

Concurrency Control

Hector Garcia-Molina and
Mahmoud SAKR

Invited Lectures

Extensible Databases - 6/10/2021



Dimitri Fontaine

CITUS BLOG AUTHOR PROFILE

PostgreSQL major contributor & author of “The Art of PostgreSQL”. Contributed extension facility & event triggers feature in Postgres. Maintains `pg_auto_failover`. Speaker at so many conferences.

Distributed Databases - 8/12/2021

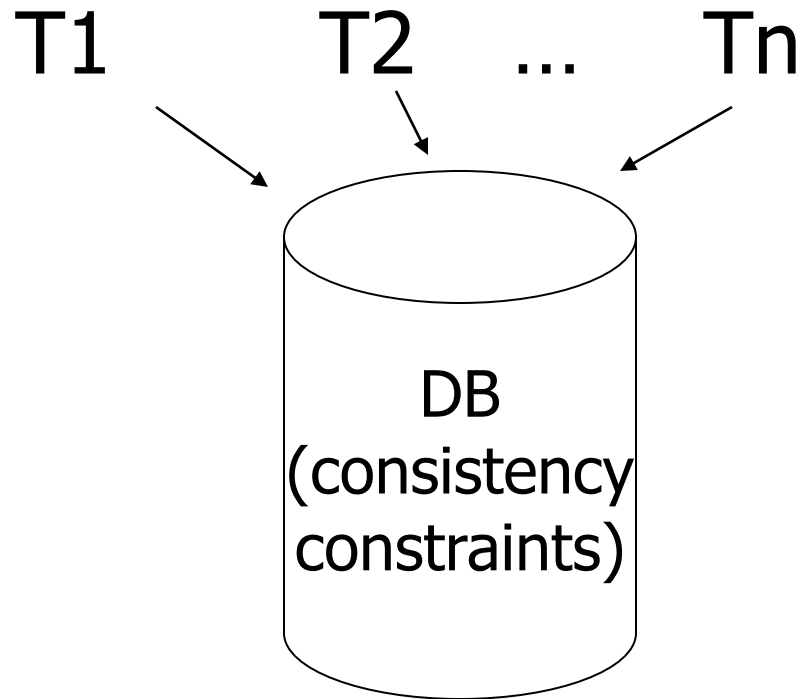


Marco Slot

CITUS BLOG AUTHOR PROFILE

Lead engineer on the Citus engine team at Microsoft. Speaker at Postgres Conf EU, PostgresOpen, pgDay Paris, Hello World, SIGMOD, & lots of meetups. PhD in distributed systems. Loves mountain hiking.

Chapter 18 [18] Concurrency Control



Example:

T1: Read(A)
A \leftarrow A+100
Write(A)
Read(B)
B \leftarrow B+100
Write(B)

T2: Read(A)
A \leftarrow A \times 2
Write(A)
Read(B)
B \leftarrow B \times 2
Write(B)

Constraint: A=B

Schedule A

T1

T2

Read(A); $A \leftarrow A + 100$

Write(A);

Read(B); $B \leftarrow B + 100$;

Write(B);

Read(A); $A \leftarrow A \times 2$;

Write(A);

Read(B); $B \leftarrow B \times 2$;

Write(B);

Schedule A

T1	T2	A	B
		25	25
Read(A); $A \leftarrow A + 100$			
Write(A);		125	
Read(B); $B \leftarrow B + 100$;			125
Write(B);			
	Read(A); $A \leftarrow A \times 2$;		
	Write(A);	250	
	Read(B); $B \leftarrow B \times 2$;		
	Write(B);		250
		250	250

Schedule B

T1

T2

Read(A); $A \leftarrow A \times 2$;

Write(A);

Read(B); $B \leftarrow B \times 2$;

Write(B);

Read(A); $A \leftarrow A + 100$

Write(A);

Read(B); $B \leftarrow B + 100$;

Write(B);

Schedule B

T1

T2

Read(A); $A \leftarrow A+100$
 Write(A);
 Read(B); $B \leftarrow B+100$;
 Write(B);

Read(A); $A \leftarrow A \times 2$;
 Write(A);
 Read(B); $B \leftarrow B \times 2$;
 Write(B);

A	B
25	25
50	
	50
150	
	150
150	150

Schedule C

T1

T2

Read(A); $A \leftarrow A + 100$

Write(A);

Read(A); $A \leftarrow A \times 2$;

Write(A);

Read(B); $B \leftarrow B + 100$;

Write(B);

Read(B); $B \leftarrow B \times 2$;

Write(B);

Schedule C

T1	T2	A	B
		25	25
Read(A); $A \leftarrow A + 100$			
Write(A);		125	
	Read(A); $A \leftarrow A \times 2$;		
	Write(A);	250	
Read(B); $B \leftarrow B + 100$;			
Write(B);			125
	Read(B); $B \leftarrow B \times 2$;		
	Write(B);		250
		250	250

Schedule D

T1

T2

Read(A); $A \leftarrow A + 100$

Write(A);

Read(A); $A \leftarrow A \times 2$;

Write(A);

Read(B); $B \leftarrow B \times 2$;

Write(B);

Read(B); $B \leftarrow B + 100$;

Write(B);

Schedule D

T1	T2	A	B
		25	25
Read(A); $A \leftarrow A + 100$			
Write(A);		125	
	Read(A); $A \leftarrow A \times 2$;		
	Write(A);	250	
	Read(B); $B \leftarrow B \times 2$;		
	Write(B);		50
Read(B); $B \leftarrow B + 100$;			
Write(B);			150
		250	150

Schedule E

Same as Schedule D
but with new T2'

T1

T2'

Read(A); $A \leftarrow A + 100$

Write(A);

Read(A); $A \leftarrow A \times 1$;

Write(A);

Read(B); $B \leftarrow B \times 1$;

Write(B);

Read(B); $B \leftarrow B + 100$;

Write(B);

Schedule E

Same as Schedule D
but with new T2'

T1	T2'	A	B
		25	25
Read(A); $A \leftarrow A + 100$			
Write(A);		125	
	Read(A); $A \leftarrow A \times 1$;		
	Write(A);	125	
	Read(B); $B \leftarrow B \times 1$;		
	Write(B);		25
Read(B); $B \leftarrow B + 100$;			
Write(B);			125
		125	125

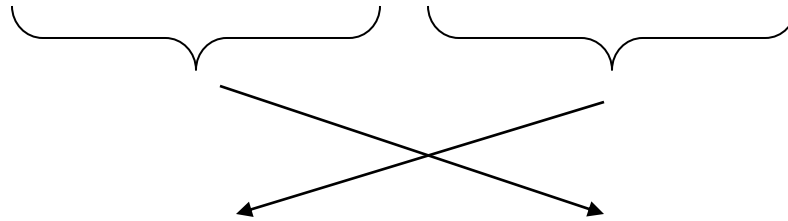
- Want schedules that are “good”, regardless of
 - initial state and
 - transaction semantics
- Only look at order of read and writes

Example:

$Sc = r_1(A)w_1(A)r_2(A)w_2(A)r_1(B)w_1(B)r_2(B)w_2(B)$

Example:

$$Sc = r_1(A)w_1(A)r_2(A)w_2(A)r_1(B)w_1(B)r_2(B)w_2(B)$$



$$Sc' = r_1(A)w_1(A) \ r_1(B)w_1(B)r_2(A)w_2(A)r_2(B)w_2(B)$$

T_1

T_2

The Transaction Game

A								
B								
T1								
T2								

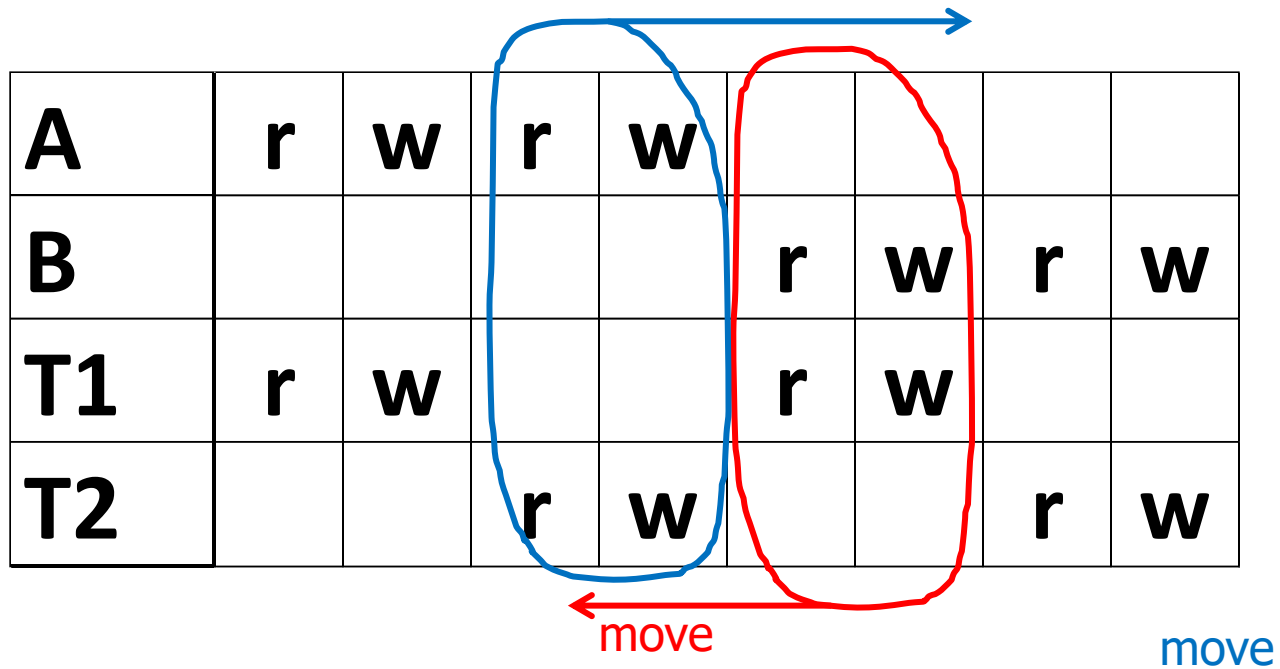
The Transaction Game

A	r	w	r	w				
B					r	w	r	w
T1	r	w			r	w		
T2			r	w			r	w

The Transaction Game

A	r	w	r	w				
B					r	w	r	w
T1	r	w			r	w		
T2			r	w			r	w

←→
can move column



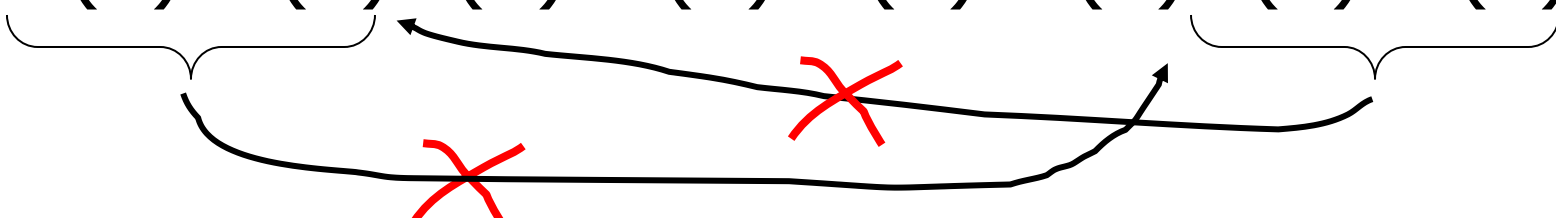
A	r	w			r	w		
B			r	w			r	w
T1	r	w	r	w				
T2					r	w	r	w

Schedule D

A	r	w	r	w				
B					r	w	r	w
T1	r	w					r	w
T2			r	w	r	w		

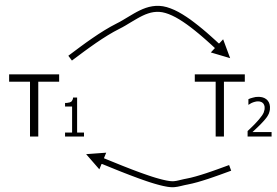
However, for S_d :

$S_d = r_1(A)w_1(A)r_2(A)w_2(A) \ r_2(B)w_2(B)r_1(B)w_1(B)$



- as a matter of fact,
 T_2 must precede T_1
in any equivalent schedule,
i.e., $T_2 \rightarrow T_1$

- $T_2 \rightarrow T_1$
- Also, $T_1 \rightarrow T_2$



Sd cannot be rearranged
into a serial schedule



Sd is not “equivalent” to
any serial schedule



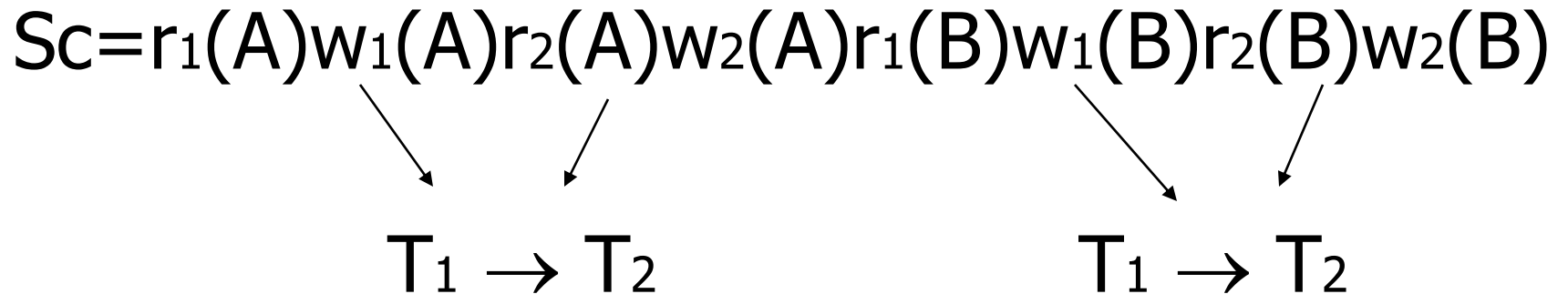
Sd is “bad”

Returning to Sc

$$\text{Sc} = r_1(A)w_1(A)r_2(A)w_2(A)r_1(B)w_1(B)r_2(B)w_2(B)$$

$T_1 \rightarrow T_2$ $T_1 \rightarrow T_2$

Returning to Sc

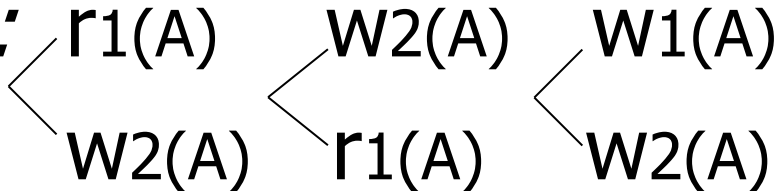


☛ no cycles \Rightarrow Sc is “equivalent” to a serial schedule
(in this case T_1, T_2)

Concepts

Transaction: sequence of $r_i(x)$, $w_i(x)$ actions

Conflicting actions:



$r_1(A) < w_2(A)$
 $w_2(A) < r_1(A)$
 $w_1(A) < w_2(A)$

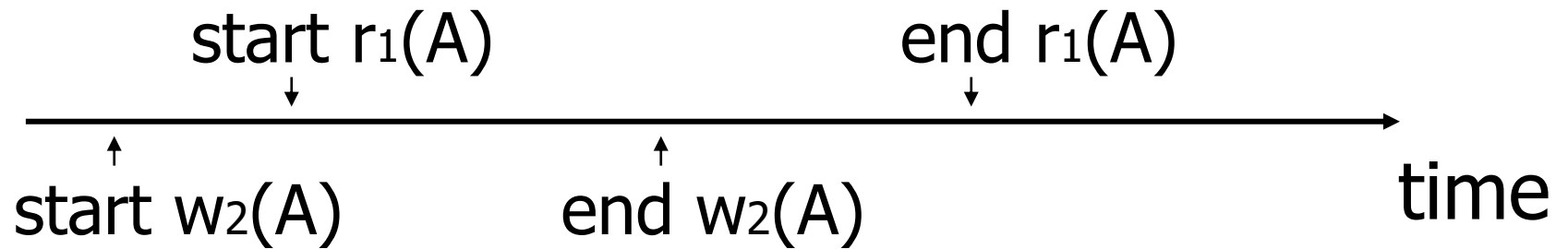
Schedule: represents chronological order
in which actions are executed

Serial schedule: no interleaving of actions
or transactions

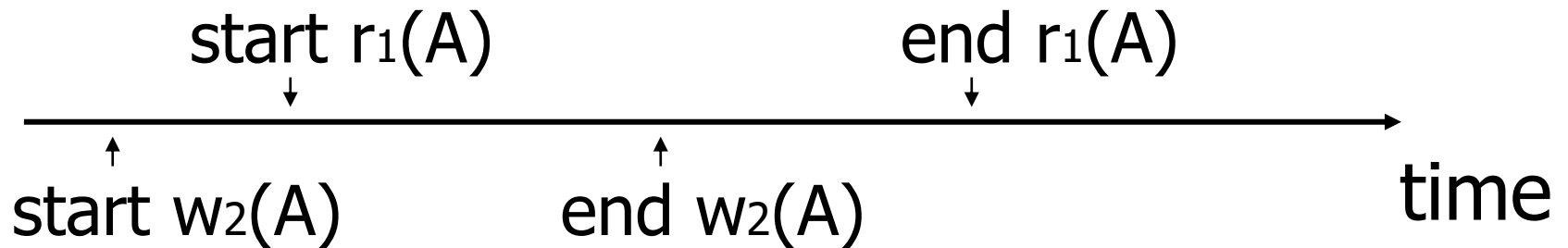
Is it OK to model reads & writes as occurring at a single point in time in a schedule?

- $S = \dots r_1(x) \dots w_2(b) \dots$

What about conflicting, concurrent actions on same object?



What about conflicting, concurrent actions on same object?



- Assume equivalent to either $r_1(A) w_2(A)$
or $w_2(A) r_1(A)$
- \Rightarrow low level synchronization mechanism
- Assumption called “atomic actions”

Definition

S_1, S_2 are conflict equivalent schedules if S_1 can be transformed into S_2 by a series of swaps on non-conflicting actions.

Definition

A schedule is conflict serializable if it is conflict equivalent to some serial schedule.

Precedence graph $P(S)$ (S is schedule)

Nodes: transactions in S

Arcs: $T_i \rightarrow T_j$ whenever

- $p_i(A), q_j(A)$ are actions in S
- $p_i(A) <_S q_j(A)$
- at least one of p_i, q_j is a write

Exercise:

- What is $P(S)$ for
 $S = w_3(A) w_2(C) r_1(A) w_1(B) r_1(C) w_2(A) r_4(A) w_4(D)$
- Is S serializable?

Another Exercise:

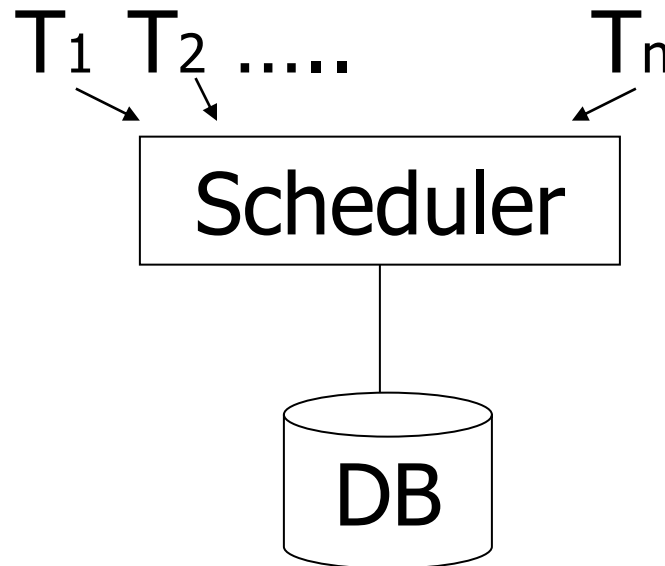
- What is $P(S)$ for
 $S = w_1(A) r_2(A) r_3(A) w_4(A) ?$

How to enforce serializable schedules?

Option 1: run system, recording $P(S)$;
at end of day, check for $P(S)$
cycles and declare if execution
was good

How to enforce serializable schedules?

Option 2: prevent P(S) cycles from occurring

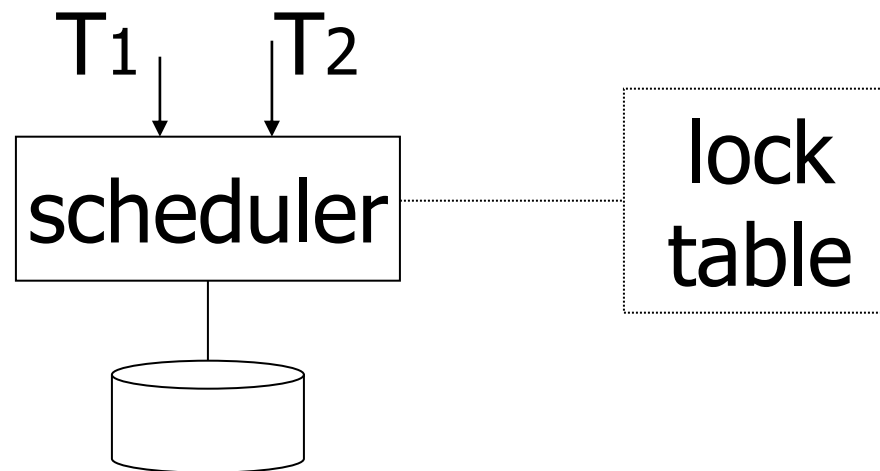


A locking protocol

Two new actions:

lock (exclusive): $li(A)$

unlock: $ui(A)$



Rule #1: Well-formed transactions

$T_i: \dots l_i(A) \dots p_i(A) \dots u_i(A) \dots$

Rule #2 Legal scheduler

$S = \dots \dots \dots l_i(A) \dots \dots \dots u_i(A) \dots \dots \dots$

\longleftrightarrow
no $l_j(A)$

Exercise:

- What schedules are legal?

What transactions are well-formed?

$S1 = l_1(A)l_1(B)r_1(A)w_1(B)l_2(B)u_1(A)u_1(B)$
 $r_2(B)w_2(B)u_2(B)l_3(B)r_3(B)u_3(B)$

$S2 = l_1(A)r_1(A)w_1(B)u_1(A)u_1(B)$
 $l_2(B)r_2(B)w_2(B)l_3(B)r_3(B)u_3(B)$

$S3 = l_1(A)r_1(A)u_1(A)l_1(B)w_1(B)u_1(B)$
 $l_2(B)r_2(B)w_2(B)u_2(B)l_3(B)r_3(B)u_3(B)$

Exercise:

- What schedules are legal?

What transactions are well-formed?

$S1 = l_1(A)l_1(B)r_1(A)w_1(B)l_2(B)u_1(A)u_1(B)$
 $r_2(B)w_2(B)u_2(B)l_3(B)r_3(B)u_3(B)$

$S2 = l_1(A)r_1(A)w_1(B)u_1(A)u_1(B)$
 $l_2(B)r_2(B)w_2(B)l_3(B)r_3(B)u_3(B)u_2(B)?$

$S3 = l_1(A)r_1(A)u_1(A)l_1(B)w_1(B)u_1(B)$
 $l_2(B)r_2(B)w_2(B)u_2(B)l_3(B)r_3(B)u_3(B)$

Schedule F

T1

$l_1(A); \text{Read}(A)$

$A \leftarrow A + 100; \text{Write}(A); u_1(A)$

$l_1(B); \text{Read}(B)$

$B \leftarrow B + 100; \text{Write}(B); u_1(B)$

T2

$l_2(A); \text{Read}(A)$

$A \leftarrow A \times 2; \text{Write}(A); u_2(A)$

$l_2(B); \text{Read}(B)$

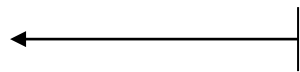
$B \leftarrow B \times 2; \text{Write}(B); u_2(B)$

Schedule F

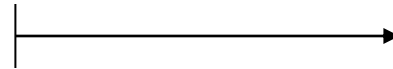
		A	B
T1	T2	25	25
l ₁ (A);Read(A)			
A ← A+100;Write(A);u ₁ (A)		125	
	l ₂ (A);Read(A)		
	A ← Ax2;Write(A);u ₂ (A)	250	
	l ₂ (B);Read(B)		
	B ← Bx2;Write(B);u ₂ (B)		50
l ₁ (B);Read(B)			
B ← B+100;Write(B);u ₁ (B)			150
		250	150

Rule #3 Two phase locking (2PL) for transactions

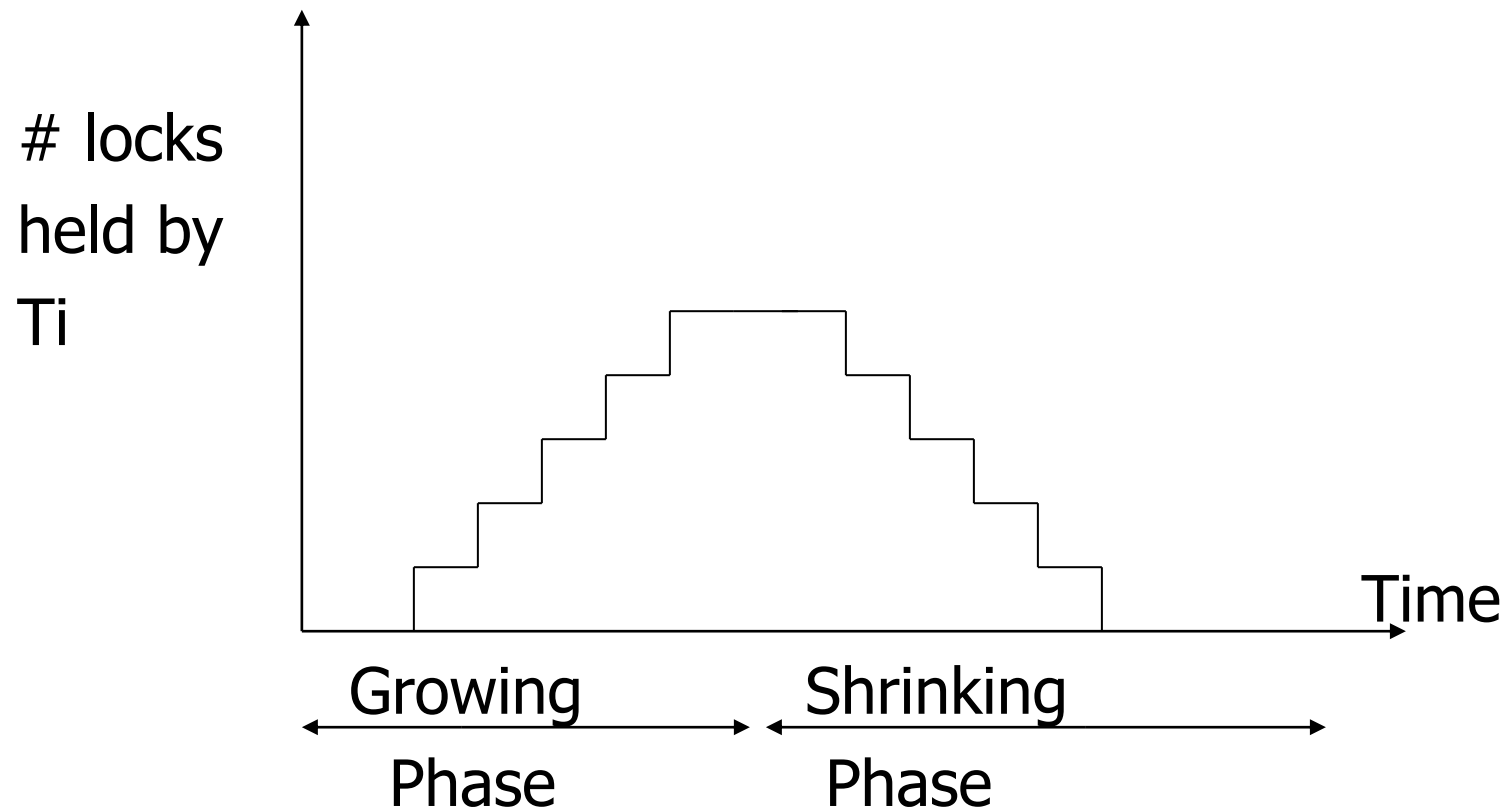
$T_i = \dots \dots \text{li}(A) \dots \dots \text{ui}(A) \dots \dots$



no unlocks



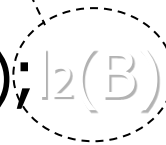
no locks



Schedule G

T1	T2
$l_1(A); \text{Read}(A)$	
$A \leftarrow A + 100; \text{Write}(A)$	
$l_1(B); u_1(A)$	
	$l_2(A); \text{Read}(A)$
	$A \leftarrow A \times 2; \text{Write}(A); l_2(B)$

delayed



Schedule G

T1	T2
$l_1(A); \text{Read}(A)$	
$A \leftarrow A + 100; \text{Write}(A)$	
$l_1(B); u_1(A)$	
	$l_2(A); \text{Read}(A)$
	$A \leftarrow A \times 2; \text{Write}(A); l_2(B)$
$\text{Read}(B); B \leftarrow B + 100$	
$\text{Write}(B); u_1(B)$	

delayed



Schedule G

T1

$l_1(A); \text{Read}(A)$

$A \leftarrow A + 100; \text{Write}(A)$

$l_1(B); u_1(A)$

$\text{Read}(B); B \leftarrow B + 100$

$\text{Write}(B); u_1(B)$

T2

$l_2(A); \text{Read}(A)$

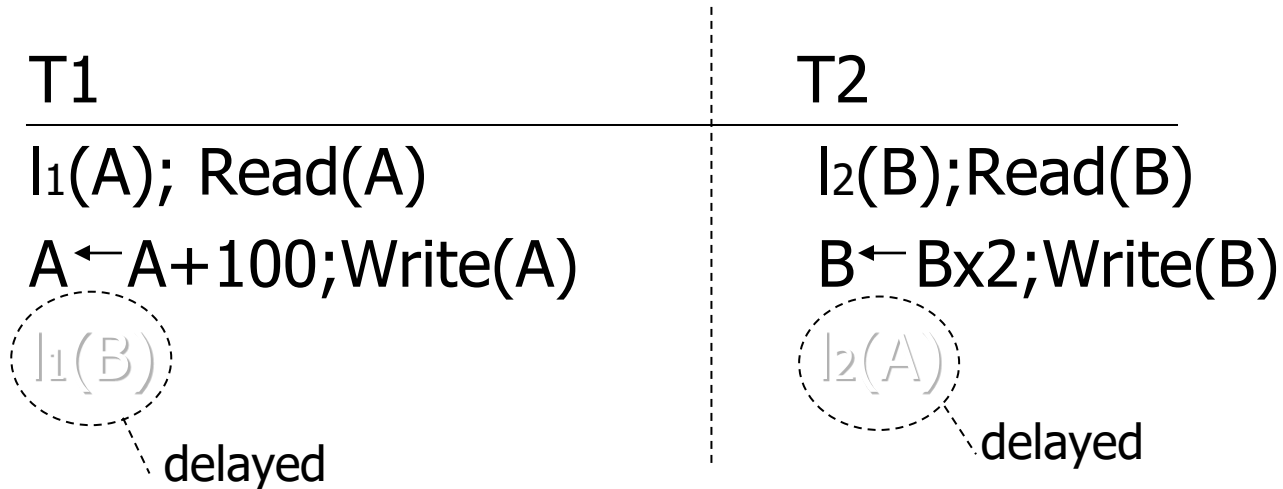
$A \leftarrow A \times 2; \text{Write}(A); l_2(B)$

delayed

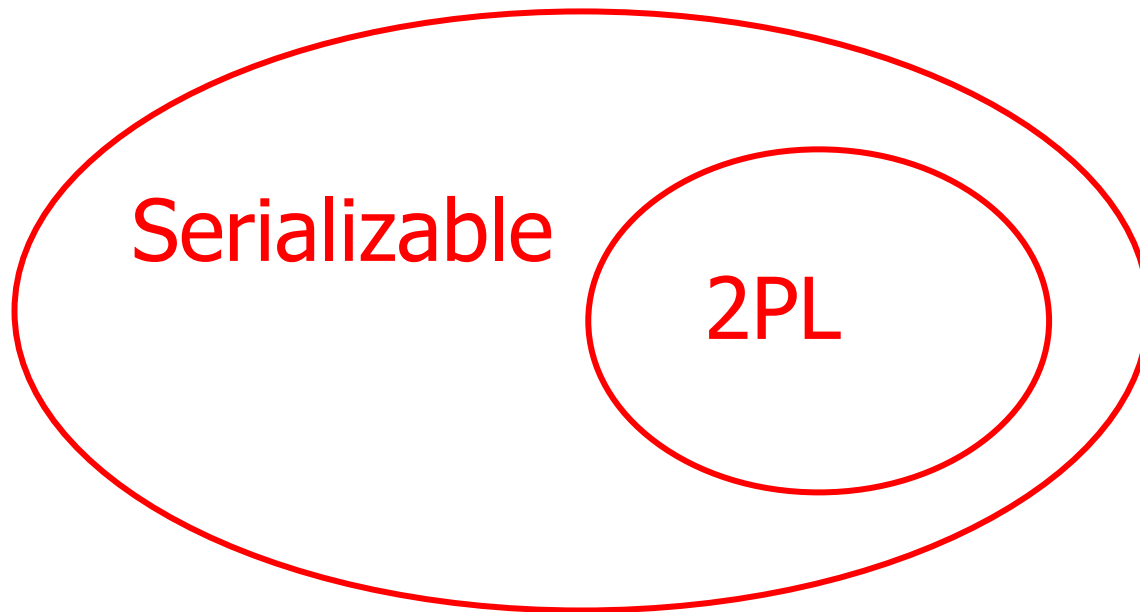
$l_2(B); u_2(A); \text{Read}(B)$

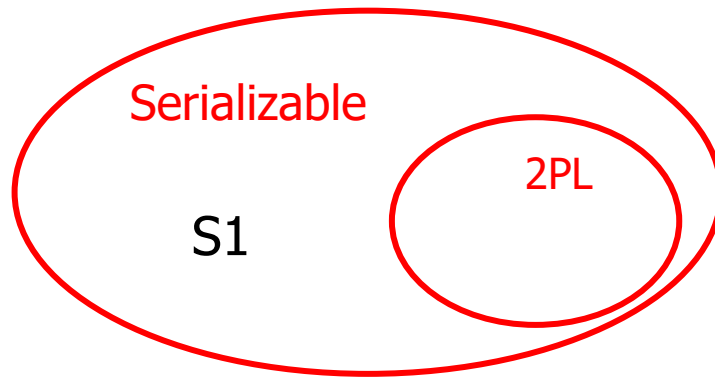
$B \leftarrow B \times 2; \text{Write}(B); u_2(B);$

Schedule H (T₂ reversed)



2PL subset of Serializable





S1: w1(x) w3(x) w2(y) w1(y)

S1: w1(x) w3(x) w2(y) w1(y)

- S1 cannot be achieved via 2PL:
The lock by T1 for y must occur after w2(y), so the unlock by T1 for x must occur after this point (and before w1(x)). Thus, w3(x) cannot occur under 2PL where shown in S1 because T1 holds the x lock at that point.
- However, S1 is serializable (equivalent to T2, T1, T3).

If you need a bit more practice:

Are our schedules S_C and S_D 2PL schedules?

S_C : $w1(A)$ $w2(A)$ $w1(B)$ $w2(B)$

S_D : $w1(A)$ $w2(A)$ $w2(B)$ $w1(B)$

- Beyond this simple 2PL protocol, it is all a matter of improving performance and allowing more concurrency....
 - Shared locks
 - Multiple granularity
 - Inserts, deletes and phantoms
 - Other types of C.C. mechanisms

Shared locks

So far:

$S = \dots l_1(A) \ r_1(A) \ u_1(A) \ \dots \ l_2(A) \ r_2(A) \ u_2(A) \ \dots$



Do not conflict

Shared locks

So far:

$S = \dots l_1(A) \ r_1(A) \ u_1(A) \ \dots \ l_2(A) \ r_2(A) \ u_2(A) \ \dots$



Do not conflict

Instead:

$S = \dots l_{s1}(A) \ r_1(A) \ l_{s2}(A) \ r_2(A) \ \dots \ u_{s1}(A) \ u_{s2}(A)$

Lock actions (Shared, Exclusive)

$l-t_i(A)$: lock A in t mode (t is S or X)

$u-t_i(A)$: unlock t mode (t is S or X)

Shorthand:

$u_i(A)$: unlock whatever modes

T_i has locked A

Rule #1 Well formed transactions

$T_i = \dots I-S_1(A) \dots r_1(A) \dots u_1(A) \dots$

$T_i = \dots I-X_1(A) \dots w_1(A) \dots u_1(A) \dots$

- What about transactions that read and write same object?

Option 1: Request exclusive lock

$T_i = \dots l-X_1(A) \dots r_1(A) \dots w_1(A) \dots u(A) \dots$

- What about transactions that read and write same object?

Option 2: Upgrade

(E.g., need to read, but don't know if will write...)

$T_i = \dots I-S_1(A) \dots r_1(A) \dots I-X_1(A) \dots w_1(A) \dots u(A) \dots$



Think of

- Get 2nd lock on A, or
- Drop S, get X lock

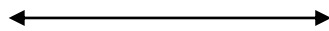
Rule #2 Legal scheduler

$$S = \dots I-S_i(A) \dots \dots u_i(A) \dots$$



no $I-X_j(A)$

$$S = \dots I-X_i(A) \dots \dots u_i(A) \dots$$



no $I-X_j(A)$

no $I-S_j(A)$

A way to summarize Rule #2

Compatibility matrix

Comp

	S	X
S	true	false
X	false	false

Rule # 3 2PL transactions

No change except for upgrades:

- (I) If upgrade gets more locks
(e.g., $S \rightarrow \{S, X\}$) then no change!
- (II) If upgrade releases read (shared)
lock (e.g., $S \rightarrow X$)
 - can be allowed in growing phase

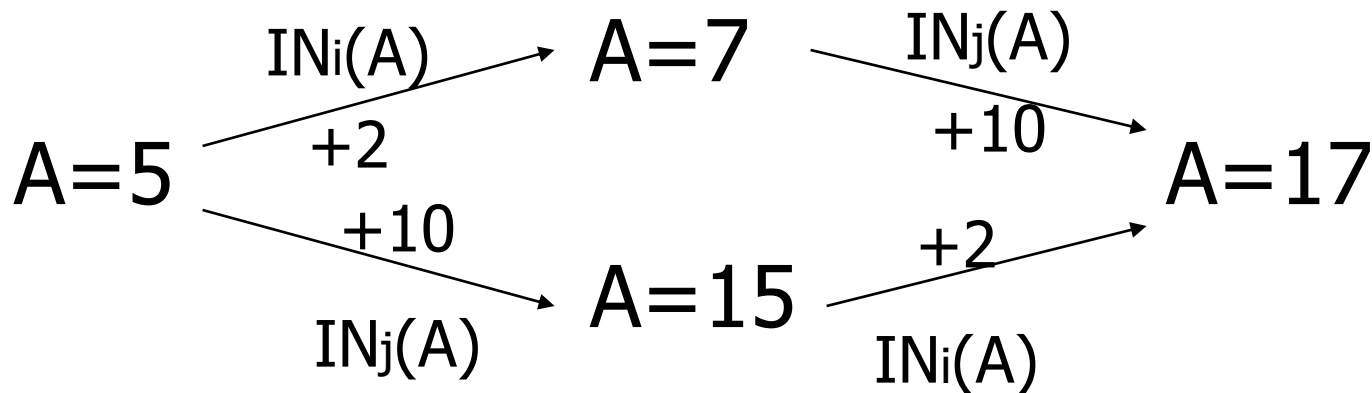
Lock types beyond S/X

Examples:

- (1) increment lock
- (2) update lock

Example (1): increment lock

- Atomic increment action: $\text{IN}_i(A)$
 $\{\text{Read}(A); A \leftarrow A+k; \text{Write}(A)\}$
- $\text{IN}_i(A), \text{IN}_j(A)$ do not conflict!



Comp

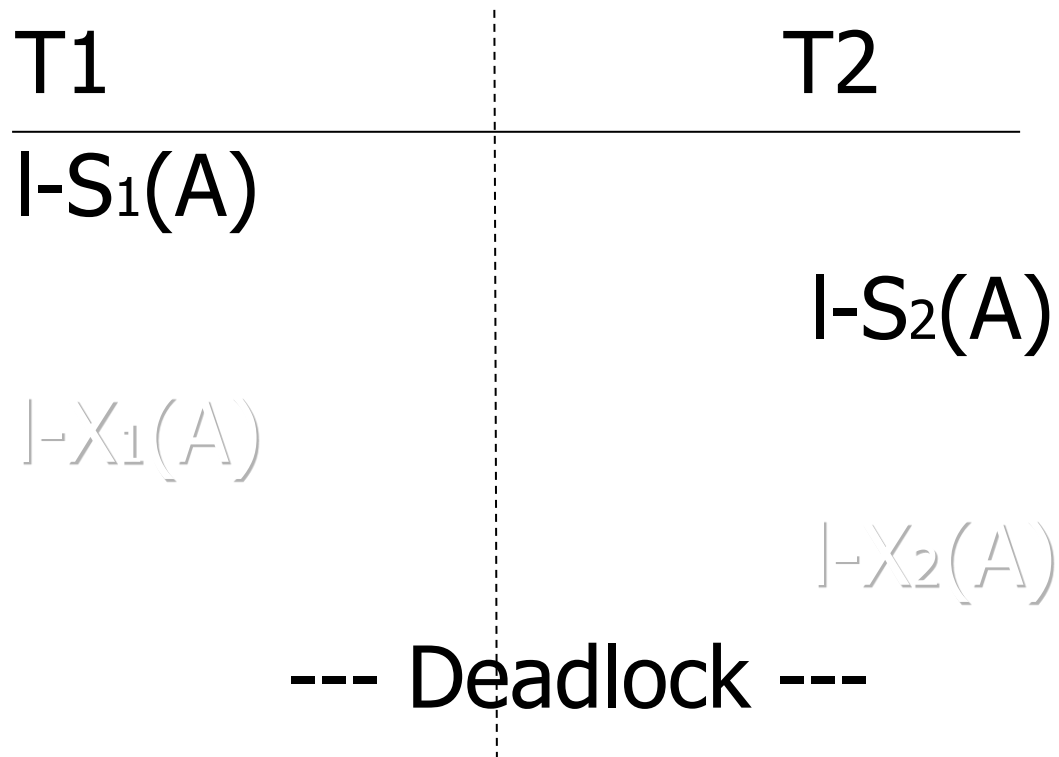
	S	X	I
S			
X			
I			

Comp

	S	X	I
S	T	F	F
X	F	F	F
I	F	F	T

Update locks

A common deadlock problem with upgrades:



Solution

If T_i wants to read A and knows it may later want to write A , it requests update lock (not shared)

New request

Comp

Lock
already
held in

	S	X	U
S			
X			
U			

New request

Comp

Lock
already
held in

	S	X	U
S	T	F	T
X	F	F	F
U	F	F	F

Note: object A may be locked in different modes at the same time...

$$S_1 = \dots I-S_1(A) \dots I-S_2(A) \dots I-U_3(A) \dots \left\{ \begin{array}{l} I-S_4(A) \dots ? \\ I-U_4(A) \dots ? \end{array} \right.$$

Note: object A may be locked in different modes at the same time...

$$S_1 = \dots I-S_1(A) \dots I-S_2(A) \dots I-U_3(A) \dots \left\{ \begin{array}{l} I-S_4(A) \dots ? \\ I-U_4(A) \dots ? \end{array} \right.$$

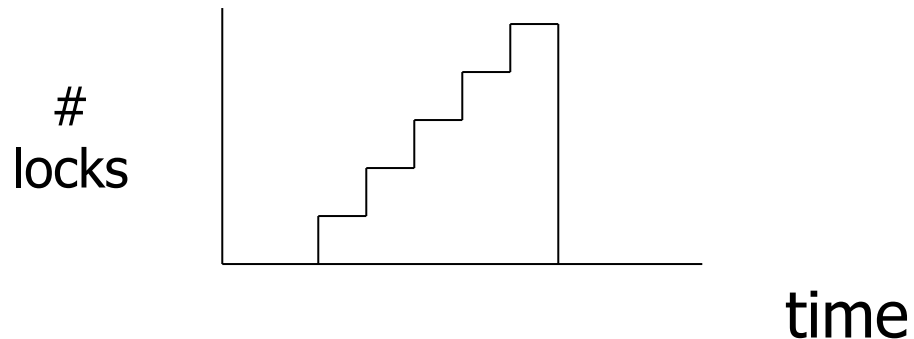
- To grant a lock in mode t, mode t must be compatible with all currently held locks on object

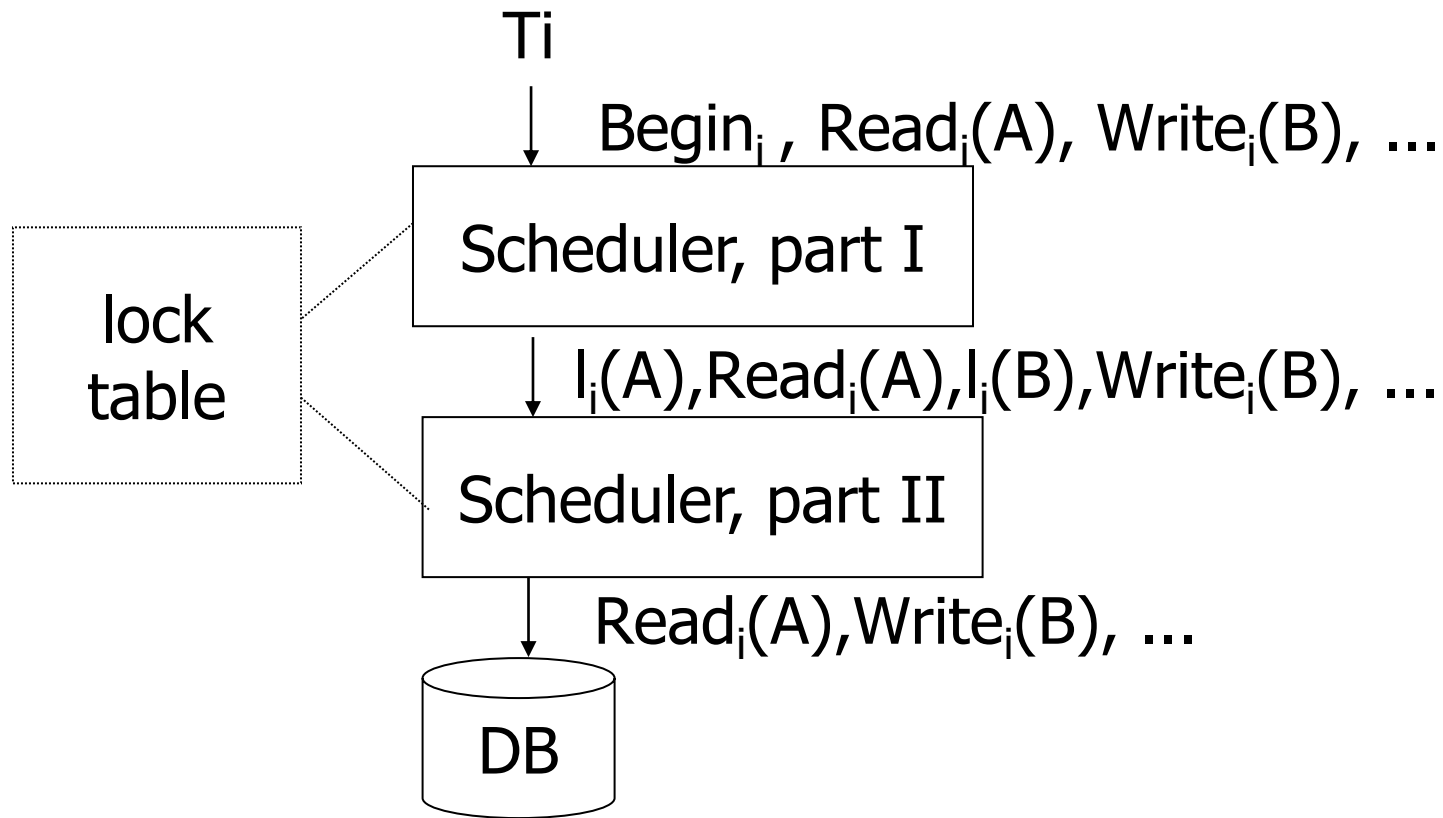
How does locking work in practice?

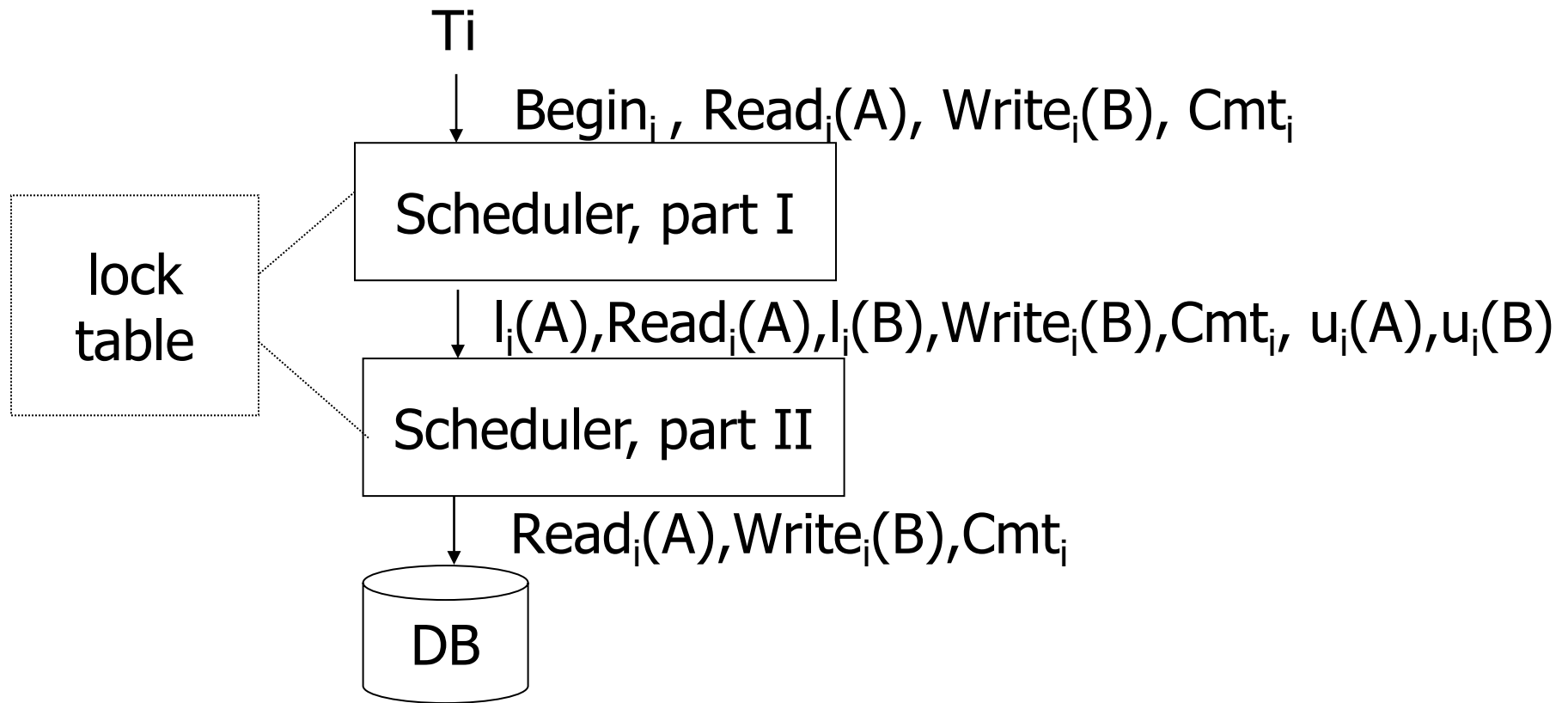
- Every system is different
(E.g., may not even provide
CONFLICT-SERIALIZABLE schedules)
- But here is one (simplified) way ...

Sample Locking System:

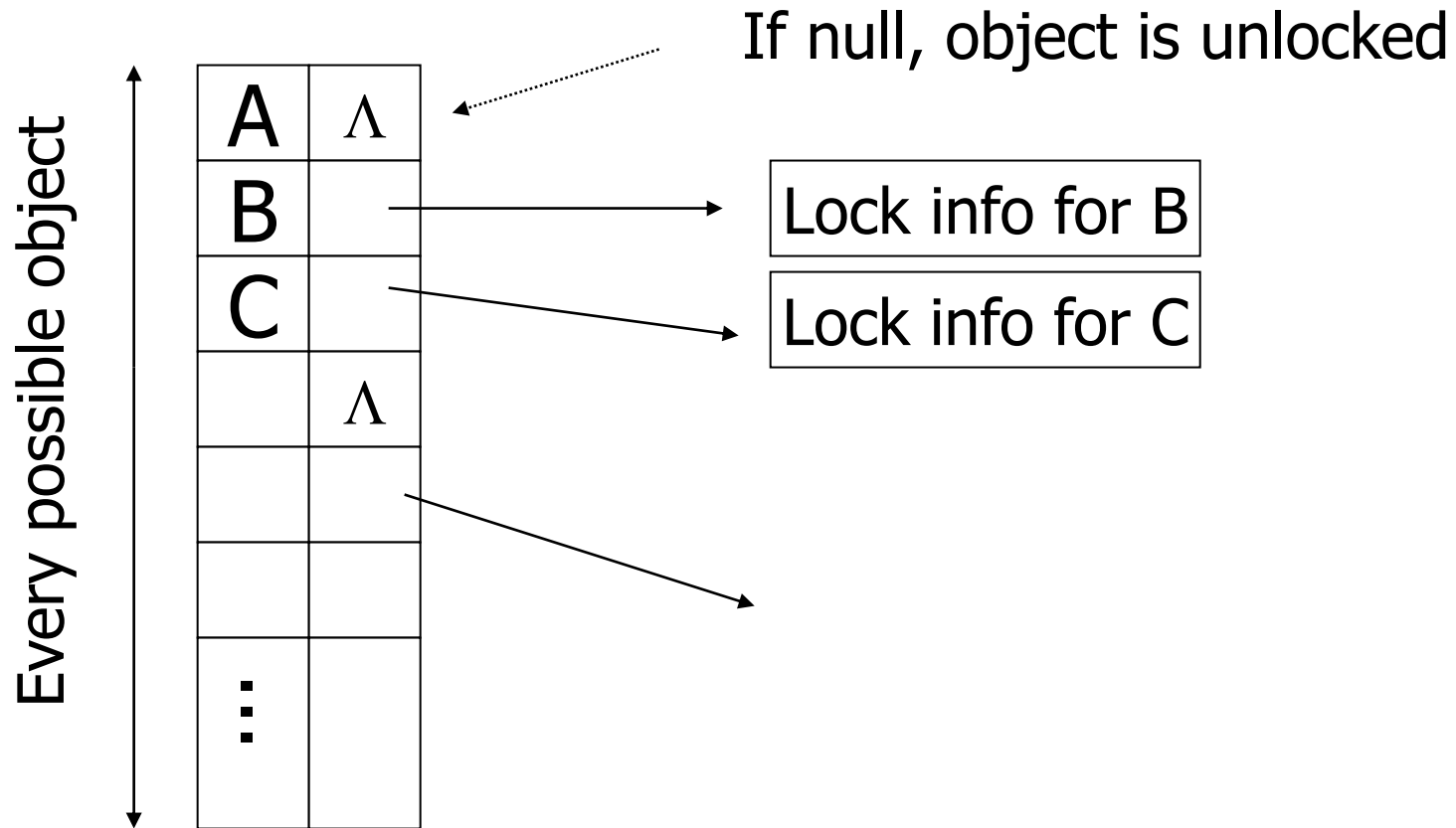
- (1) Don't trust transactions to request/release locks
- (2) Hold all locks until transaction commits



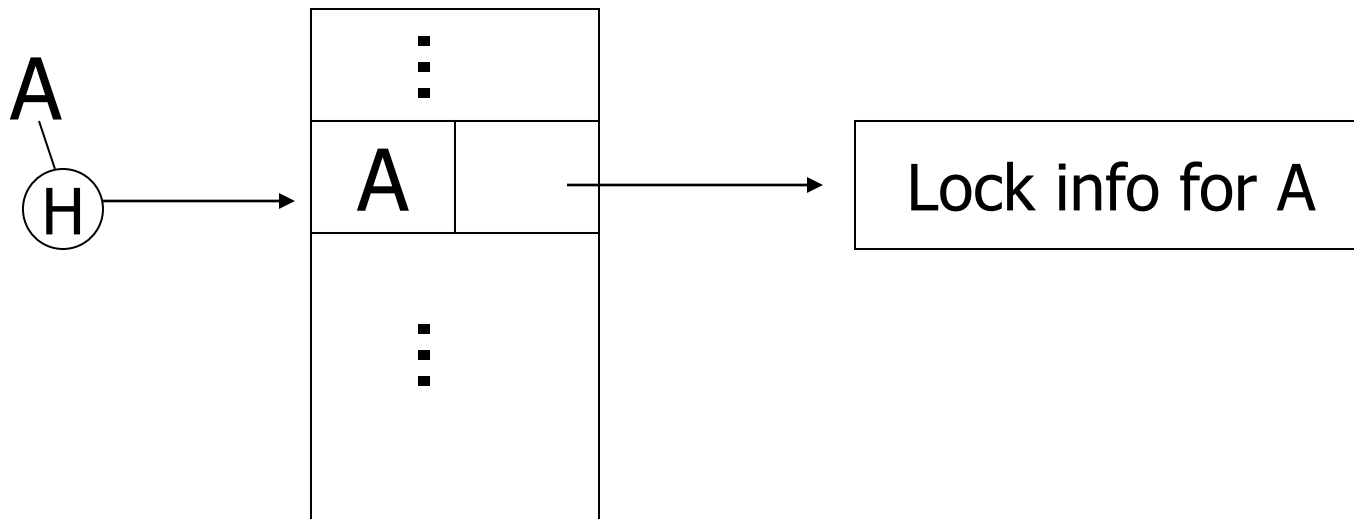




Lock table Conceptually

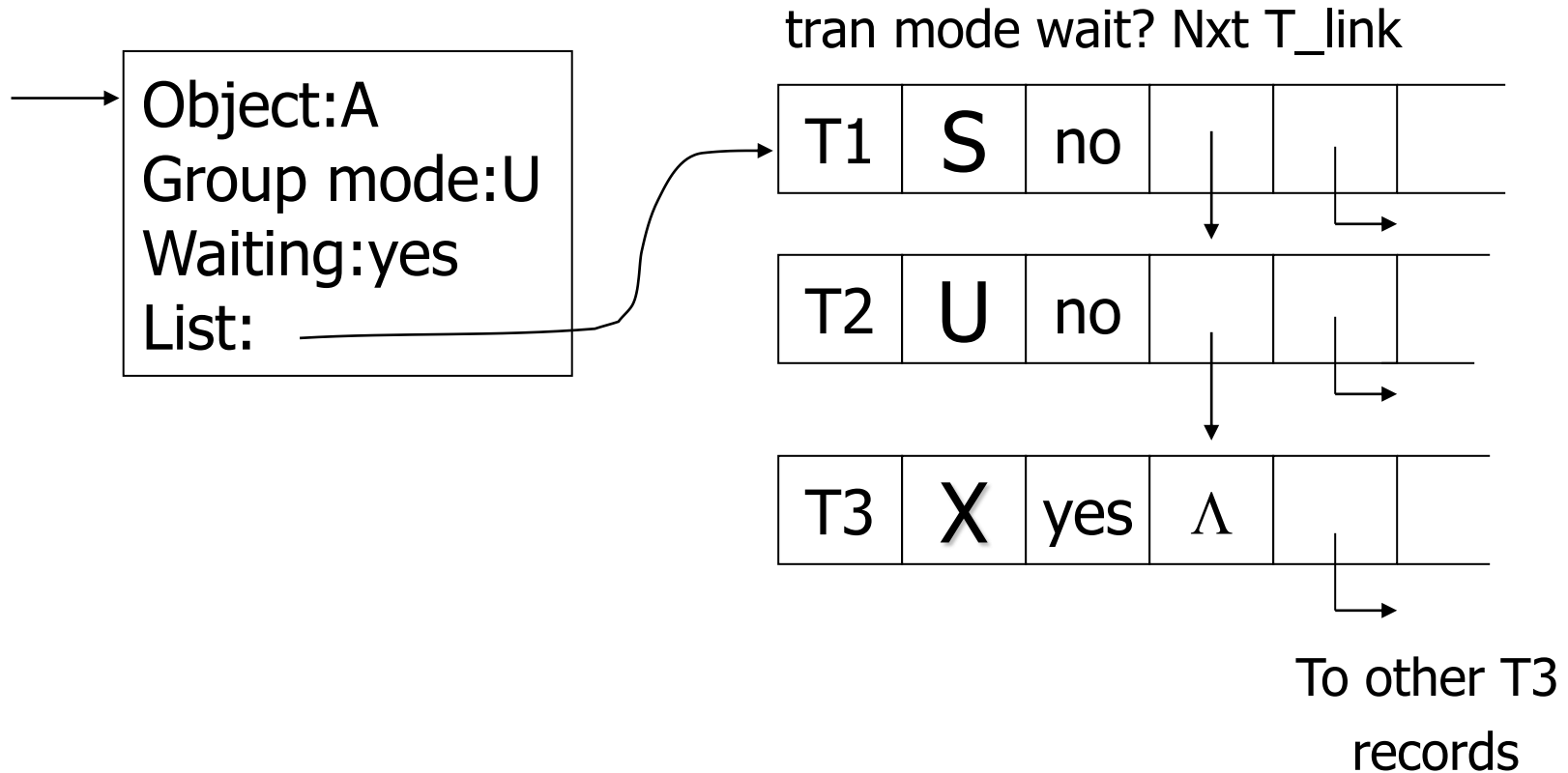


But use hash table:

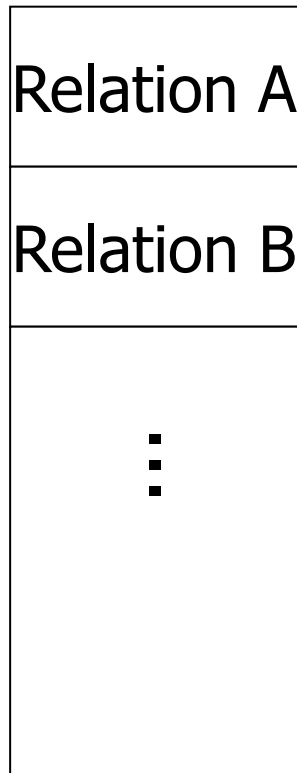


If object not found in hash table, it is unlocked

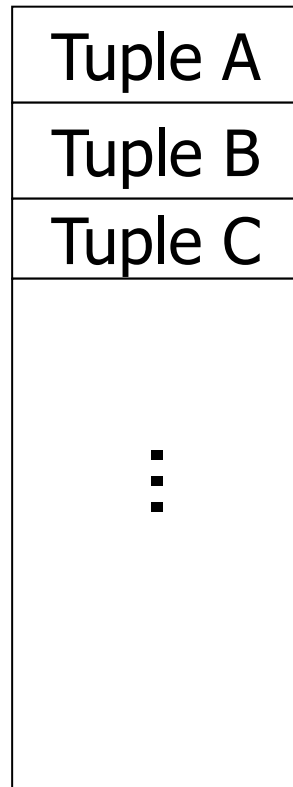
Lock info for A - example



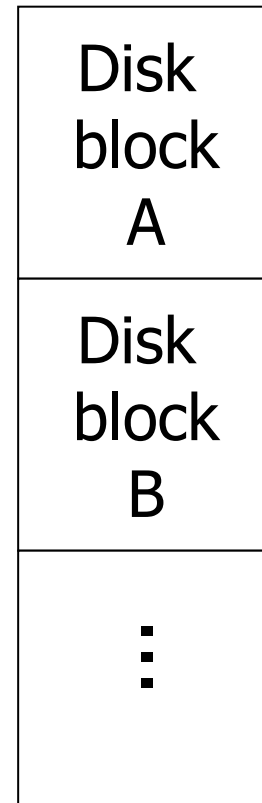
What are the objects we lock?



DB



DB



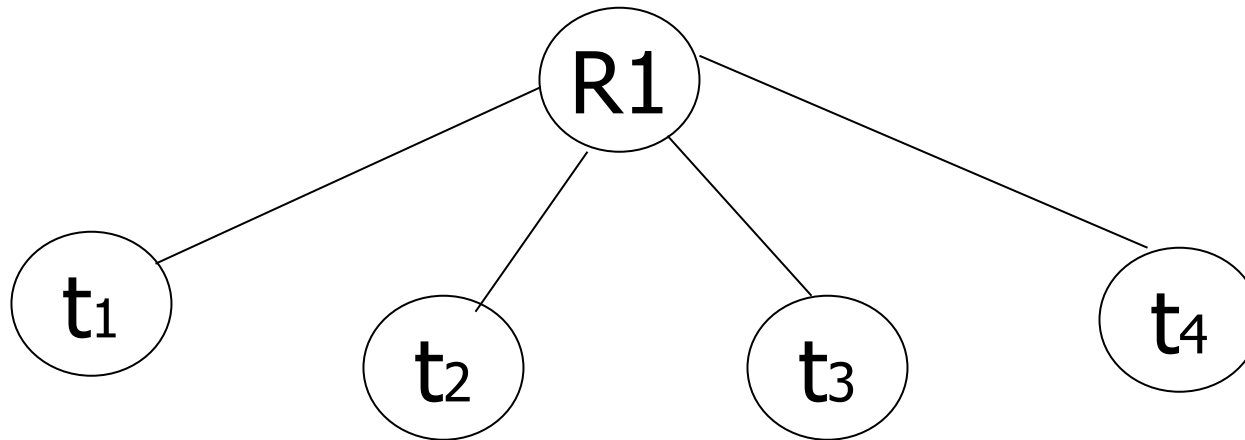
DB

?

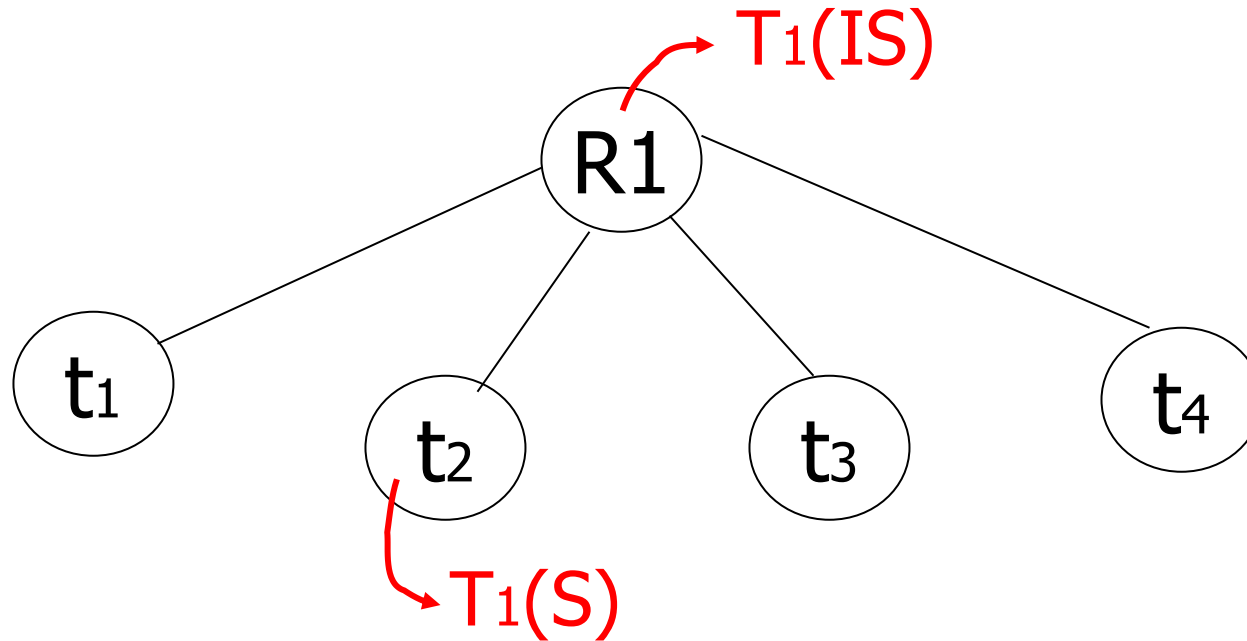
- Locking works in any case, but should we choose small or large objects?

- Locking works in any case, but should we choose small or large objects?
- If we lock large objects (e.g., Relations)
 - Need few locks
 - Low concurrency
- If we lock small objects (e.g., tuples, fields)
 - Need more locks
 - More concurrency

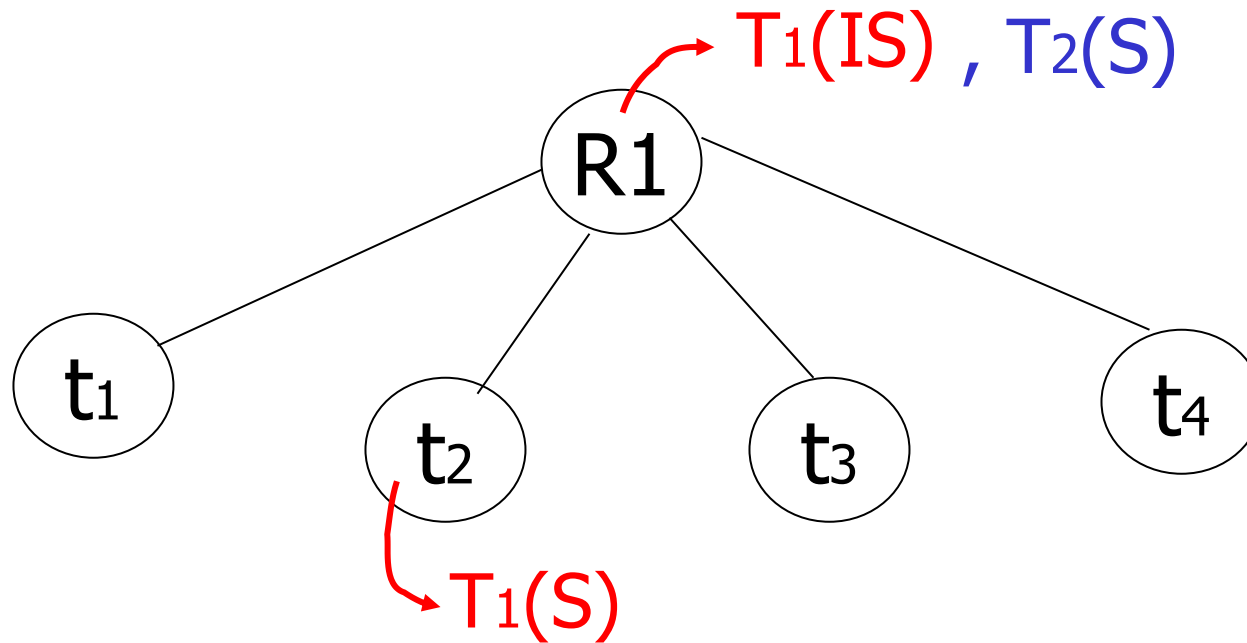
Example



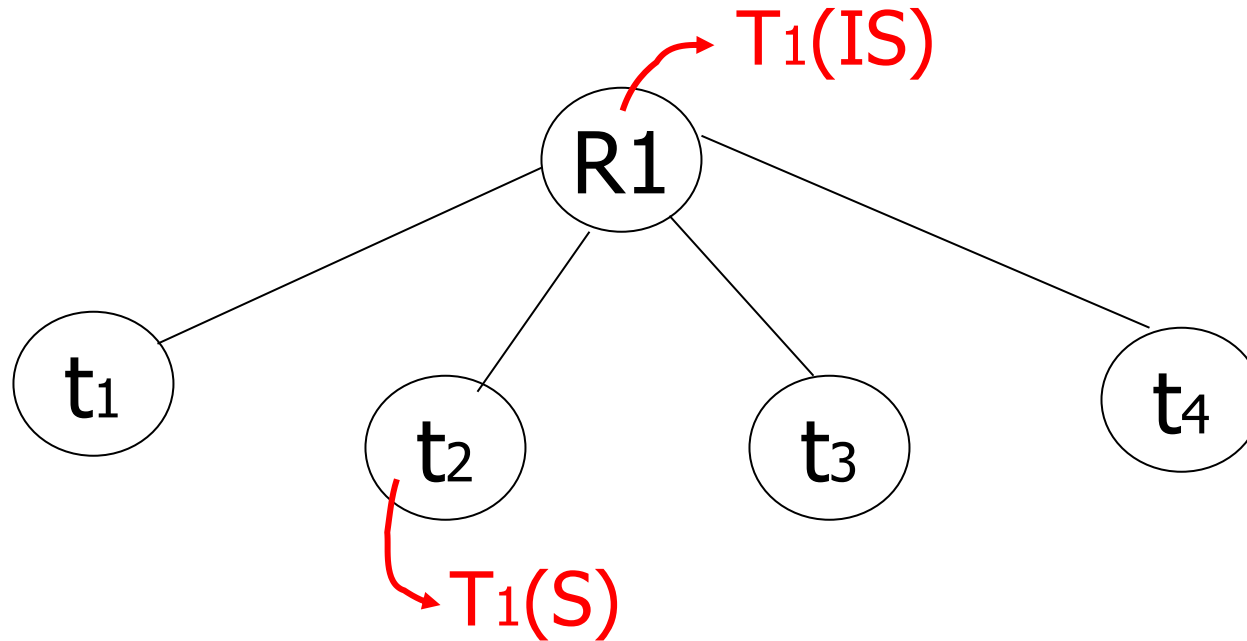
Example



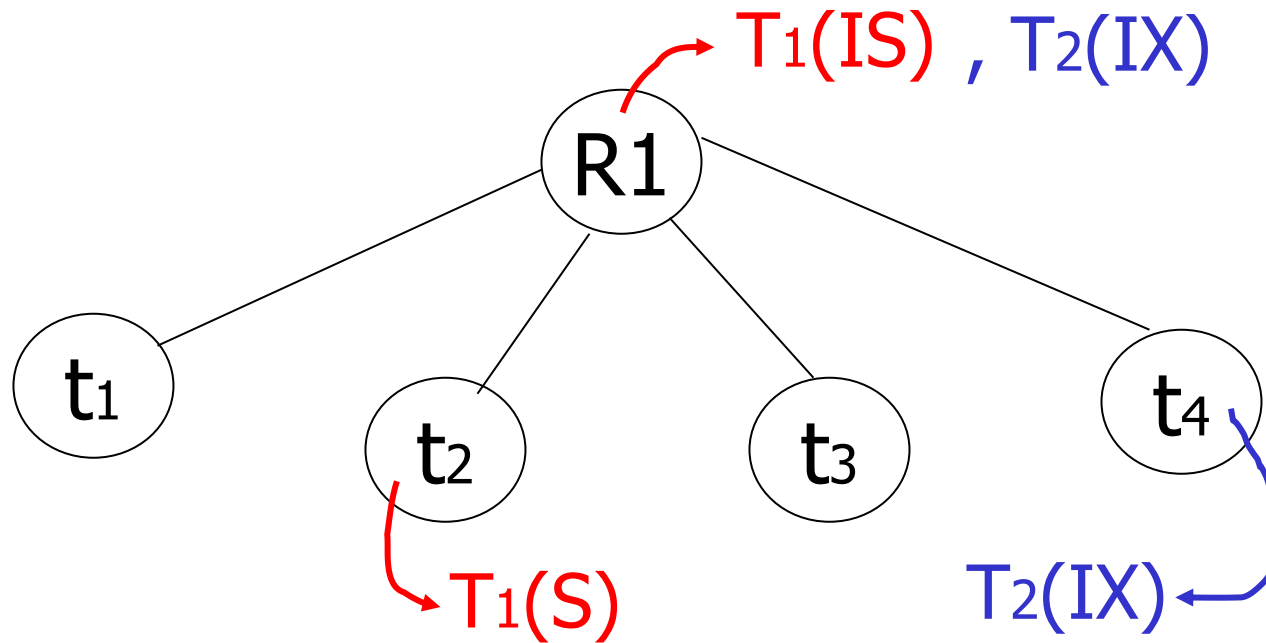
Example



Example (b)



Example



Multiple granularity

Comp

Requestor

		IS	IX	S	SIX	X
Holder	IS					
	IX					
	S					
	SIX					
	X					

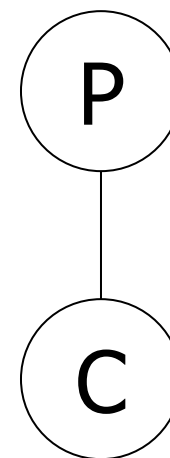
Multiple granularity

Comp

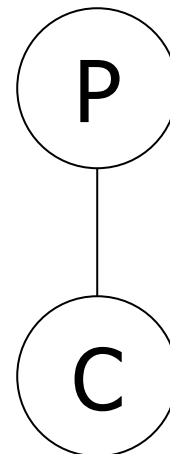
Requestor

		IS	IX	S	SIX	X
Holder	IS	T	T	T	T	F
	IX	T	T	F	F	F
	S	T	F	T	F	F
	SIX	T	F	F	F	F
	X	F	F	F	F	F

Parent locked in	Child can be locked in
IS	
IX	
S	
SIX	
X	



Parent locked in	Child can be locked by same transaction in
IS	IS, S
IX	IS, S, IX, X, SIX
S	none
SIX	X, IX, [SIX]
X	none



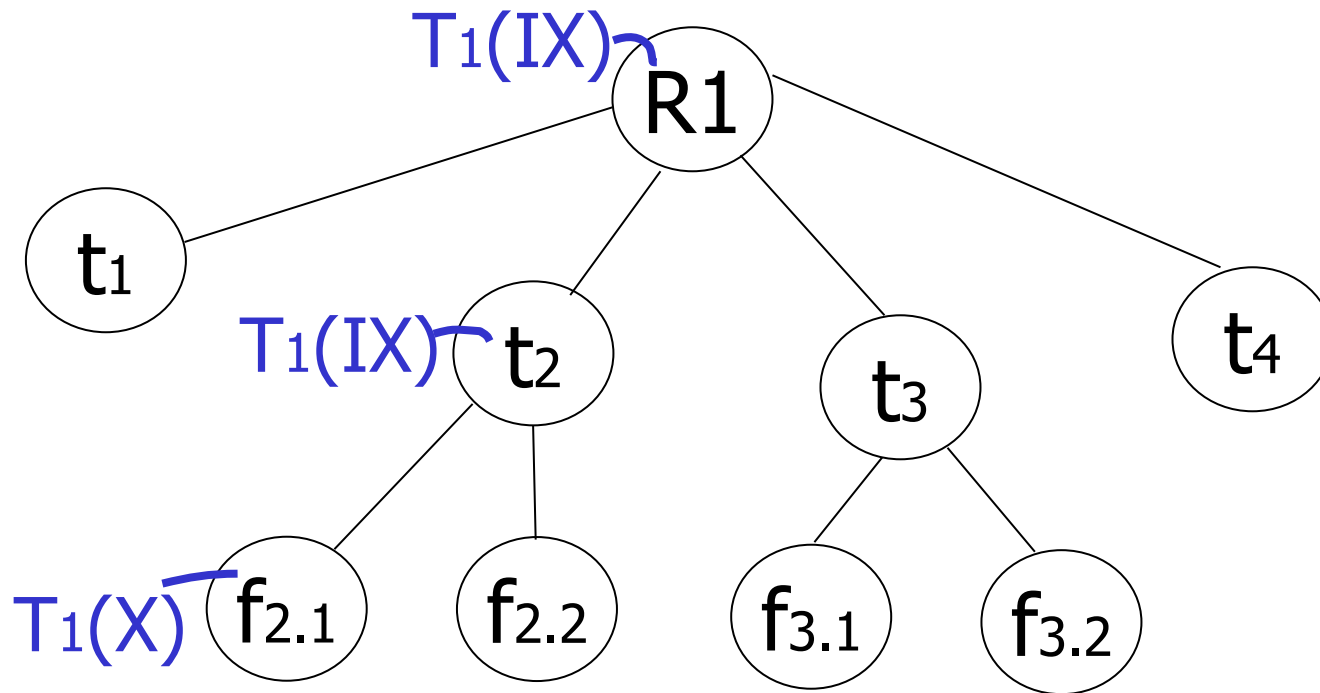
not necessary

Rules

- (1) Follow multiple granularity comp function
- (2) Lock root of tree first, any mode
- (3) Node Q can be locked by Ti in S or IS only if
parent(Q) locked by Ti in IX or IS
- (4) Node Q can be locked by Ti in X,SIX,IX only
if parent(Q) locked by Ti in IX,SIX
- (5) Ti is two-phase
- (6) Ti can unlock node Q only if none of Q's
children are locked by Ti

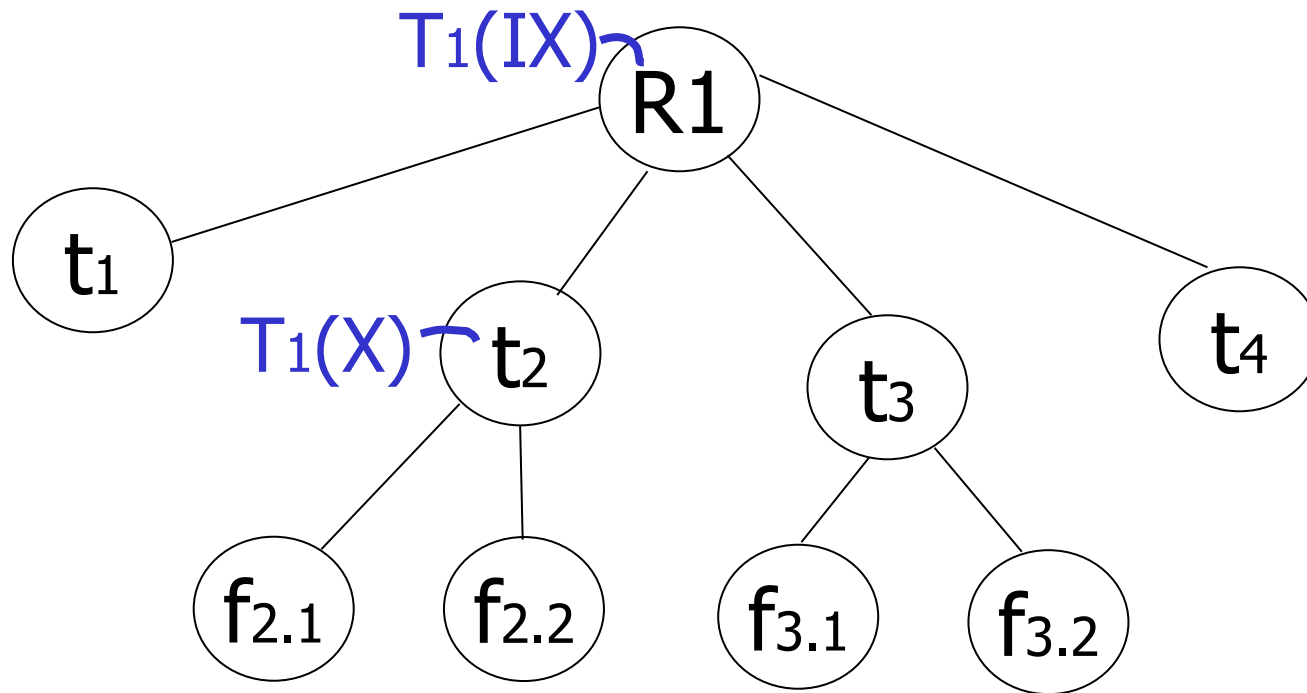
Exercise:

- Can T2 access object f2.2 in X mode?
What locks will T2 get?



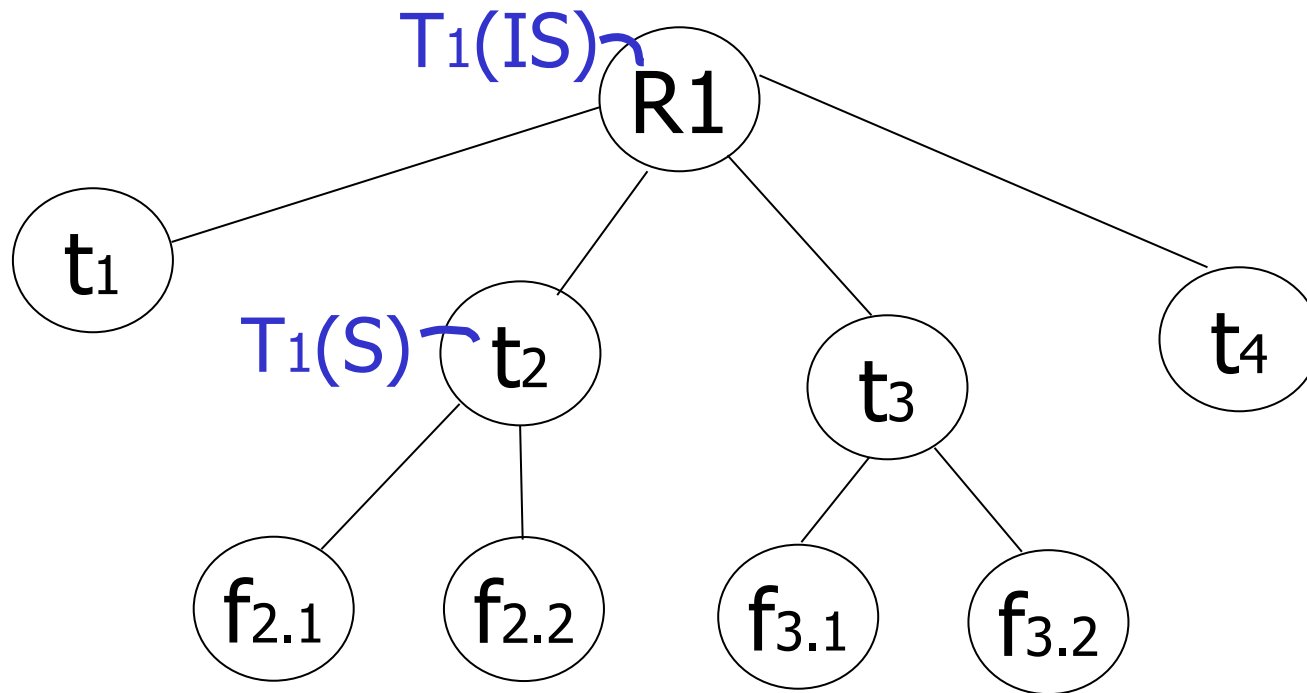
Exercise:

- Can T2 access object f2.2 in X mode?
What locks will T2 get?



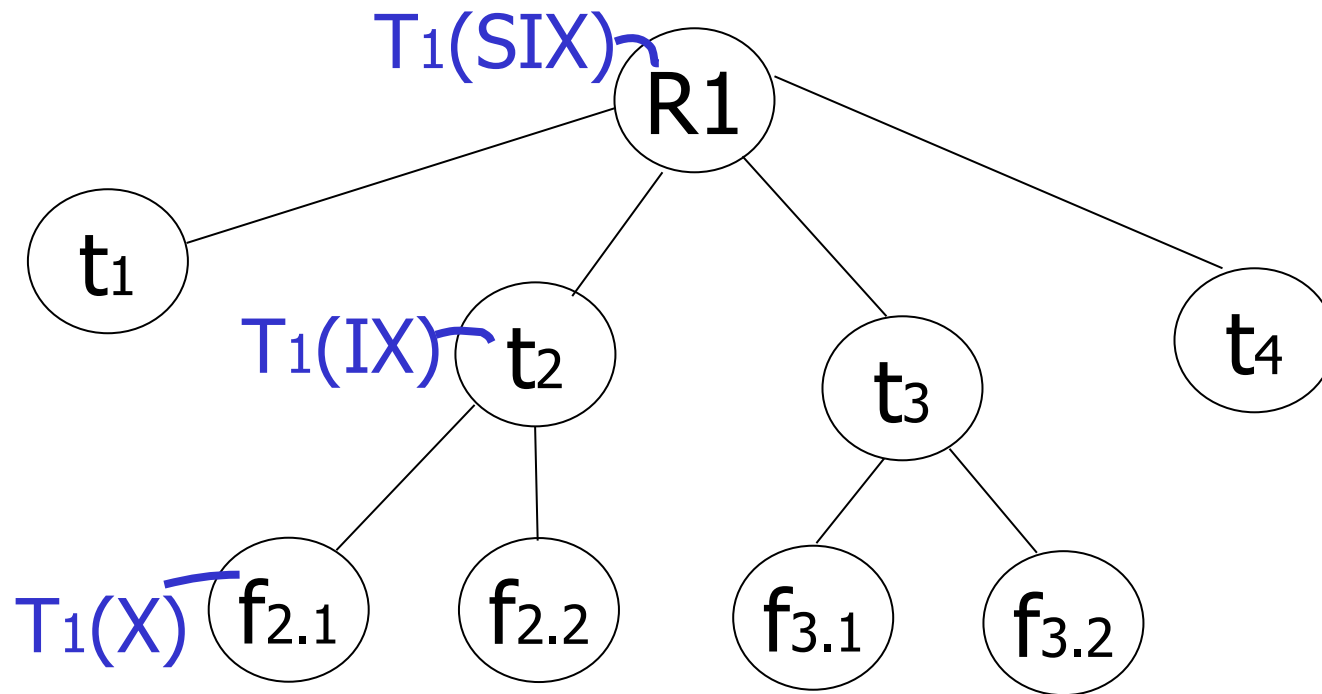
Exercise:

- Can T2 access object f3.1 in X mode?
What locks will T2 get?



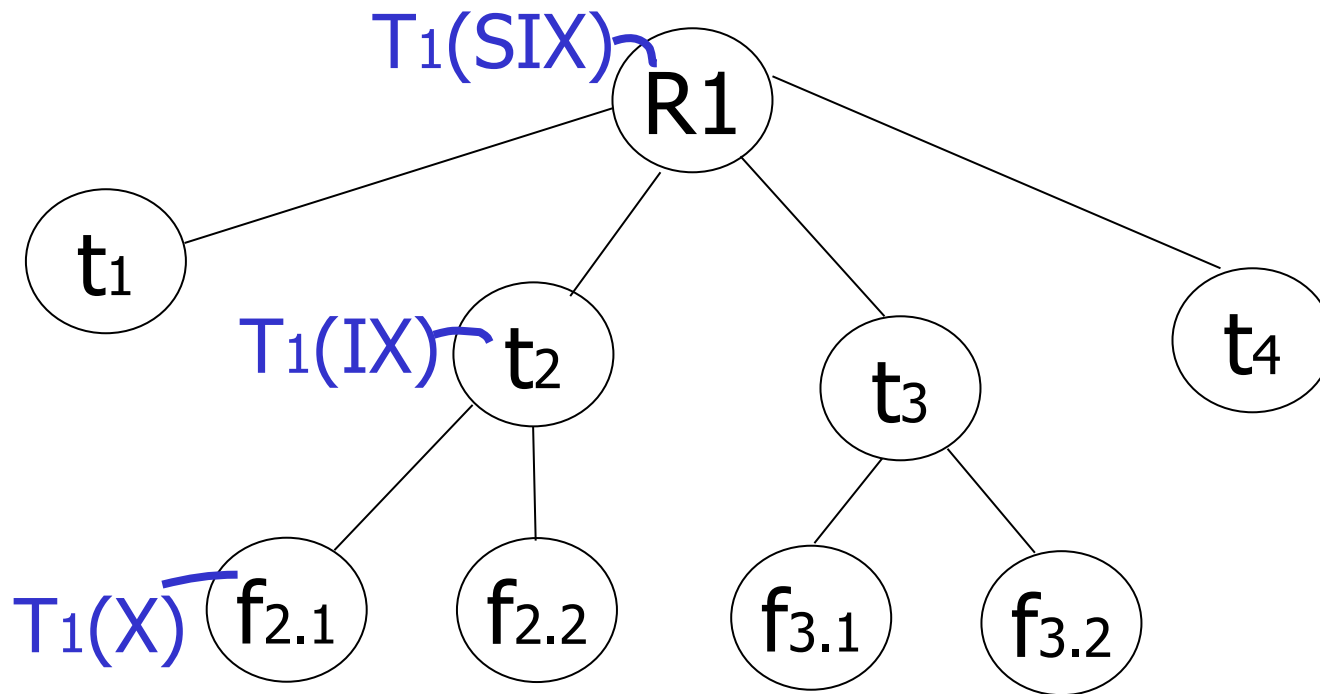
Exercise:

- Can T2 access object f2.2 in S mode?
What locks will T2 get?



Exercise:

- Can T2 access object f2.2 in X mode?
What locks will T2 get?



Reading

Ch18 Concurrency Control

Héctor García-Molina, Jeffrey Ullman, and
Jennifer Widom. Database Systems:
The Complete Book.