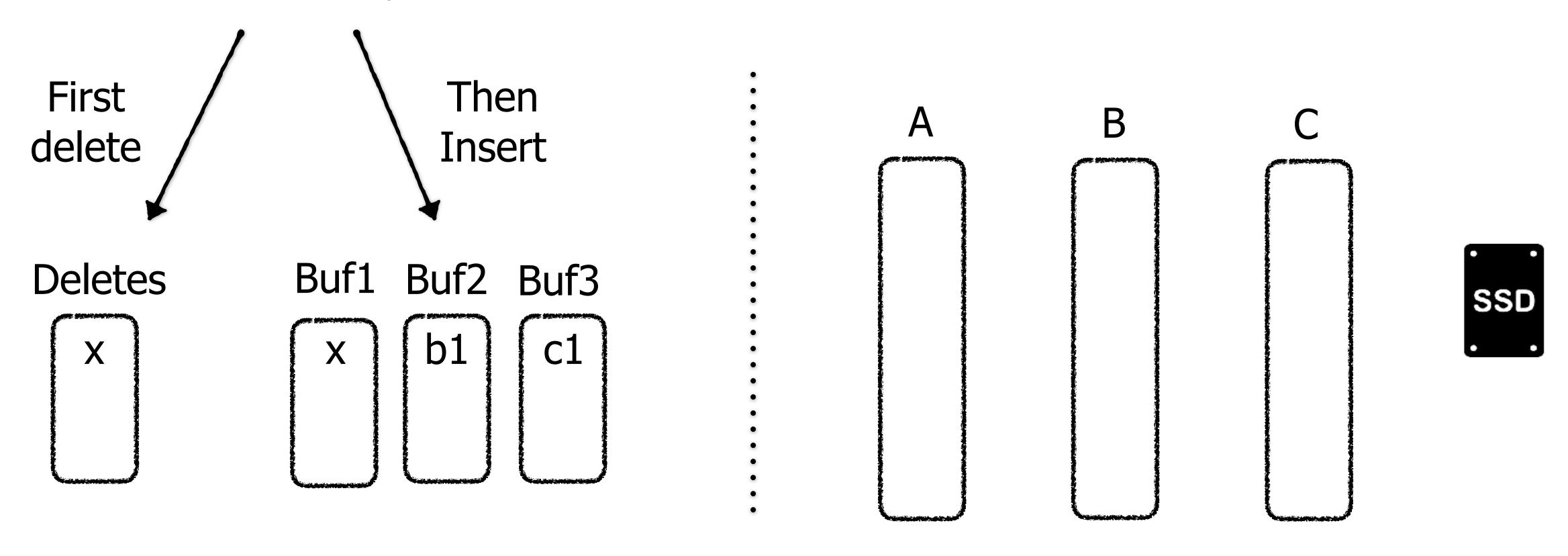
Can we combine our solutions for insertions & deletes to efficiently update?



update table set B="b1", C="c1" where A="x"

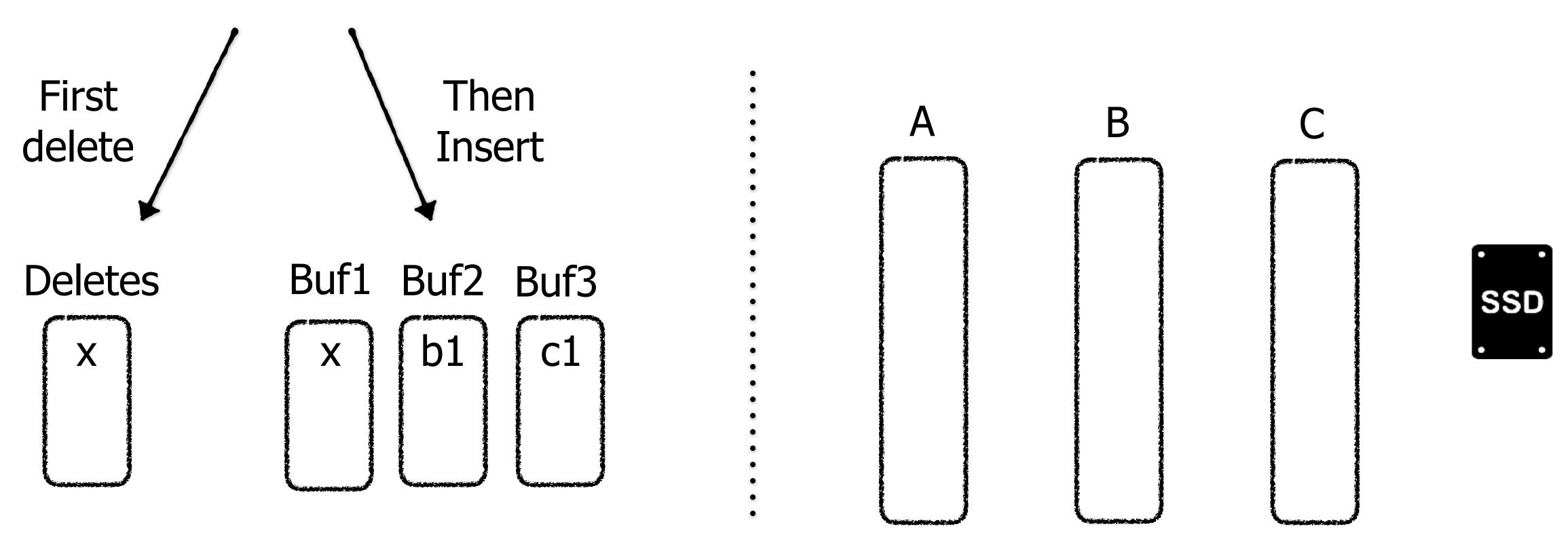


update table set B="b1", C="c1" where A="x"



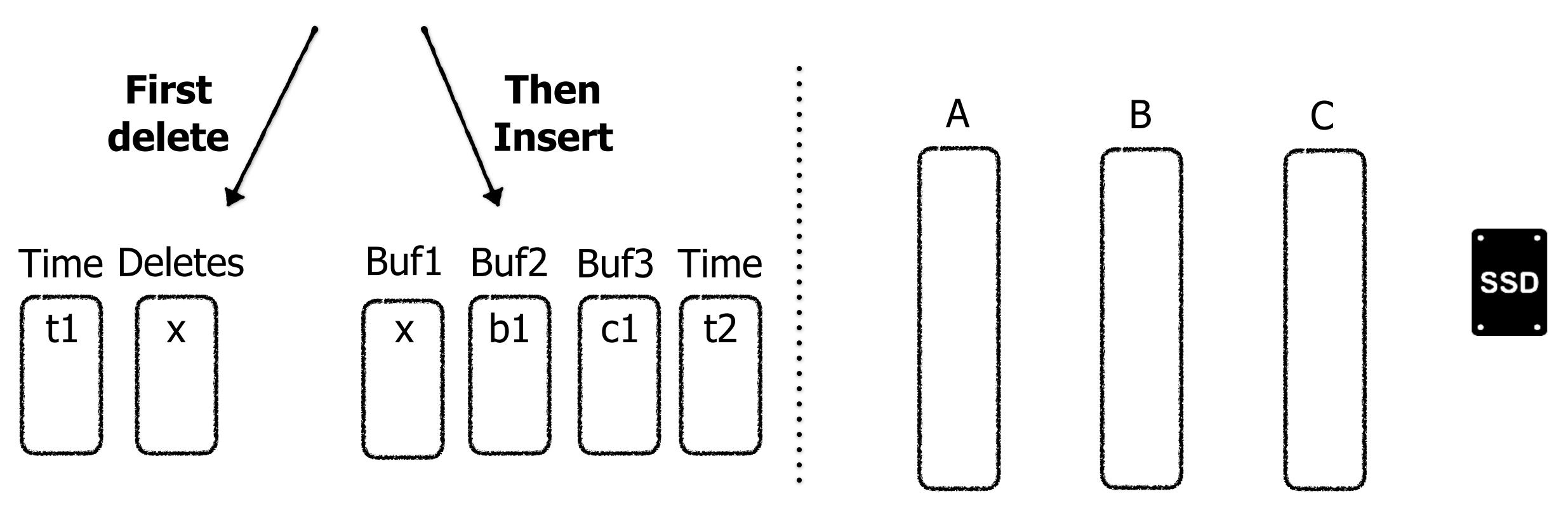
**Problem?** 

update table set B="b1", C="c1" where A="x"



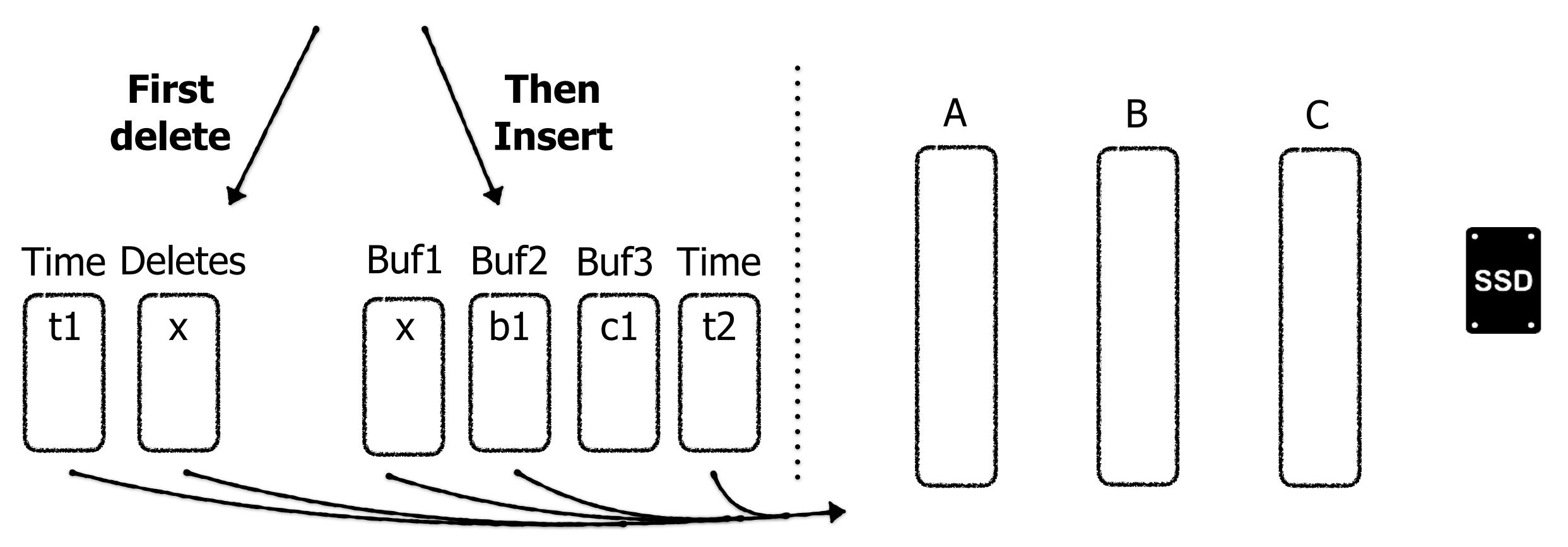
**Problem: Which came first?** 

update table set B="b1", C="c1" where A="x"



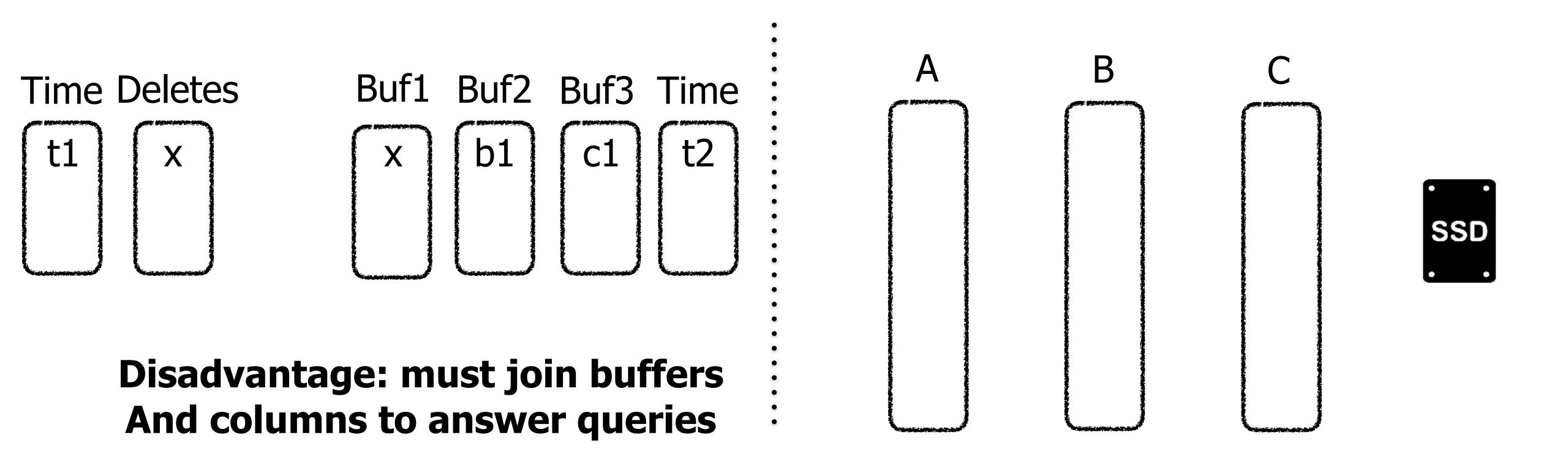
Problem: Which came first? Fix using timestamps

update table set B="b1", C="c1" where A="x"

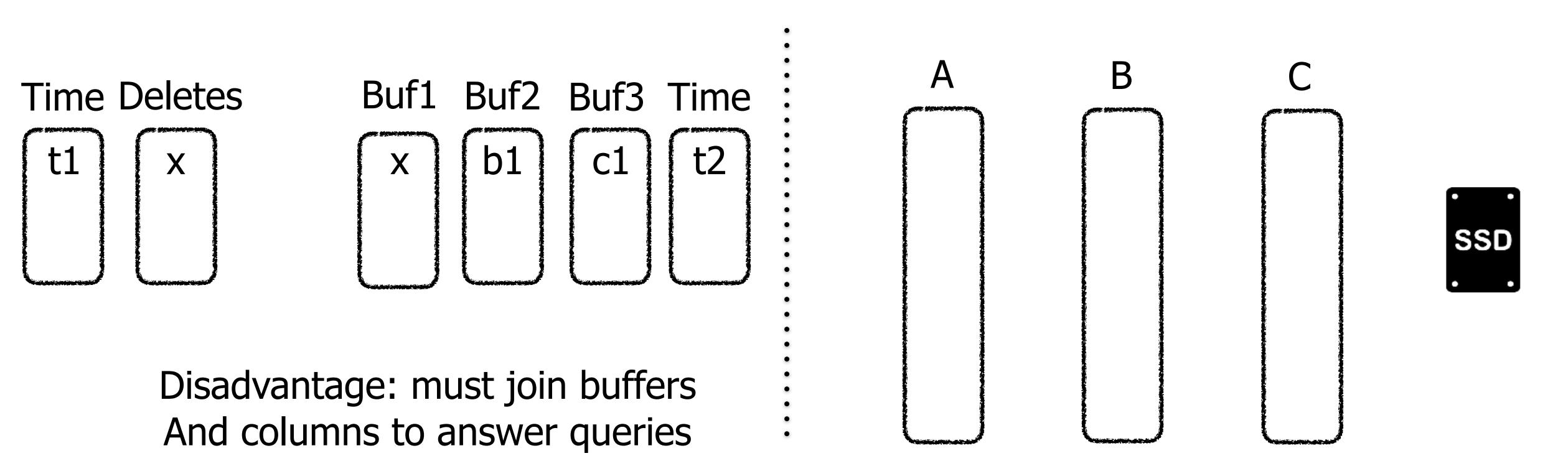


**Eventually merge into columns** 

update table set B="b1", C="c1" where A="x"



update table set B="b1", C="c1" where A="x"

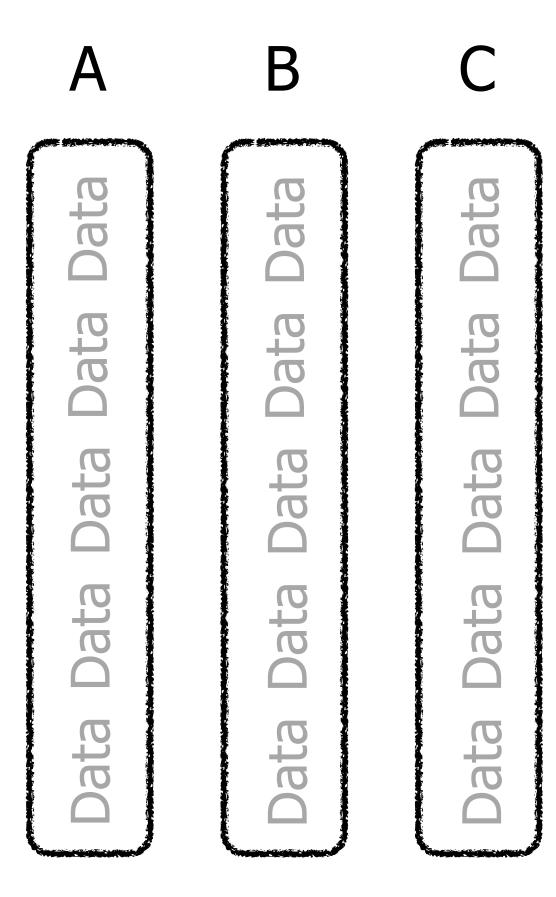


Hence, modifications are best done in batch and offline (e.g., overnight)

#### Columns are immutable and their values inside are mostly fixed-sized

Columns are immutable and their values inside are mostly fixed-sized

#### This allows for further optimizations



# Tight processing loops with no pointer-chasing or function calls

Select from table where A > v



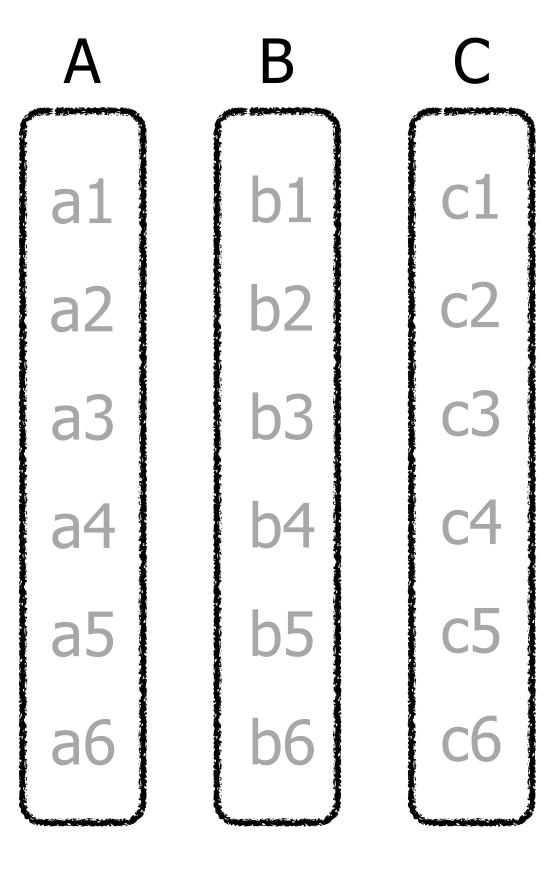
for (i = 0; i < size; i++)
 if A[i] > v
 qualifying[j++] = i

B Data ata Data Data Data ata Data Data Data

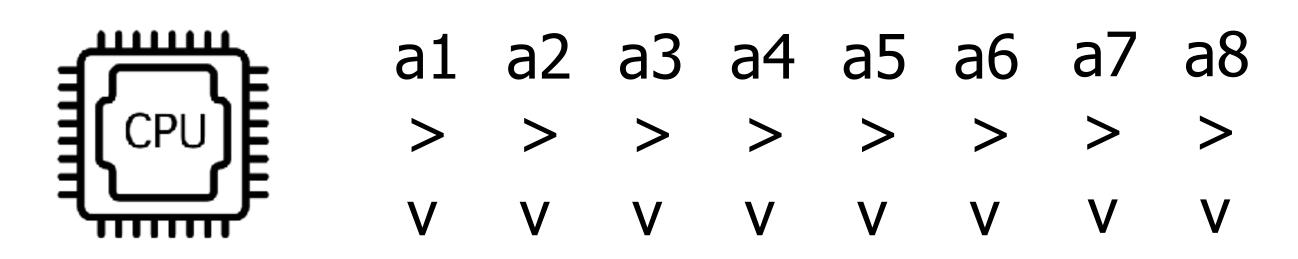
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)

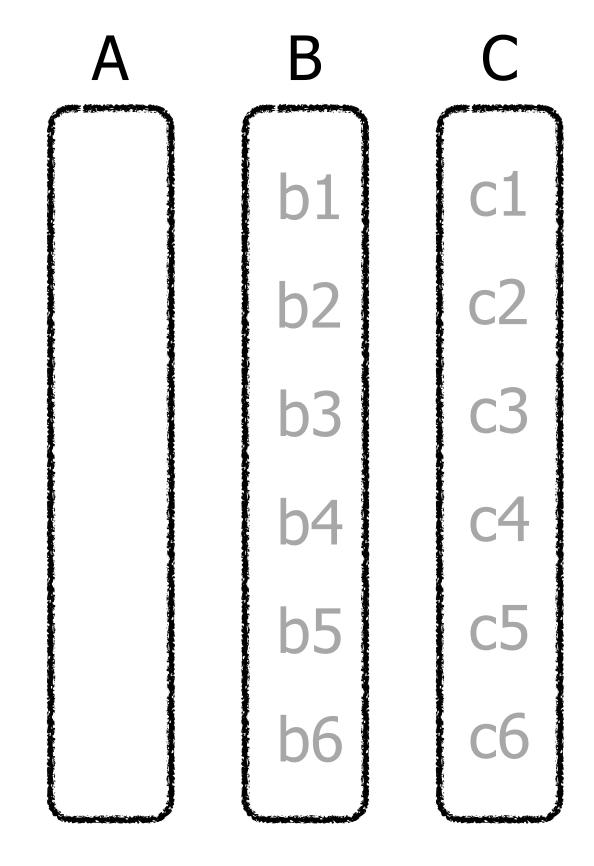
В	C
<b>b1</b>	<b>c1</b>
b2	<b>c</b> 2
<b>b</b> 3	<b>c</b> 3
b4	<b>c4</b>
b5	<b>c</b> 5
b6	с6
	b1 b2 b3 b4 b5

Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)

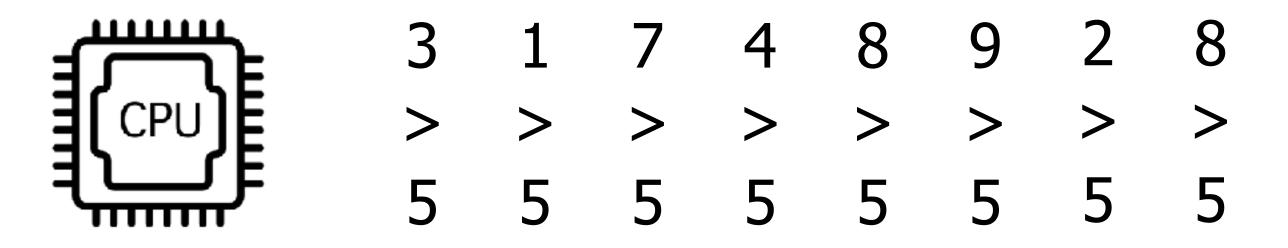


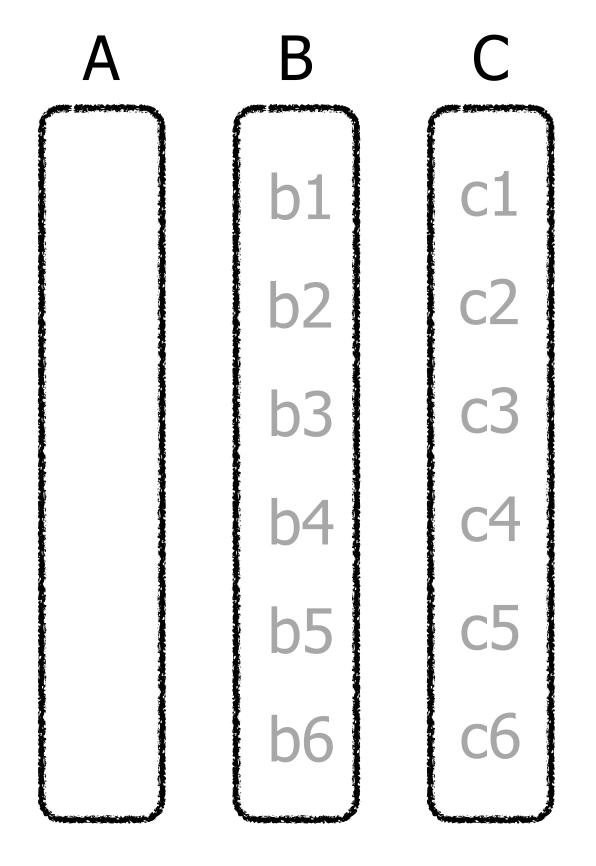
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)



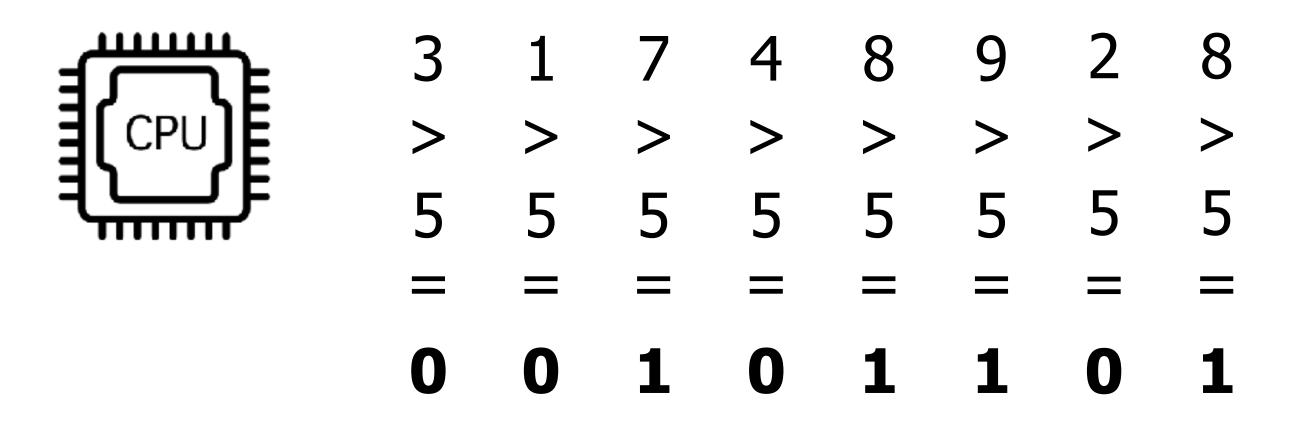


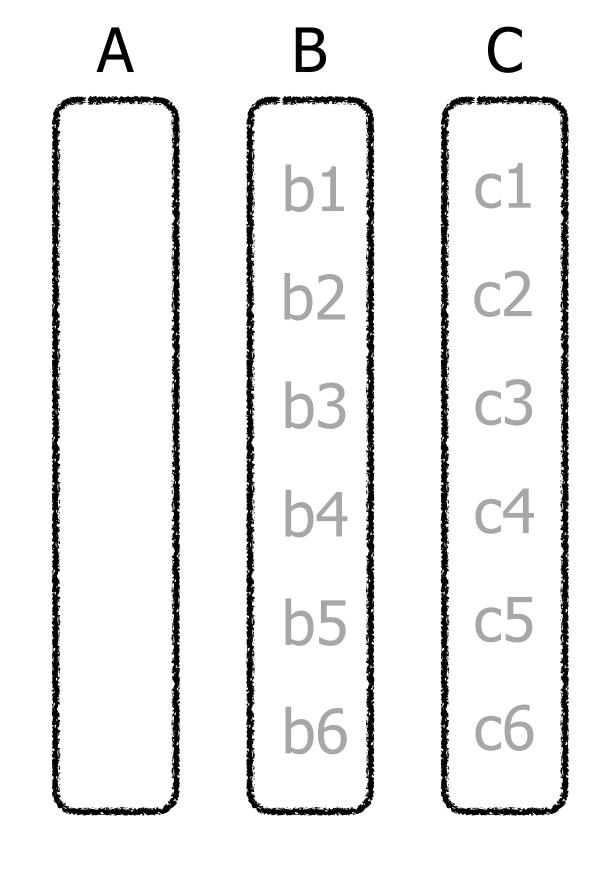
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)



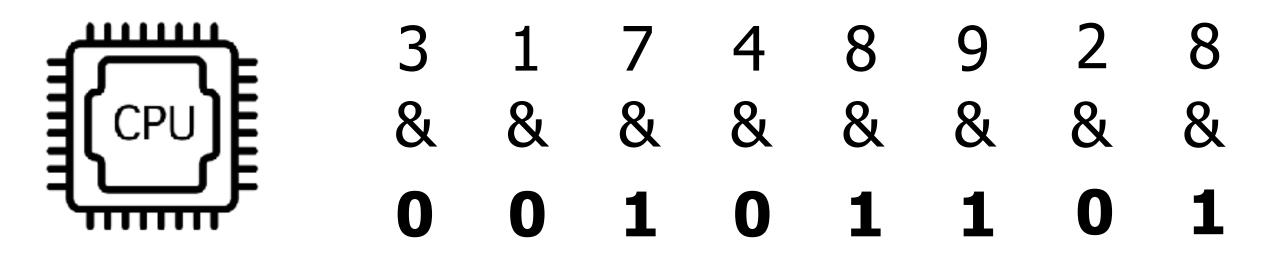


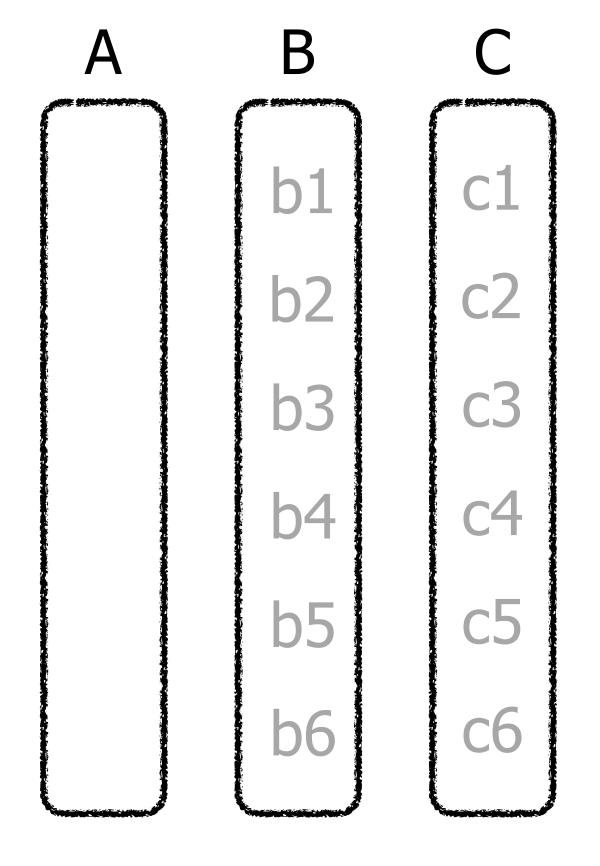
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)



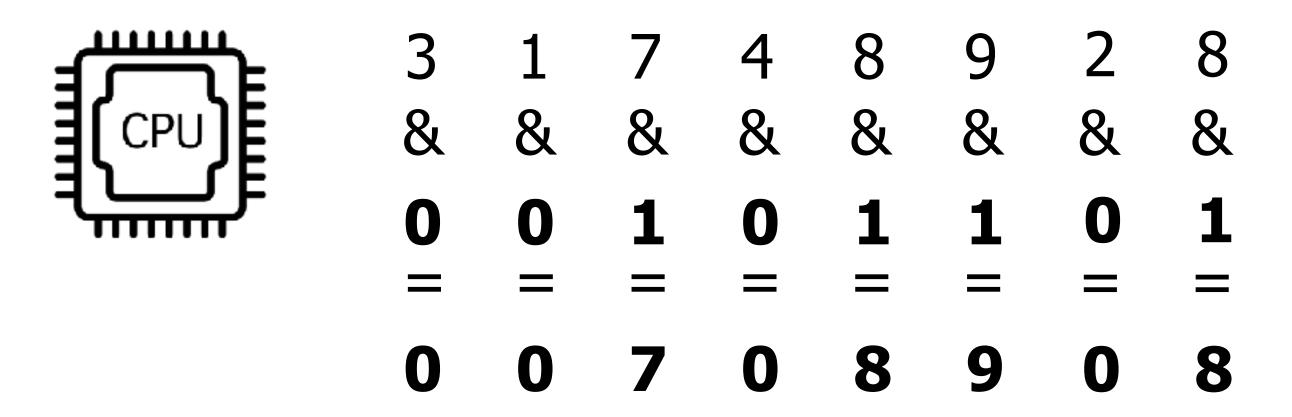


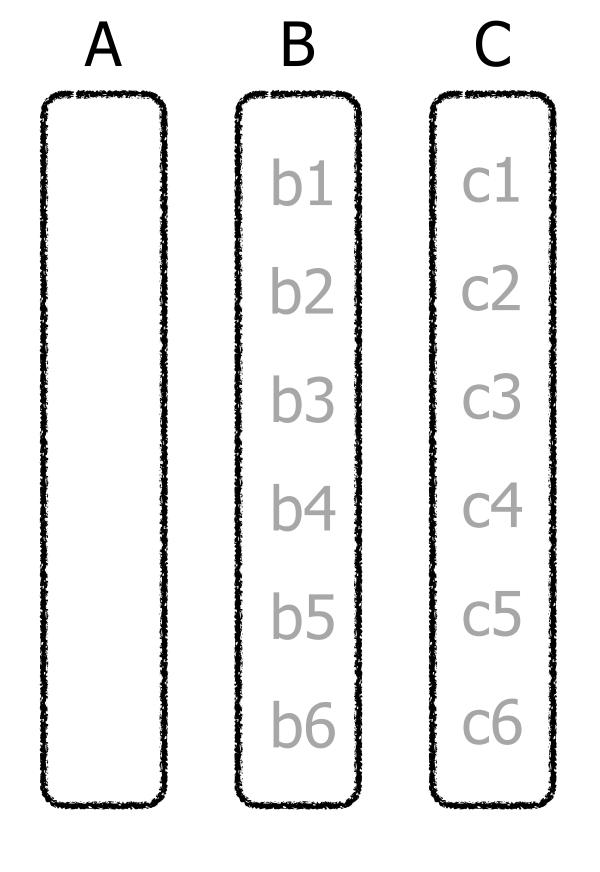
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)



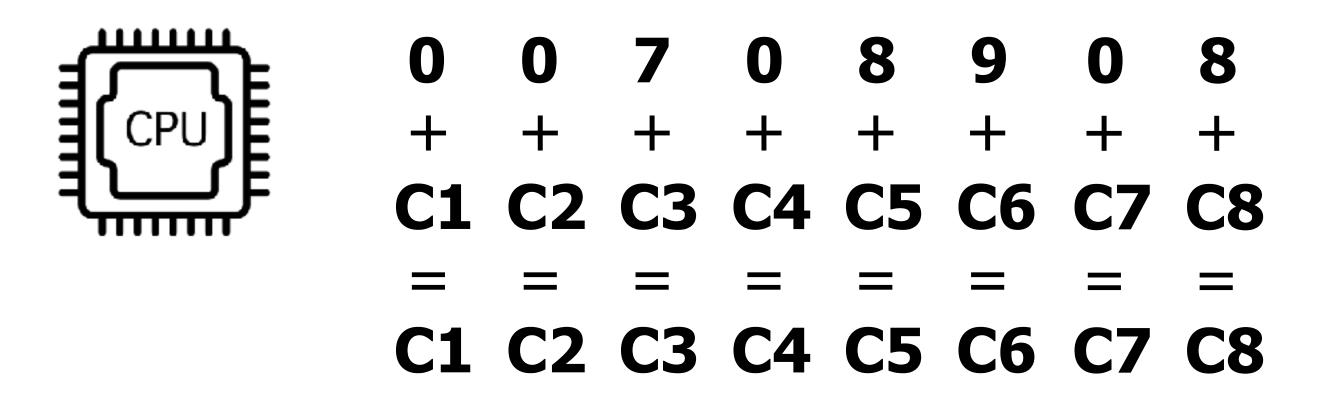


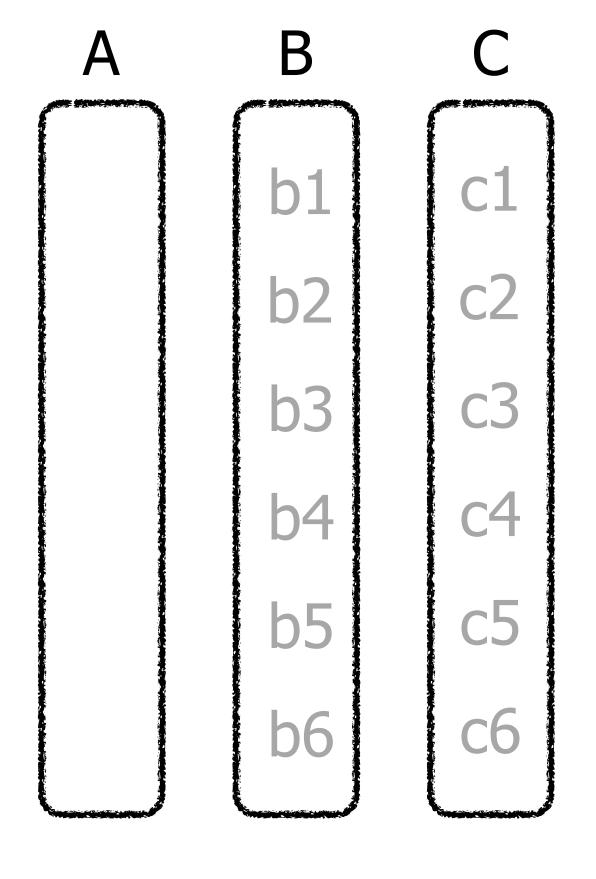
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)



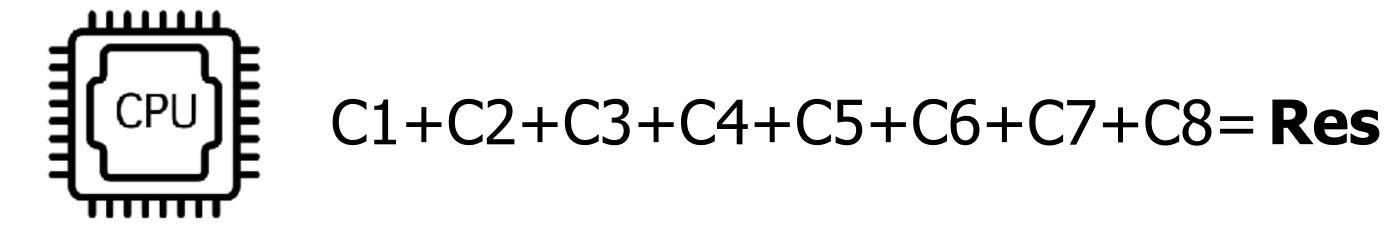


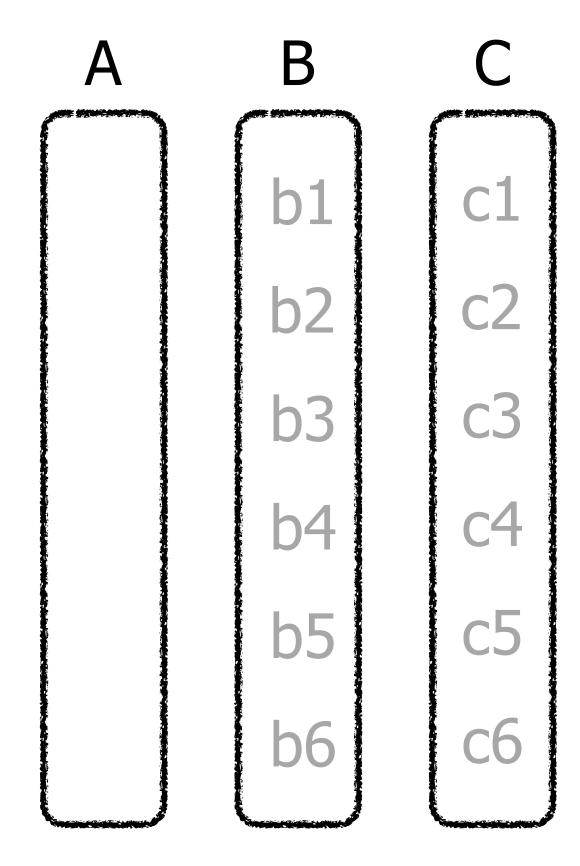
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)





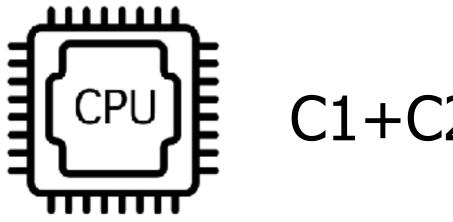
Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)



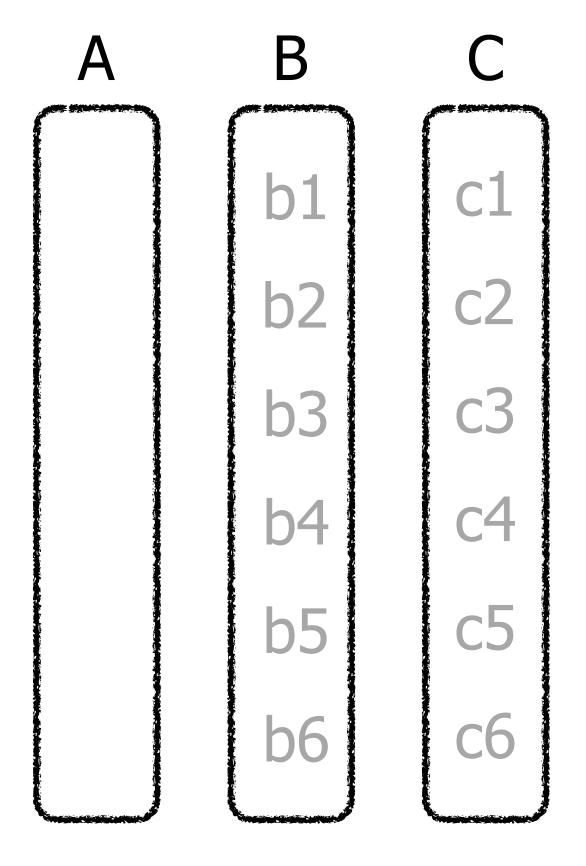


Apply one instruction in parallel to multiple values within one cache line (e.g., 128-256 bits at a time)

Select sum(A) from table where A > 5

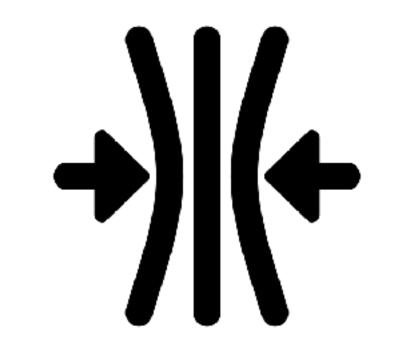


$$C1+C2+C3+C4+C5+C6+C7+C8=Res$$



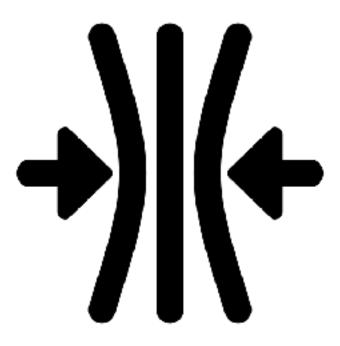
Works best with late materialization:)

## Compression



## Compression

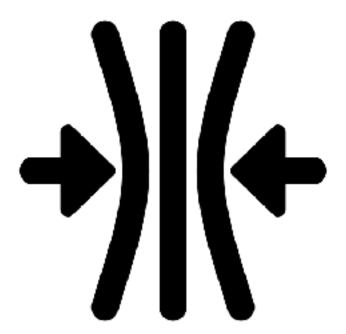
#### crucial for column-stores



### Compression

crucial for column-stores

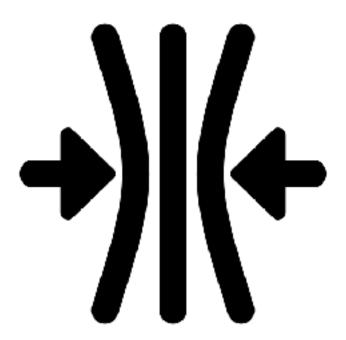
The reason is not only to save space but to improve performance. How?



### Compression

crucial for column-stores

The reason is not only to save space but to improve performance. How?



CPU cost of compression and decompression

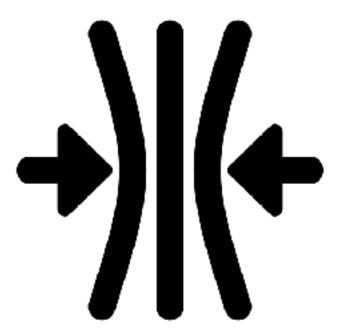
cost savings of moving data across the memory hierarchy

### Compression

crucial for column-stores

The reason is not only to save space but to improve performance. How?

Important for compressed values to be fixed-size to support positional ID lookups



### (1) Bit-Vector Encoding

Employ one bit string for each possible value indicating if the entry has the given value

Pet

Cat

Dog

Cat

Cat

Dog

Cat

Horse

Cat

Dog

### (1) Bit-Vector Encoding

Employ one bit string for each possible value indicating if the entry has the given value

Cat	Dog	Pet
1	0	Cat
0	1	Dog
1	0	Cat
1	0	Cat
0	1	Dog
1	0	Cat
0	0	Horse
1	0	Cat
0	1	Dog
0	1	Dog
	1 0 1 0 1 0	<ul> <li>0</li> <li>1</li> <li>1</li> <li>0</li> <li>1</li> <li>0</li> <li>0</li> <li>1</li> <li>0</li> <li>0</li> <li>1</li> <li>0</li> <li>1</li> <li>0</li> <li>1</li> <li>0</li> <li>1</li> <li>1</li> <li>0</li> <li>1</li> <li>1&lt;</li></ul>

### (1) Bit-Vector Encoding

Employ one bit string for each possible value indicating if the entry has the given value

e.g., select \* from table where specie = "Cat"

**Pros and cons?** 



Horse Cat	Dog
-----------	-----

**1** 0

1 0

1 (

0 1

1 (

0

1 (

0 1

0 1

Pet

Cat

Dog

Cat

Cat

Dog

Cat

Horse

Cat

Dog

### (1) Bit-Vector Encoding Pet Horse Cat Dog Cat 0 Employ one bit string for each possible value indicating if the entry has the given value Dog Cat 0 0 e.g., select \* from table where specie = "Cat" Cat 0 Dog 0 0 Cat 0 Pro: Fast to read & compare 1 bit per entry. Horse 0 Con: Only applicable if there are very few values. Cat 0 Dog

Employ a dictionary with smaller strings to represent larger ones

Pet

Cat

Dog

Cat

Cat

Dog

Cat

Horse

Cat

Parrot

Employ a dictionary with smaller strings to represent larger ones

Dictionary

00 Cat

01 Dog

10 Horse

11 Parrot

Pet

Cat

Dog

Cat

Cat

Dog

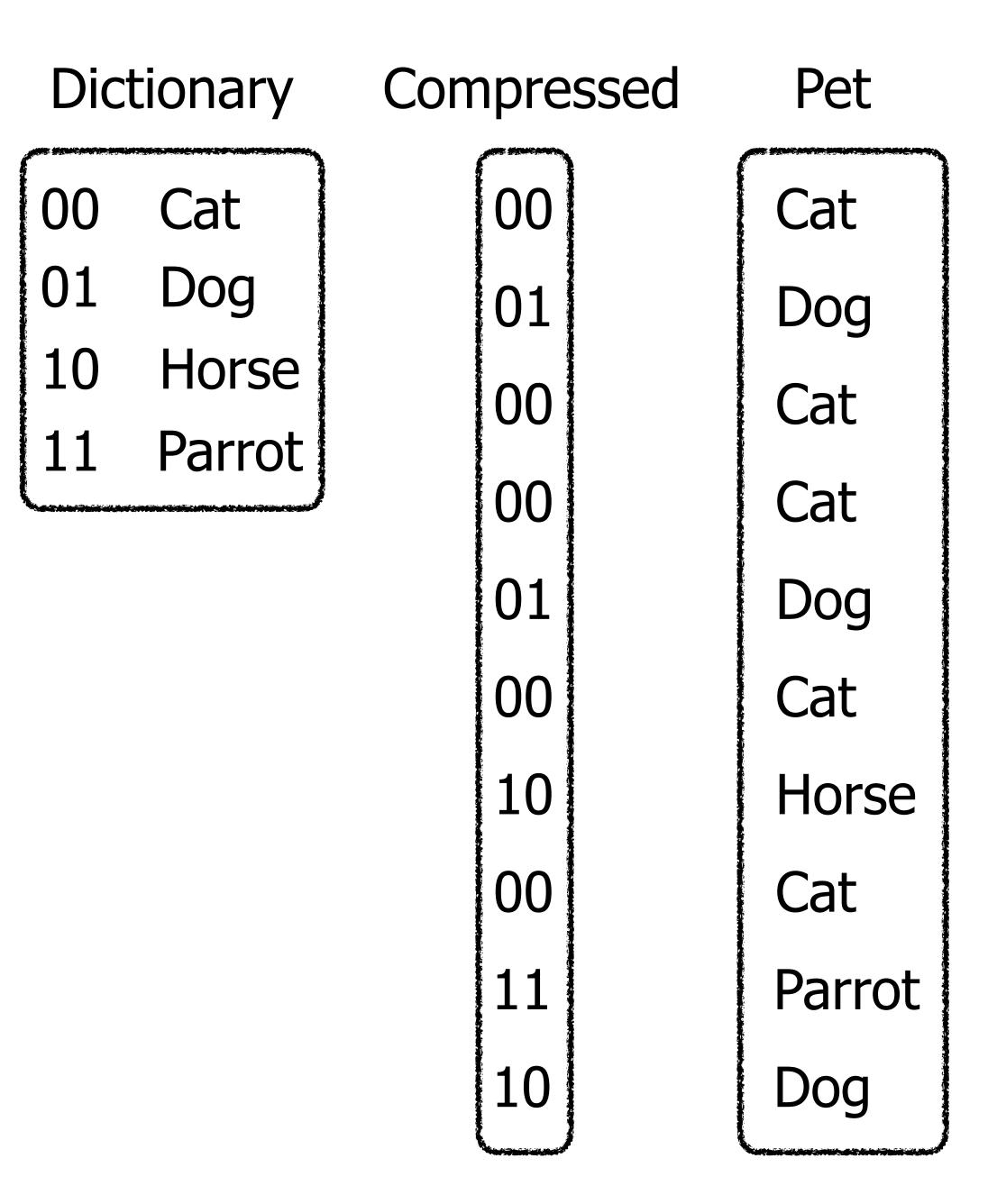
Cat

Horse

Cat

**Parrot** 

Employ a dictionary with smaller strings to represent larger ones



Employ a dictionary with smaller strings to represent larger ones

Dictionary

Compressed

Pet

00

01

Cat

Dog

10 Horse

11 Parrot

00

01

00

01

00

00

10

00

11

10

Dog

Cat

Cat

Cat

Dog

Cat

Horse

Cat

Parrot

Dog

Compare to bit-vector encoding from before



#### (2) Dictionary Encoding Dictionary Compressed Pet Cat 00 00 Cat Employ a dictionary with smaller strings to represent larger ones 01 Dog 01 Dog 10 Horse 00 Cat Parrot 00 Cat 01 Dog 00 Cat Pro: applicable across larger value spaces 10 Horse Con: slower as we need to read more bits 00 Cat **Parrot** Dog

Represent repeating values using one entry

Pet

Cat

Dog

Cat

Cat

Cat

Cat

Horse

Cat

Cat

Represent repeating values using one entry

Compressed

Cat

Dog

**Cat · 4** 

Horse

**Cat · 2** 

Dog

Pet

Cat

Dog

Cat

Cat

Cat

Cat

Horse

Cat

Cat

Represent repeating values using one entry

Compatible with dictionary encoding

Compressed

Cat

Dog

**Cat** • 4

Horse

**Cat · 2** 

Dog

Pet

Cat

Dog

Cat

Cat

Cat

Cat

Horse

Cat

Cat

Represent repeating values using one entry

Compatible with dictionary encoding

Cons?



Compressed

Cat

Dog

**Cat** • 4

Horse

**Cat · 2** 

Dog

Pet

Cat

Dog

Cat

Cat

Cat

Cat

Horse

Cat

Cat

Represent repeating values using one entry

Compatible with dictionary encoding

Cons: must scan column to get entry at a given offset

Compressed

Cat

Dog

**Cat** • 4

Horse

**Cat · 2** 

Dog

Pet

Cat

Dog

Cat

Cat

Cat

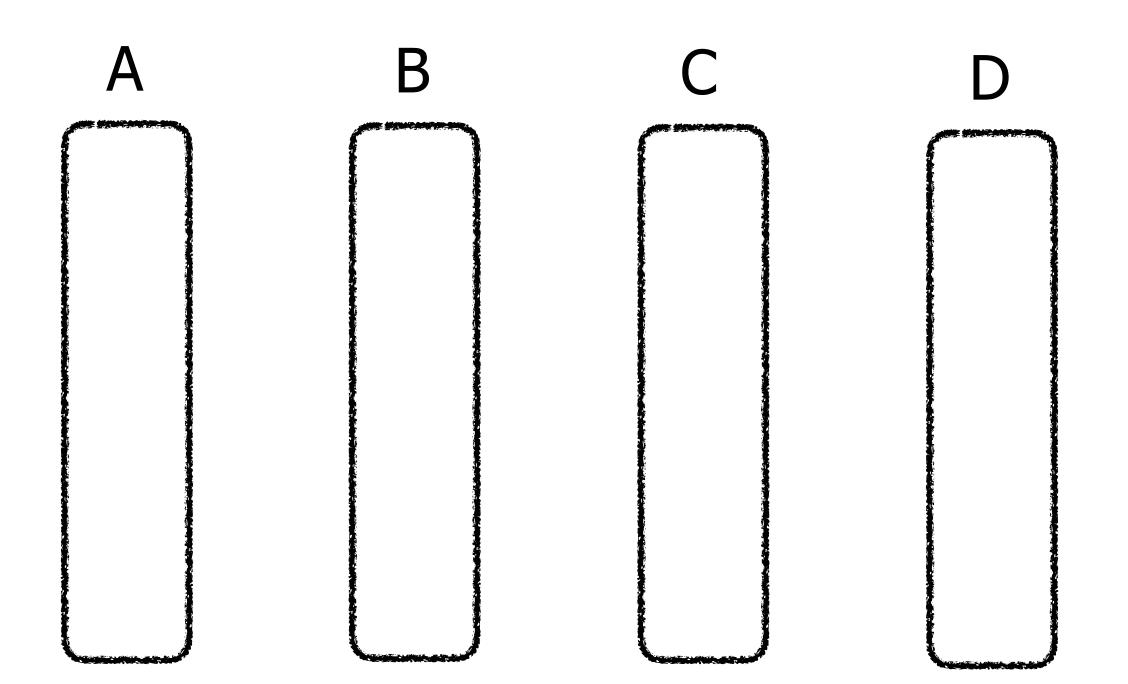
Cat

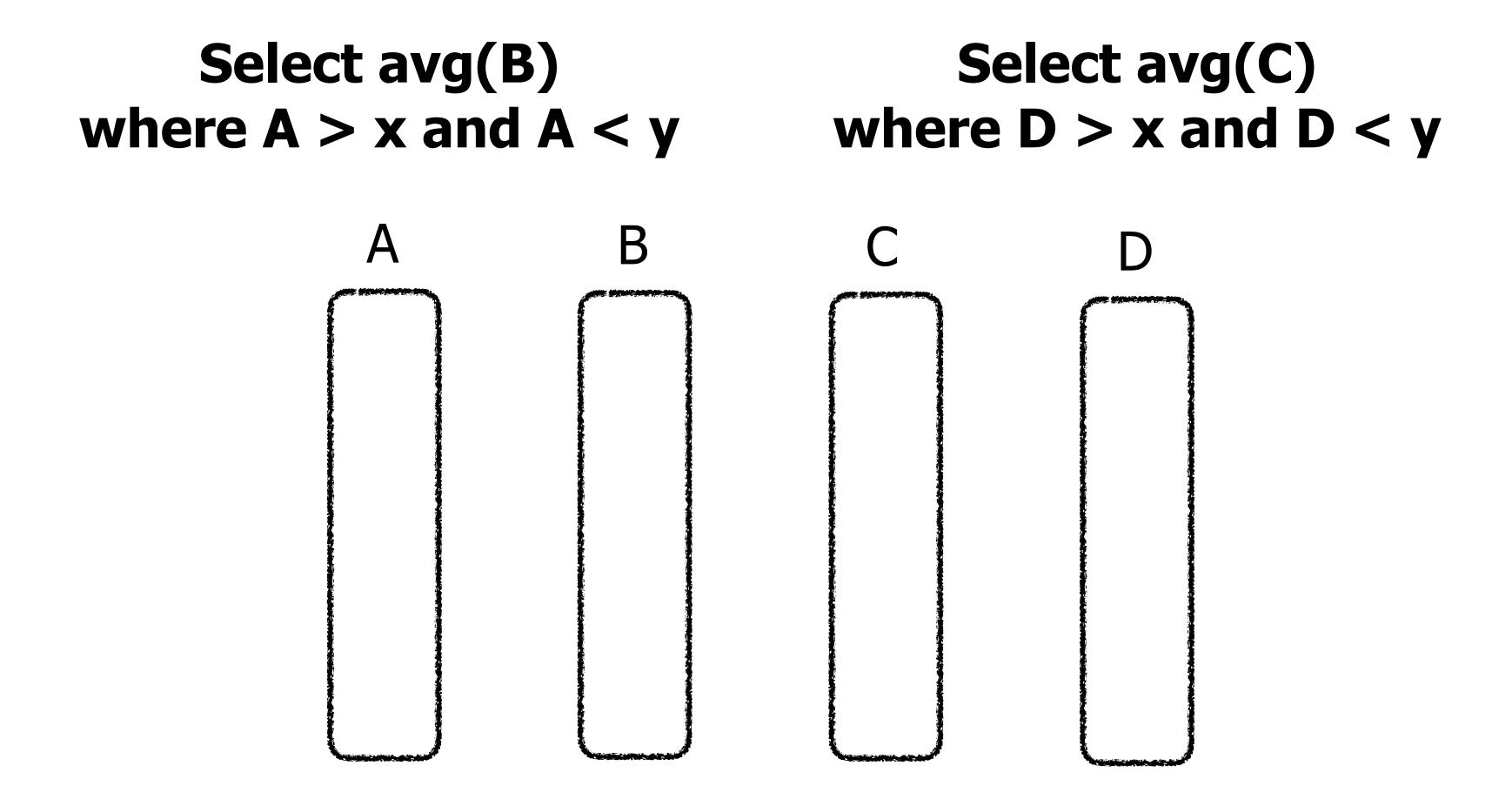
Horse

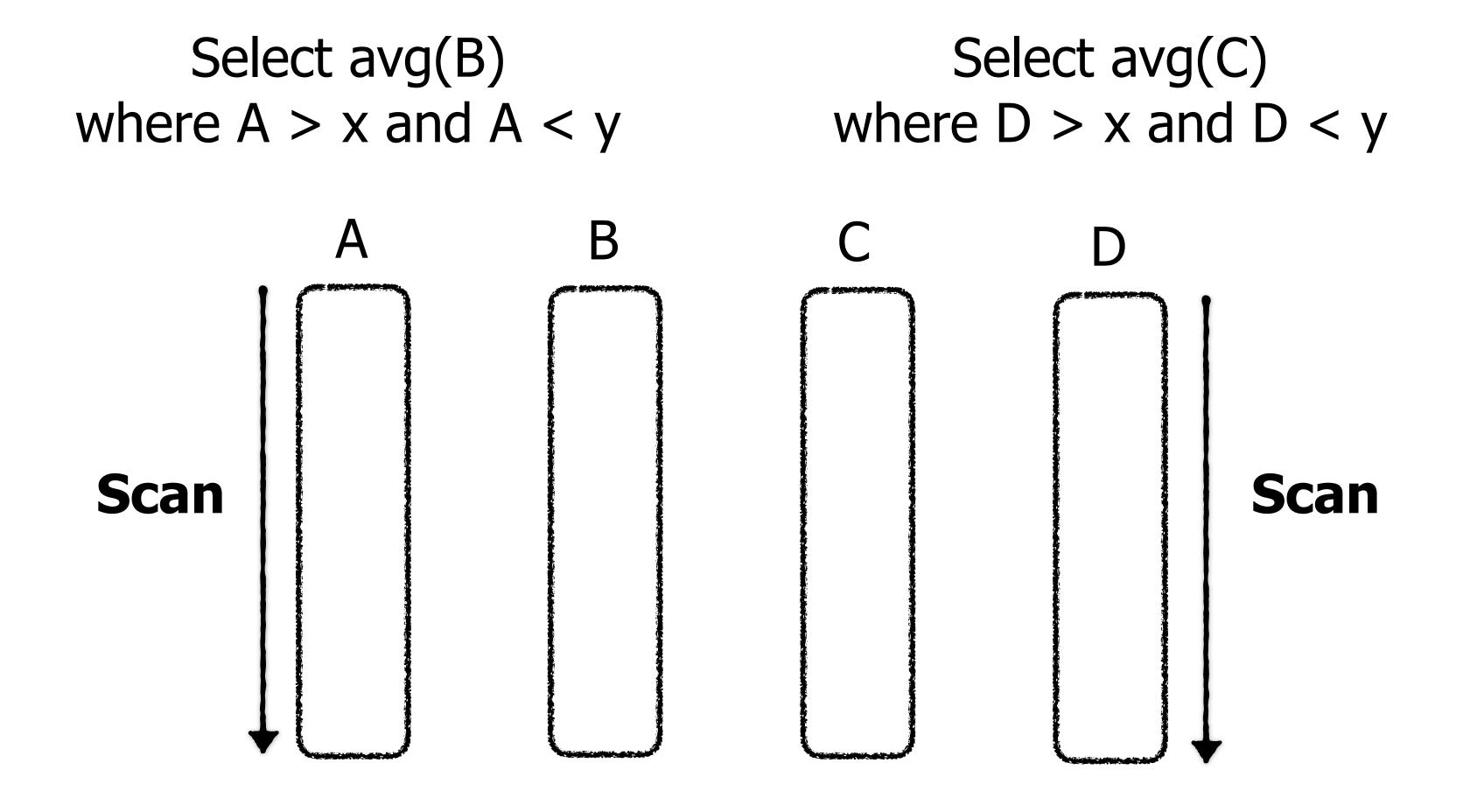
Cat

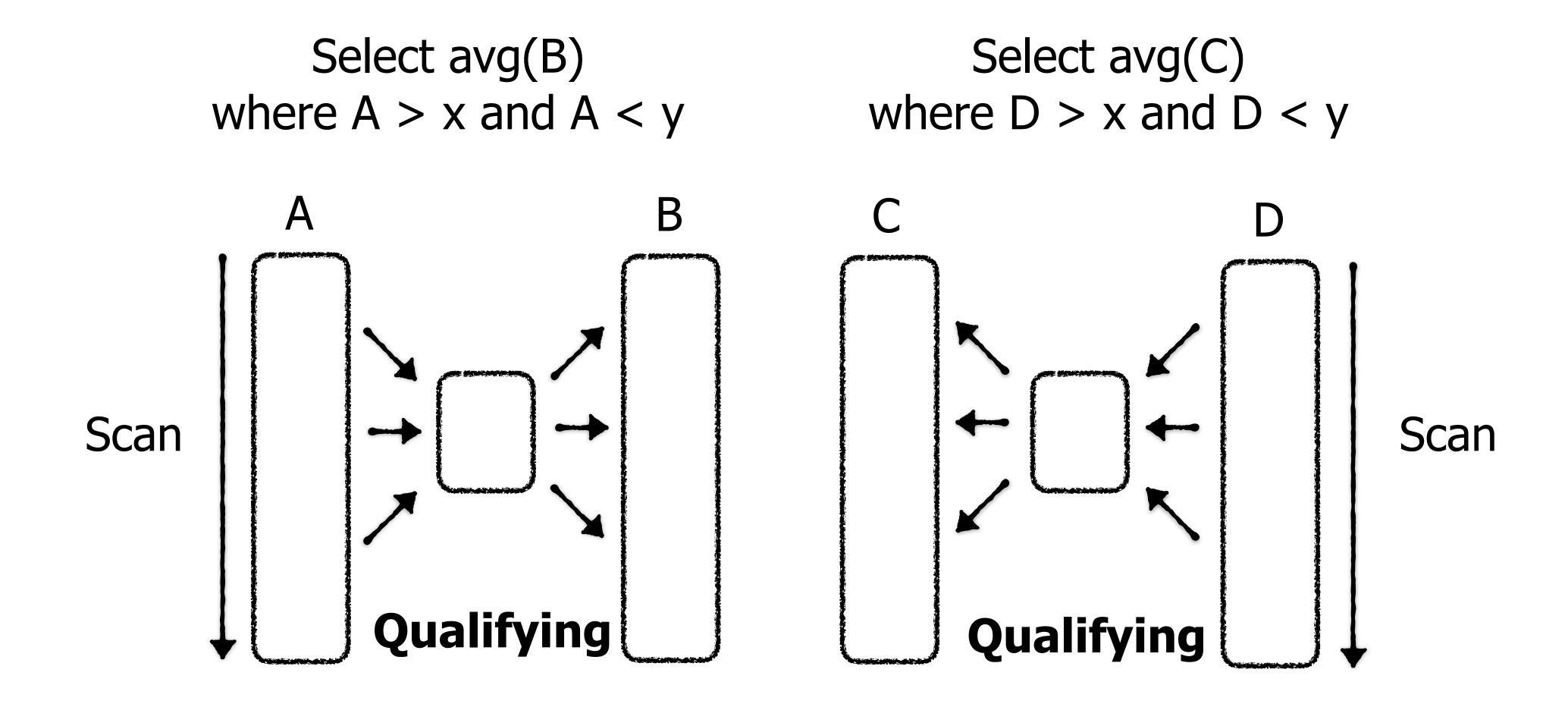
Cat

**Indexing in Column Stores** 

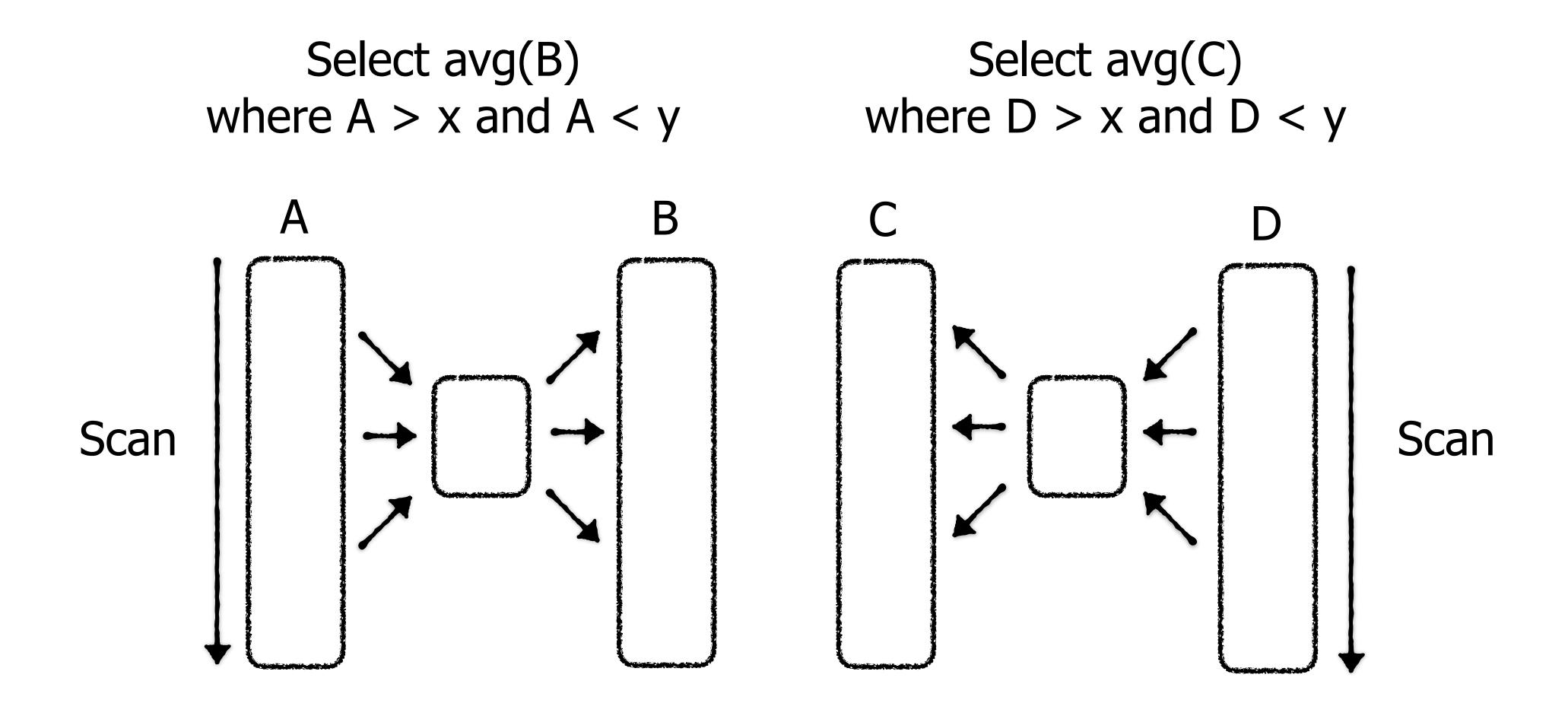




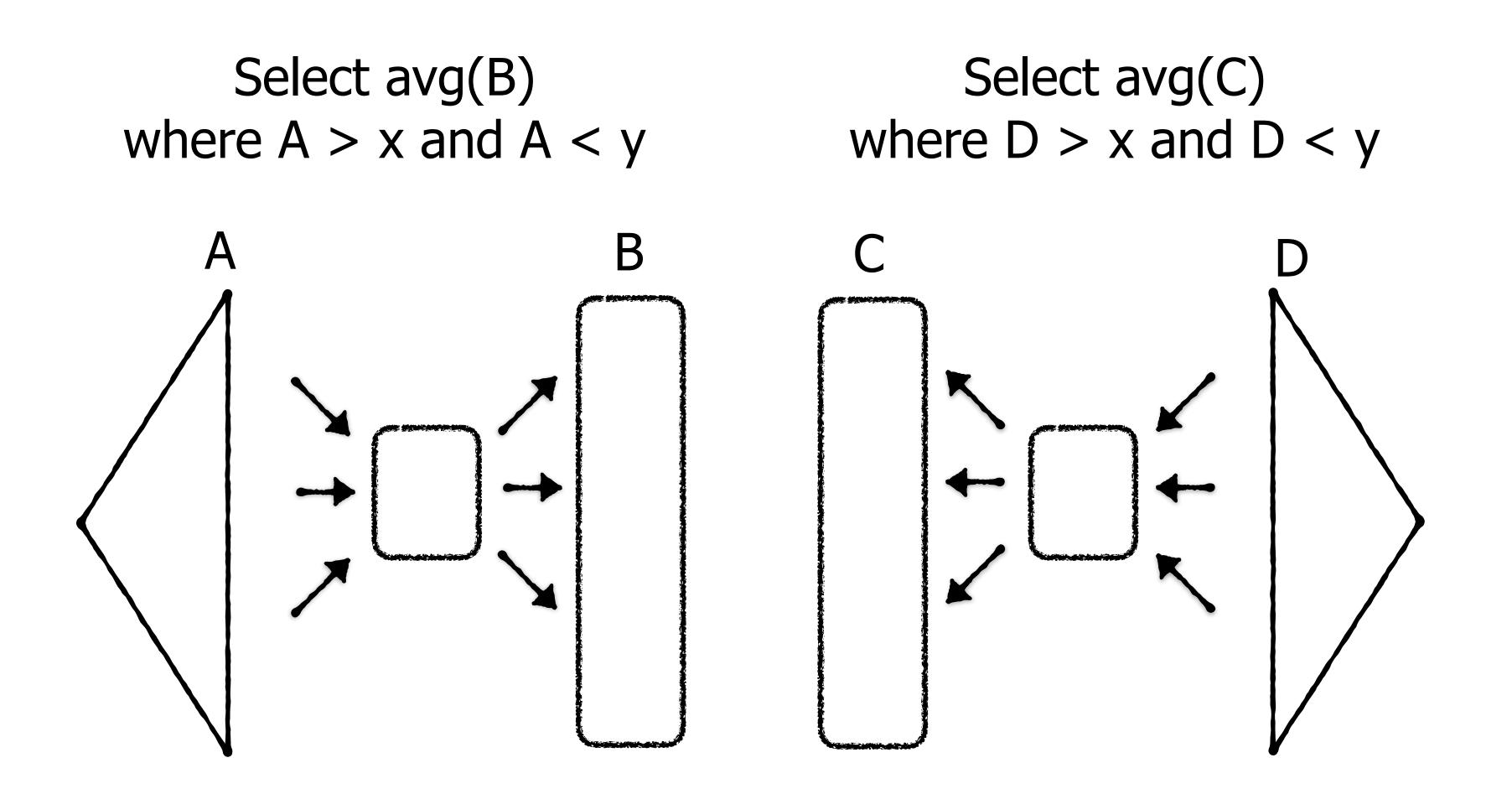




We can answer both queries with full scan and late materialization

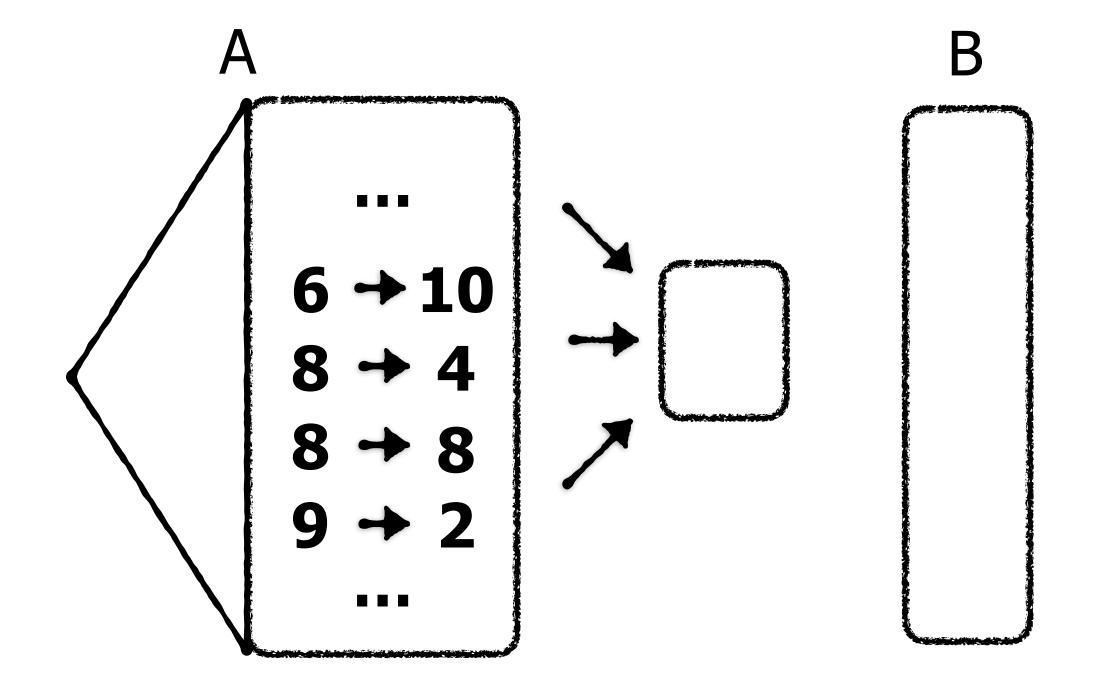


But can we avoid the full scans?



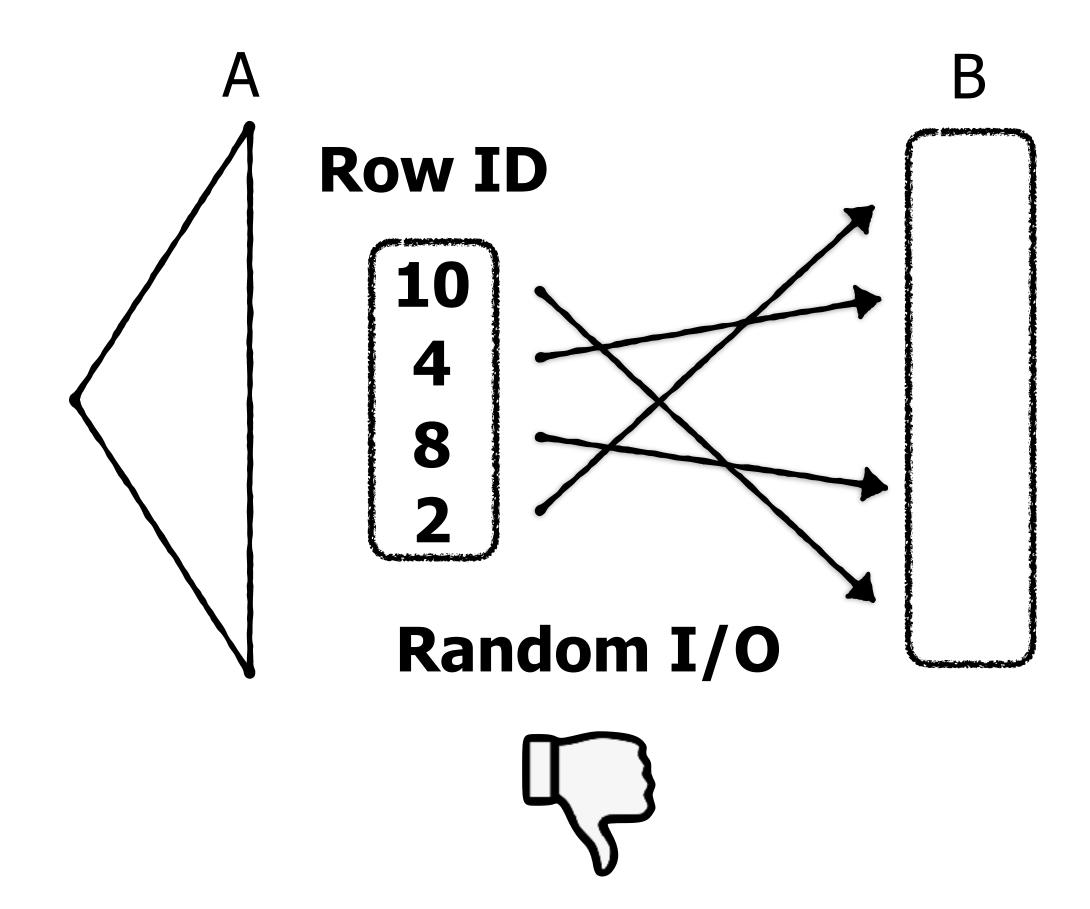
We can employ indexes mapping from key to positional ID

Select avg(B) where A > 5 and A < 10

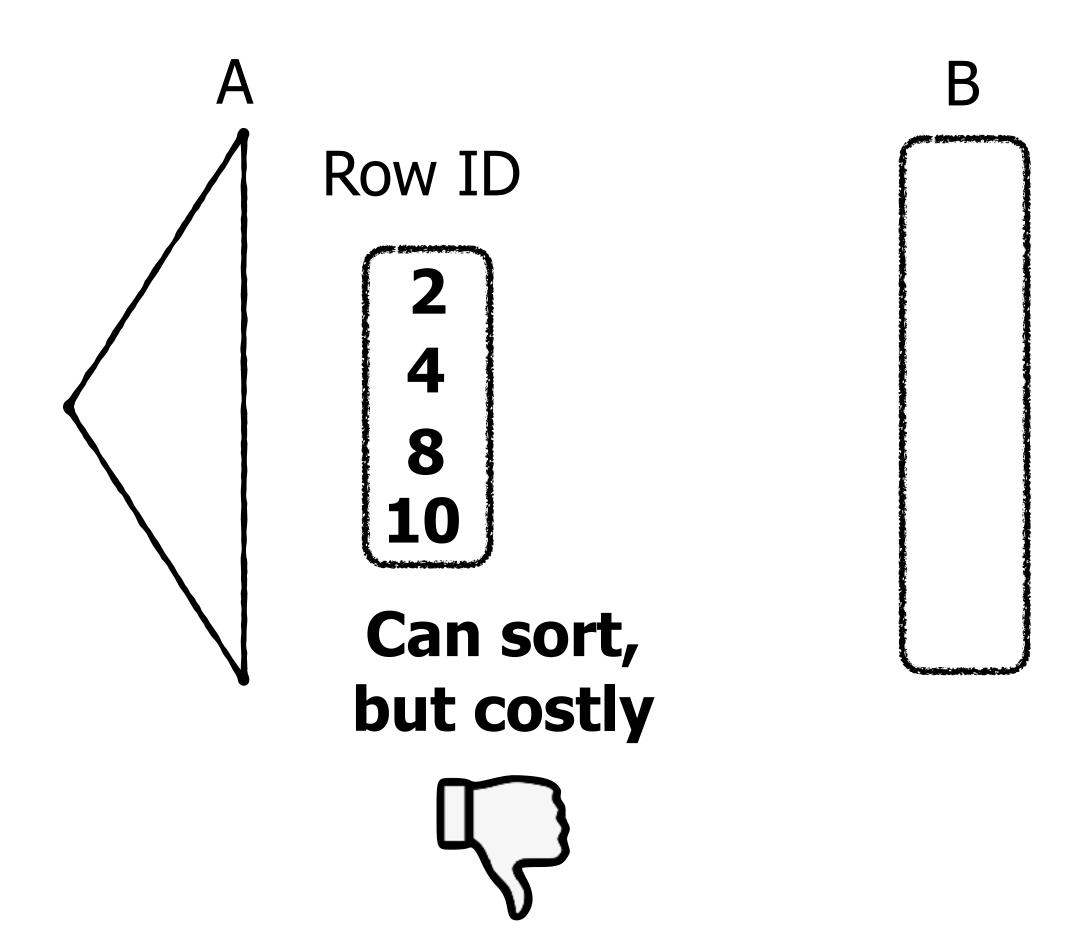


We can employ indexes mapping from key to positional ID

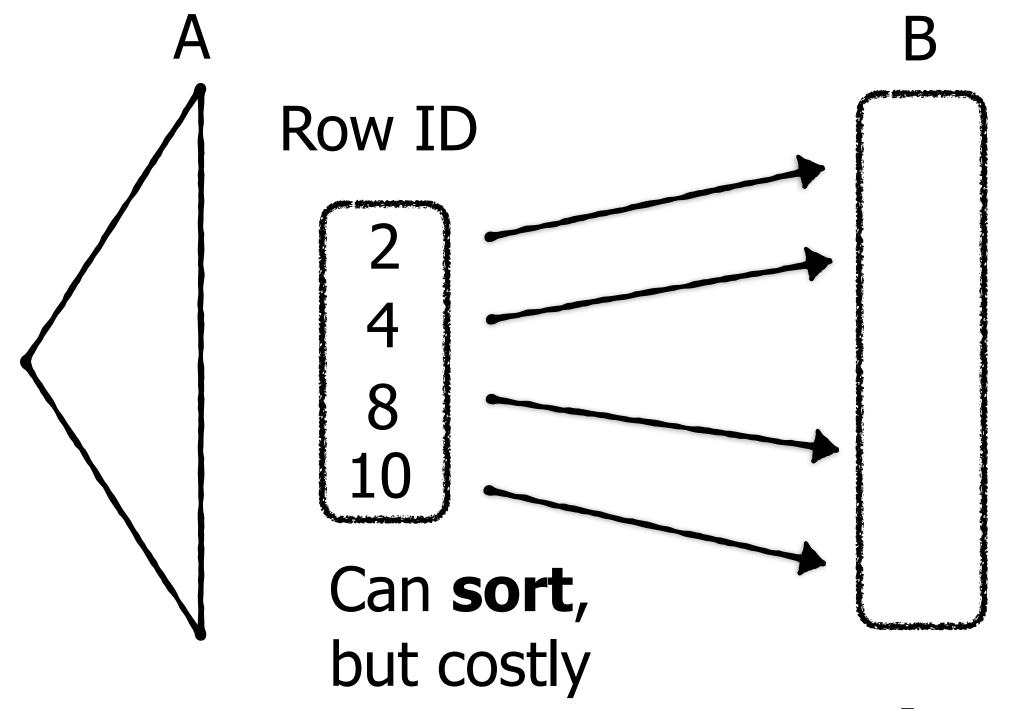
Select avg(B) where A > 5 and A < 10



# Select avg(B) where A > 5 and A < 10

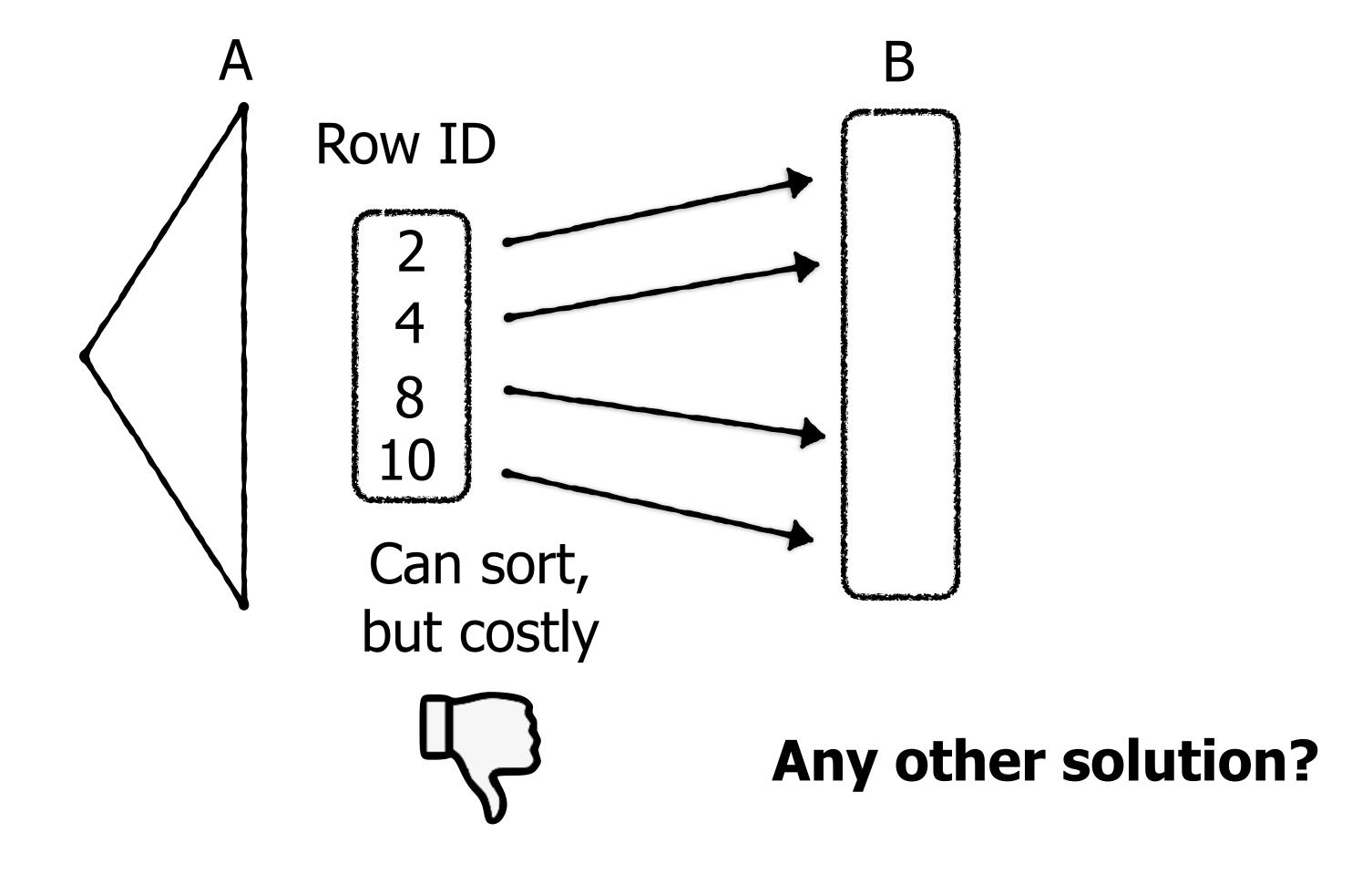


Select avg(B) where A > 5 and A < 10

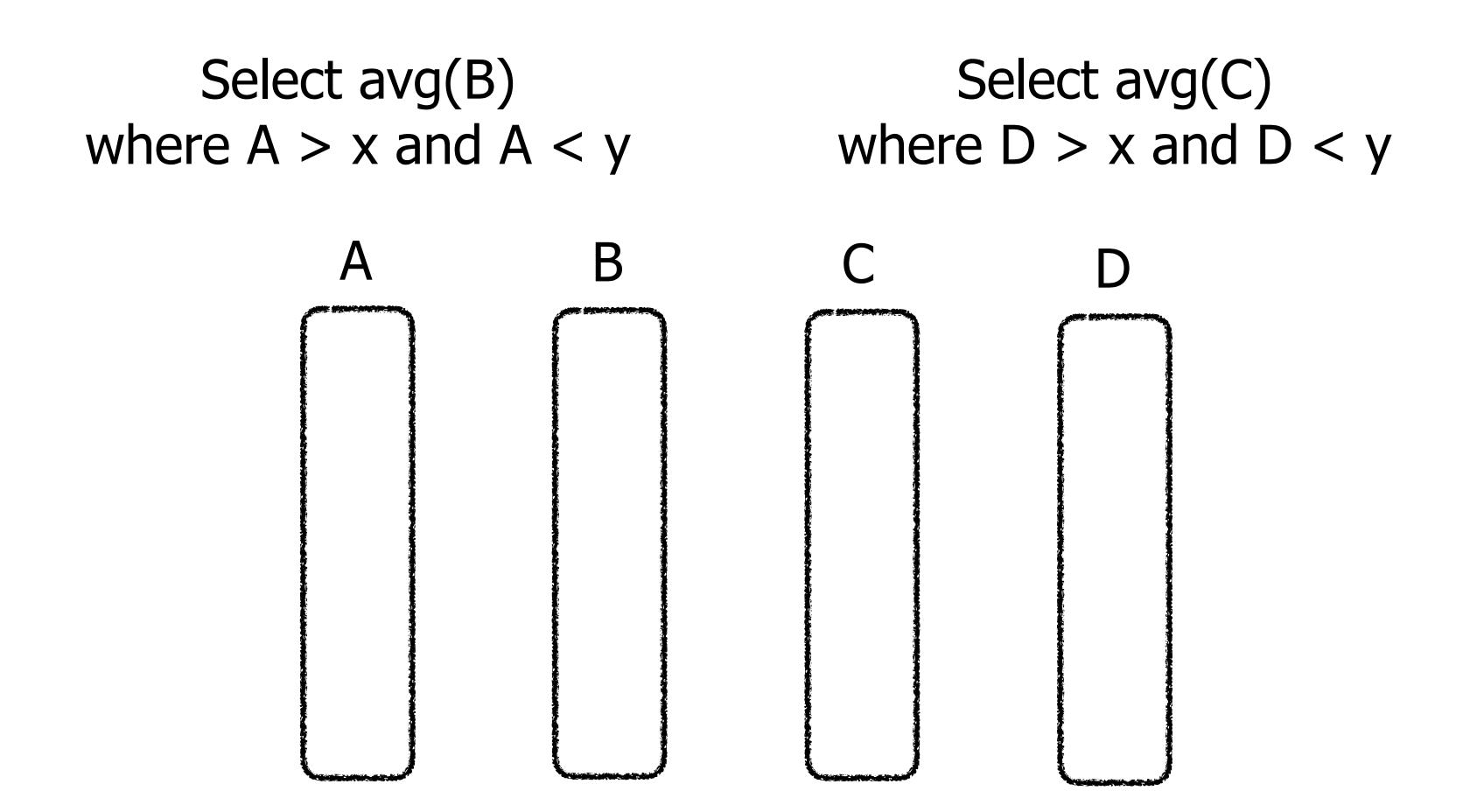


Access to B becomes "skipsequential" rather than random, but the sorting can be an issue

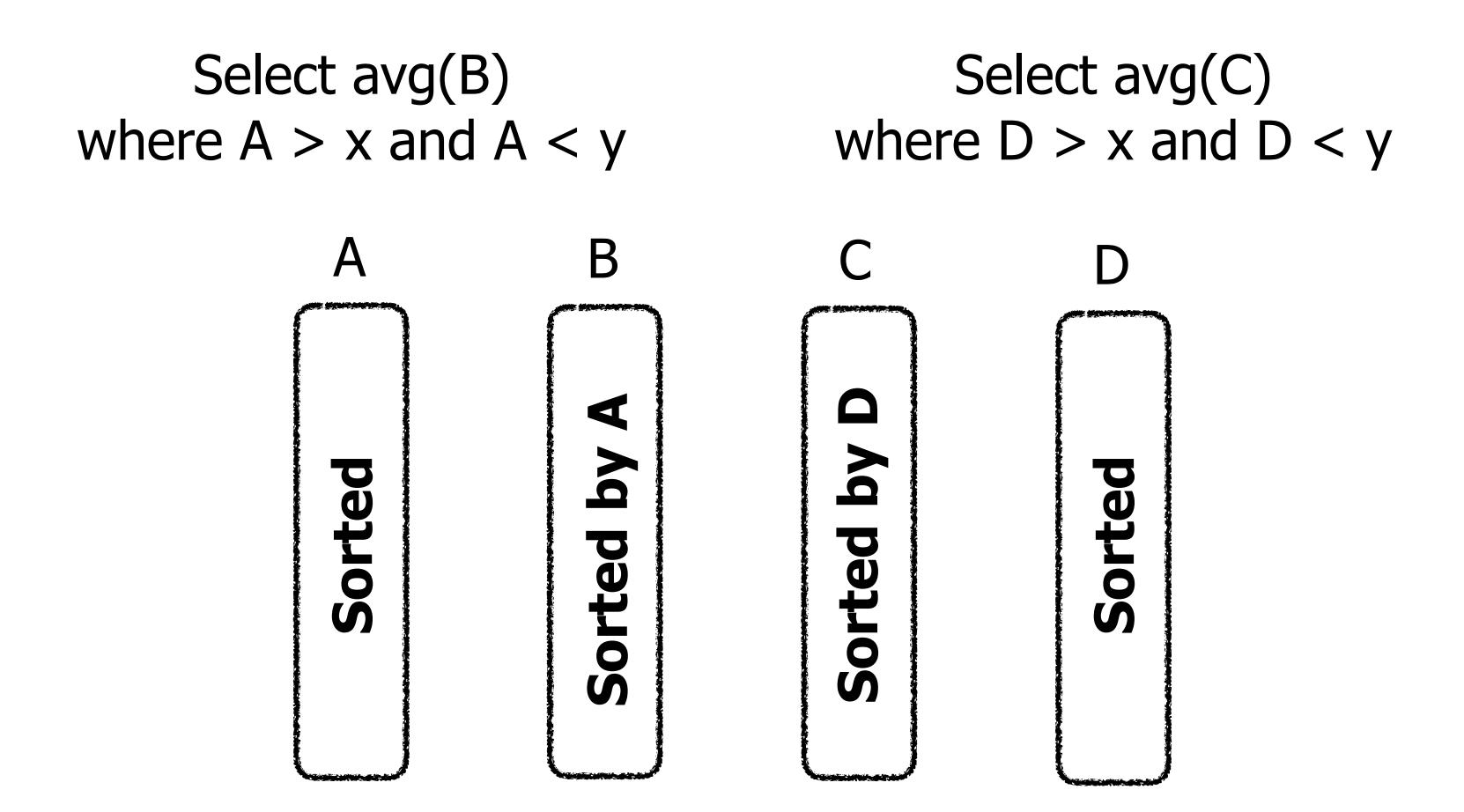
Select avg(B) where A > 5 and A < 10



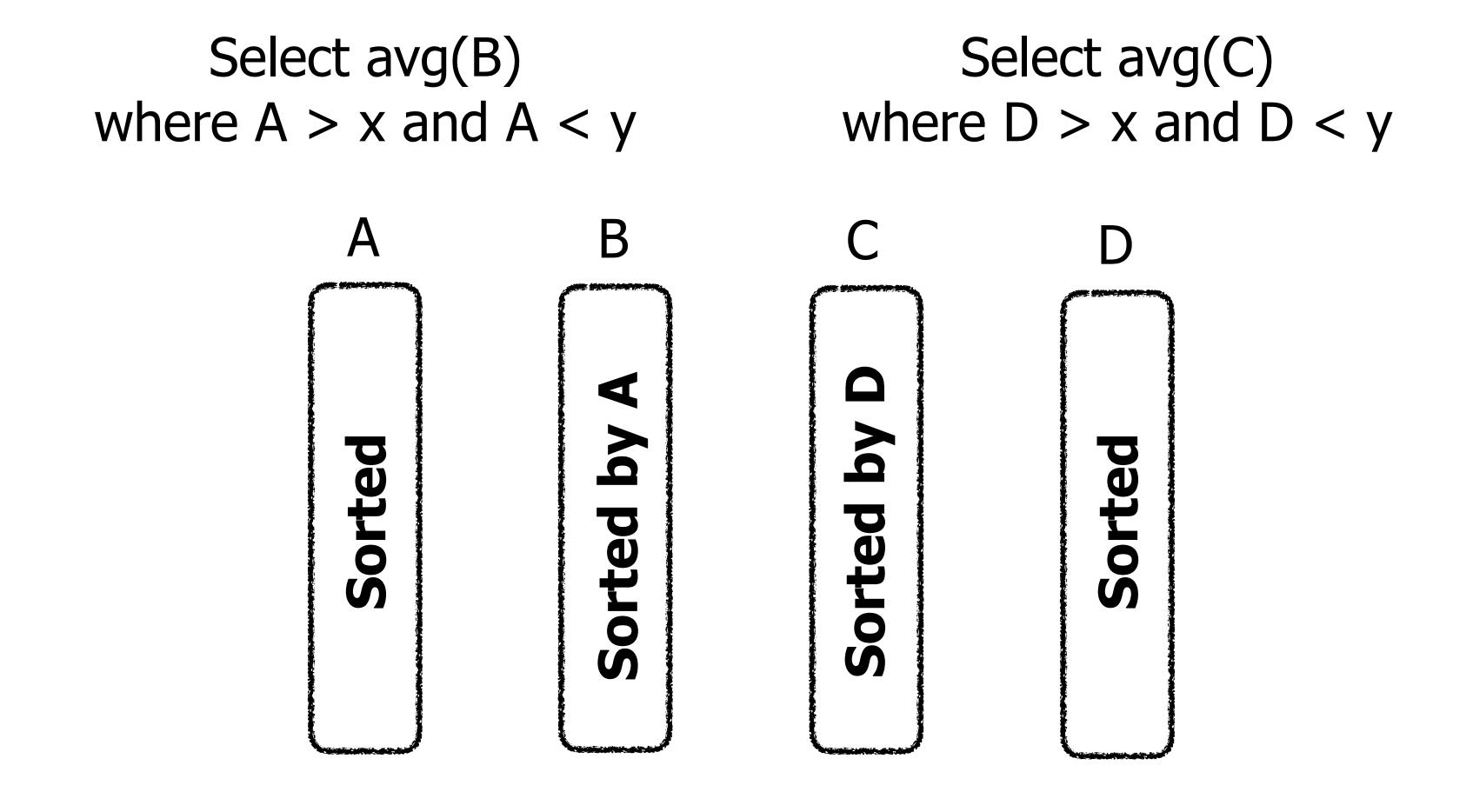
### Column Projections: sort subset of columns by one column



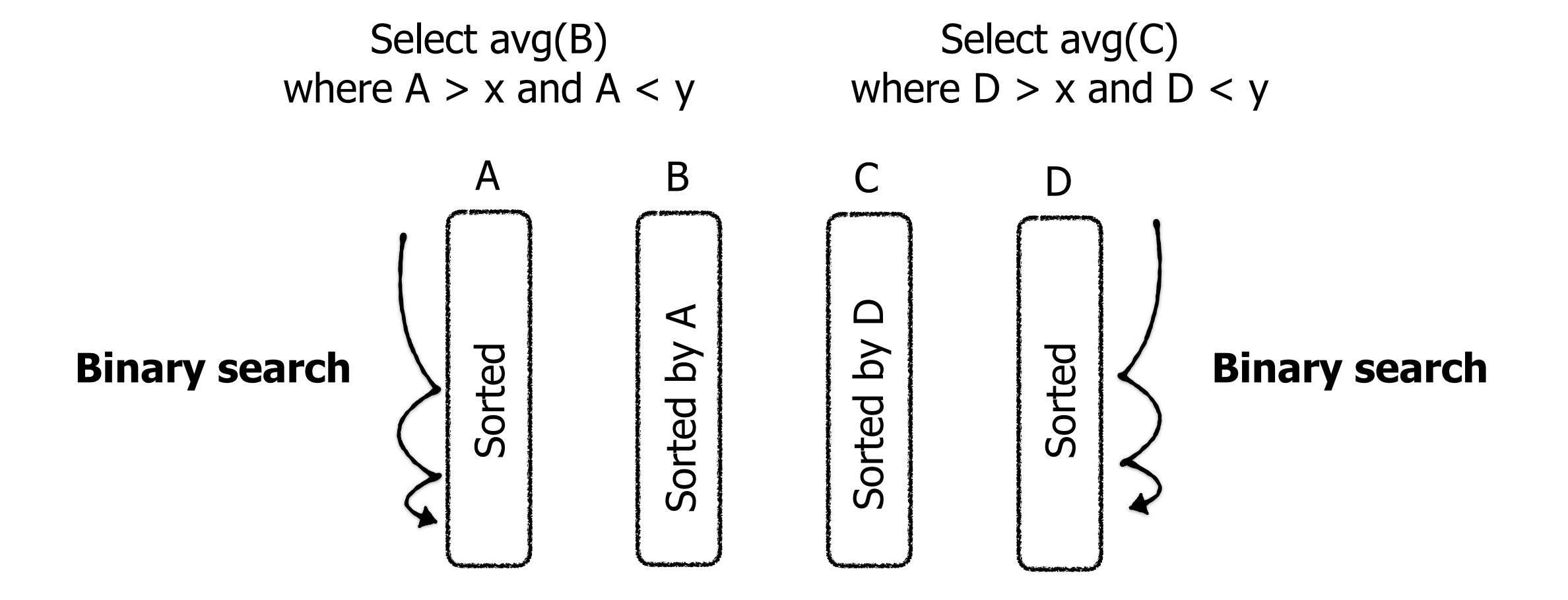
### Column Projections: sort subset of columns by one column

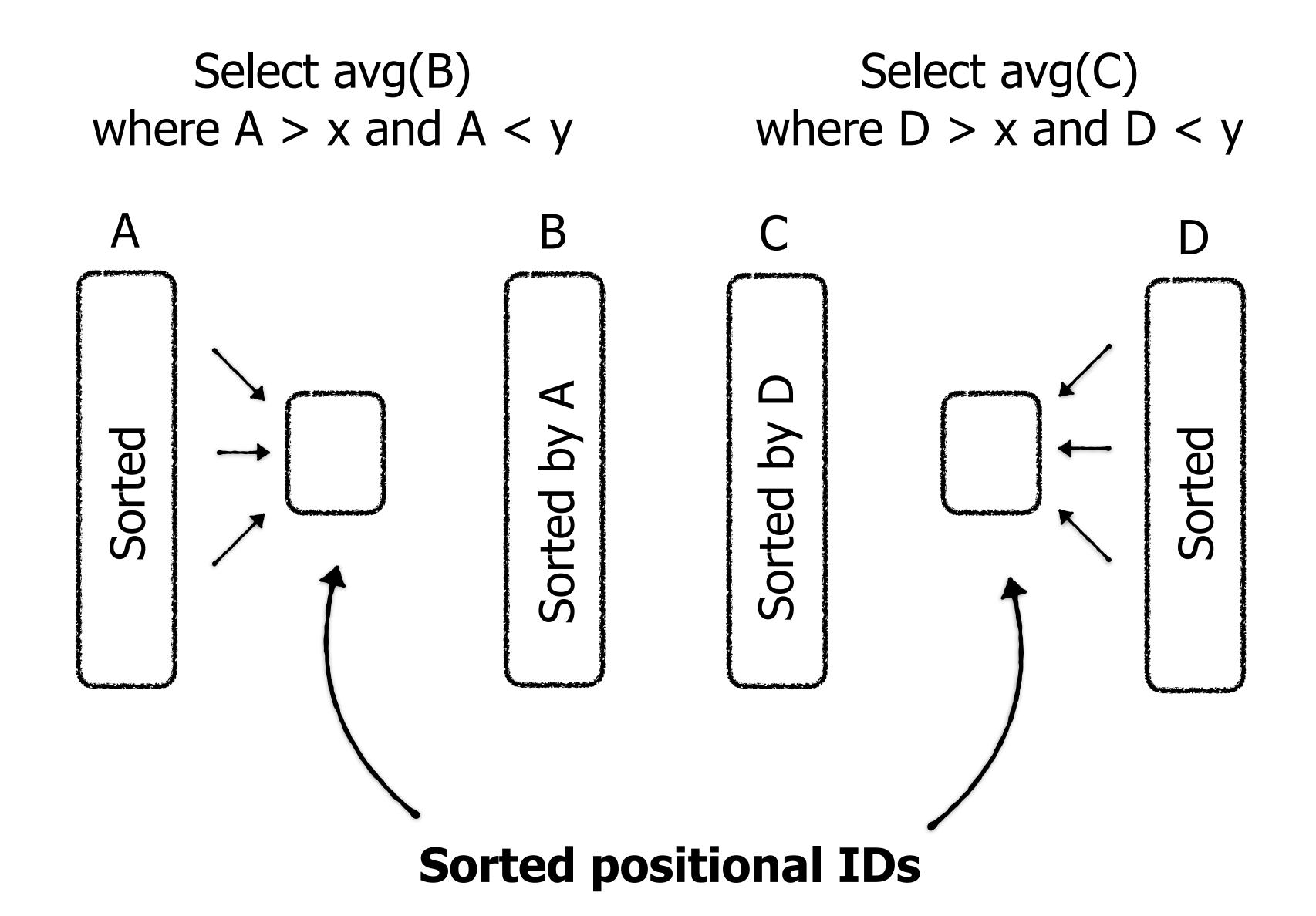


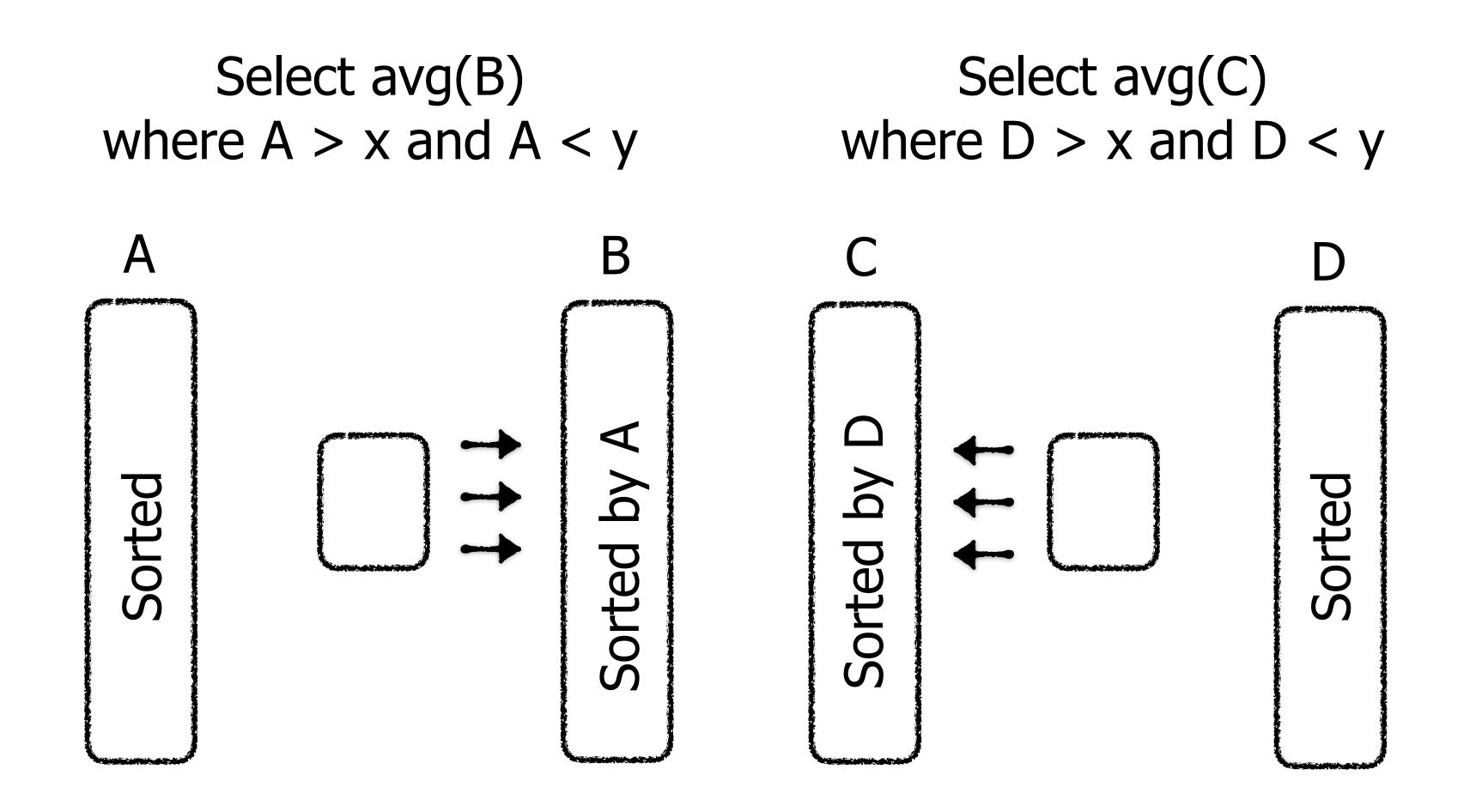
Column Projections: sort subset of columns by one column



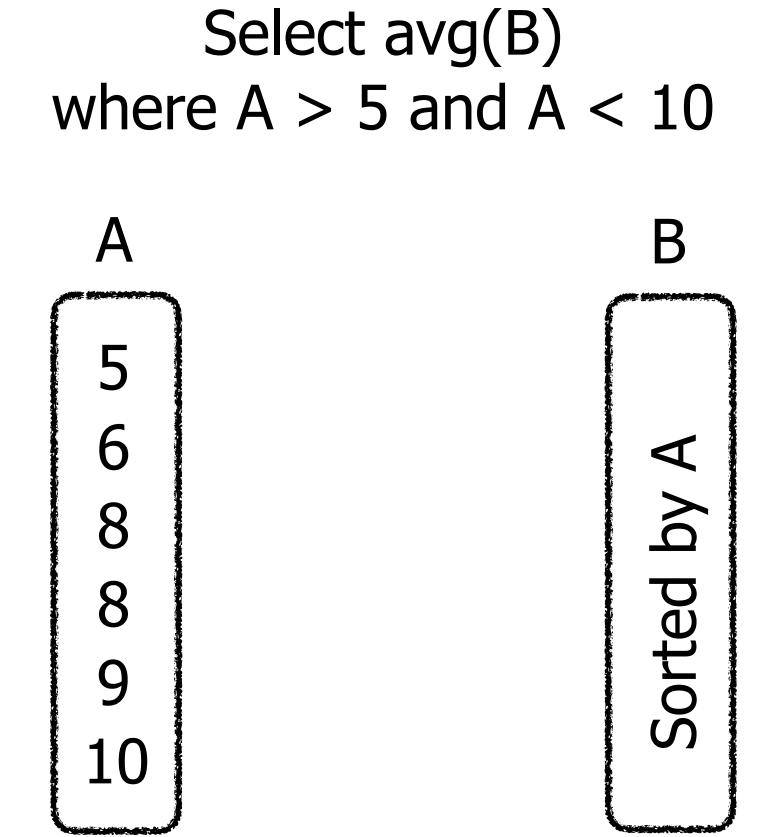
Like having multiple clustered indexes on subset of columns







**Enables sequential access over B and C** 



For example

Select avg(B) where A > 5 and A < 10

Binary search

V

D

O

D

Sorted by A

Binary search

P

Sorted by A

Binary search

Binary search

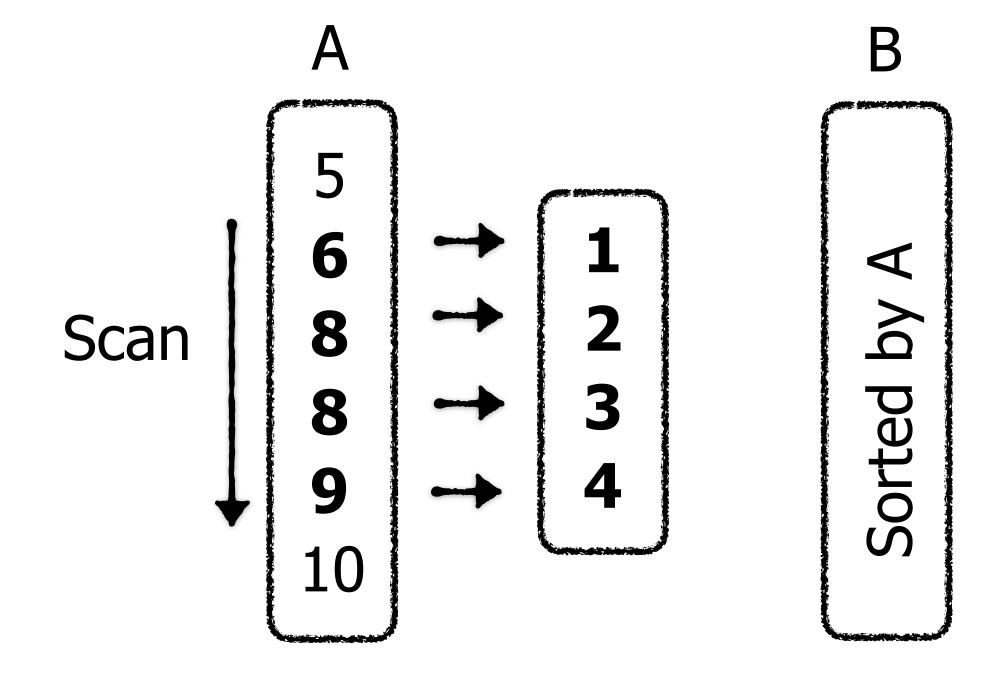
A

Binary search

Binary

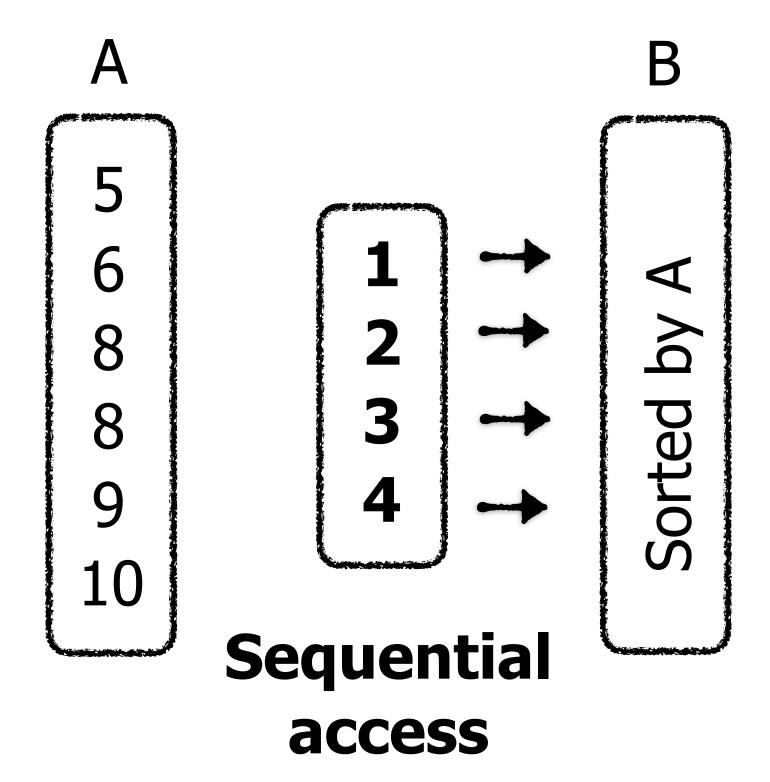
For example

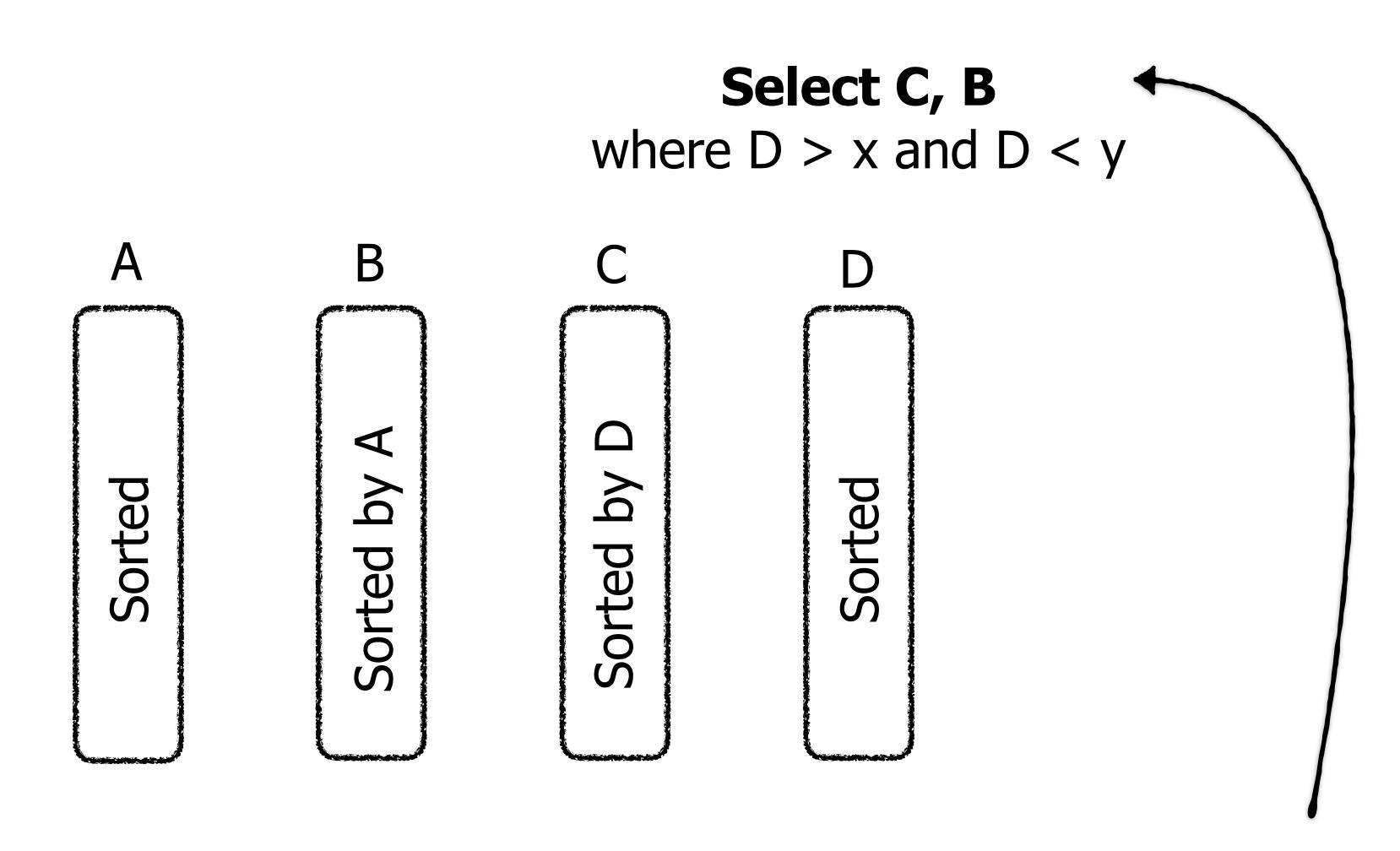
Select avg(B) where A > 5 and A < 10



For example

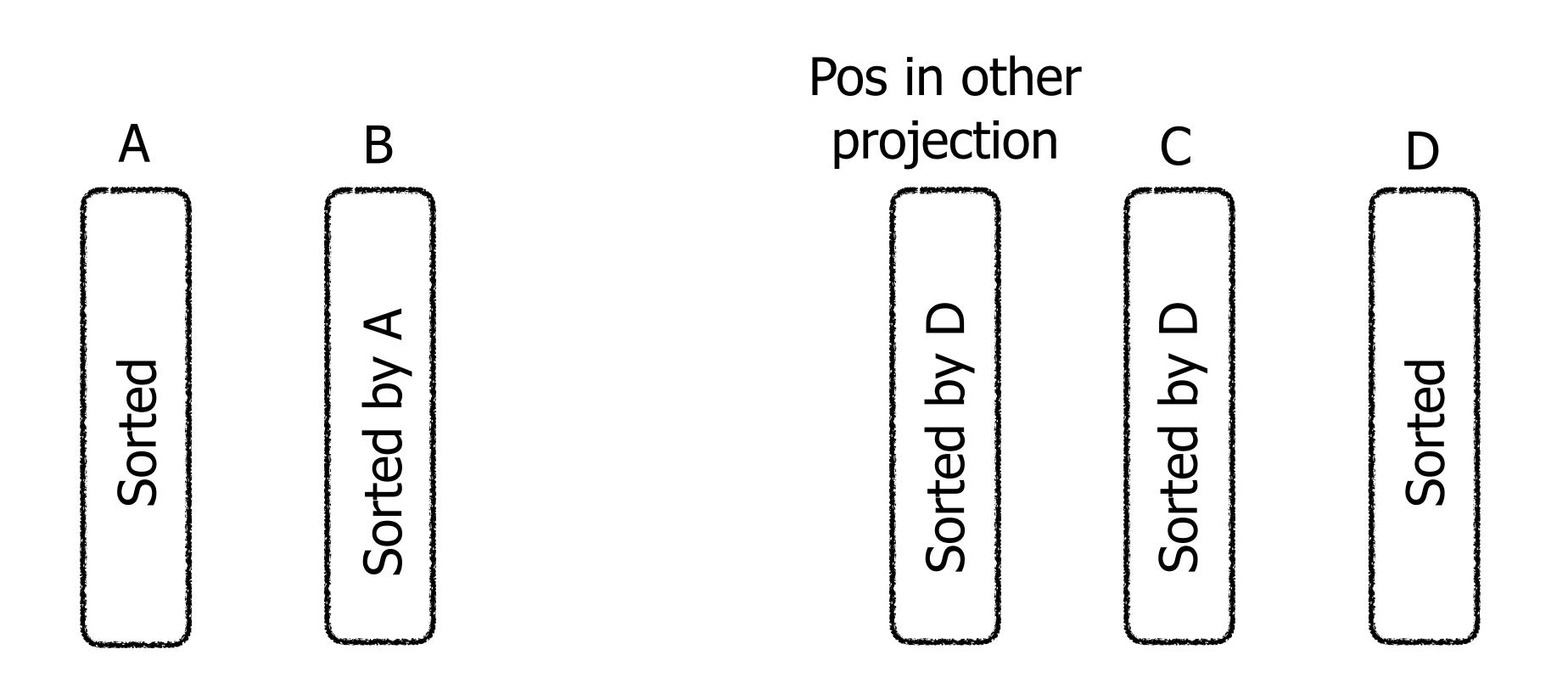
Select avg(B) where A > 5 and A < 10





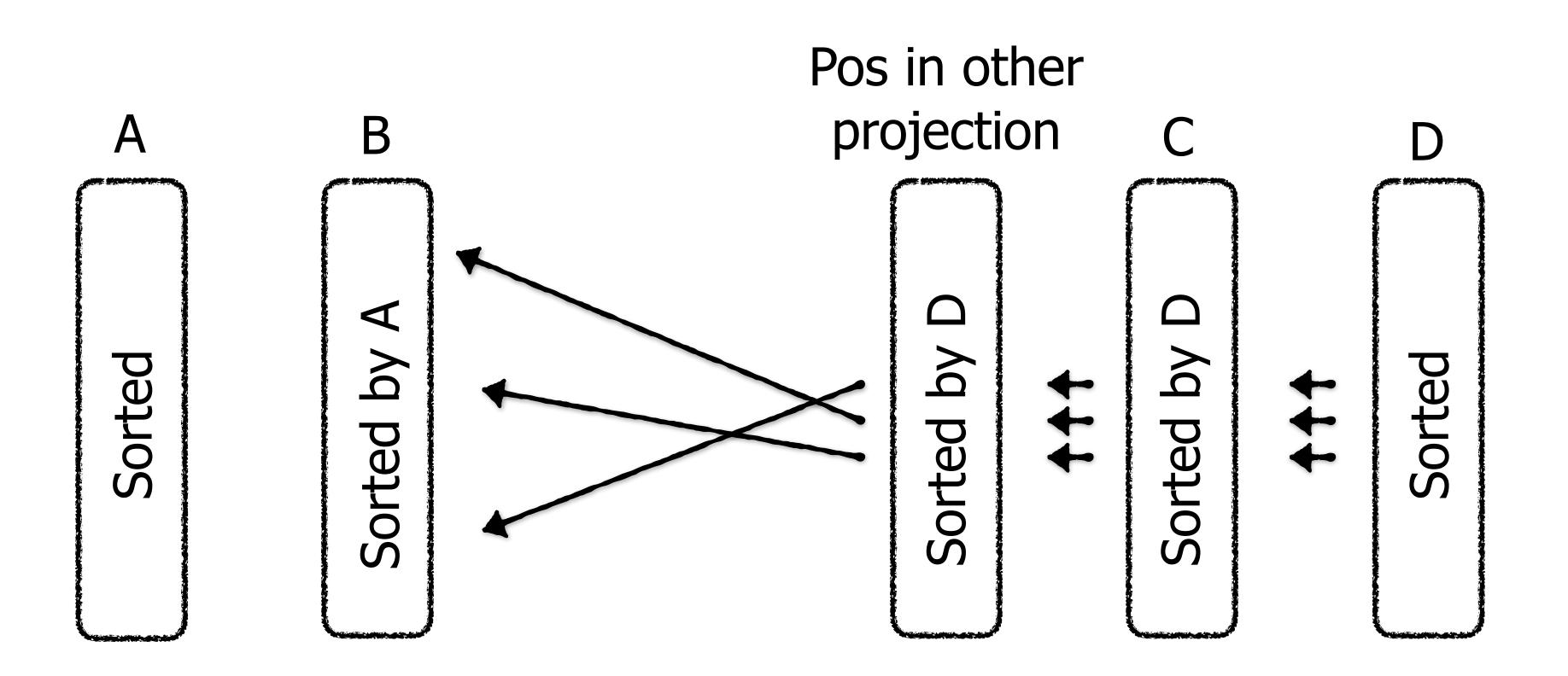
What if the queries do not target mutually exclusive columns?

Select C, B where D > x and D < y

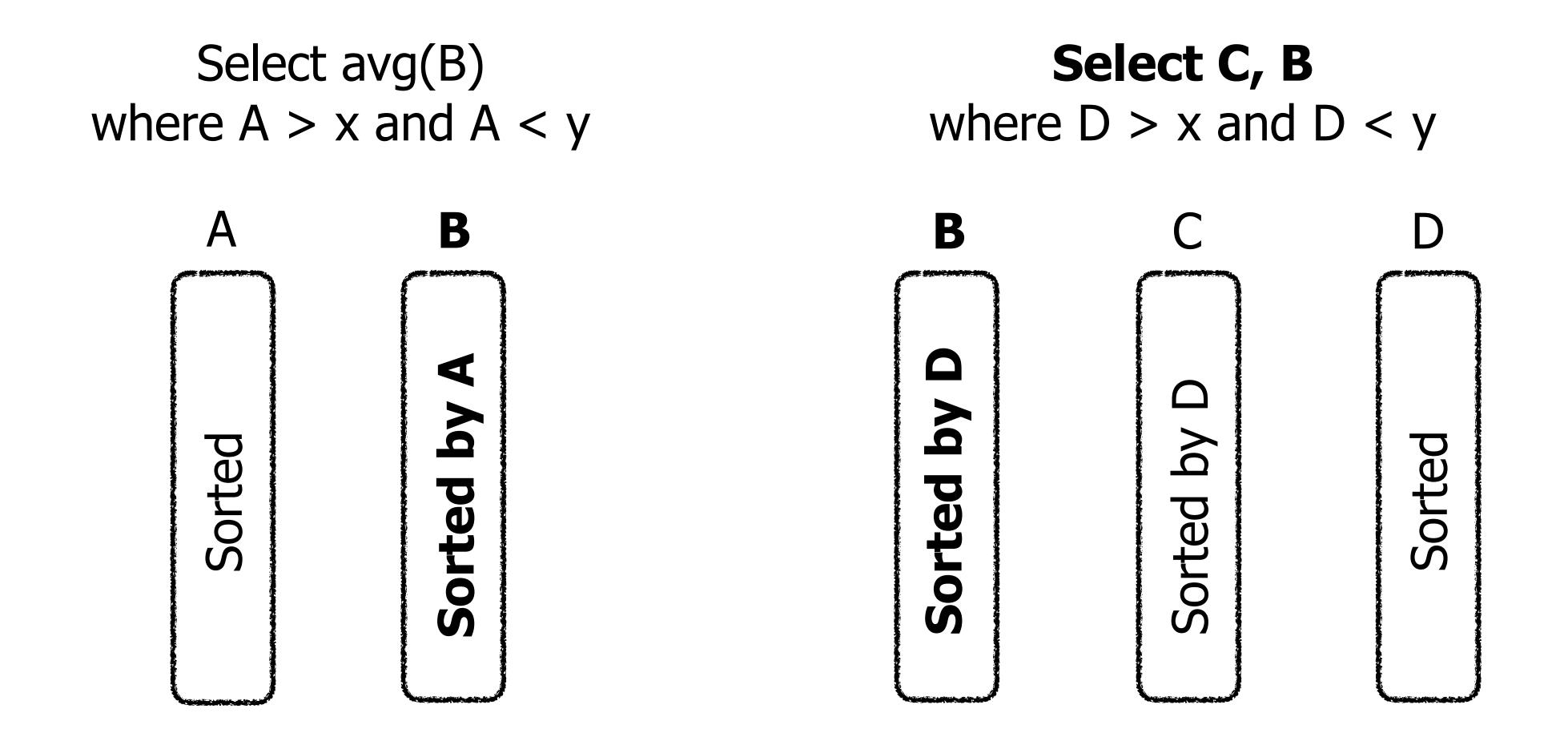


What if the queries do not target mutually exclusive columns?

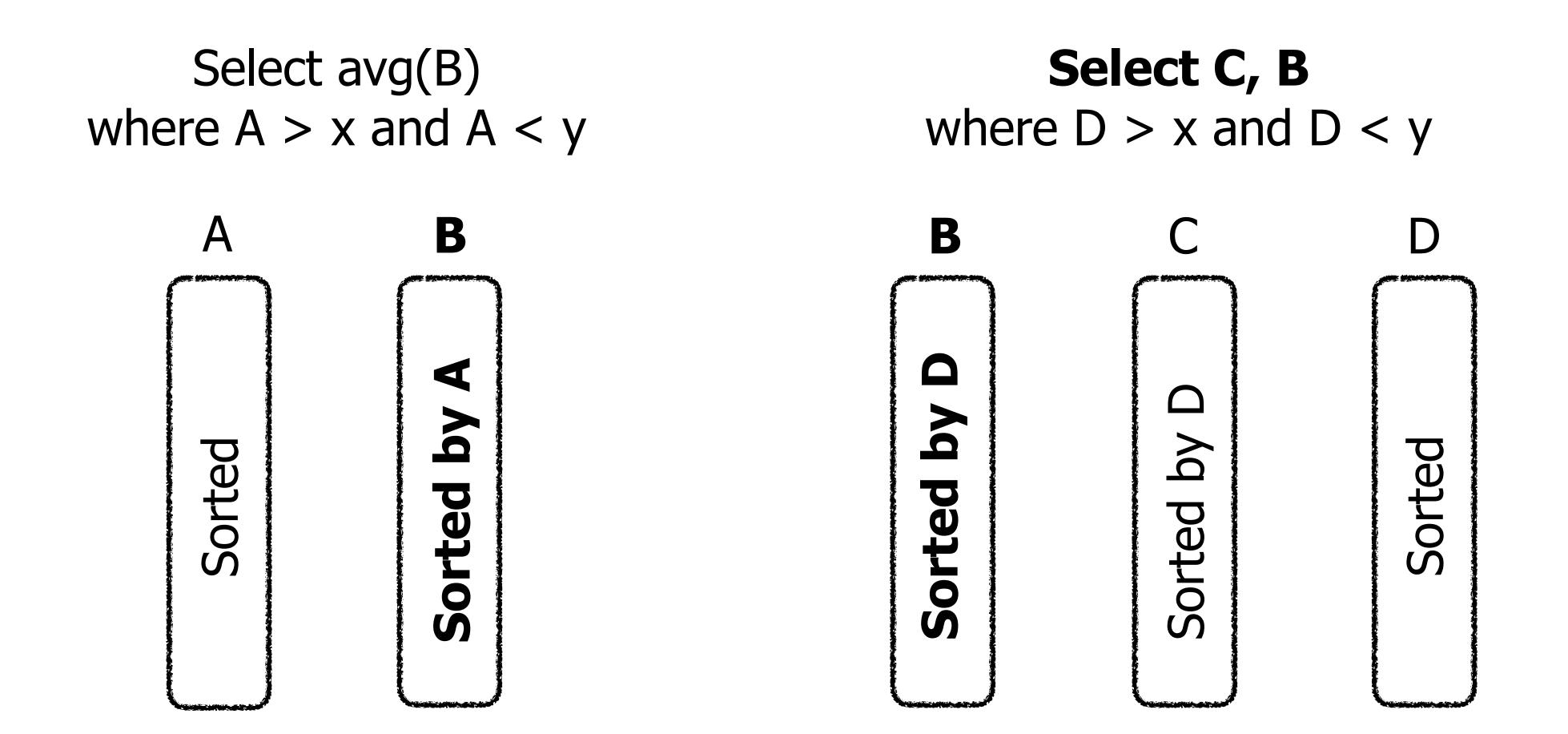
Select C, B where D > x and D < y



What if the queries do not target mutually exclusive columns?

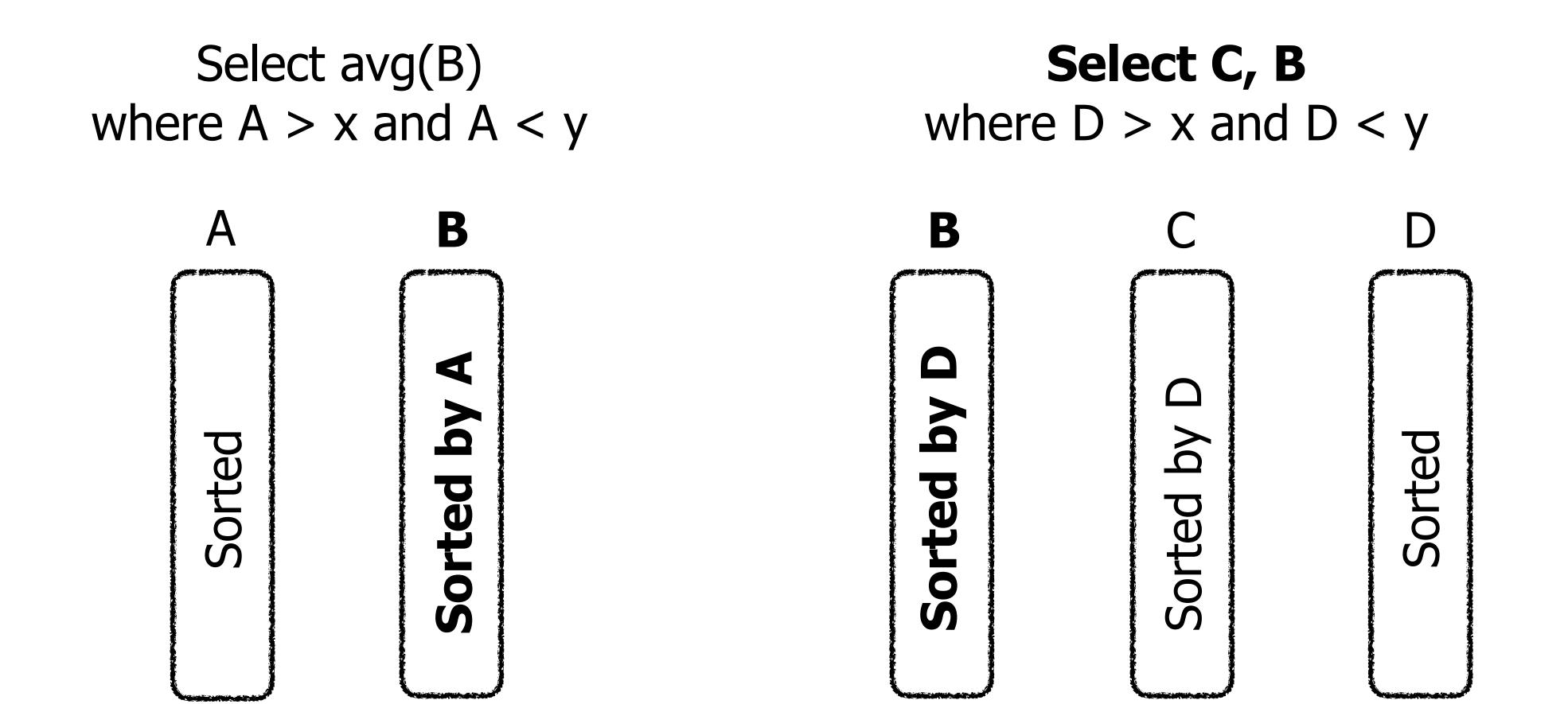


Or we can duplicate some columns and sort them in different orders

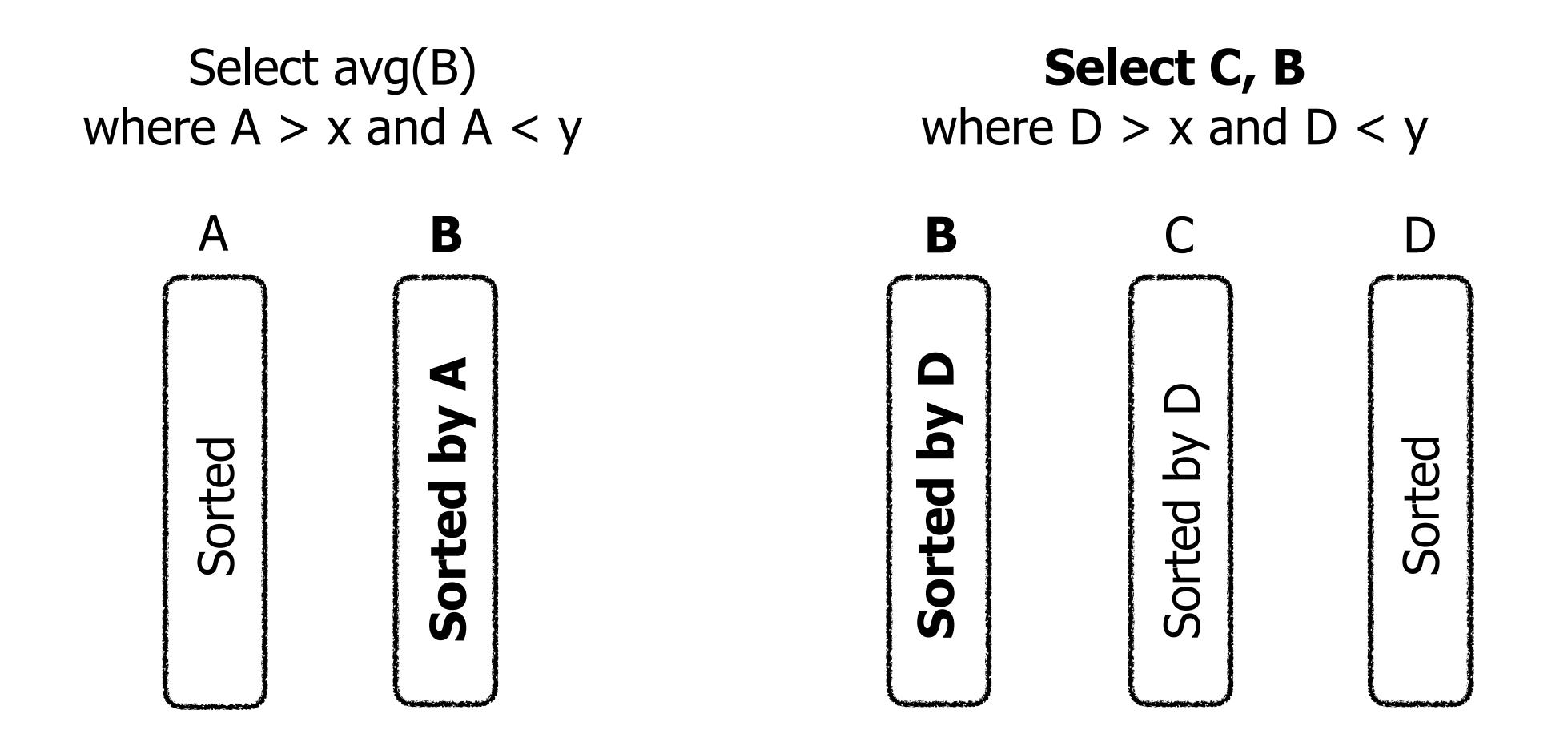


Or we can duplicate some columns and sort them in different orders

#### **Downsides?**



- Downsides: (1) More space
  - (2) Inserts are more expensive
  - (3) Construction can take a long time with many projections

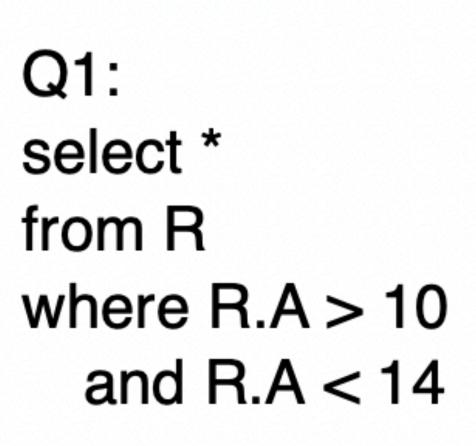


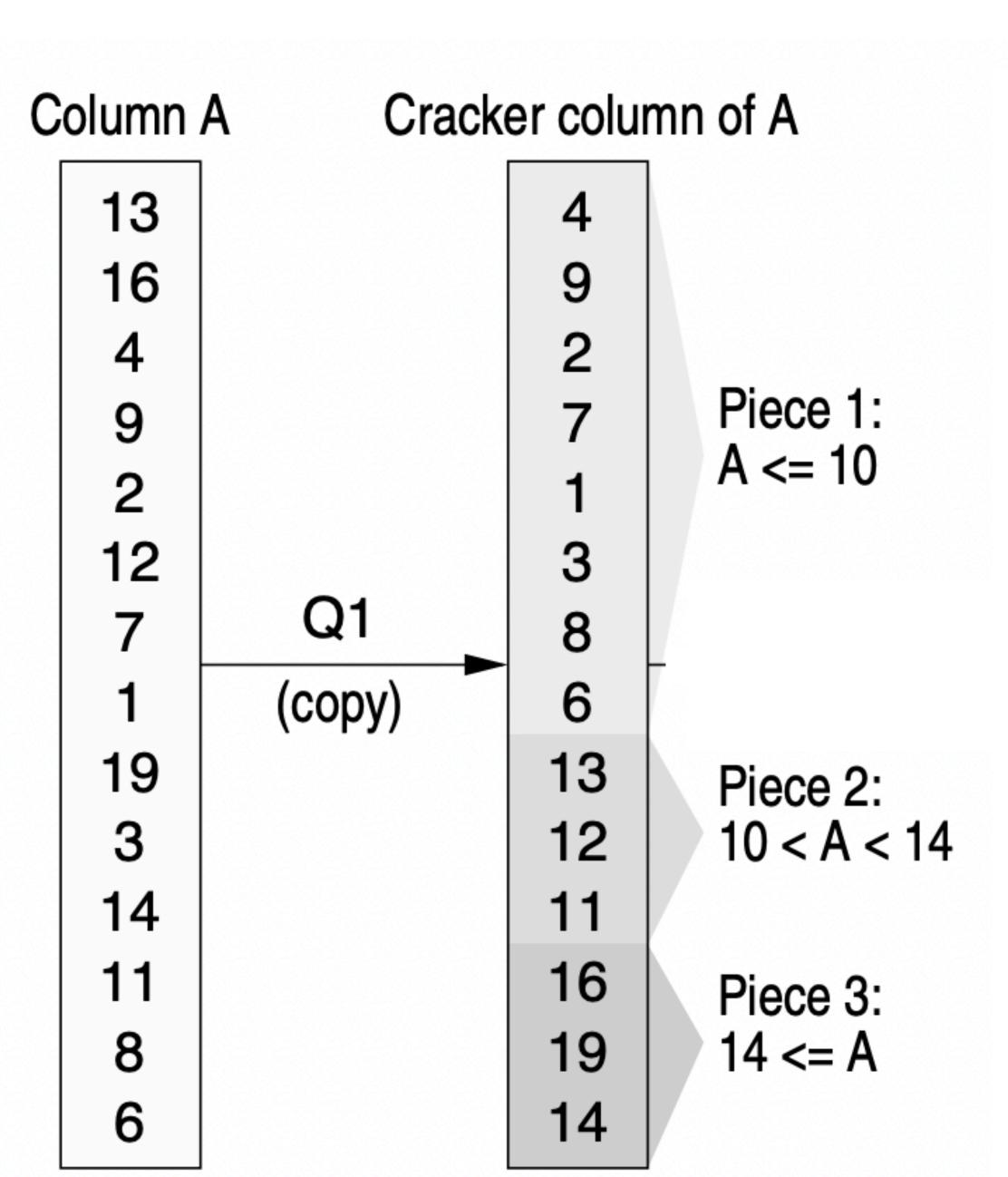
- (1) More space
- (2) Inserts are more expensive
- Tackle this: —— (3) Construction can take a long time with many projections

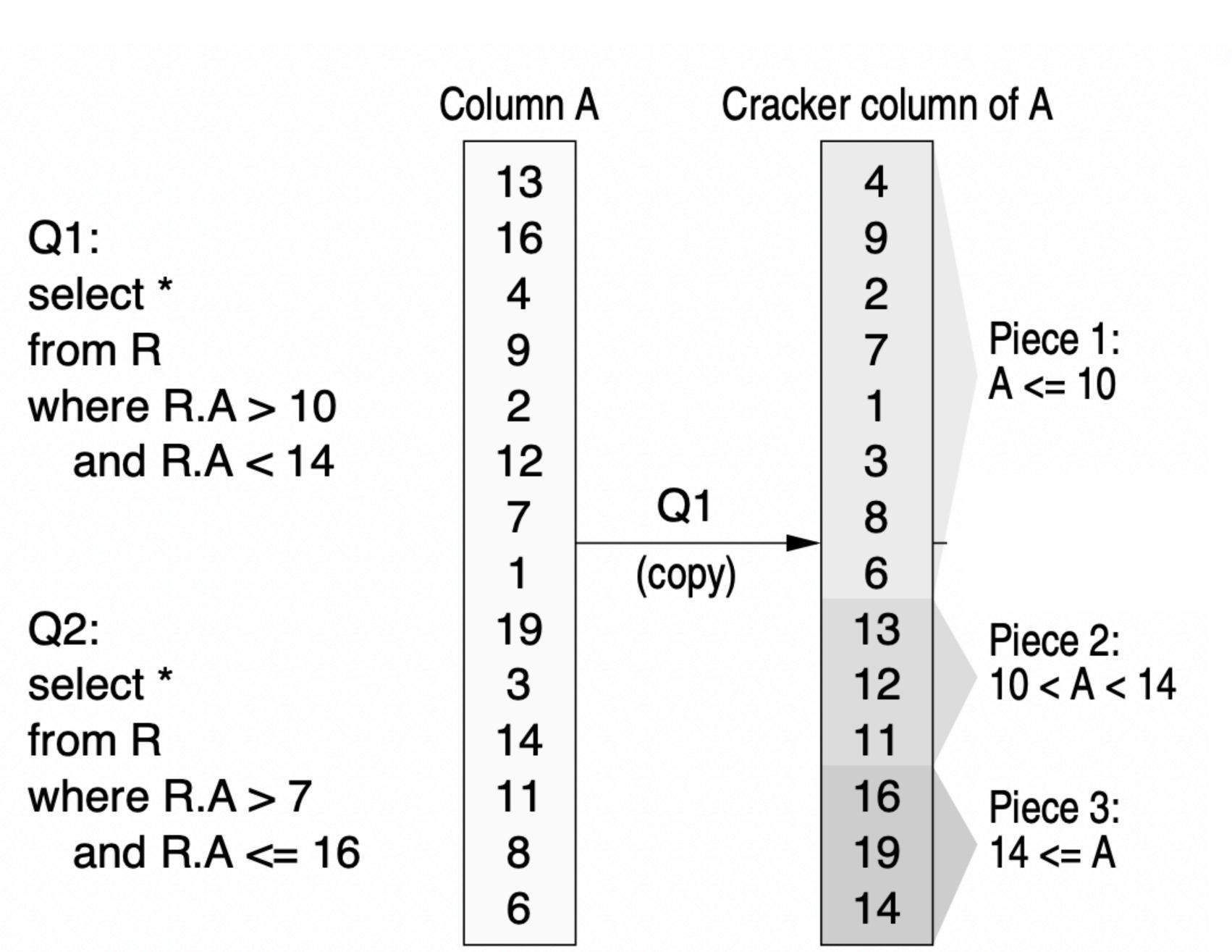
#### Column A

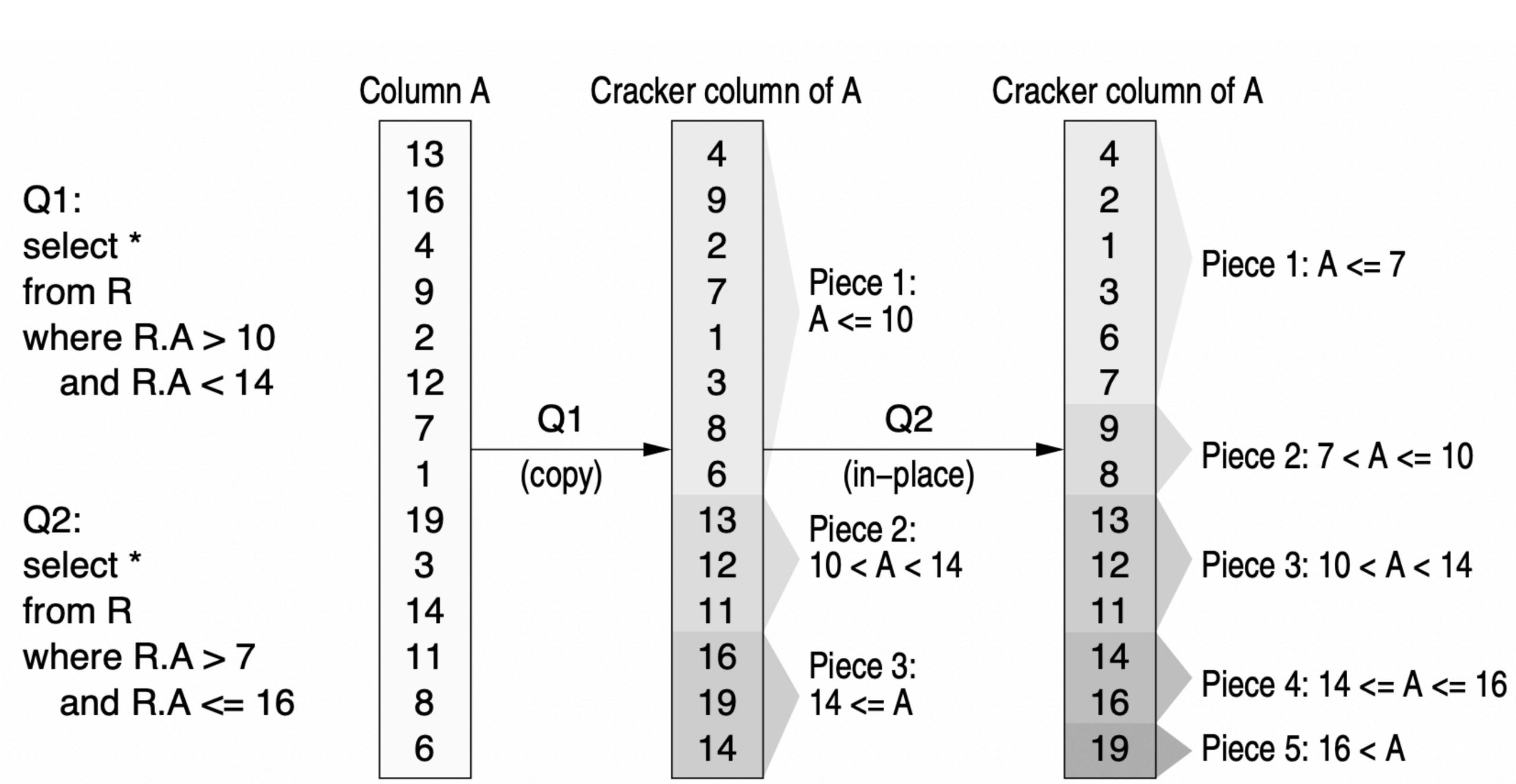
#### Column A

Q1:
select \*
from R
where R.A > 10
and R.A < 14

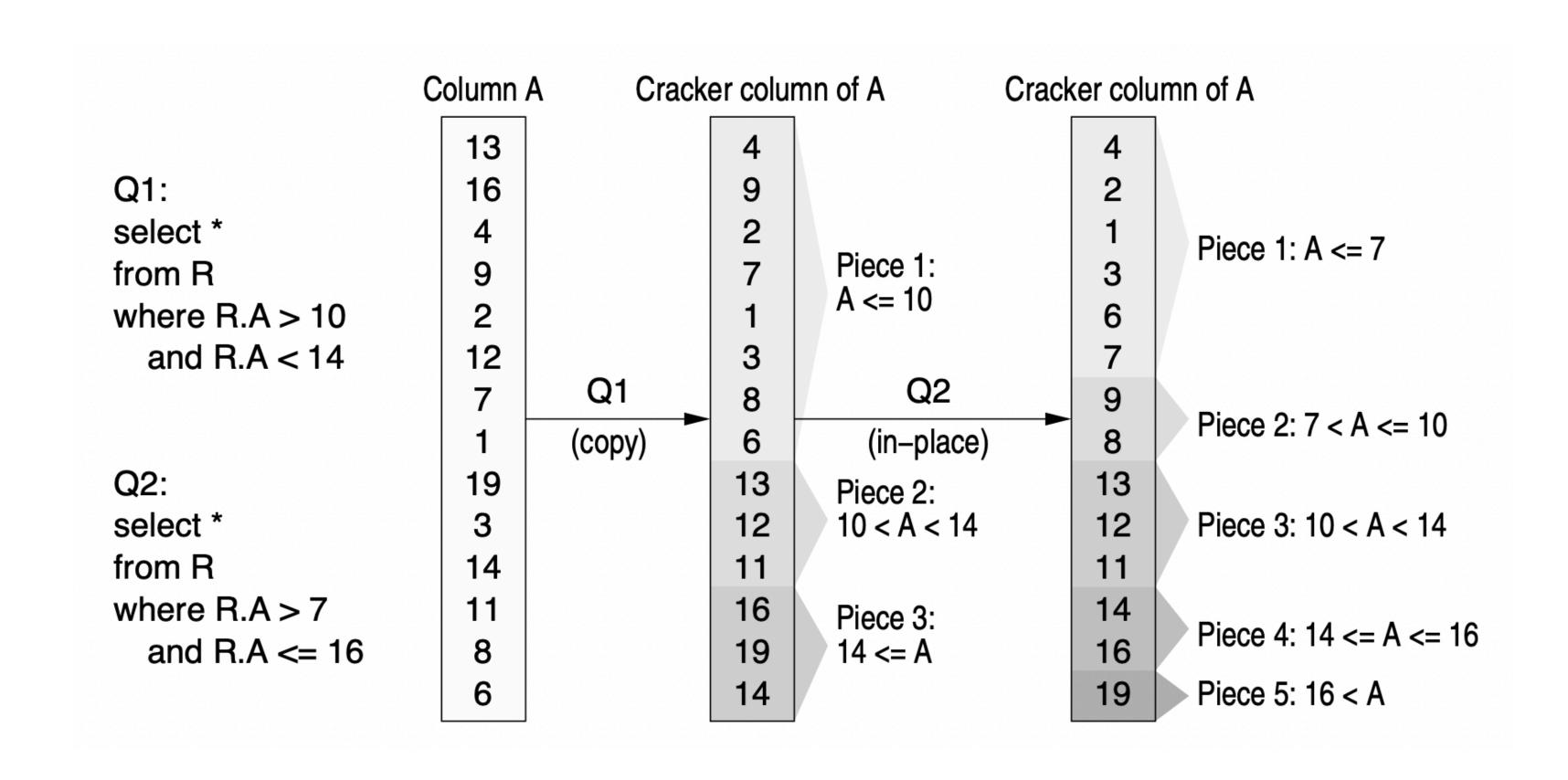






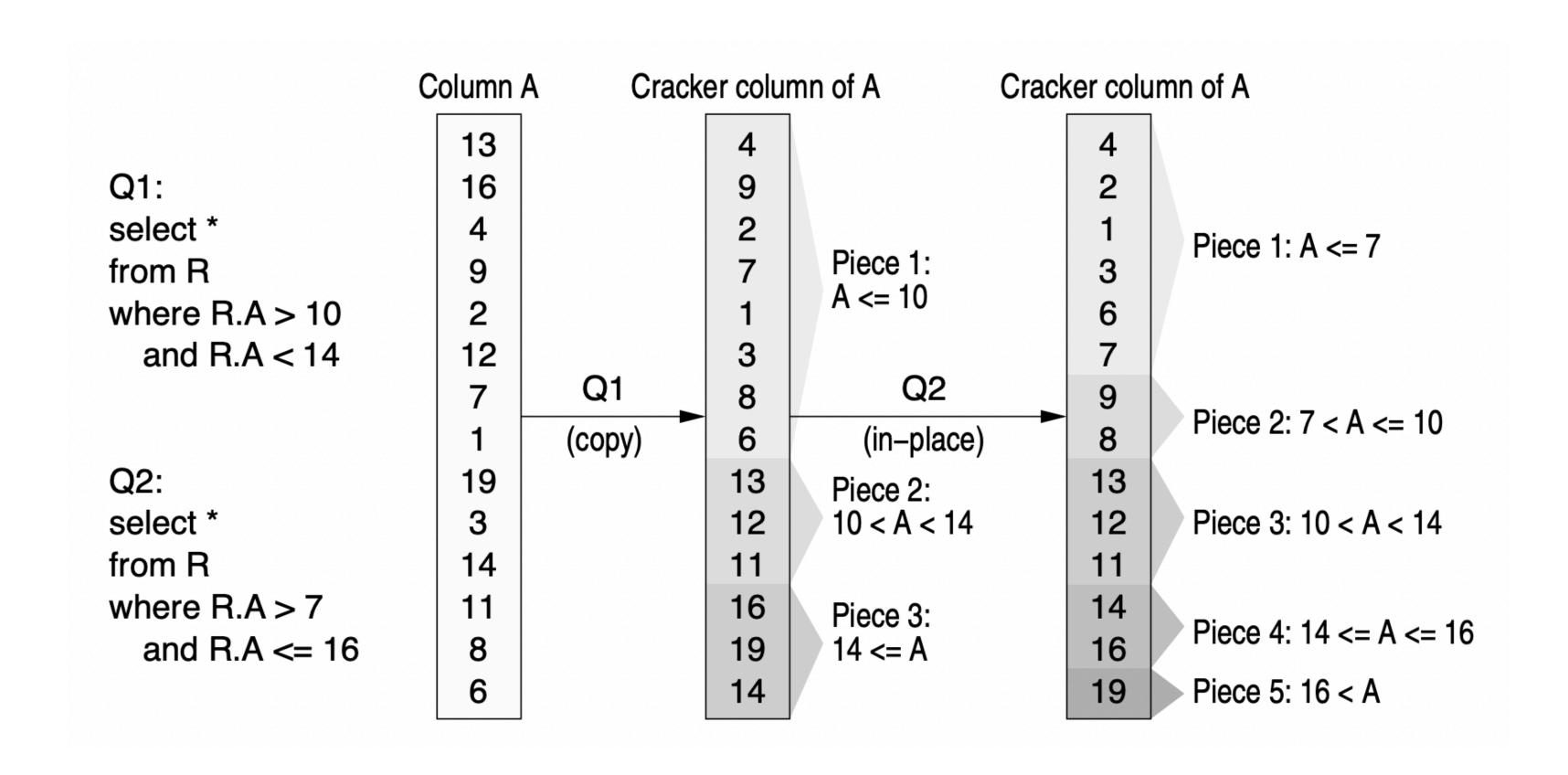


### Queries adaptively refine the sorting and speed up subsequent queries



Queries adaptively refine the sorting and speed up subsequent queries

# We can begin querying the data immediately without creating too many projections from the onset, which would take a while



And now: TA office hours