



COMP90050 ADBS

---

# COMP90050 ADBS

授课老师: Rita





# 1 – Query Processing & Optimisation





## 基本概念

- 8 continuous bits = 1 byte
- 4000/ 8000 continuous bytes = 1 block
- Bit – b; byte – B; block – A





record < block < file

## How data is stored in disk?

record  $\Rightarrow$  行

→ 多个连续 record.

o Files – A database is mapped into different files. A file is a sequence of records.

o Data blocks  $\rightarrow$  多行(连续) – Each file is mapped into fixed length storage units, called data blocks (also called logical blocks, or pages)

e.g., size of each record: 55 byte; fixed size of 1 data block: 4096 byte

$$\frac{4096}{55} \approx 74.47 \rightarrow 74 \text{ 行} \text{ 记录} \quad \text{一条数据不能被拆成两个 blocks 中}$$

record





# COMP90050 ADBS

Table – “Employee”

ID	Name	Age
1001	A	25
1002	B	32
1003	C	19
1004	D	27
1005	E	40
1006	F	36

- Attribute (column) *31*
- Record (row)
- Query
  - Please tell me the name of all employees who is over 25 *filter*
  - please tell me the name of all managers *join on ID →*
  - Please tell me the name of all managers located in Australia
- Join

Table – “Manager”

ID	Department
1001	a
1003	b
1004	c
1006	d

Table – “Location”

ID	Country	City
1001	China	Beijing
1002	Australia	Melbourne
1004	China	Shanghai
1006	Australia	Sydney





```
Select Name  
from (Employee inner join Manager)  
On Employee.ID = Manager.ID;
```

```
Select Name  
from (A inner join location)  
On A.ID = Location.ID  
Where Location.country == 'Australia';
```





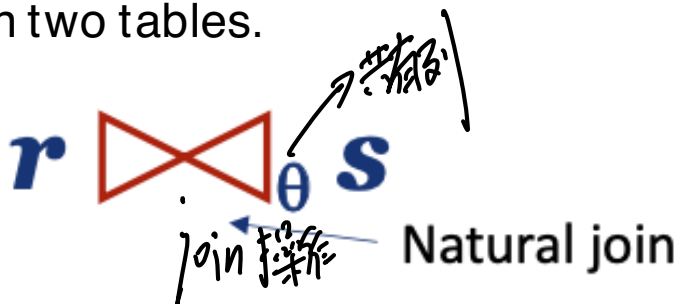
## Join

速度慢  $\rightarrow$  中间生成物  
空间占用大

Join is a very common and very expensive operation

Different types of join: inner join, outer join... (we focus on inner join)

**Natural join:** a join operation that can be performed of a column that is common in two tables.



r: outer relation

s: inner relation

r and s are two tables

Theta is the common column.





## Relational Algebra Expressions

*Select \* from T1  
inner join T2  
on T1.a = T2.b*

*all attributes*

```
SELECT *  
FROM Employees  
INNER JOIN Managers  
ON Employees.ID = Managers.ID;
```

.... Is same as the following in relational algebra expression:

$\Pi_{*}(\sigma_{\text{Employees.ID=Managers.ID}}(\text{Employees} \text{ join Managers}))$

*all attribute*      *common attribute*







## Relational Algebra Expressions

*select* A1, A2, ..., An  
*from* r1, r2, ..., rm  
*where* P

**Select** salary  
**From** Employees  
**Where** salary < 60000

.... Is same as the following in relational algebra expression:

由里向外执行

$$\Pi_{A_1, A_2, \dots, A_n} (\sigma_P (r_1 \times r_2 \times \dots \times r_m))$$

**Equivalent  
Expressions**

Can also be  
written as:

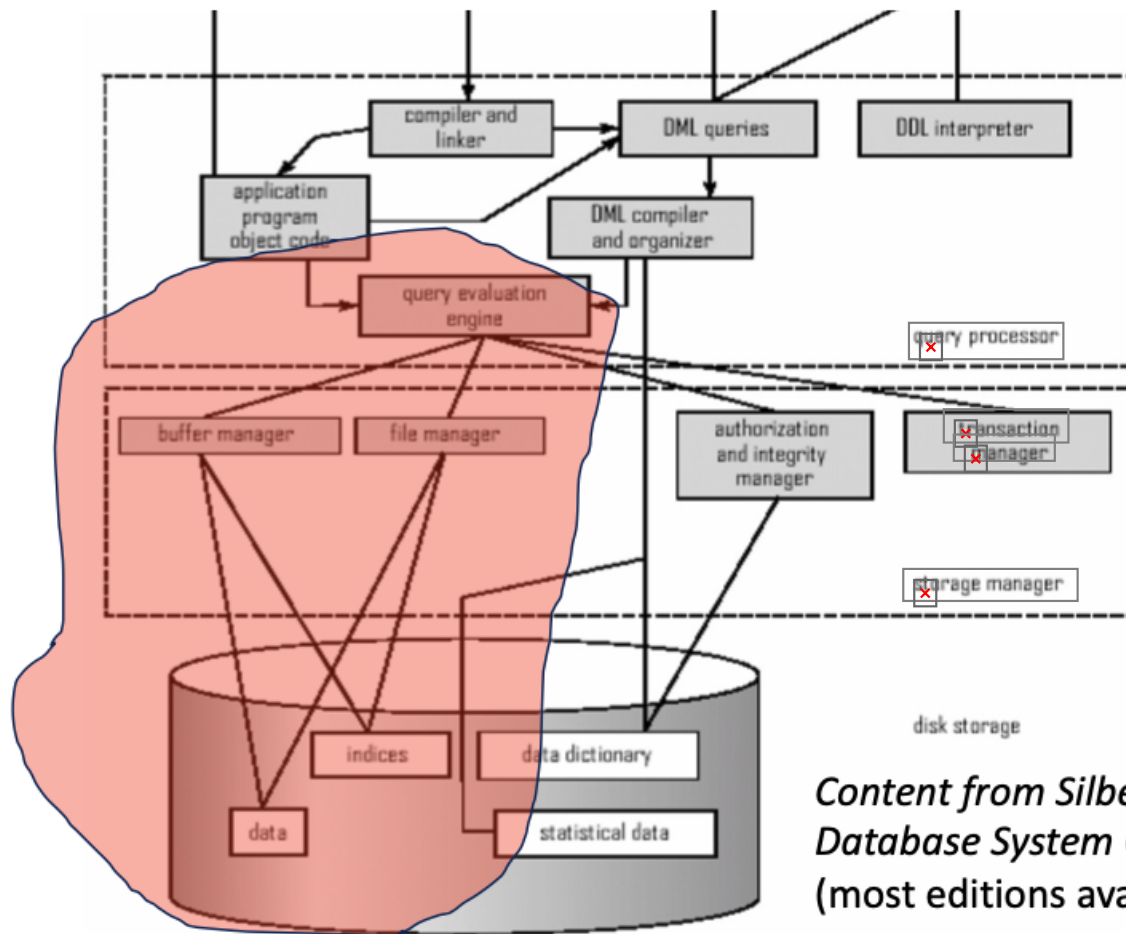
$$\left\{ \begin{array}{l} \Pi_{salary} (\sigma_{salary < 60000} (Employees)) \\ \sigma_{salary < 60000} (\Pi_{salary} (Employees)) \end{array} \right.$$

①: 筛 salary < 60000 ②: 取 salary 值  
 ①: 取 salary = 7 ②: 筛 salary < 60000





## DBEngine



*Content from Silberschatz et al  
Database System Concepts  
(most editions available online)*





- **Hardware**

- – The speed of the processor
- – Number of processors
- – Number of disk drives and I/ O bandwidth
- – Size of main memory
- – Communication network
- – Type of architecture

DB Engine

- **Software**

- – Type of database technology used for a given application

- **Database tuning, crash recovery**

- – Indexing parameters
- – Data duplication
- – Sharing data, etc

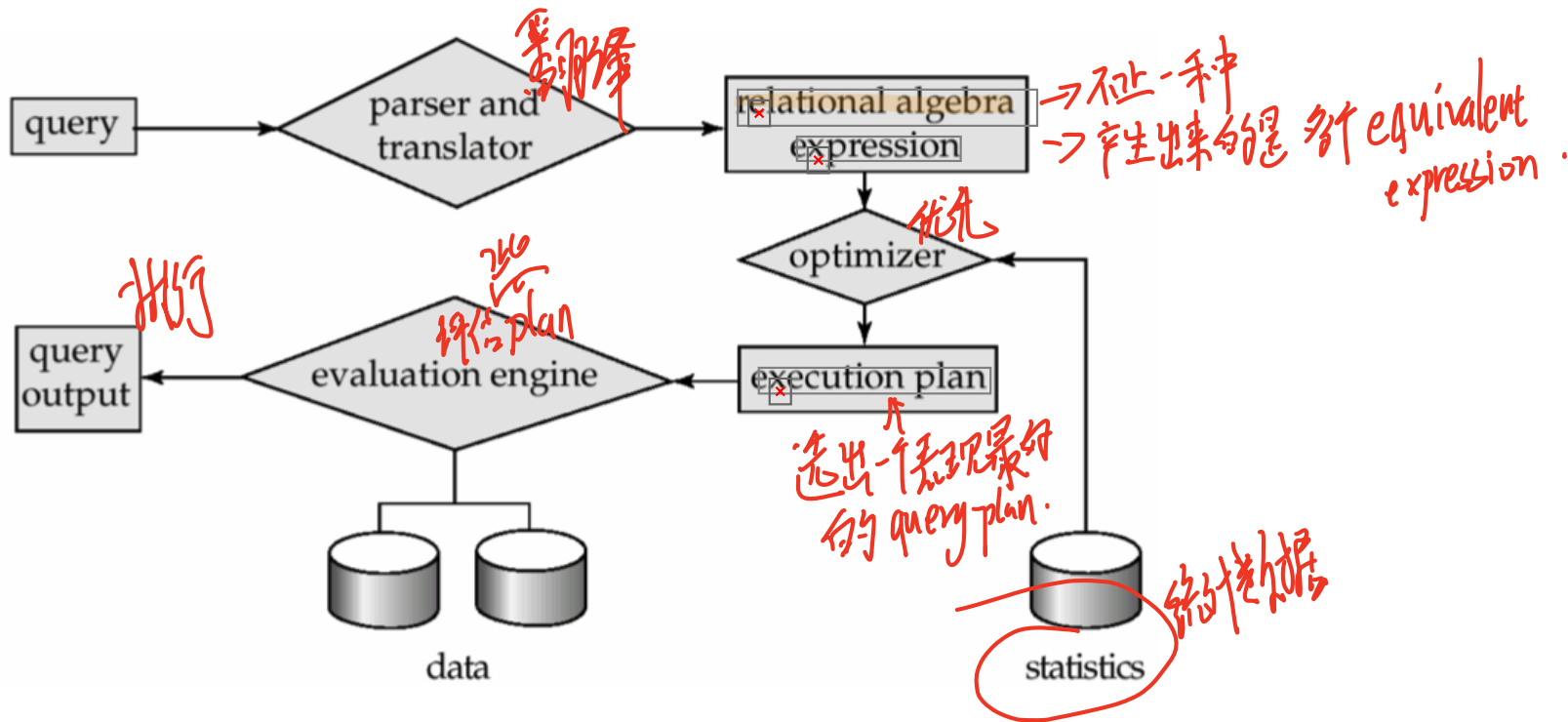




\* 一个 query 会有多个 plan 产生  
\* 优化后会选取一个最优的 plan.

## Query Processing Steps

E.g., select Salary From Employees Where Salary < 60000



- Translate query to relational algebra expression
- Make execution plan

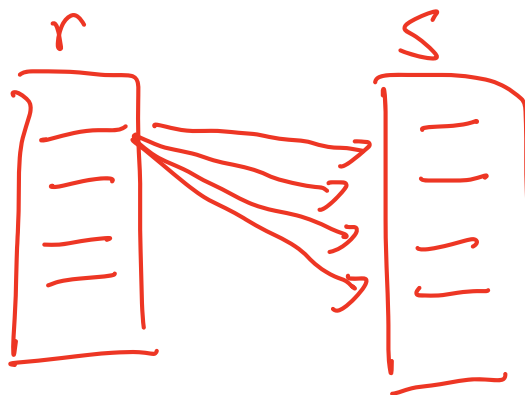




## Nested-loop Join (不用index)

for each tuple  $t_r$  in  $r$  do begin  
  for each tuple  $t_s$  in  $s$  do begin  
    test pair  $(t_r, t_s)$  to see if they satisfy the join condition theta ( $\theta$ )  
    if they do, add  $t_r \cdot t_s$  to the result.  
  end  
end

$r$   $\bowtie$   $s$   
outer  $\theta$  inner  
condition



- ? Requires no indices because it checks everything in  $r$  against everything in  $s$ .
- ? Expensive since it examines every pair of tuples in the two relations.
- ? Could be cheap if you do it on two small tables where they fit to main memory (disk brings the whole tables with first block access).

两小表容易

→ 目的(方便快速查找)

record 每个值





# COMP90050 ADBS

definition of cost: how many blocks are read from disk?

In the worst case, if there is enough memory only to hold one block of each table, the estimated cost is

$n$ : record

$b$ : block

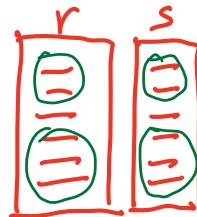
access time:  $n_r * b_s + b_r$  block transfers, and  
 $n_r + b_r$  seeks (查找次数)

→ 只有外层表块数, 内层表加载好了不用 access time: <sup>块数</sup>





*fk nested loop join 简单*



*1. 2 blocks 互相比较匹配*

## Block Nested-loop Join (Page-Oriented Nested-loop Join)

```

for each block  $B_r$  of  $r$  do begin
  for each block  $B_s$  of  $s$  do begin
    for each tuple  $t_r$  in  $B_r$  do begin
      for each tuple  $t_s$  in  $B_s$  do begin
        Check if  $(t_r, t_s)$  satisfy the join condition
        if they do, add  $t_r \bullet t_s$  to the result.
      end
    end
  end
end

```

? Variant of nested-loop join in which every block of inner relation is paired with every block of outer relation.





把 nested loop join  $n_r$  换进  $b_r$  即可

$b_r * b_s + b_r$  block transfers +  $2 * b_r$  seeks

$n_r \rightarrow b_r$







e.g.,

- Number of records of *customer*: 10,000    *depositor*: 5000
- Number of blocks of *customer*: 400    *depositor*: 100





1 seek time = 12ms nested loop join

## Nested-loop Join

不同 outer 结果不同

Depositor as the outer relation

$R \bowtie S$

$R = \text{depositor}$

$$\text{access time} : n_r \times b_s + b_r = 5000 \times 400 + 1000 = 2001000$$

$$\text{seek time} : n_r + b_r = 5000 + 100 = 5100$$

$$5100 \times 12 \approx$$

Customer as the outer relation

$R = \text{Customer}$

$$n_r \times b_s + b_r = 10000 \times 100 + 400 = 1000000 + 400 = 1000400$$

$$n_r + b_r = 10000 + 400 = 10400$$

$$10400 \times 12 \approx$$





block transfer:  $b_r \times b_s + b_r$

access time:  $2 \times b_r$

## Block Nested-loop Join (Page-Oriented Nested-loop Join)

Depositor as the outer relation

$\left\{ \begin{array}{l} M \\ R \end{array} \right\}$

$$b_{\text{depositor}} \times b_{\text{customer}} + b_{\text{depositor}}$$

$$= 100 \times 400 + 100$$

Customer as the outer relation

seek:  $2 \times 100$

$$400 \times 100 + 400$$

$$2 \times 400$$

Customer 作为 outer relation 时  
查询更快 (谁少谁更快)





## Query Optimization

Query optimization is about the right choices and annotations

What to optimize? join 的顺序优化 相对量的变化

? e.g., I have 3 tables, which two can I join first?

? Join algorithms





## Query Plans

Alternatives 备选

Execution Plan (最终执行方案)

Final decision



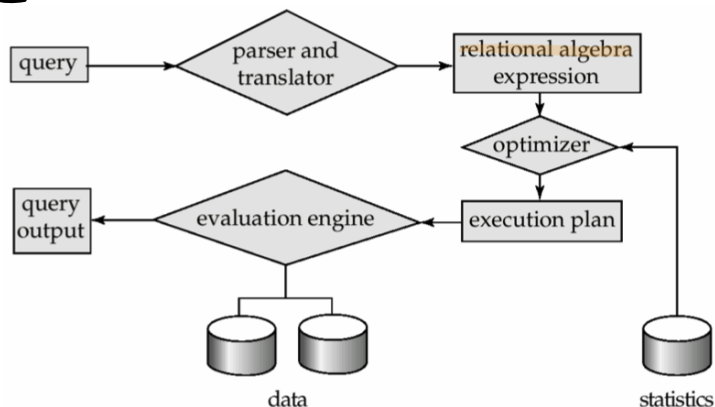


## How do we make the choices?

steps in cost-based query optimization

Cost: how many blocks are read from disk

1. Generate logically equivalent expressions of the query
2. Annotate resultant expressions to get alternative query plans
3. Choose the cheapest plan based on estimated cost





选择 plan 依据

Estimation of plan cost based on:

- ① Statistical information about tables. Example:  $\Rightarrow$   
number of distinct values for an attribute

- Statistics estimation for intermediate results to compute cost of complex expressions

$A+B+C$  (AB) 中间结果

- Cost formulae for algorithms, computed using statistics again

算法 cost

ID	Name	Age
1	A	18
2	B	21
3	C	19
4	D	20, 18, 21, 20

distinct.





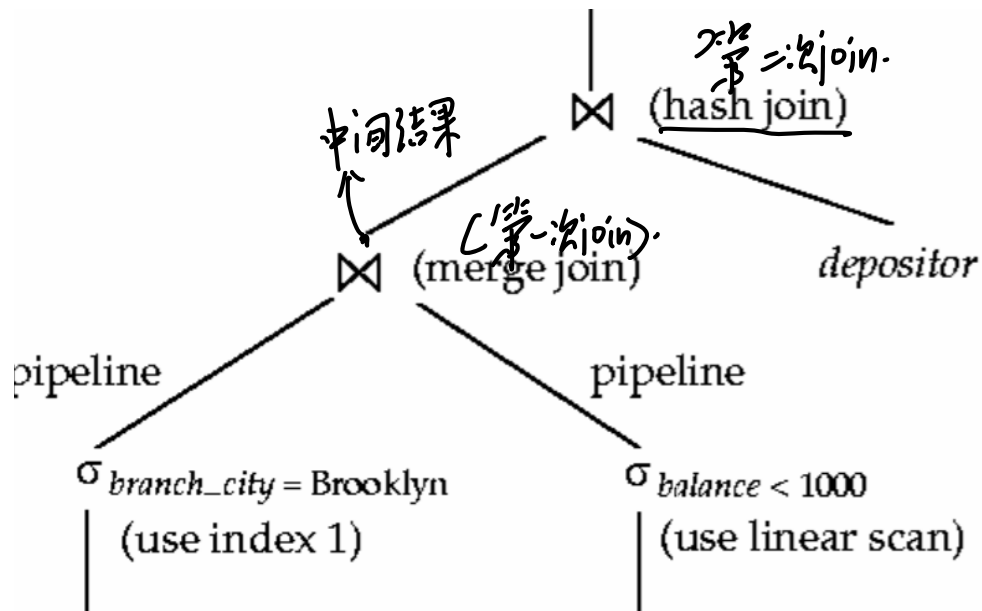
## 2 – Query Optimization in Real Life





**Query Optimization** is about the right choice on a graph

(不仅在上)





## How to generate alternatives?

- Query optimizers use equivalence rules to systematically generate expressions equivalent to the given expression
- One can generate all equivalent expressions exhaustively (所有 plan 都出来)
- The above approach is very **expensive** in space and time though ,  
(In query optimizer, some expressions are not generated if they are for sure very complex)  
若有预见可能





But

- Must consider the interaction of evaluation techniques when choosing evaluation plans
- Choosing the cheapest algorithm for each operation independently may not yield best overall algorithm (每步最好的, 不一定overall最好)

e.g., merge-join may be costlier than hash-join, but may provide a sorted output which could be useful later (the sorted result may be benefit to the later operation)

局部最优 ≠ 全局最优





## In Real Life

- Practical query optimizers incorporate elements of the following two broad approaches:
  1. Search all the plans and choose the best plan in a cost- based fashion.
  2. Uses **heuristics** to choose a plan. <sup>history</sup>  
<sub>history</sub> <sub>history plan</sub> 有的plan很差 → 减少选择 [总结下来直接跳过以]
- Systems may use heuristics to **reduce the number of choices** that must be made in a cost- based fashion (because cost- based optimization is expensive)





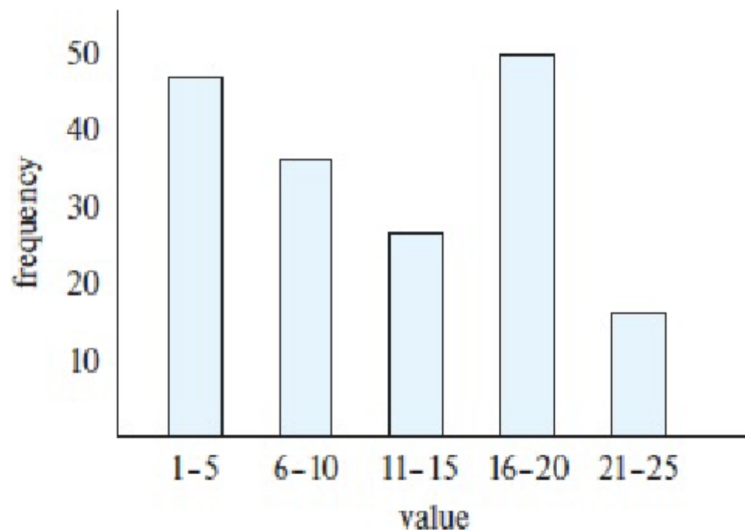
- Heuristic optimization transforms the query- tree by using a set of rules that typically (but not in all cases) improve execution performance:
  1. **Perform selections early** (reduces the number of tuples) *select A from I where C*
  2. **Perform projections early** (reduces the number of attributes) *列少的更好*
  3. **Perform most restrictive selection and join operations** (i.e. with smallest result size) before other similar operations *(中间表 size 少的好)*
- Some systems use only heuristics, others combine heuristics with cost- based optimization
- Optimizers often use simple heuristics for very cheap queries, and perform exhaustive enumeration for more expensive queries

对 cheap query: 用 simple heuristic plan  
对 贵的 query: exhaustive enumeration.





- Further Optimisations
  - Sampling
  - Histogram



- Adaptive Plans - Wait for one/ some parts of a plan to execute first, then choose the next best alternative  
先看看到底如何  $\Rightarrow$  好的话再继续下一个阶段  
过程中运行





## 3 – Query Cost in Practice





## Troubleshooting to manage costs (e.g., Query Store in Microsoft)

监控 query

(Query store: SQL server management studio for monitoring)

识别变差的 query 表现

- Identify 'regressed queries' - Pinpoint the queries for which execution metrics have recently regressed (for example, changed to worse).
- Track specific queries - Track the execution of the most important queries in real time. (e.g., most frequently asked queries)

跟踪重要的 query 的执行







When you identify a query with suboptimal performance  
(没有达到最好的performance)

- Force a query plan instead of the plan chosen by the optimizer  
*对表现差的 → 找一个表现好的强制* force execution
- Do we need an **index**? - - - quickly find the data in the query
- Enforce statistic recompilation (重新编译)
- Rewrite query? (with parameters)

*next page*





## Parameter in query

给一类 query plan 设定为同一个参数。

Query rewriting with parameters for execution plan reuse

```
SELECT *  
FROM Product  
WHERE categoryID = 1;
```

```
SELECT *  
FROM Product  
WHERE categoryID = 4;
```

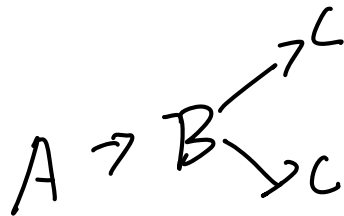
They use the same plan!

We expect the optimizer to generate essentially the same plan and reuse the plans - **parameterize**

```
DECLARE @MyIntParm INT  
SET @MyIntParm = 1  
EXEC sp_executesql  
    N'SELECT -  
    FROM Product  
    WHERE categoryID = @Parm',  
    N'@Parm INT',  
    @MyIntParm
```

↓  
"1" & "4"





To further lower query cost

## • Store derived data

- When you frequently need derived values
- Original data do not change frequently

## • Use pre-joined tables

- When tables need to be joined frequently
- Regularly check and update pre-joined table for updates in the original table
- May still return some 'outdated' result (pre-joined tables are not updated)

original table 不能 update or change frequently  $\Rightarrow$  中间表  $\left\{ \begin{array}{l} \text{pre-joined table} \\ \text{derived table} \end{array} \right.$  会频繁更新  
 $\rightarrow$  开销更大

$\rightarrow$  存储经常使用的中间值 要保证 A 不能频繁更新  
 频繁更新  $\Rightarrow$  开销更大

