

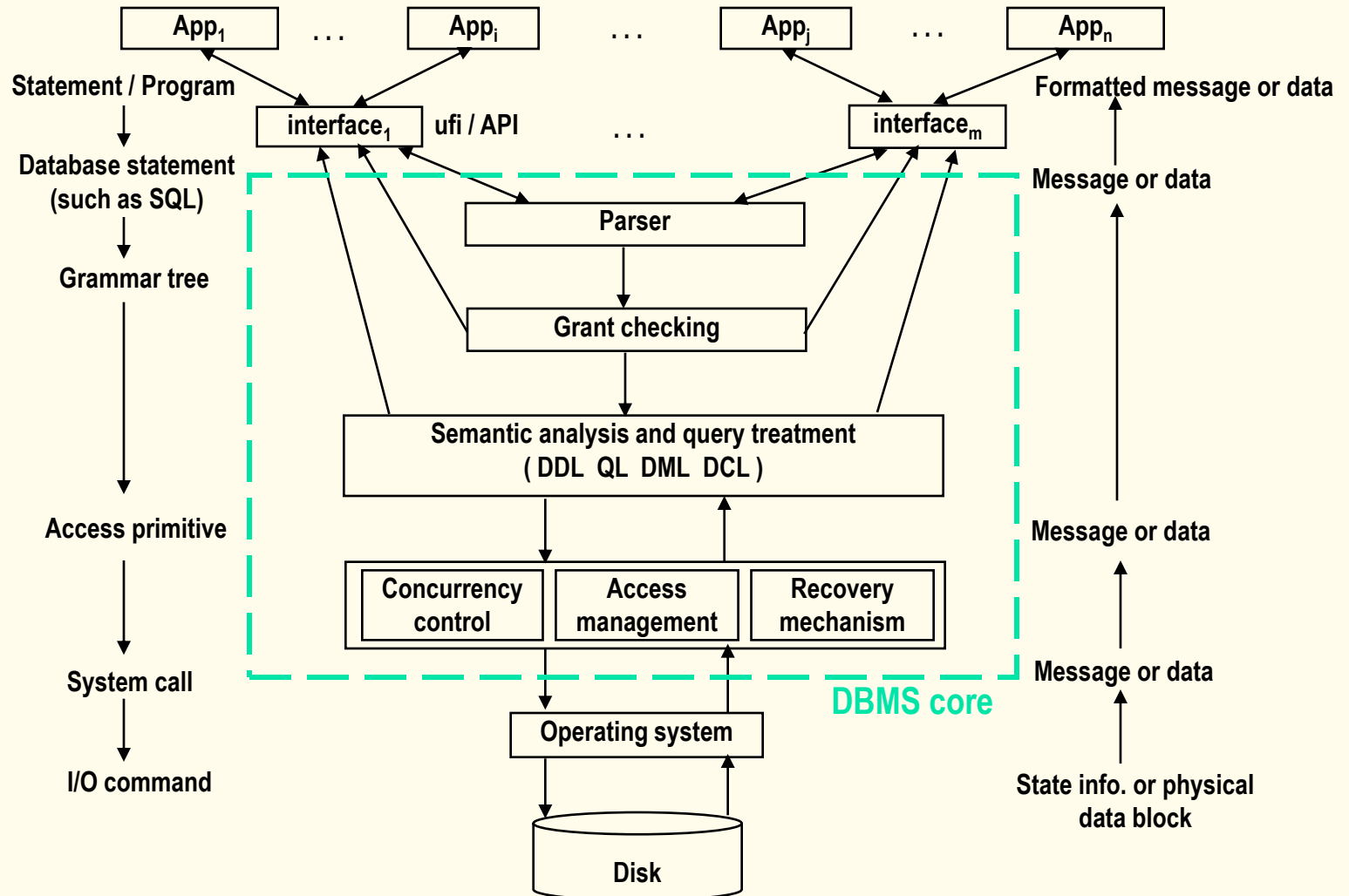
4. Database Management Systems*



Contents

- The Architecture of DBMS
 - The components of DBMS core
 - The process structure of DBMS
- Database Access Management
- Query Optimization
- Transaction Management
 - Recovery
 - Concurrent Control

4.1 The Components of DBMS Core





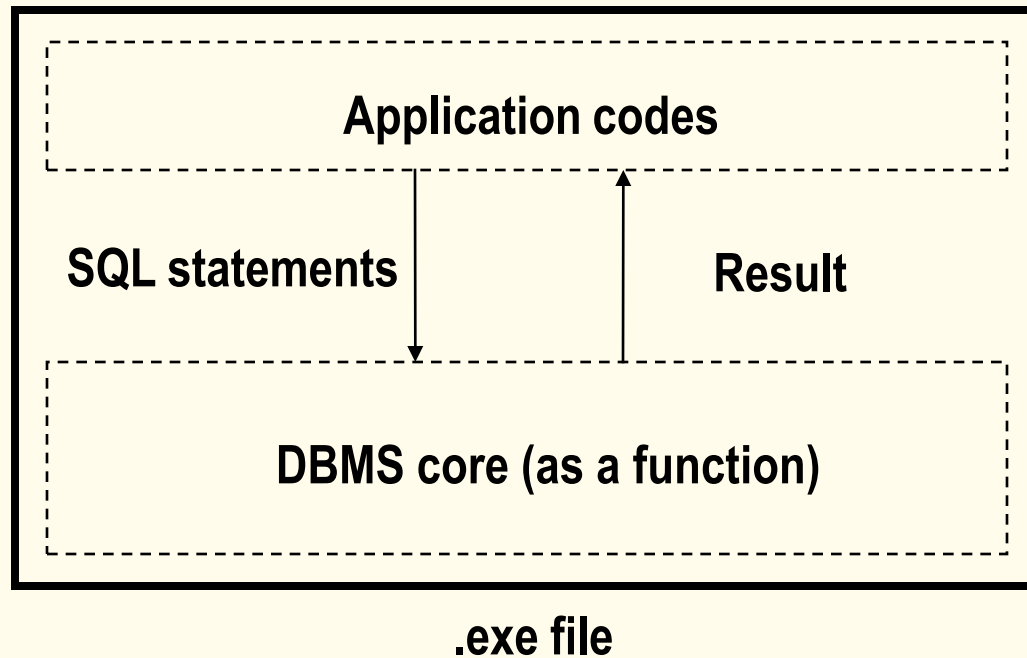
4.2 The Process Structure of DBMS

- Single process structure
- Multi processes structure
- Multi threads structure
- Communication protocols between processes / threads



Single process structure

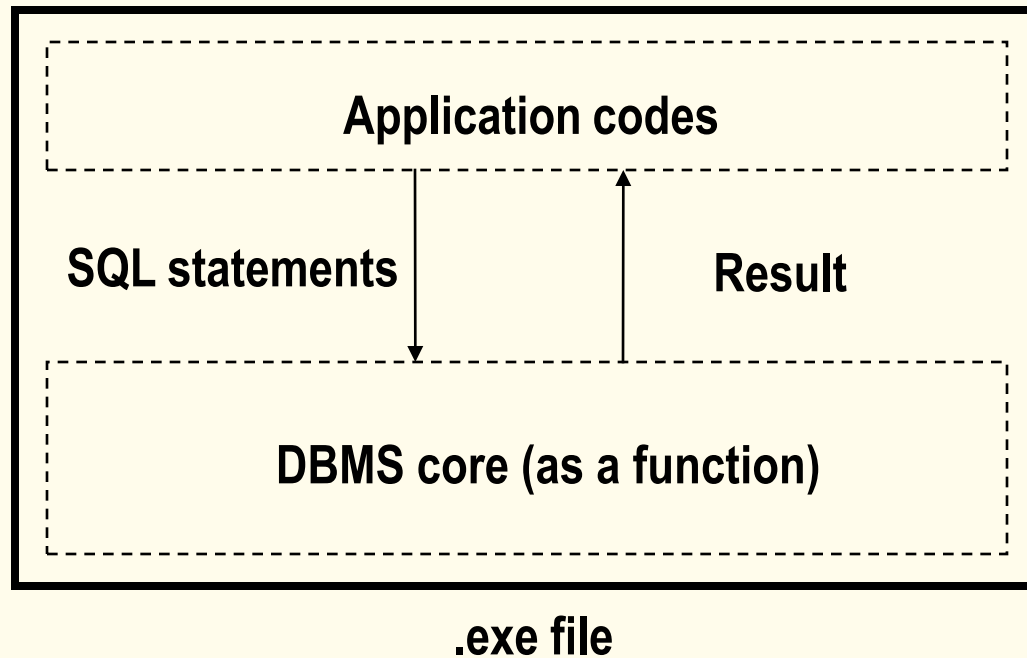
- The application program is compiled with DBMS core as a single .exe file, running as a single process.





Single process structure

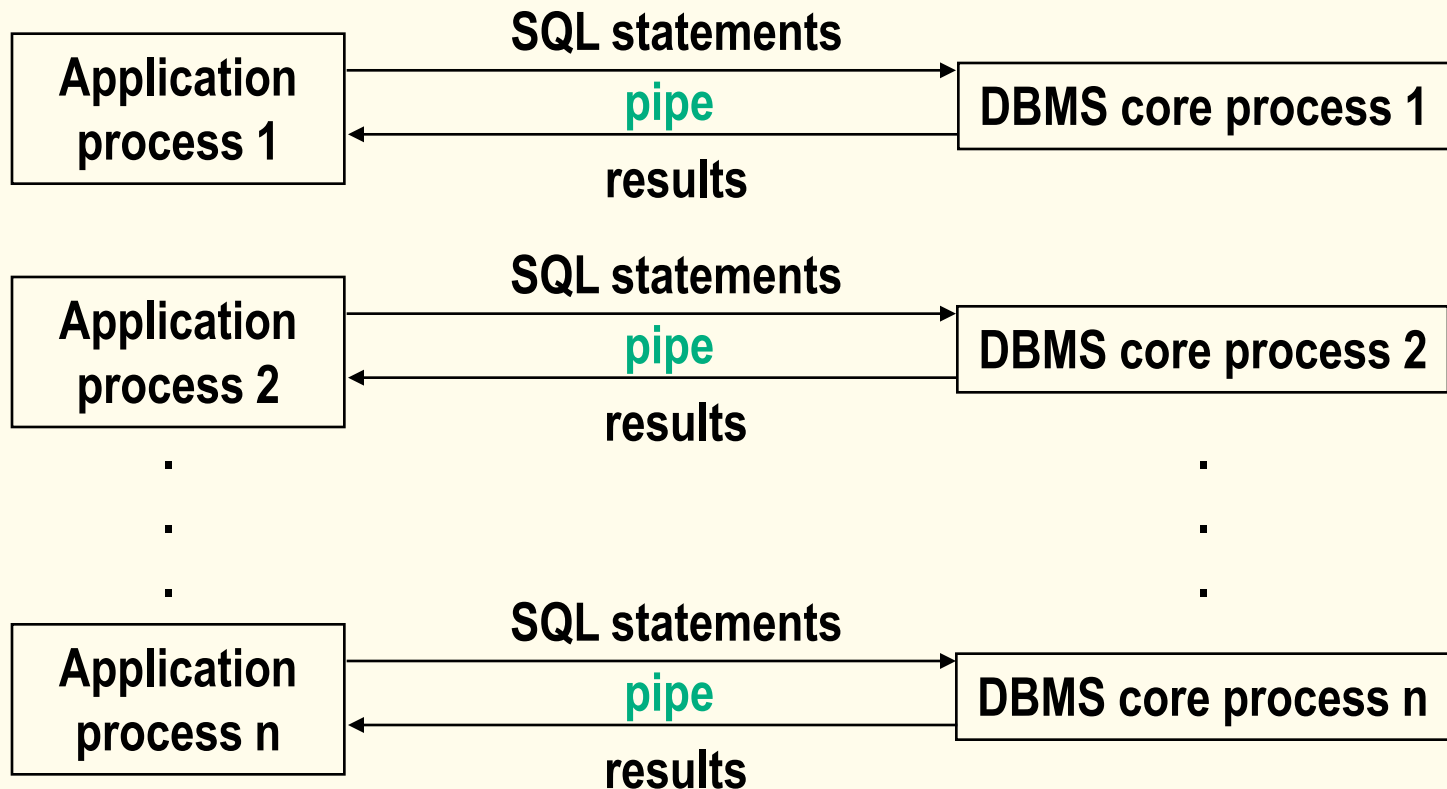
- The application program is compiled with DBMS core as a single .exe file, running as a single process.





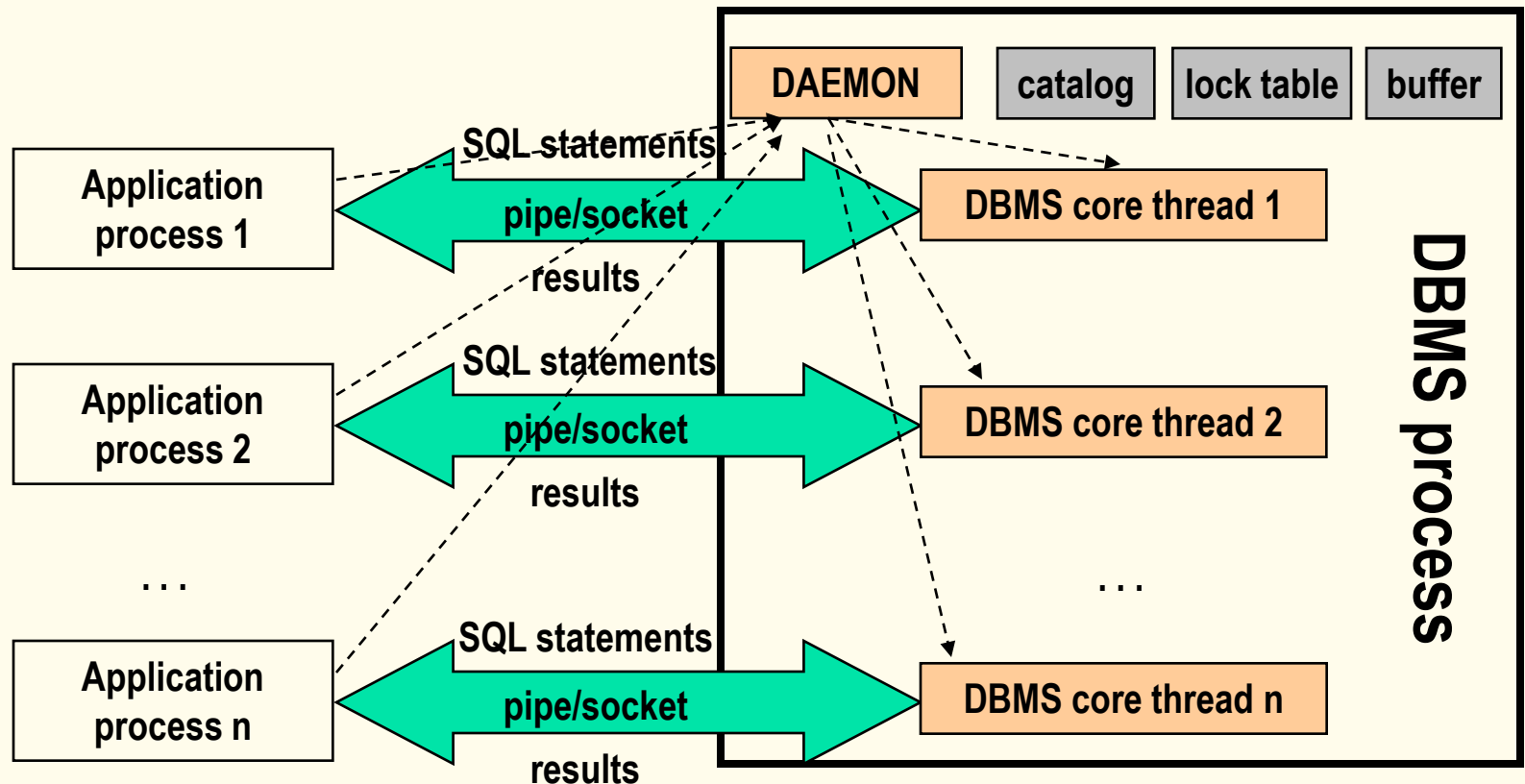
Multi processes structure

- One application process corresponding to one DBMS core process



Multi threads structure

- Only one DBMS process, every application process corresponding to a DBMS core thread.





4.3 Database Access Management

The access to database is transferred to the operations on files (of OS) eventually. The file structure and access route offered on it will affect the speed of data access directly. It is impossible that one kind of file structure will be effective for all kinds of data access.

- Access types
- File organization
- Index technique
- Access primitives



Access Types

- Query all or most records of a file ($>15\%$)
- Query some special record
- Query some records ($<15\%$)
- Scope query
- Update



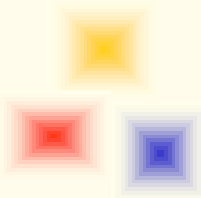
Access Types

- Query all or most records of a file ($>15\%$)
- Query some special record
- Query some records ($<15\%$)
- Scope query
- Update



File Organization

- Heap file: records stored according to their inserted order, and retrieved sequentially. This is the most basic and general form of file organization.
- Direct file: the record address is mapped through hash function according to some attribute's value.
- Indexed file: index + heap file/cluster
- Dynamic hashing: p115
- Grid structure file: p118 (suitable for multi attributes queries)
- Raw disk (notice the difference between the logical block and physical block of file. You can control physical blocks in OS by using raw disk)



File Organization

- Heap file: records stored according to their inserted order, and retrieved sequentially. This is the most basic and general form of file organization.
- Direct file: the record address is mapped through hash function according to some attribute's value.
- Indexed file: index + heap file/cluster
- Dynamic hashing: p115
- Grid structure file: p118 (suitable for multi attributes queries)
- Raw disk (notice the difference between the logical block and physical block of file. You can control physical blocks in OS by using raw disk)



Index Technique

- B+ Tree (✓ ✓)
- Clustering index (✓)
- Inverted file
- Dynamic hashing
- Grid structure file and partitioned hash function
- Bitmap index (used in data warehouse)
- Others



4.4 Query Optimization

“Rewrite” the query statements submitted by user first, and then decide the most effective operating method and steps to get the result. The goal is to gain the result of user’s query with the lowest cost and in shortest time.

4.4.1 Summary of Query Optimization

- **Algebra Optimization**
- **Operation Optimization**



Example

S(SNUM, SNAME, CITY)

SP(SNUM, PNUM, QUAN)

P(PNUM, PNAME, WEIGHT, SIZE)

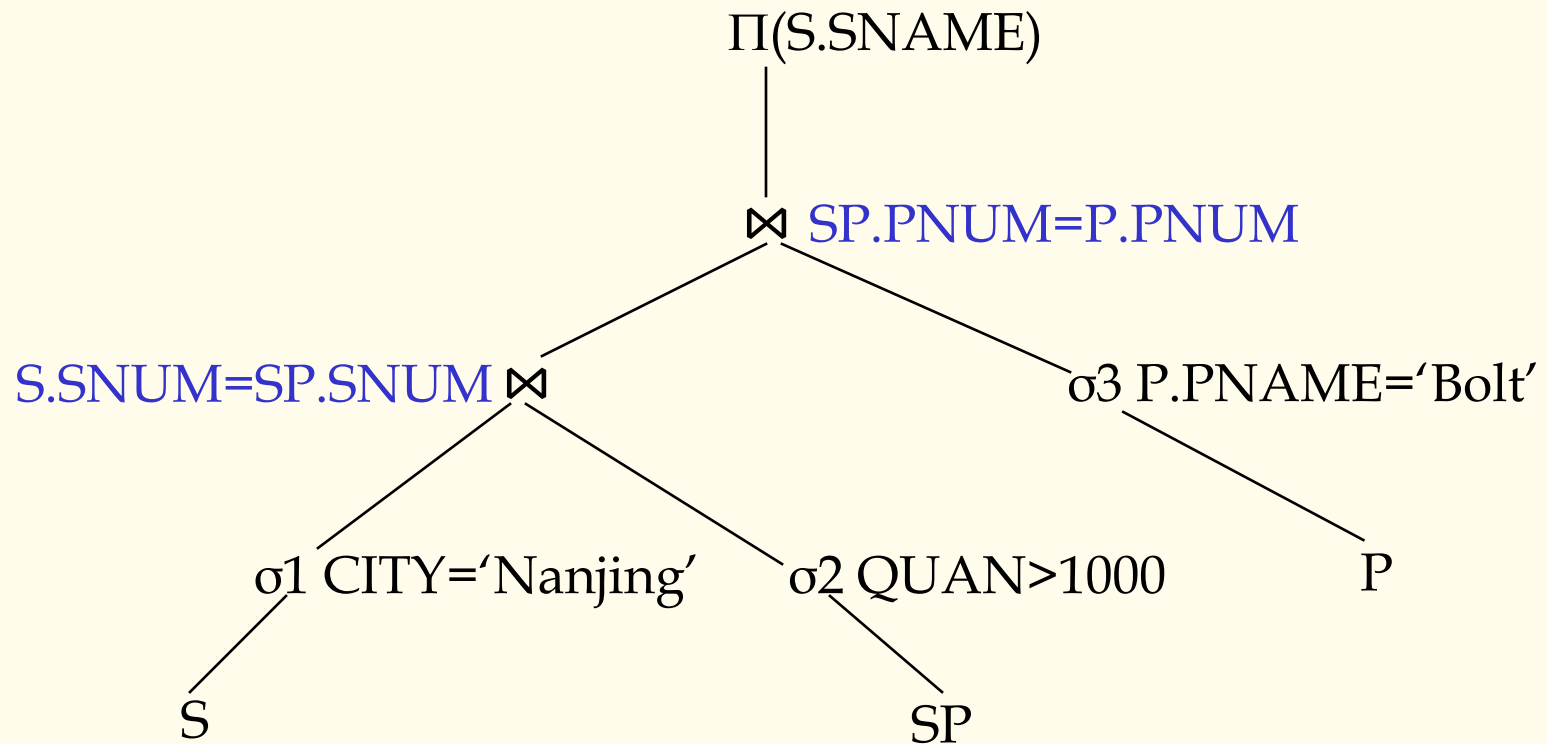
SELECT SNAME

FROM S, SP, P

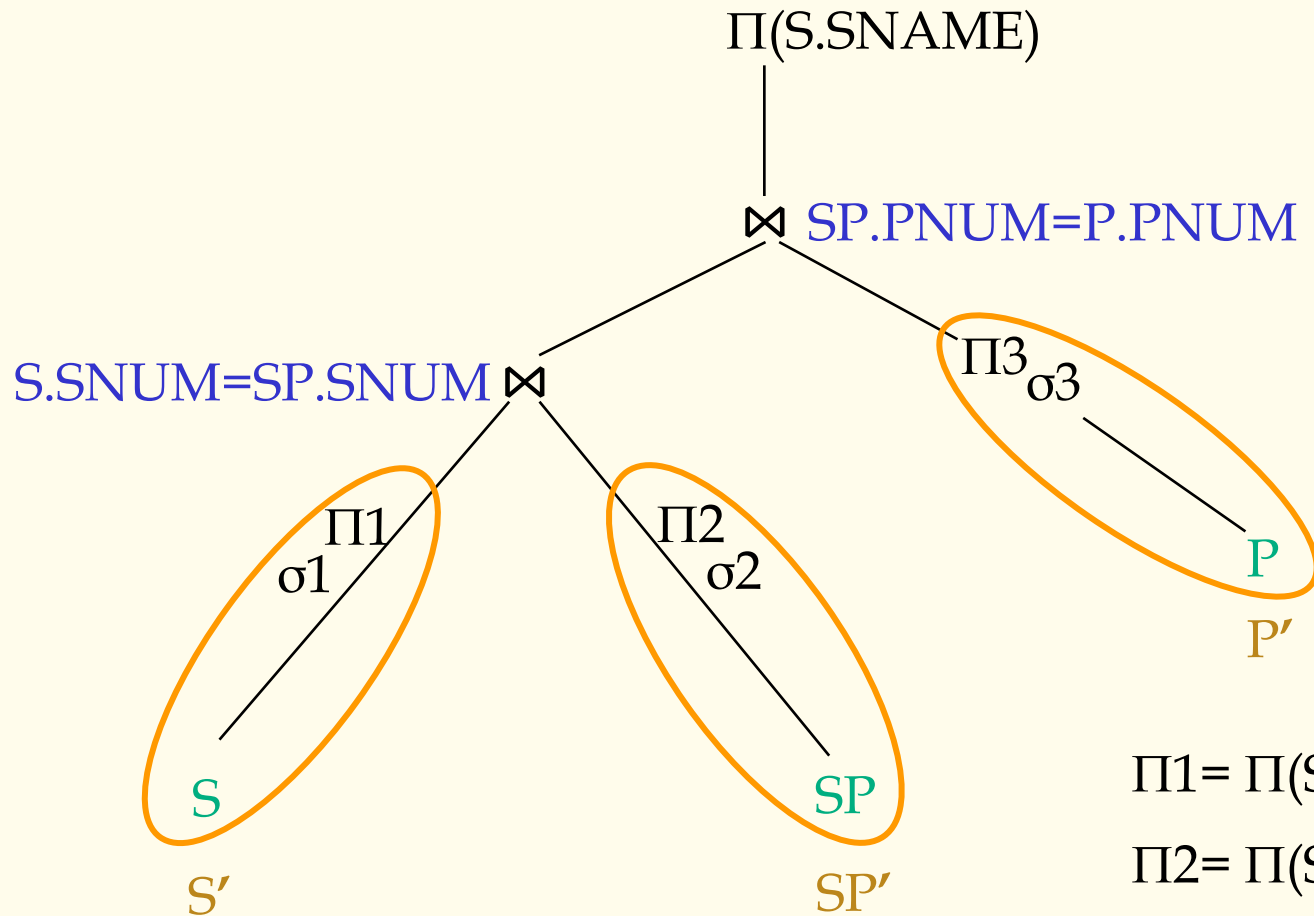
WHERE S.SNUM=SP.SNUM AND
SP.PNUM=P.PNUM AND
S.CITY='Nanjing' AND
P.PNAME='Bolt' AND
SP.QUAN>1000;



Query tree



After equivalent transform (Algebra optimization) :

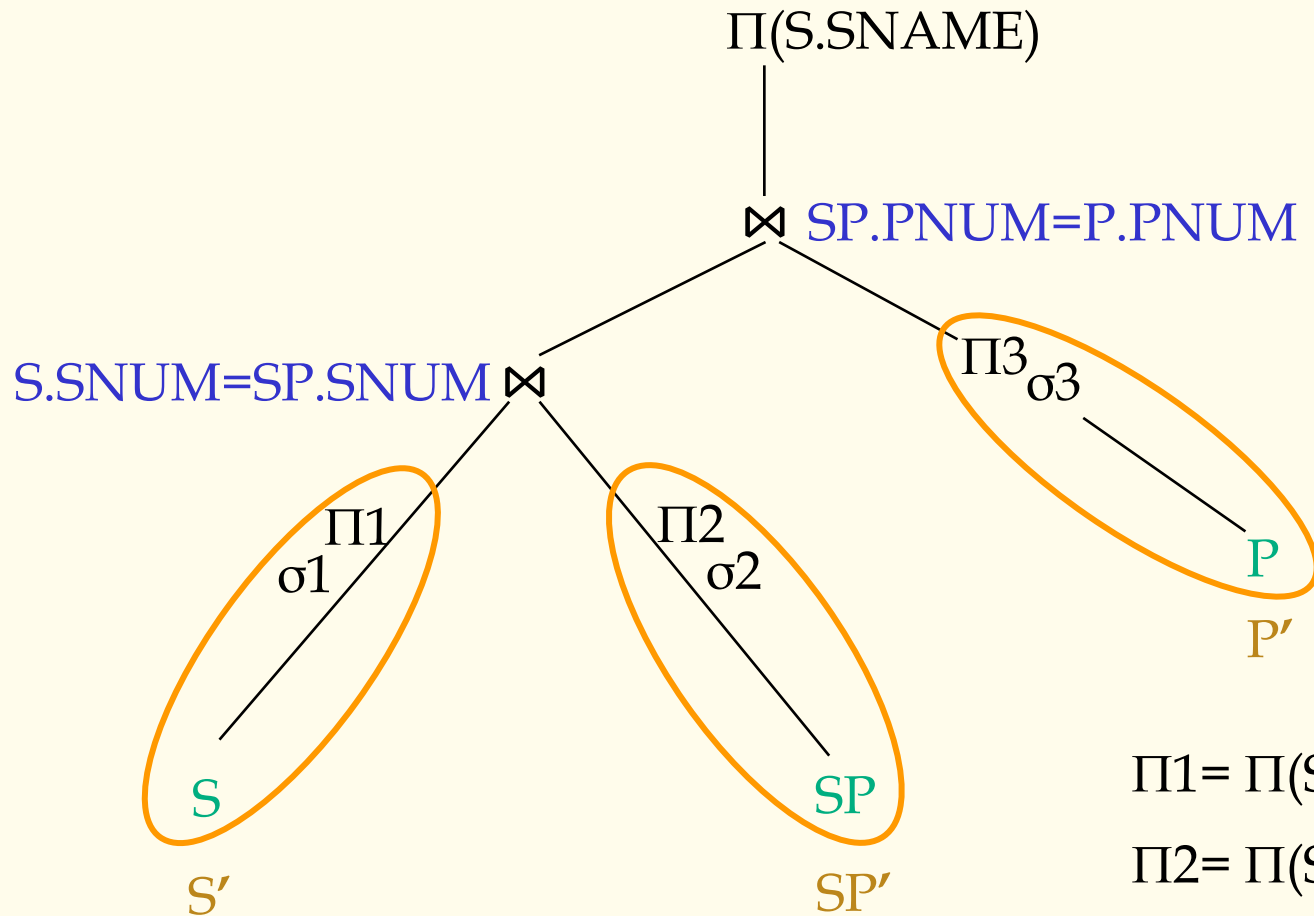


$\Pi_1 = \Pi(S.SNUM, S.SNAME)$

$\Pi_2 = \Pi(SP.SNUM, SP.PNUM)$

$\Pi_3 = \Pi(P.PNUM)$

After equivalent transform (Algebra optimization) :

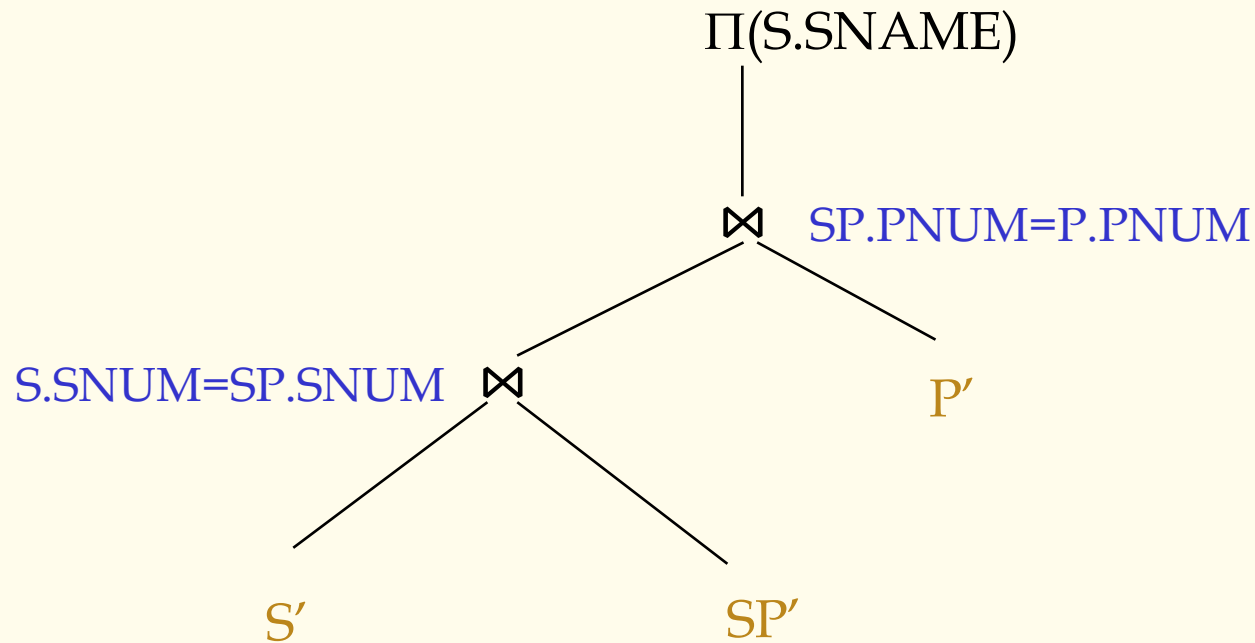


$\Pi_1 = \Pi(S.SNUM, S.SNAME)$

$\Pi_2 = \Pi(SP.SNUM, SP.PNUM)$

$\Pi_3 = \Pi(P.PNUM)$

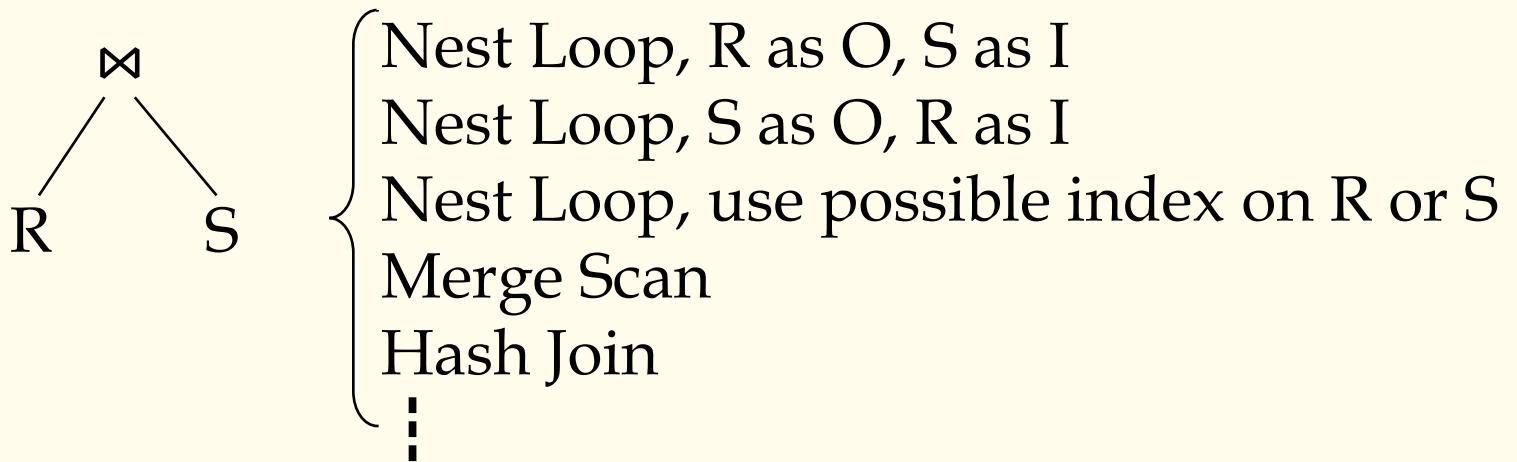
The result of equivalent transform





The operation optimization of the tree :

- Decide the order of two joins
- For every join operation, there are many computing method:



The goal of query optimization is to select a “good” solution from so many possible execution strategies. So it is a complex task.

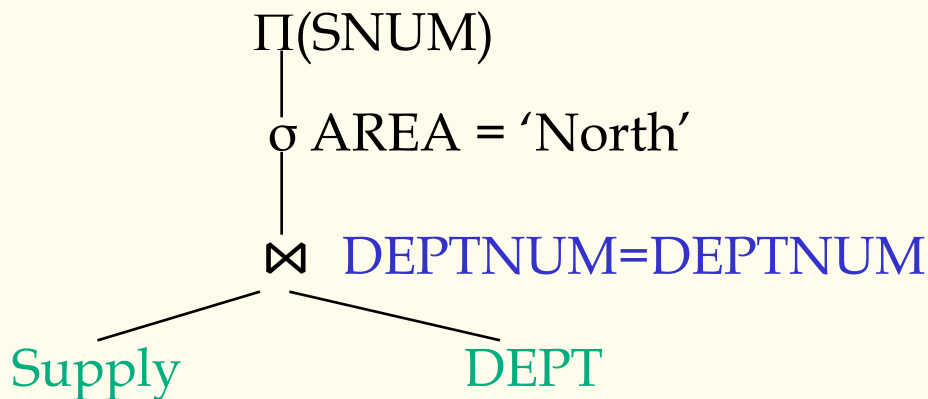
4.4.2 The Equivalent Transform of a Query

That is so called algebra optimization. It takes a series of transform on original query expression, and transform it into an equivalent, most effective form to be executed.

For example: $\Pi_{\text{NAME,DEPT}} \sigma_{\text{DEPT}=15}(\text{EMP}) \equiv \sigma_{\text{DEPT}=15} \Pi_{\text{NAME,DEPT}}(\text{EMP})$

(1) Query tree

For example: $\Pi_{\text{SNUM}} \sigma_{\text{AREA}=\text{'NORTH'}}(\text{SUPPLY} \bowtie_{\text{DEPTNUM}} \text{DEPT})$



Leaves: relations


Middle nodes: unary/binary operations

Leaves \rightarrow root: the executing order of operations



(2) The equivalent transform rules of relational algebra

- 1) Exchange rule of \bowtie/\times : $E1 \times E2 \equiv E2 \times E1$
- 2) Combination rule of \bowtie/\times : $E1 \times (E2 \times E3) \equiv (E1 \times E2) \times E3$
- 3) Cluster rule of Π : $\Pi_{A1 \dots An}(\Pi_{B1 \dots Bm}(E)) \equiv \Pi_{A1 \dots An}(E)$,
legal when $A_1 \dots A_n$ is the sub set of $\{B_1 \dots B_m\}$
- 4) Cluster rule of σ : $\sigma_{F1}(\sigma_{F2}(E)) \equiv \sigma_{F1 \wedge F2}(E)$
- 5) Exchange rule of σ and Π : $\sigma_F(\Pi_{A1 \dots An}(E)) \equiv \Pi_{A1 \dots An}(\sigma_F(E))$
if F includes attributes $B_1 \dots B_m$ which don't belong to $A_1 \dots A_n$, then $\Pi_{A1 \dots An}(\sigma_F(E)) \equiv \Pi_{A1 \dots An} \sigma_F(\Pi_{A1 \dots An, B1 \dots Bm}(E))$
- 6) If the attributes in F are all the attributes in $E1$, then
 $\sigma_F(E1 \times E2) \equiv \sigma_F(E1) \times E2$



if F in the form of $F1 \wedge F2$, and there are only $E1$'s attributes in $F1$, and there are only $E2$'s attributes in $F2$, then $\sigma_F(E1 \times E2) \equiv \sigma_{F1}(E1) \times \sigma_{F2}(E2)$

if F in the form of $F1 \wedge F2$, and there are only $E1$'s attributes in $F1$, while $F2$ includes the attributes both in $E1$ and $E2$, then $\sigma_F(E1 \times E2) \equiv \sigma_{F2}(\sigma_{F1}(E1) \times E2)$


7) $\sigma_F(E1 \cup E2) \equiv \sigma_F(E1) \cup \sigma_F(E2)$

8) $\sigma_F(E1 - E2) \equiv \sigma_F(E1) - \sigma_F(E2)$

9) Suppose $A_1 \dots A_n$ is a set of attributes, in which $B_1 \dots B_m$ are $E1$'s attributes, and $C_1 \dots C_k$ are $E2$'s attributes, then

$$\Pi_{A1 \dots An}(E1 \times E2) \equiv \Pi_{B1 \dots Bm}(E1) \times \Pi_{C1 \dots Ck}(E2)$$

10) $\Pi_{A1 \dots An}(E1 \cup E2) \equiv \Pi_{A1 \dots An}(E1) \cup \Pi_{A1 \dots An}(E2)$



if F in the form of $F1 \wedge F2$, and there are only $E1$'s attributes in $F1$, and there are only $E2$'s attributes in $F2$, then $\sigma_F(E1 \times E2) \equiv \sigma_{F1}(E1) \times \sigma_{F2}(E2)$

if F in the form of $F1 \wedge F2$, and there are only $E1$'s attributes in $F1$, while $F2$ includes the attributes both in $E1$ and $E2$, then $\sigma_F(E1 \times E2) \equiv \sigma_{F2}(\sigma_{F1}(E1) \times E2)$


7) $\sigma_F(E1 \cup E2) \equiv \sigma_F(E1) \cup \sigma_F(E2)$

8) $\sigma_F(E1 - E2) \equiv \sigma_F(E1) - \sigma_F(E2)$


9) Suppose $A_1 \dots A_n$ is a set of attributes, in which $B_1 \dots B_m$ are $E1$'s attributes, and $C_1 \dots C_k$ are $E2$'s attributes, then


$$\Pi_{A1 \dots An}(E1 \times E2) \equiv \Pi_{B1 \dots Bm}(E1) \times \Pi_{C1 \dots Ck}(E2)$$

10) $\Pi_{A1 \dots An}(E1 \cup E2) \equiv \Pi_{A1 \dots An}(E1) \cup \Pi_{A1 \dots An}(E2)$



```
SELECT S.sname
FROM   Sailors S
WHERE  EXISTS (SELECT *
               FROM   Reserves R
               WHERE  R.bid=103 AND S.sid=R.sid)
```

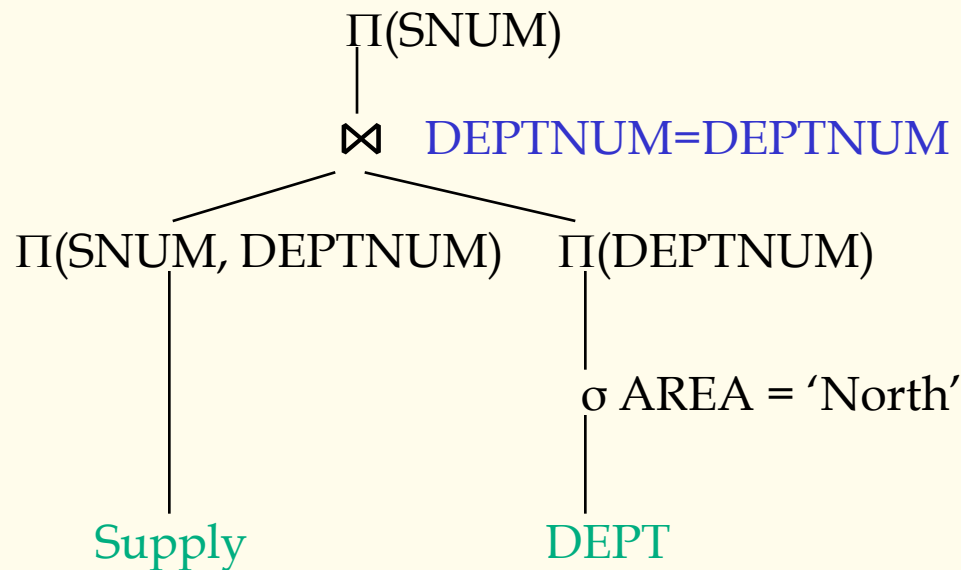




(3) Basic principles

The target of algebra optimization is to make the scale of the operands which involved in binary operations be as small as possible :

- ✓ Push down the unary operations as low as possible
- ✓ Look for and combine the common sub-expression





4.4.3 The Operation Optimization

How to find a “good” access strategy to compute the query improved by algebra optimization is introduced in this section:

- Optimization of select operation
- Optimization of project operation
- Optimization of set operation
- Optimization of join operation
- Optimization of combined operations



Optimization of join operation

- Nested loop: one relation acts as outer loop relation (O), the other acts as inner loop relation (I). For every tuple in O, scan I one time to check join condition.

Because the relation is accessed from disk in the unit of block, we can use block buffer to improve efficiency. For $R \bowtie S$, if let R as O, S as I, b_R is physical block number of R, b_S is physical block number of S, there are n_B block buffers in system ($n_B \geq 2$), and $n_B - 1$ buffers used for O, one buffer used for I, then the total disk access times needed to compute $R \bowtie S$ is:

$$b_R + \lceil b_R / (n_B - 1) \rceil \times b_S$$

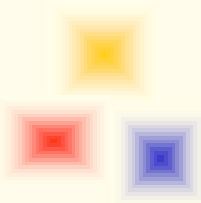


Optimization of join operation

- Nested loop: one relation acts as outer loop relation (O), the other acts as inner loop relation (I). For every tuple in O, scan I one time to check join condition.

Because the relation is accessed from disk in the unit of block, we can use block buffer to improve efficiency. For $R \bowtie S$, if let R as O, S as I, b_R is physical block number of R, b_S is physical block number of S, there are n_B block buffers in system ($n_B \geq 2$), and $n_B - 1$ buffers used for O, one buffer used for I, then the total disk access times needed to compute $R \bowtie S$ is:

$$b_R + \lceil b_R / (n_B - 1) \rceil \times b_S$$



Optimization of join operation

- Merge scan: order the relation R and S on disk in ahead, then we can compare their tuples in order, and both relation only need to scan one time. If R and S have not ordered in ahead, must consider the ordering cost to see if it is worth to use this method (p122)
- Using index or hash to look for mapping tuples: in nested loop method, if there is suitable access route on I (say B+ tree index), it can be used to substitute sequence scan. It is best when there is cluster index or hash on join attributes.
- Hash join: because the join attributes of R and S have the same domain, R and S can be hashed into the same hash file using the same hash function, then $R \bowtie S$ can be computed based on the hash file.



Optimization of join operation

- Merge scan: order the relation R and S on disk in ahead, then we can compare their tuples in order, and both relation only need to scan one time. If R and S have not ordered in ahead, must consider the ordering cost to see if it is worth to use this method (p122)
- Using index or hash to look for mapping tuples: in nested loop method, if there is suitable access route on I (say B+ tree index), it can be used to substitute sequence scan. It is best when there is cluster index or hash on join attributes.
- Hash join: because the join attributes of R and S have the same domain, R and S can be hashed into the same hash file using the same hash function, then $R \bowtie S$ can be computed based on the hash file.



4.5 Recovery

4.5.1 Introduction

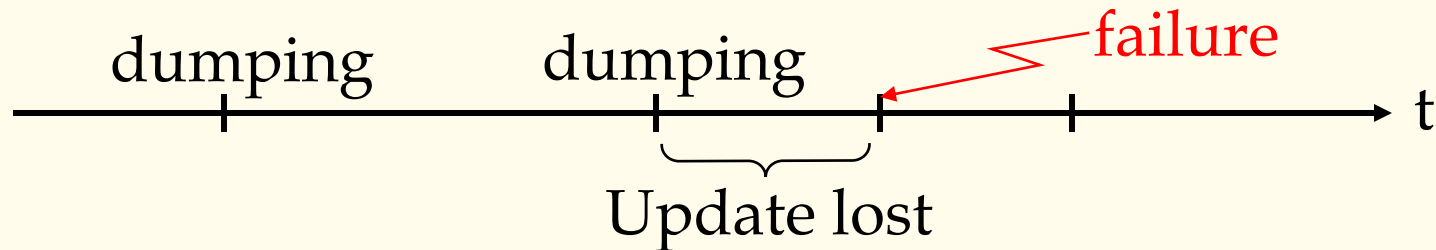
The main roles of recovery mechanism in DBMS are:

- (1) Reducing the likelihood of failures (prevention)
- (2) Recover from failures (solving)

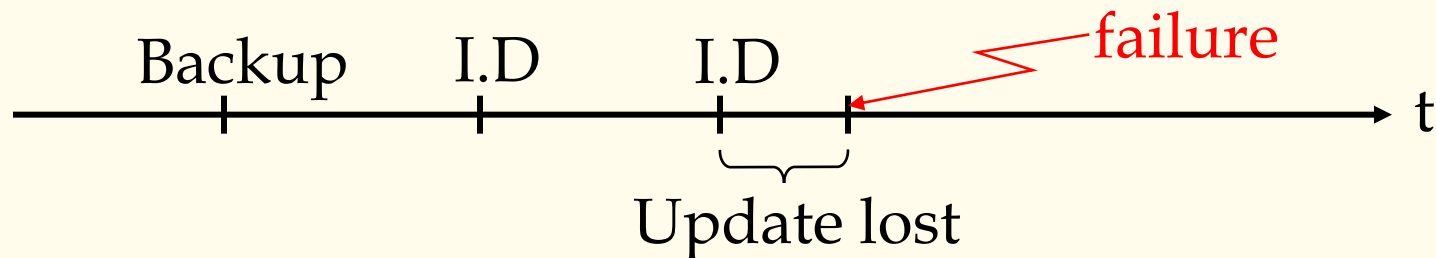
Restore DB to a consistent state after some failures.

- Redundancy is necessary.
- Should inspect **all possible** failures.
- General method:

1) Periodical dumping



- Variation : Backup + Incremental dumping
I.D --- updated parts of DB



This method is easy to be implemented and the overhead is low, but the update maybe lost after failure occurring. So it is often used in file system or small DBMS.

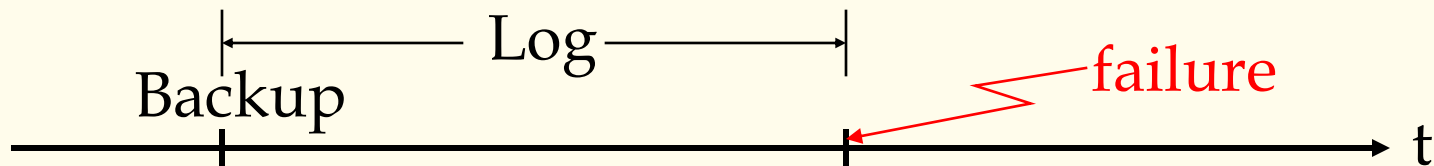


2) Backup + Log

Log : record of **all** changes on DB since the last backup copy was made.

Change: $\left\{ \begin{array}{l} \text{Old value (before image --- B.I)} \\ \text{New value (after image --- A.I)} \end{array} \right\}$ Recorded into Log

| | | | |
|-----|--------------|------|------|
| For | update op. : | B.I | A.I |
| | insert op. : | ---- | A.I |
| | delete op. : | B.I | ---- |





While recovering:

- Some transactions maybe half done, should undo them with B.I recorded in Log.
- Some transactions have finished but the results have not been written into DB in time, should redo them with A.I recorded in Log. (finish writing into DB)

It is possible to recover DB to the **most recent** consistent state with Log.



4.5.2 Transaction

A transaction T is a finite sequence of actions on DB exhibiting the following effects:

- ✓ **A**tomic action: Nothing or All.
- ✓ **C**onsistency preservation: consistency state of DB \rightarrow another consistency state of DB.
- ✓ **I**solation: concurrent transactions should run as if they are independent each other.
- ✓ **D**urability: The effects of a successfully completed transaction are permanently reflected in DB and recoverable even failure occurs later.



Example: transfer money s from account A to account B

Begin transaction

read A

$A := A - s$

if $A < 0$ then Display “insufficient fund”

Rollback /*undo and terminate */

else $B := B + s$

Display “transfer complete”

Commit /*commit the update and terminate */

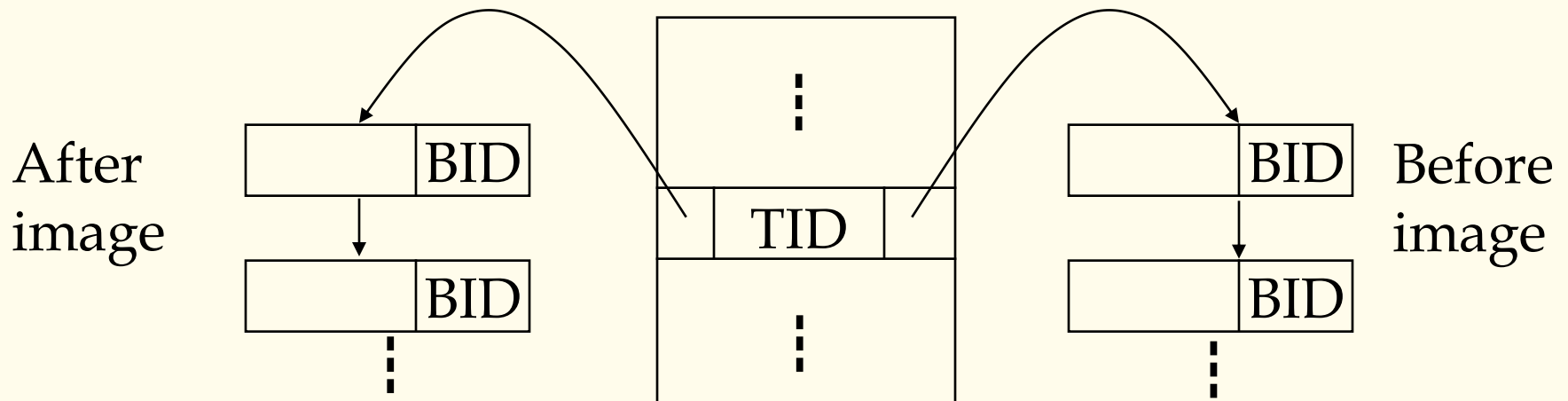
Rollback --- abnormal termination. (Nothing)

Commit --- normal termination. (All)

4.5.3 Some Structures to Support Recovery

Recovery information (such as Log) should be stored in nonvolatile storage. The following information need to be stored in order to support recovery:

- 1) Commit list : list of TID which have been committed.
- 2) Active list : list of TID which is in progress.
- 3) Log :





4.5.4 Commit Rule and Log Ahead Rule

1) Commit Rule

A.I must be written to nonvolatile storage before commit of the transaction.

2) Log Ahead Rule

If A.I is written to DB before commit then B.I must first written to log.

3) Recovery strategies

(1) The features of undo and redo (are idempotent) :

$\text{undo}(\text{undo}(\text{undo} \dots \text{undo}(x) \dots)) = \text{undo}(x)$

$\text{redo}(\text{redo}(\text{redo} \dots \text{redo}(x) \dots)) = \text{redo}(x)$



4.5.4 Commit Rule and Log Ahead Rule

1) Commit Rule

A.I must be written to nonvolatile storage before commit of the transaction.

2) Log Ahead Rule

If A.I is written to DB before commit then B.I must first written to log.

3) Recovery strategies

(1) The features of undo and redo (are idempotent) :

$\text{undo}(\text{undo}(\text{undo} \dots \text{undo}(x) \dots)) = \text{undo}(x)$

$\text{redo}(\text{redo}(\text{redo} \dots \text{redo}(x) \dots)) = \text{redo}(x)$



(2) Three kinds of update strategy

a) $A.I \rightarrow DB$ before commit

$TID \rightarrow \text{active list}$

$\left[\begin{array}{l} \downarrow \\ \end{array} \right\} \begin{array}{l} B.I \rightarrow \text{Log} \\ A.I \rightarrow DB \\ \vdots \end{array} \quad (\text{Log Ahead Rule})$

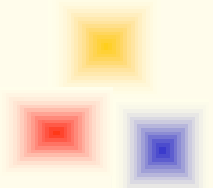
commit $\left\{ \begin{array}{l} TID \rightarrow \text{commit list} \\ \text{delete TID from active list} \end{array} \right.$



The recovery after failure in this situation

Check two lists for every TID while restarting after failure:

| Commit list | Active list | |
|-------------|-------------|-----------------------------------|
| | ✓ | Undo, delete TID from active list |
| ✓ | ✓ | delete TID from active list |
| ✓ | | nothing to do |



b) $A.I \rightarrow DB$ after commit

$TID \rightarrow \text{active list}$

$\boxed{\downarrow} \left\{ \begin{array}{l} A.I \rightarrow \text{Log} \\ \vdots \end{array} \right. \quad (\text{Commit Rule})$

commit $\left\{ \begin{array}{l} TID \rightarrow \text{commit list} \\ \boxed{\downarrow} A.I \rightarrow DB \\ \text{delete TID from active list} \end{array} \right.$



The recovery after failure in this situation

Check two lists for every TID while restarting after failure:

| Commit list | Active list | |
|-------------|-------------|-----------------------------------|
| | ✓ | delete TID from active list |
| ✓ | ✓ | redo, delete TID from active list |
| ✓ | | nothing to do |



c) $A.I \rightarrow DB$ concurrently with commit

$TID \rightarrow \text{active list}$

$\left[\begin{array}{l} \downarrow \\ \downarrow \end{array} \right] \left\{ \begin{array}{l} A.I, B.I \rightarrow \text{Log} \\ A.I \rightarrow DB \end{array} \right. \quad \begin{array}{l} \text{(Two Rules)} \\ \text{(partially done)} \end{array}$

\vdots

$\text{commit} \left\{ \begin{array}{l} TID \rightarrow \text{commit list} \\ \left[\begin{array}{l} \downarrow \\ \downarrow \end{array} \right] A.I \rightarrow DB \\ \text{delete TID from active list} \end{array} \right. \quad \begin{array}{l} \\ \text{(completed)} \end{array}$



c) $A.I \rightarrow DB$ concurrently with commit

$TID \rightarrow \text{active list}$

$\left[\begin{array}{l} \downarrow \\ \downarrow \end{array} \right] \left\{ \begin{array}{l} A.I, B.I \rightarrow \text{Log} \\ A.I \rightarrow DB \end{array} \right. \quad \begin{array}{l} \text{(Two Rules)} \\ \text{(partially done)} \end{array}$

\vdots

$\text{commit} \left\{ \begin{array}{l} TID \rightarrow \text{commit list} \\ \left[\begin{array}{l} \downarrow \end{array} \right] A.I \rightarrow DB \\ \text{delete TID from active list} \end{array} \right. \quad \begin{array}{l} \\ \text{(completed)} \end{array}$



The recovery after failure in this situation

Check two lists for every TID while restarting after failure:

| Commit list | Active list | |
|-------------|-------------|-----------------------------------|
| | ✓ | Undo, delete TID from active list |
| ✓ | ✓ | redo, delete TID from active list |
| ✓ | | nothing to do |



Conclusion :

| | redo | undo |
|----|------|------|
| a) | ✗ | ✓ |
| b) | ✓ | ✗ |
| c) | ✓ | ✓ |
| d) | ✗ | ✗ |

?

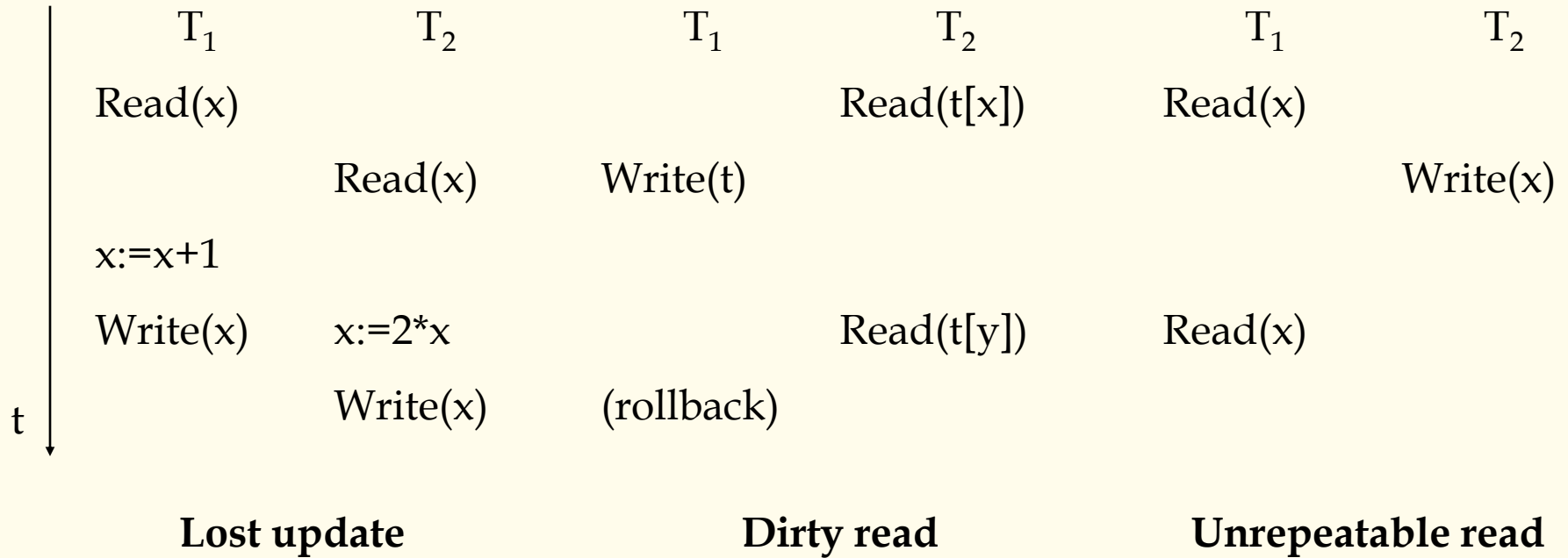


4.6 Concurrency Control

4.6.1 Introduction

In multi users DBMS, permit multi transaction access the database concurrently.

- Why concurrency?
 - 1) Improving system utilization & response time.
 - 2) Different transaction may access to different parts of database.
- Problems arise from concurrent executions



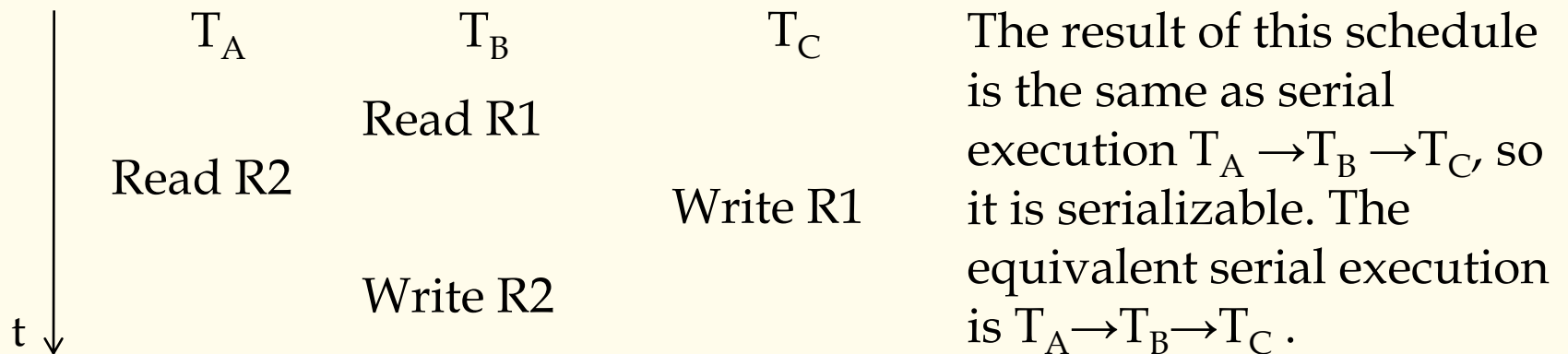
So there may be three kinds of conflict when transactions execute concurrently. They are write – write, write – read, and read – write conflicts. Write – write conflict must be avoided anytime. Write – read and read – write conflicts should be avoided generally, but they are endurable in some applications.



4.6.2 Serialization --- the criterion for concurrency consistency

Definition: suppose $\{T_1, T_2, \dots, T_n\}$ is a set of transactions executing concurrently. If a schedule of $\{T_1, T_2, \dots, T_n\}$ produces the same effect on database as some serial execution of this set of transactions, then the schedule is serializable.

Problem: different schedule \rightarrow different equivalent serial execution \rightarrow different result? (yes, $n!$)



4.6.3 Locking Protocol

Locking method is the most basic concurrency control method. There maybe many kinds of locking protocols.

(1) X locks

Only one type of lock, for both read and write.

Compatibility matrix : NL—no lock X—X lock
 Y —compatible N—incompatible

| | NL | X |
|----|----|---|
| NL | Y | Y |
| X | Y | N |

T_A
 X_lock R
 Update R
 ⋮
 X_unlock R
 EOT

T_B
 X_lock R
 wait
 ↓
 X_lock R
 Read R
 ⋮



*Two Phase Locking

- *Definiton1:* In a transaction, if all locks precede all unlocks, then the transaction is called two phase transaction. This restriction is called two phase locking protocol.
- *Definition2:* In a transaction, if it first acquires a lock on the object before operating on it, it is called well-formed.

- *Theorem:* If S is any schedule of well-formed and two phase transactions, then S is serializable. (proving is on p151)

| | T ₁ | T ₂ |
|-----------------|----------------------------------|--|
| Growing phase | Lock A Lock B Lock C ⋮ | Lock A Lock B Unlock A Unlock B |
| Shrinking phase | Unlock A Unlock B Unlock C | Lock C ⋮ Unlock C |
| | 2PL | not 2PL |



Conclusions :

- 1) Well-formed + 2PL : serializable
- 2) Well-formed + 2PL + unlock update at EOT: serializable and recoverable. (without domino phenomena)
- 3) Well-formed + 2PL + holding all locks to EOT: strict two phase locking transaction.

(2) (S,X) locks

S lock --- if read access is intended.

X lock --- if update access is intended.

| | NL | S | X |
|----|----|---|---|
| NL | Y | Y | Y |
| S | Y | Y | N |
| X | Y | N | N |

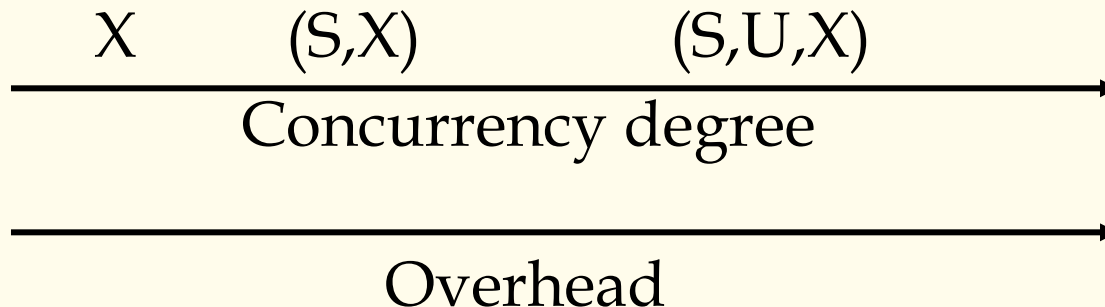


(3) (S,U,X) locks

U lock --- update lock. For an update access the transaction first acquires a U-lock and then promote it to X-lock.

Purpose: shorten the time of exclusion, so as to boost concurrency degree, and reduce deadlock.

| | NL | S | U | X |
|----|----|---|---|---|
| NL | Y | Y | Y | Y |
| S | Y | Y | Y | N |
| U | Y | Y | N | N |
| X | Y | N | N | N |



4.6.4 Deadlock & Live Lock

Dead lock: wait in cycle, no transaction can obtain all of resources needed to complete.

Live lock: although other transactions release their resource in limited time, some transaction can not get the resources needed for a very long time.

T_A
X_lock R1

⋮

X_lock R2

wait

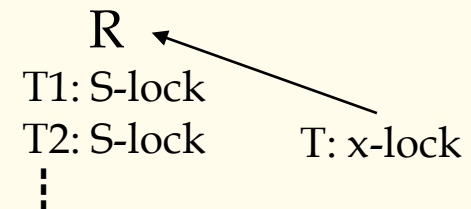


T_B
X_lock R2

⋮

X_lock R1

wait




- Live lock is simpler, only need to adjust schedule strategy, such as FIFO
- Deadlock: (1) Prevention(don't let it occur); (2) Solving(permit it occurs, but can solve it)




(1) Deadlock Detection

- 1) Timeout: If a transaction waits for some specified time then deadlock is **assumed** and the transaction should be aborted.
- 2) Detect deadlock by wait-for graph $G = \langle V, E \rangle$
 V : set of transactions $\{T_i \mid T_i \text{ is a transaction in DBS } (i=1,2,\dots,n)\}$
 E : $\{ \langle T_i, T_j \rangle \mid T_i \text{ waits for } T_j (i \neq j) \}$
 - If there is cycle in the graph, the deadlock occurs.
 - When to detect?
 - 1) whenever one transaction waits.
 - 2) periodically

- 
- What to do when detected?
 - 1) Pick a victim (youngest, minimum abort cost, ...)
 - 2) Abort the victim and release its locks and resources
 - 3) Grant a waiter
 - 4) Restart the victim (automatically or manually)


(2) Deadlock avoidance

- 1) Requesting all locks at initial time of transaction.
- 2) Requesting locks in a specified order of resource.
- 3) Abort once conflicted.
- 4) Transaction Retry

- 
- What to do when detected?
 - 1) Pick a victim (youngest, minimum abort cost, ...)
 - 2) Abort the victim and release its locks and resources
 - 3) Grant a waiter
 - 4) Restart the victim (automatically or manually)

(2) Deadlock avoidance

- 1) Requesting all locks at initial time of transaction.
- 2) Requesting locks in a specified order of resource.
- 3) Abort once conflicted.
- 4) Transaction Retry



Every transaction is uniquely time stamped. If T_A requires a lock on a data object that is already locked by T_B , one of the following methods is used:

- a) Wait-die: T_A waits if it is older than T_B , otherwise it “dies”, i.e. it is aborted and automatically retried with original timestamp.
- b) Wound-wait: T_A waits if it is younger than T_B , otherwise it “wound” T_B , i.e. T_B is aborted and automatically retried with original timestamp.

In above, both have only one direction wait, either older \rightarrow younger or younger \rightarrow older. It is impossible to occur wait in cycle, so the dead lock is avoided.