# BIG DATA MANAGEMENT

CZ/CE4123

# Tutorial 5:
# Column Stores

Given column store table T as follow.

(1) Give the flow chart (the flow graph presented in the lecture slides) using "column at a time" for the query
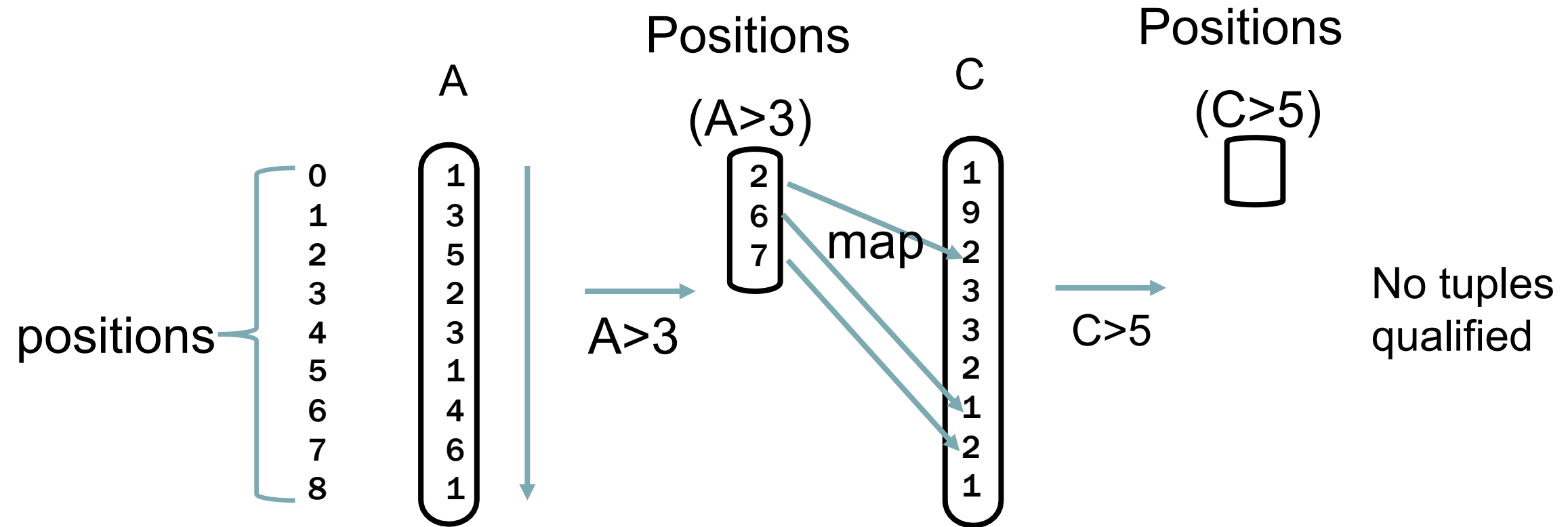
"**SELECT** *min(B)* **FROM** *T* **WHERE** *A>3 and C>5*"

(2) Give the flow chart using "column at a time" for the query

"**SELECT** *sum(C)* **FROM** *T* **WHERE** *A<10 and B>3*"

| A | B | C |
|---|---|---|
| 1 | 4 | 1 |
| 3 | 3 | 9 |
| 5 | 5 | 2 |
| 2 | 5 | 3 |
| 3 | 7 | 3 |
| 1 | 1 | 2 |
| 4 | 2 | 1 |
| 6 | 7 | 2 |
| 1 | 1 | 1 |

# QUESTION 2

Redo Question 1 using "vector at a time". Assume that the vector size is 3.

Positions

A

C

Positions

(A>3)

(C>5)

positions
0
1
2

1
3
5

A>3

2

1
9
2

map

C>5

Batch1

No tuples
qualified

positions
3
4
5

2
3
1

A>3

3
3
2

Batch2

positions
6
7
8

4
6
1

A>3

6
7

map

1
2
1

C>5

Batch3

Suppose we are querying a **Student** information table with three columns **Name**, **Email**, **Age**. Given a query of the following form:

"**SELECT** *Name* **FROM** *Student* **WHERE** *predicate (Email) and predicate (Age)*".

Here, a predicate applied on a column is a filtering function (e.g., **Email** ending with *ntu.edu.sg,* **Age**>19). We define the selectivity of a predicate by the percentage of the qualified results in the corresponding column. Assume that the selectivity of predicate(**Email**) is $p$ and the selectivity of predicate(**Age**) is $q$, where $0<p<1$, and $0<q<1$. Let page size be $P$. We assume each column width is less than $P$ and each value in a column is contained in a page. Consider two options in scanning columns: scanning **Email** first and scanning **Age** first.
(1) If the column widths are the same (denoted by $w$), please analyze which is better.
(2) If the width of **Email** is $2w$, and the widths for **Name** and **Age** are $w$, then which option is better?

Cost for column store (number of page access):

$Zw/P+2result(A)*4/P+result(A)+2result(AB)*4/P+result(AB)$

$=Zw/P+result(A)*(8/P+1)+result(AB)*(8/P+1)$

$\approx Zw/P+result(A)+result(AB)$

# SOLUTION TO (1)

**Answer**:

Use the formula we learnt in the lectures.

Let Email be *A*, and Age be *B*. Let *Z* be the length of a column

The cost of option 1 (Scanning Email first) is approximately

**Zw/P+result(A)+result(AB)**

The cost of option 2 (Scanning Age first) is approximately

**Zw/P+result(B)+result(AB)**

# SOLUTION TO (1)

**Answer**:

Use the formula we learnt in the lectures.

Let Email be *A*, and Age be *B*. Let *Z* be the length of a column

The cost of option 1 (Scanning Email first) is approximately

**Zw/P+result(A)+result(AB)**

The cost of option 2 (Scanning Age first) is approximately

**Zw/P+result(B)+result(AB)**

Note that **result(A)=Zp,  result(B)=Zq**

Then, if p<q, then scanning Email first is better; if p=q, equally good; if p>q, then scanning Age first is better.

Recap the formula we learnt in the lectures.

**Zw/P+result(A)+result(AB)**

Width of 1st accessed column

# SOLUTION TO (2)

**Answer**:

Let Email be Column *A*, and Age be Column *B*. Let width of Email be 2*w*. Then the width of Age is *w*.

Revise the formula we learnt in the lectures.

The cost of option 1 (Scanning Email first) is approximately

**2Zw/P+result(A)+result(AB)**

The cost of option 2 (Scanning Age first) is approximately

**Zw/P+result(B)+result(AB)**

Note that **result(A)=Zp,  result(B)=Zq**

**Answer**:

Then, option 1 is worse (or option 2 is better) when

$2Zw/P+Zp+\text{result}(AB) > Zw/P+Zq+\text{result}(AB)$

$\rightarrow Zw/P+Zp > Zq$

$\rightarrow w/P+p>q$

(Note: it is also okay to use more refined formulas, i.e., considering the cost of positions.)

# In practice, how do we know p and q?