



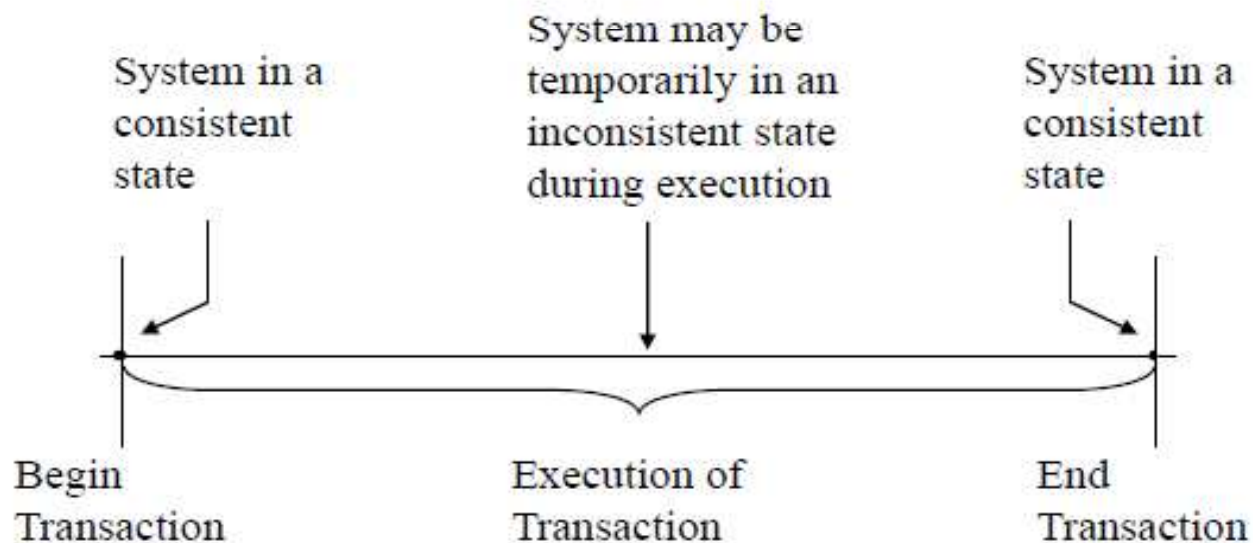
# Distributed Transactions

From  
George  
Coulouris  
Material

# Transaction

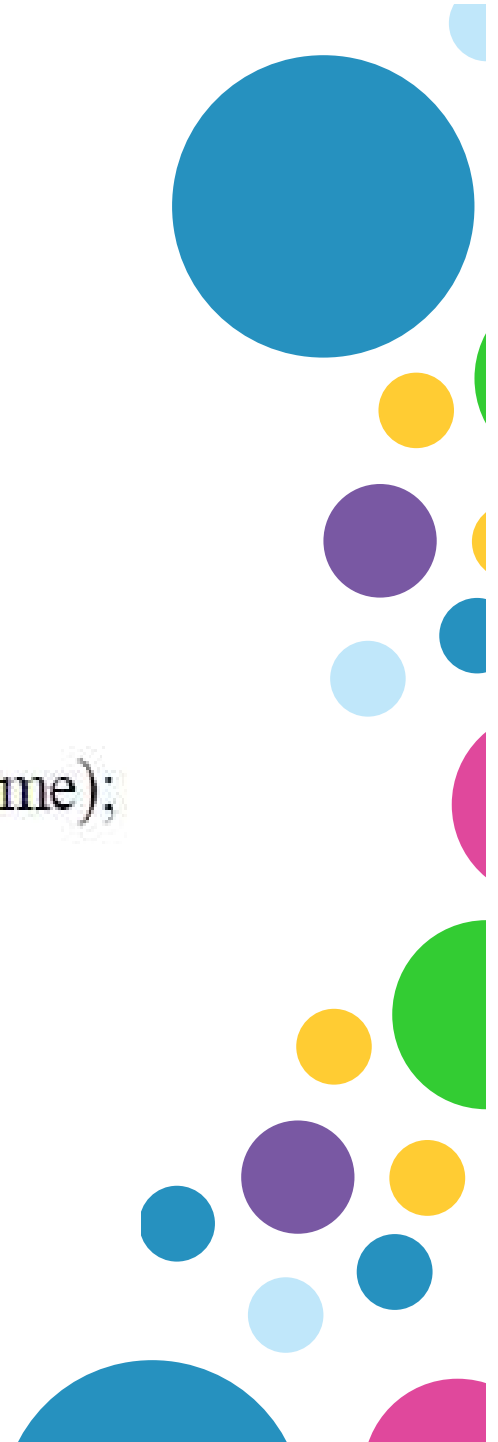
A transaction is a collection of actions that make consistent transformations of system states while preserving system consistency.

- concurrency transparency
- failure transparency



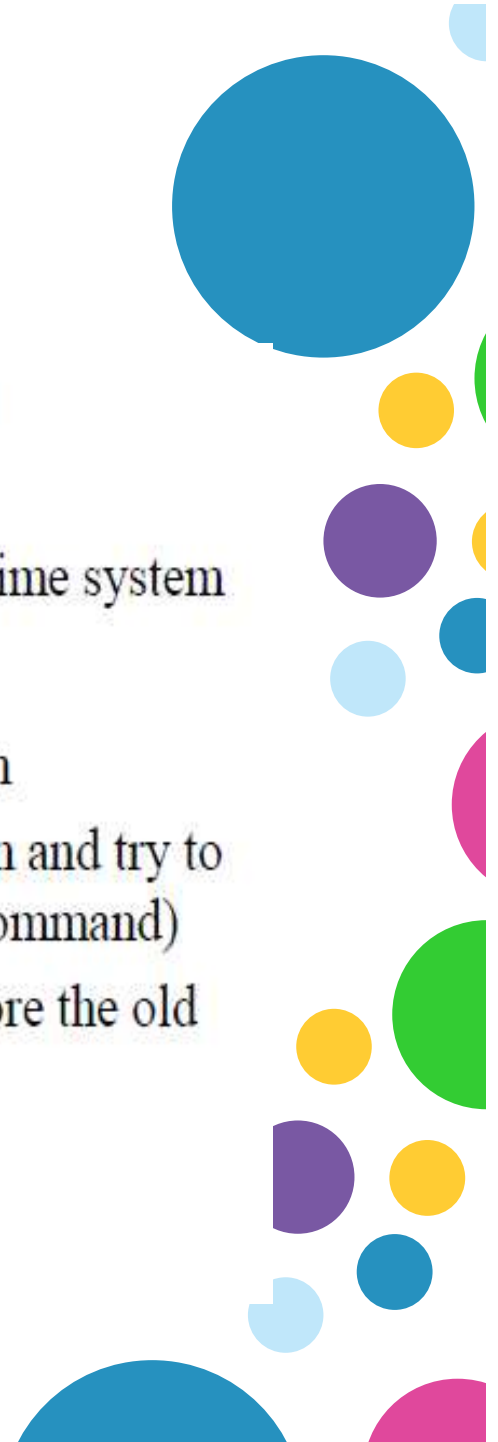
# Example of Transaction

```
begin
  input(flight_no, date, customer_name);
  Begin_transaction Reservation
  begin
    Write(flight(date).stsold++);
    Write(flight(date).cname, customer_name);
    Commit
  end. {Reservation}
  output("reservation completed")
...
end
```



# Transaction Primitives

- Special primitives required for programming using transactions
  - Supplied by the operating system or by the language runtime system
- Examples of transaction primitives:
  - **BEGIN\_TRANSACTION**: Mark the start of a transaction
  - **END\_TRANSACTION (EOT)**: Terminate the transaction and try to commit (there may or may not be a separate **COMMIT** command)
  - **ABORT\_TRANSACTION**: Kill the transaction and restore the old values
  - **READ**: Read data from a file (or other object)
  - **WRITE**: Write data to a file (or other object)



# Operations of the *Account* interface

***deposit(amount)***

deposit amount in the account

***withdraw(amount)***

withdraw amount from the account

***getBalance()* -> amount**

return the balance of the account

***setBalance(amount)***

set the balance of the account to amount

---

## Operations of the **Branch** interface

***create(name)* -> account**

create a new account with a given name

***lookUp(name)* -> account**

return a reference to the account with the given name

***branchTotal()* -> amount**

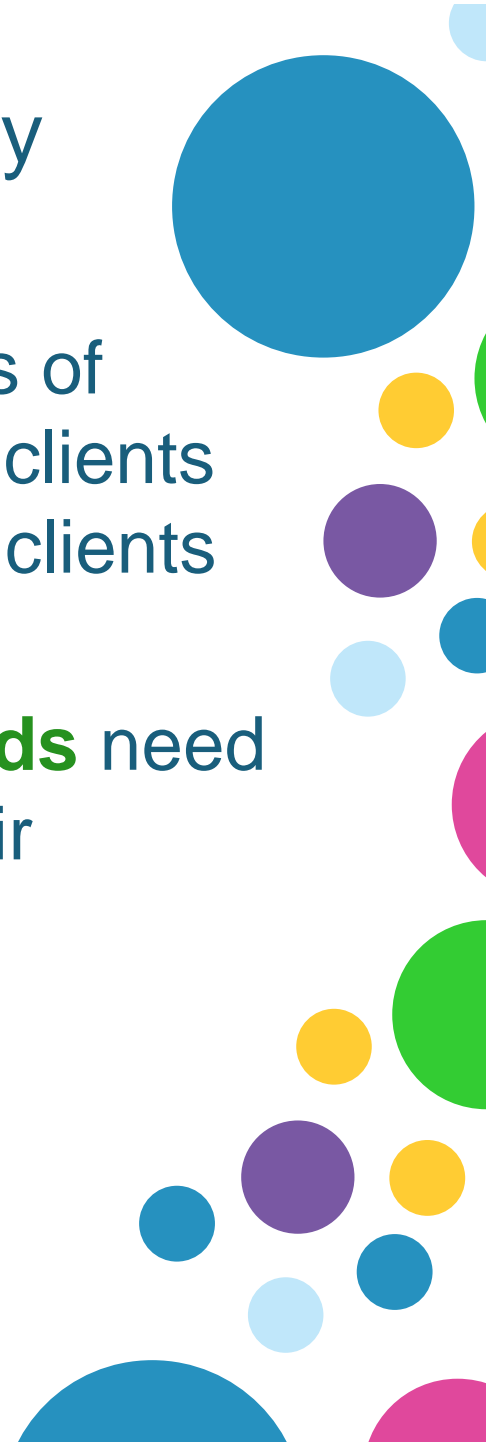
return the total of all the balances at the branch

---

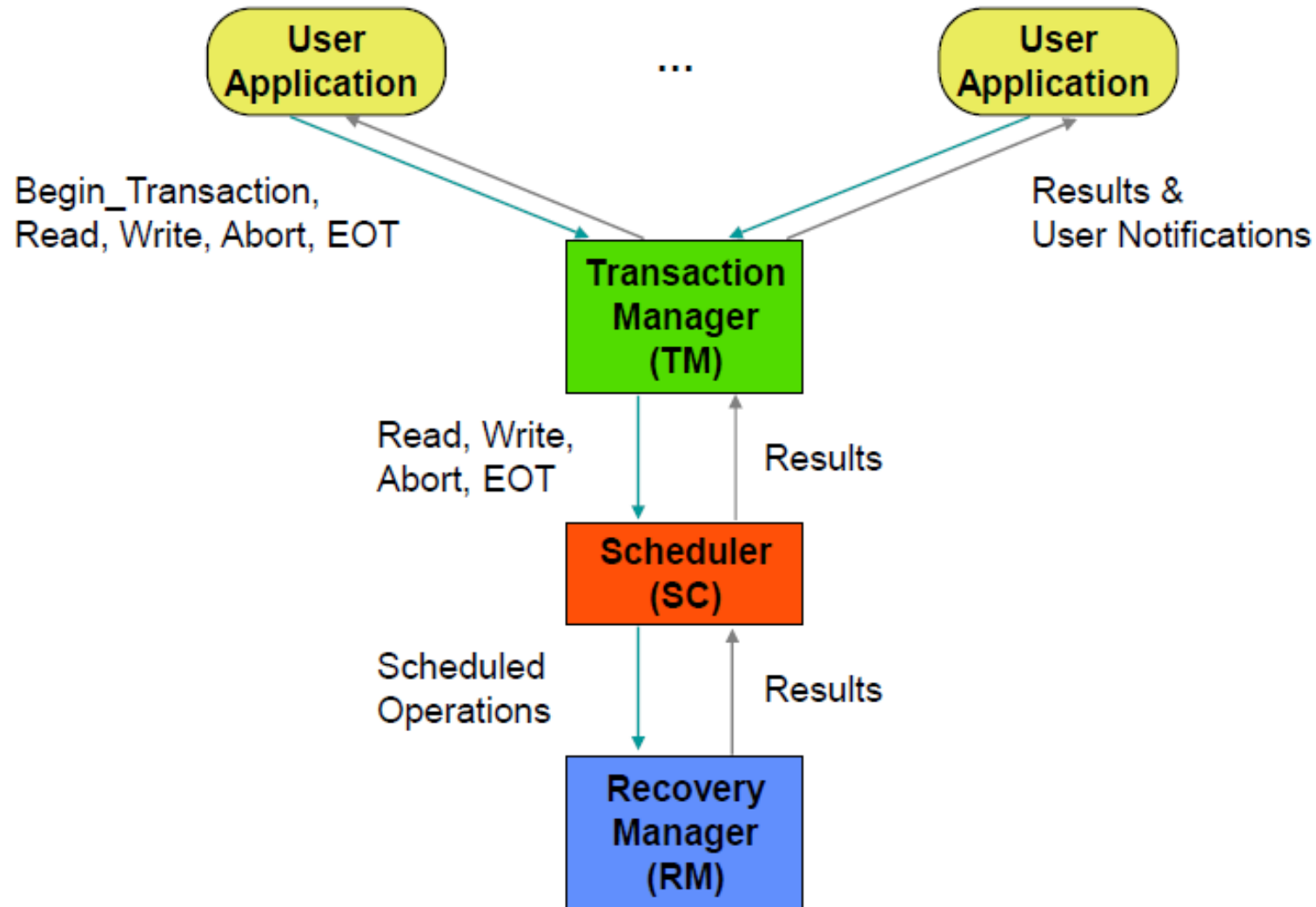


# Enhancing Client Cooperation by Signaling

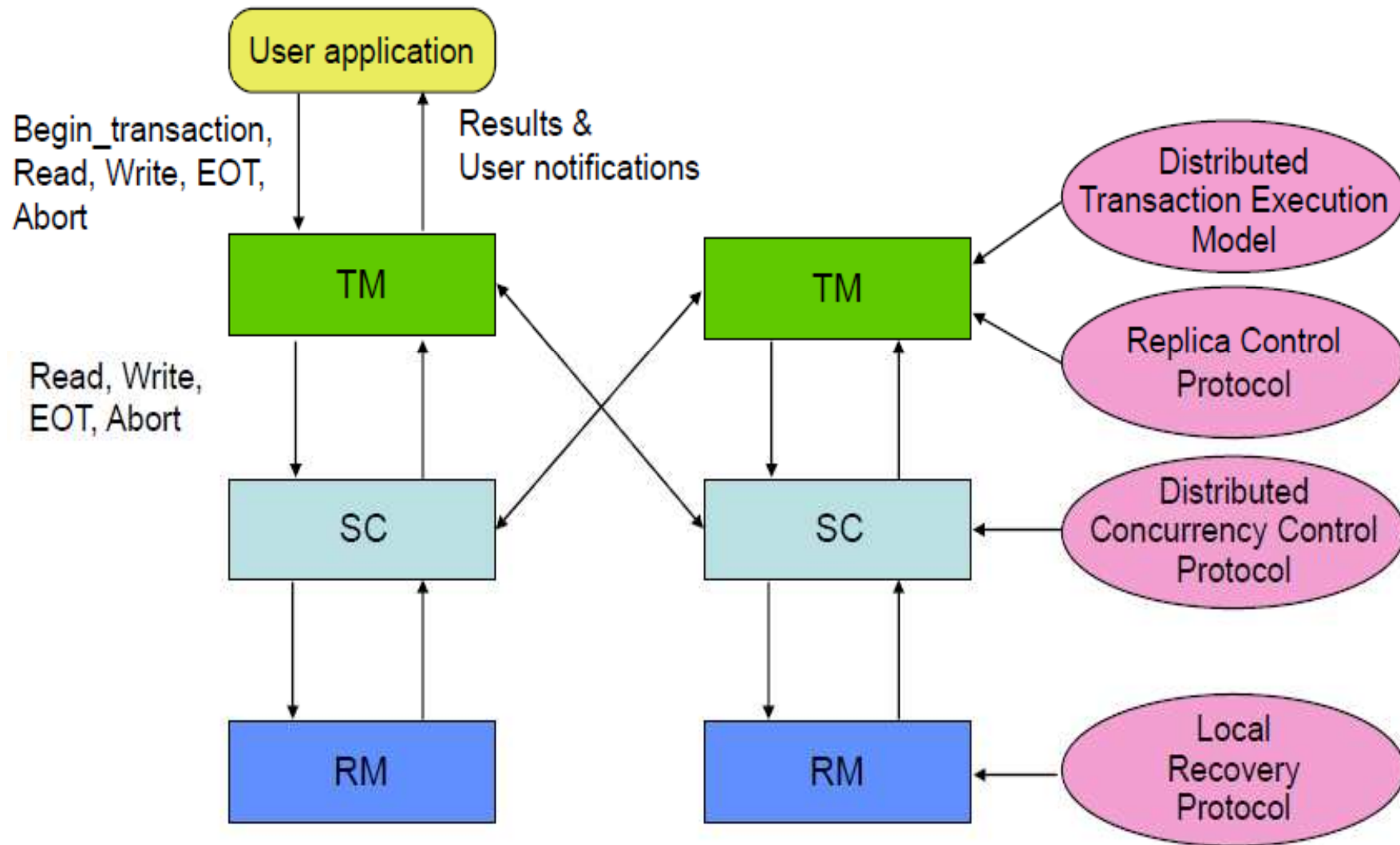
- **Clients** may use a **server** as a means of **sharing** some **resources**. E.g. some clients update the server's objects and other clients access them.
- However, in some applications, **threads** need to **communicate** and **coordinate** their actions.
- **Producer** and **Consumer** problem.
  - **Wait** and **Notify** actions.



# Centralized Transaction Execution



# Distributed Transaction Execution





# Properties of Transactions

## **A** TOMICITY

- All or nothing
- Multiple operations combined as an atomic transaction

## **C**ONSISTENCY

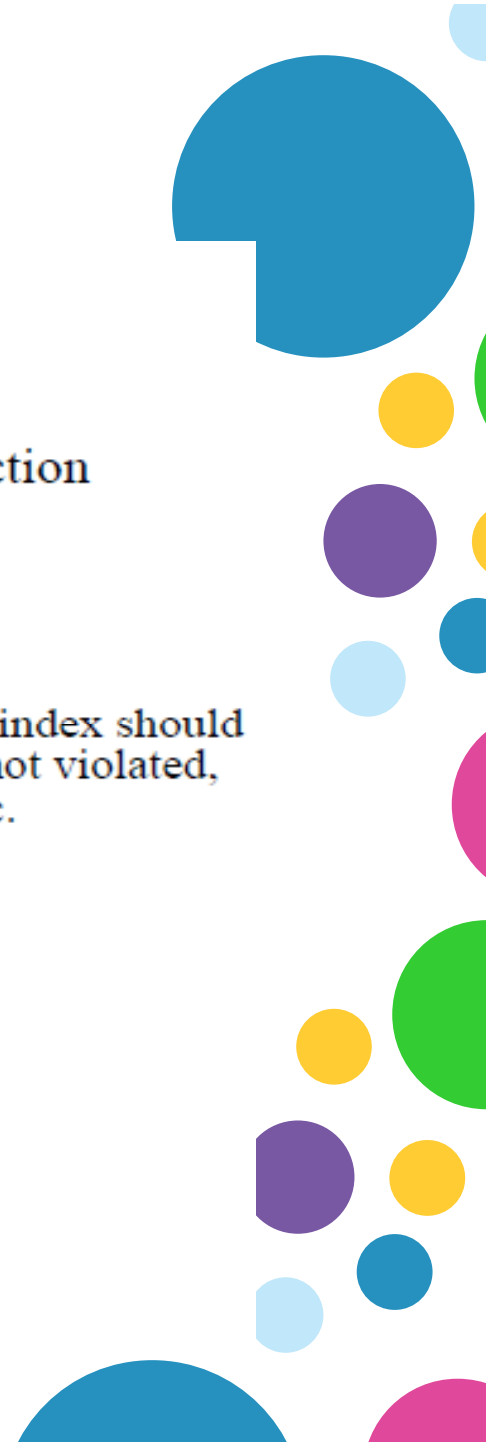
- No violation of integrity constraints
  - System specific rules. In a distributed database, the index should always reflect the data, foreign key constraints are not violated, triggers are issued, replicas have the same value, etc.
- Transactions are correct programs

## **I**SOLATION $\Leftarrow$ Our focus in this module

- Concurrent changes invisible  $\rightarrow$  serializable

## **D**URABILITY

- Committed updates persist
- Database recovery



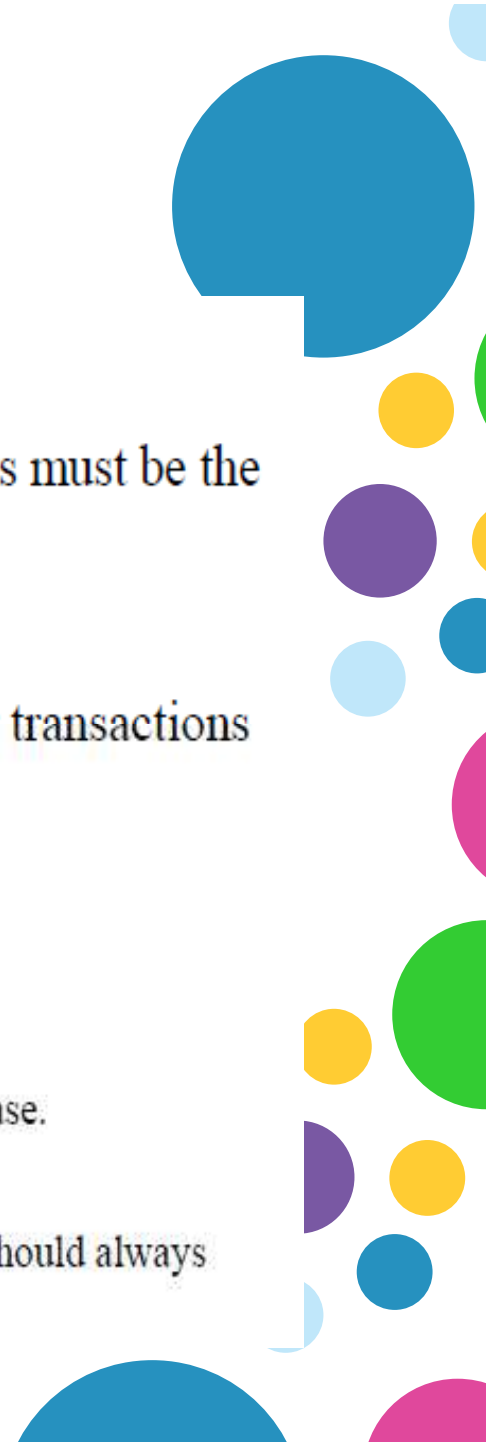
# Serializability of Transactions

## ■ Serializability

- If several transactions are executed concurrently, the results must be the same as if they were executed serially in some order.

## ■ Incomplete results

- An incomplete transaction cannot reveal its results to other transactions before its commitment.
- Necessary to avoid cascading aborts.
- Anomalies:
  - Lost updates
    - The effects of some transactions are not reflected on the database.
  - Inconsistent retrievals
    - E.g. a transaction, if it reads the same object more than once, should always read the same value.



# Concurrency Control – Lost Update Problem

Initial values: A=100, B=200, C=300

Transaction <i>T</i> :	Transaction <i>U</i> :
<i>balance</i> = <i>b.getBalance</i> (); <i>b.setBalance</i> ( <i>balance</i> *1.1); <i>a.withdraw</i> ( <i>balance</i> /10)	<i>balance</i> = <i>b.getBalance</i> (); <i>b.setBalance</i> ( <i>balance</i> *1.1); <i>c.withdraw</i> ( <i>balance</i> /10)
<i>balance</i> = <i>b.getBalance</i> ();      \$200	<i>balance</i> = <i>b.getBalance</i> ();      \$200
<i>b.setBalance</i> ( <i>balance</i> *1.1);      \$220	<i>b.setBalance</i> ( <i>balance</i> *1.1);      \$220
<i>a.withdraw</i> ( <i>balance</i> /10)      \$80	<i>c.withdraw</i> ( <i>balance</i> /10)      \$280

a, b and c initially have bank account balance are: 100, 200, and 300.  
T transfers an amount from a to b. U transfers an amount from c to b.  
b is increased by 10% on its balance in each. Totally 20 % hike

# Concurrency Control – Inconsistent Retrieval Problem

Initial values: A=200, B=200

Transaction <i>V</i> :		Transaction <i>W</i> :	
<i>a.withdraw(100)</i> <i>b.deposit(100)</i>		<i>aBranch.branchTotal()</i>	
<i>a.withdraw(100);</i>	\$100	<i>total = a.getBalance()</i>	\$100
		<i>total = total+b.getBalance()</i>	\$300
		<i>total = total+c.getBalance()</i>	
<i>b.deposit(100)</i>	\$300	⋮	

# Serially Equivalent

- If these transactions are done one at a time in some order, then the final result will be correct.
- If we do not want to sacrifice the concurrency, an **interleaving** of the **operations** of **transactions** may lead to the **same effect** as if the **transactions** had been **performed one** at a **time** in **some order**.
- We say it is a **serially equivalent interleaving**.
- The use of **serial equivalence** is a criterion for **correct concurrent execution** to prevent lost updates and inconsistent retrievals.



# A serially equivalent interleaving of $T$ and $U$

Transaction $T$ :		Transaction $U$ :	
<i>balance = b.getBalance()</i>		<i>balance = b.getBalance()</i>	
<i>b.setBalance(balance*1.1)</i>		<i>b.setBalance(balance*1.1)</i>	
<i>a.withdraw(balance/10)</i>		<i>c.withdraw(balance/10)</i>	
<i>balance = b.getBalance()</i>	\$200	<i>balance = b.getBalance()</i>	\$220
<i>b.setBalance(balance*1.1)</i>	\$220	<i>b.setBalance(balance*1.1)</i>	\$242
<i>a.withdraw(balance/10)</i>	\$80	<i>c.withdraw(balance/10)</i>	\$278

# A serially equivalent interleaving of *V* and *W*

TransactionV:		TransactionW:
<i>a.withdraw(100);</i> <i>b.deposit(100)</i>		<i>aBranch.branchTotal()</i>
<i>a.withdraw(100);</i>	\$100	
<i>b.deposit(100)</i>	\$300	
		<i>total = a.getBalance()</i> \$100
		<i>total = total+b.getBalance()</i> \$400
		<i>total = total+c.getBalance()</i>
		...

# *Read and write operation conflict rules*

Operations of different transactions		Conflict	Reason
read	read	No	Because the effect of a pair of read operations does not depend on the order in which they are executed
read	write	Yes	Because the effect of a read and a write operation depends on the order of their execution
write	write	Yes	Because the effect of a pair of write operations depends on the order of their execution



# A non-serially equivalent interleaving of operations of transactions $T$ and $U$

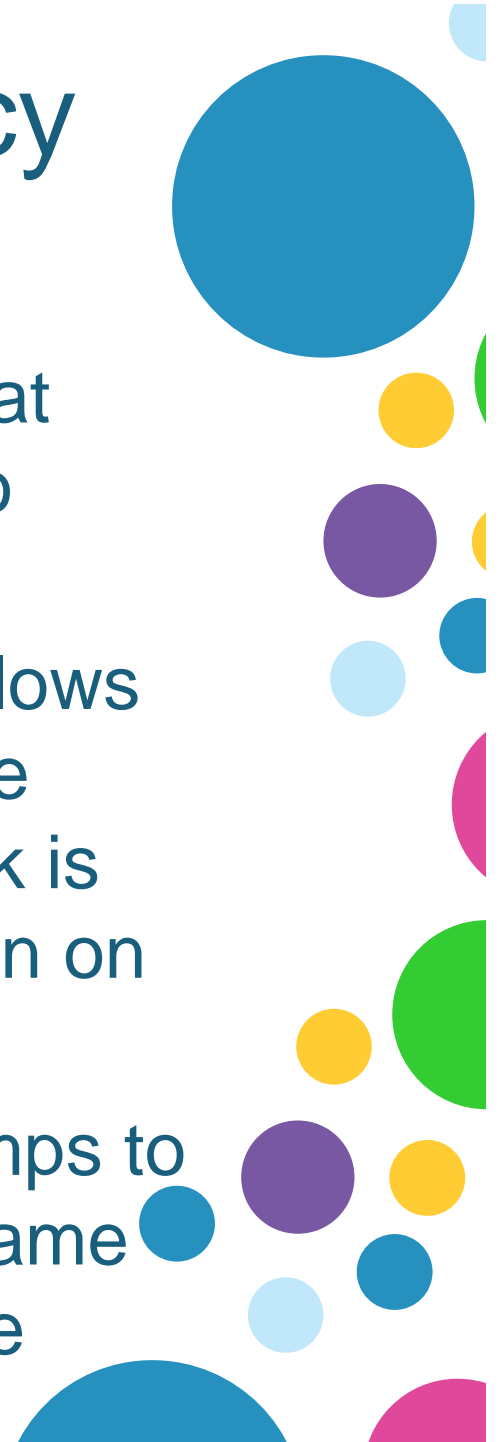
Transaction $T$ :	Transaction $U$ :
$x = \text{read}(i)$ $\text{write}(i, 10)$	$y = \text{read}(j)$ $\text{write}(j, 30)$
$\text{write}(j, 20)$	$z = \text{read}(i)$

Access to objects  $i$  &  $j$  are serial, but transactions are not serially equivalent



# Solutions to Concurrency control problems

- **Locks** used to order transactions that access the same object according to request order.
- **Optimistic concurrency control** allows transactions to proceed until they are ready to commit, whereupon a check is made to see any conflicting operation on objects.
- **Timestamp ordering** uses timestamps to order transactions that access the same object according to their starting time



# A dirty read when transaction *T* aborts

Transaction *T*:

*a.getBalance()*

*a.setBalance(balance + 10)*

*balance = a.getBalance()*      \$100

*a.setBalance(balance + 10)*      \$110

*abort transaction*

Transaction *U*:

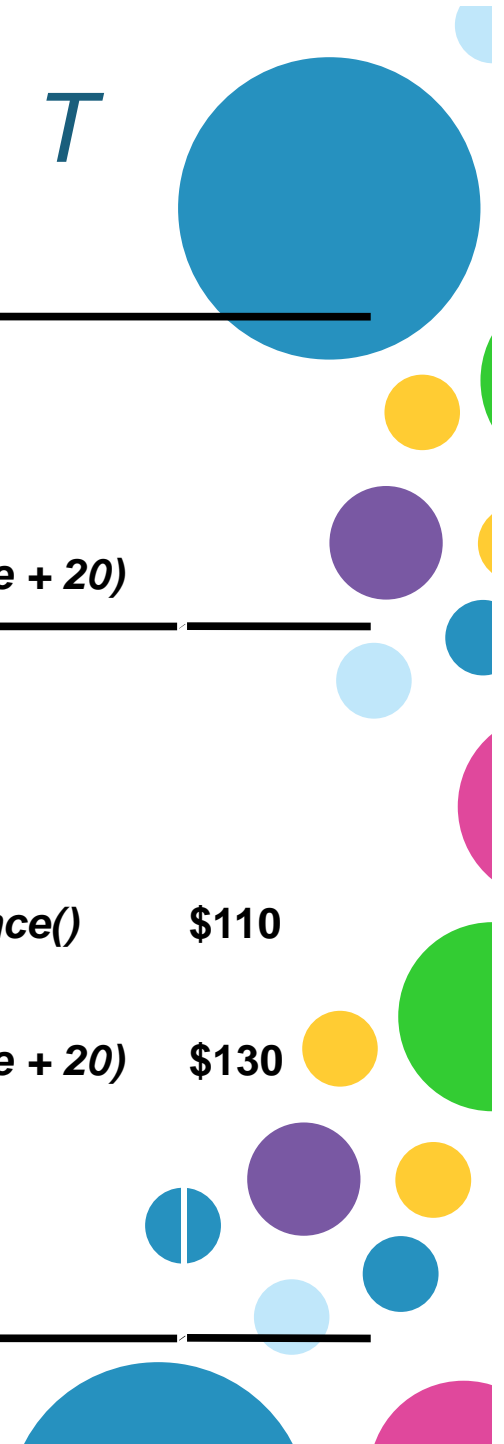
*a.getBalance()*

*a.setBalance(balance + 20)*

*balance = a.getBalance()*      \$110

*a.setBalance(balance + 20)*      \$130

*commit transaction*



# Recoverability of Transactions

- The strategy for recoverability is to **delay commits** until after the **commitment** of any **other transaction** whose **uncommitted** state has been observed.
- In our example, *U delays its commit until after T commits.*
- *In the case that T aborts, then U must abort as well*



# Cascading aborts

- Abort of one transaction will cause other transactions to abort.
- Transactions are **only allowed to read objects** that were **written** by **committed transactions**.
- To ensure that this is the case, any *read operation must be **delayed until** other **transactions** that applied a **write operation** to the **same object** have **committed** or **aborted**.*



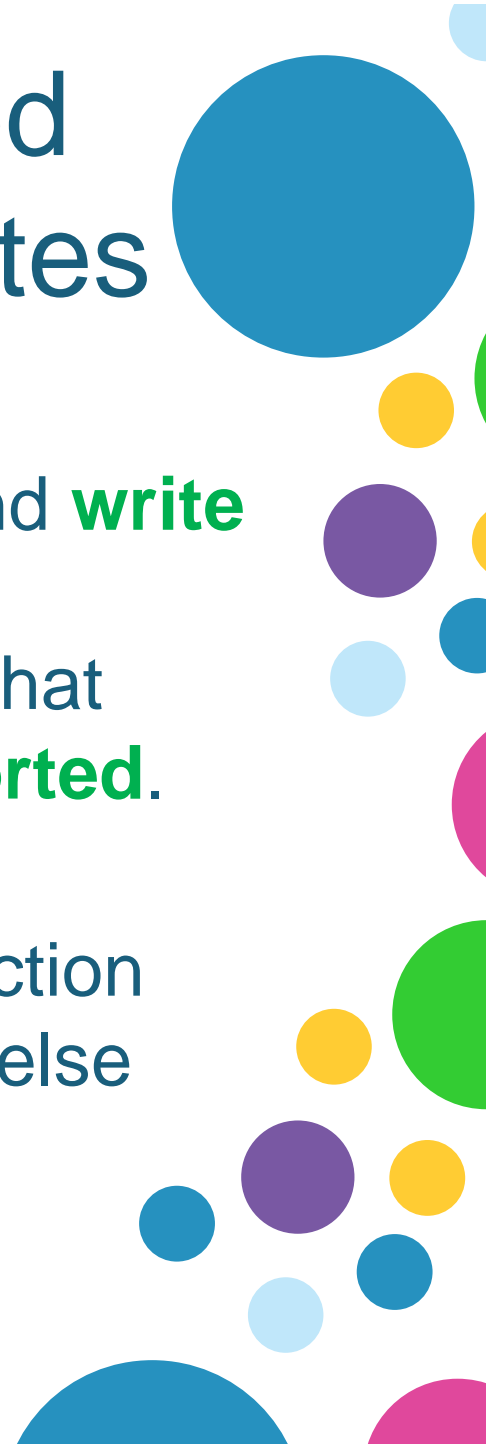


# Overwriting uncommitted values or Pre-mature Writes

- **Strict Execution of Transactions:**

If the transaction **delays** both **read** and **write** operations on an **object** until all **transactions** that **previously wrote** that **object** have either **committed** or **aborted**.

- **Tentative Version:** Make changes to tentative versions of objects. If transaction commits, transfer updates to objects, else delete tentative version.



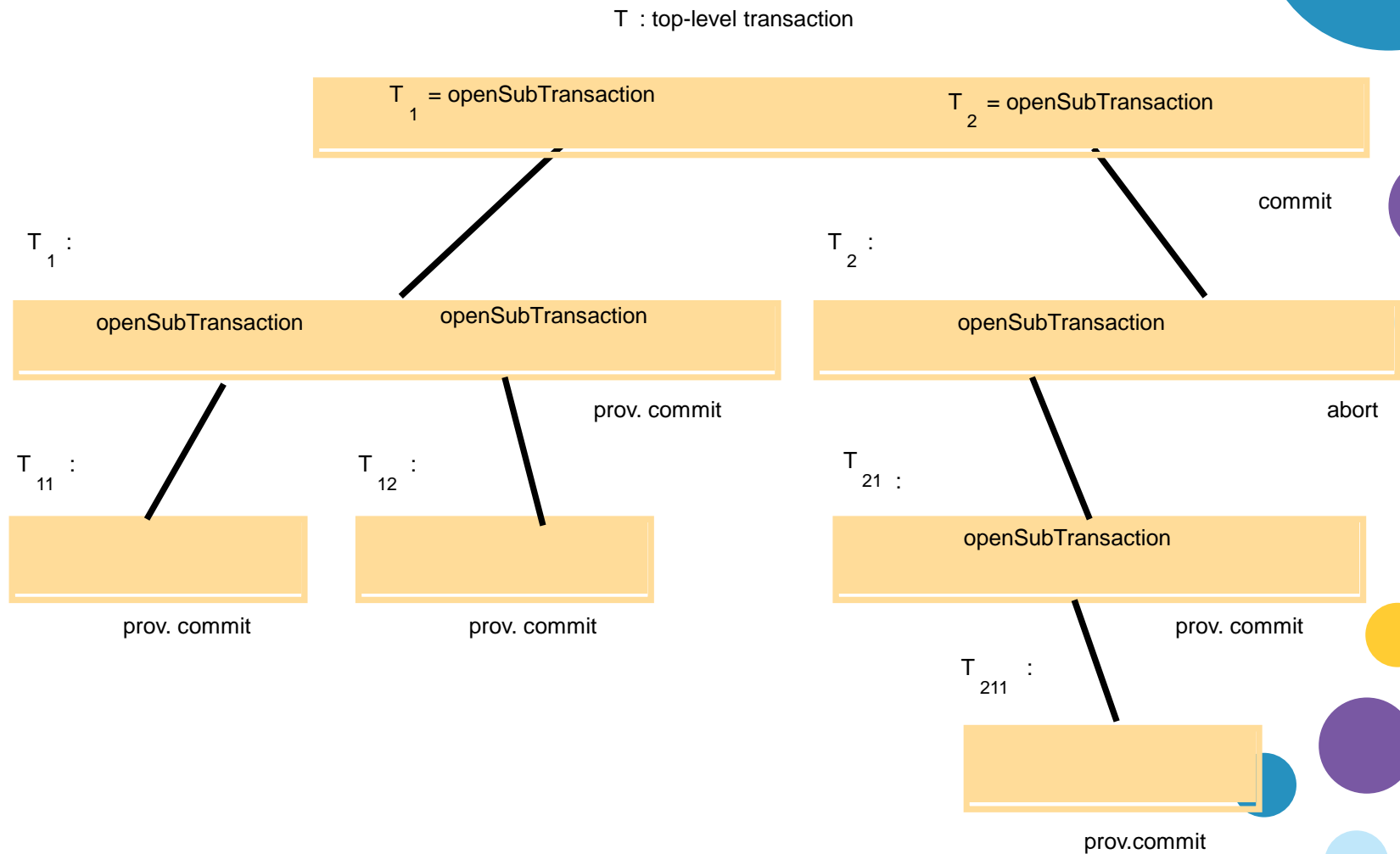
# Nested transactions

- **Nested transactions** extend the above transaction model by allowing **transactions** to be **composed** of **other transactions**.
- Thus **several transactions** may be **started from within a transaction**.





# Nested transactions



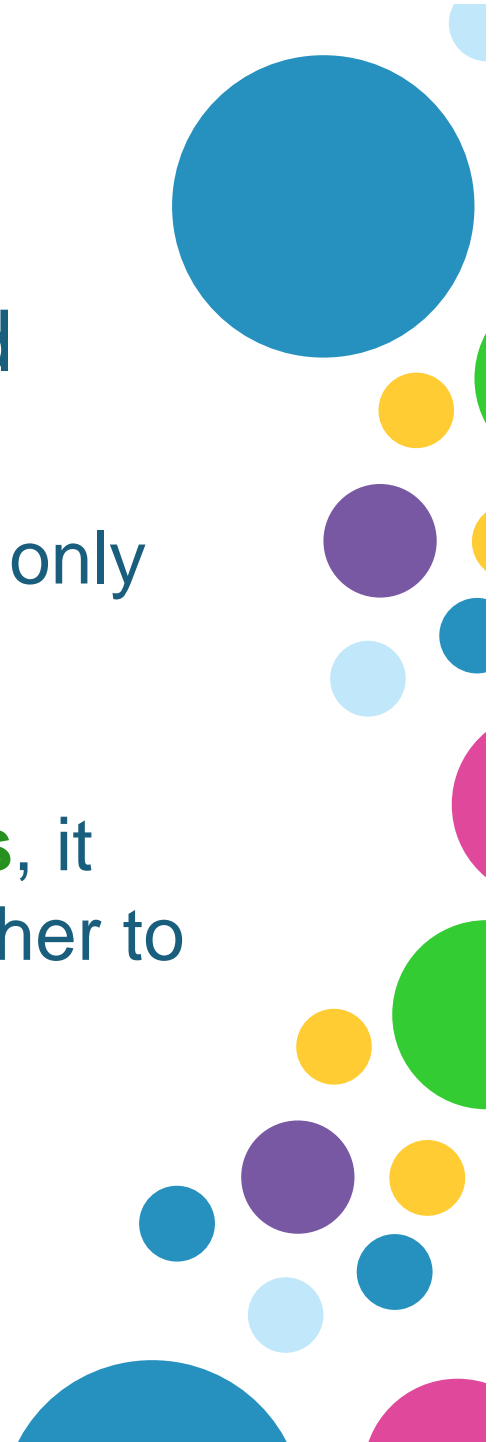
# Nested transactions

- **Subtransactions** at one **level** (and their descendants) may **run concurrently** with other **subtransactions** at the same level in the hierarchy.
- **Subtransactions** can **commit** or **abort independently**.



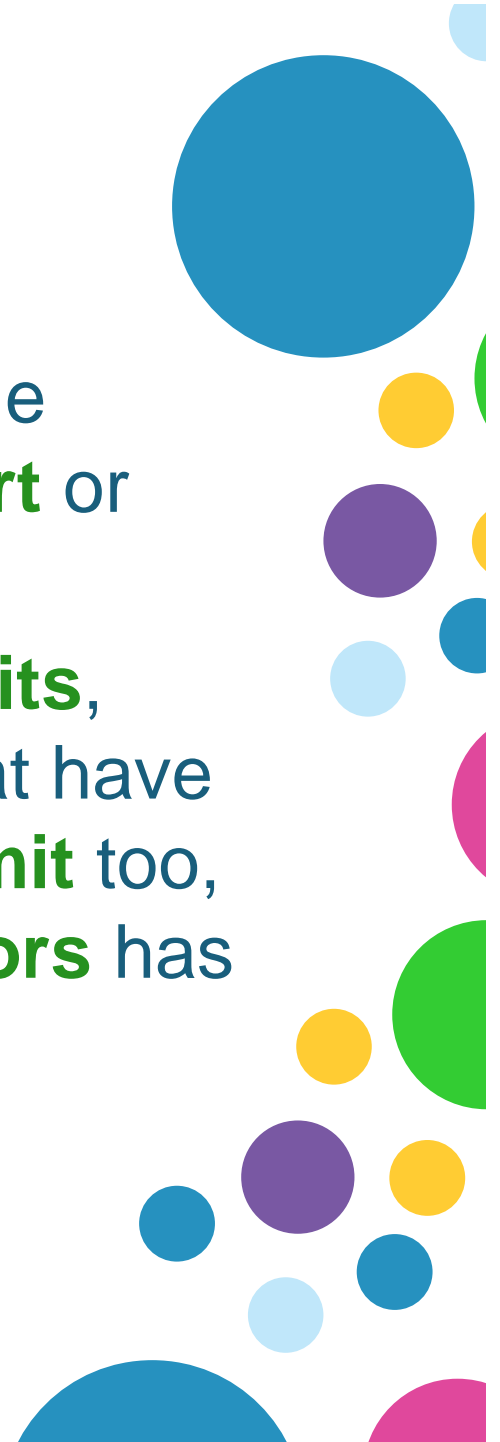
# Nested transactions

- The rules for committing of nested transactions
  - A transaction may **commit** or **abort** only after its **child** transactions have completed.
  - When a **subtransaction completes**, it makes an **independent** decision either to commit provisionally or to abort. Its decision to abort is final.
  - When a **parent aborts**, **all** of its **subtransactions** are aborted.



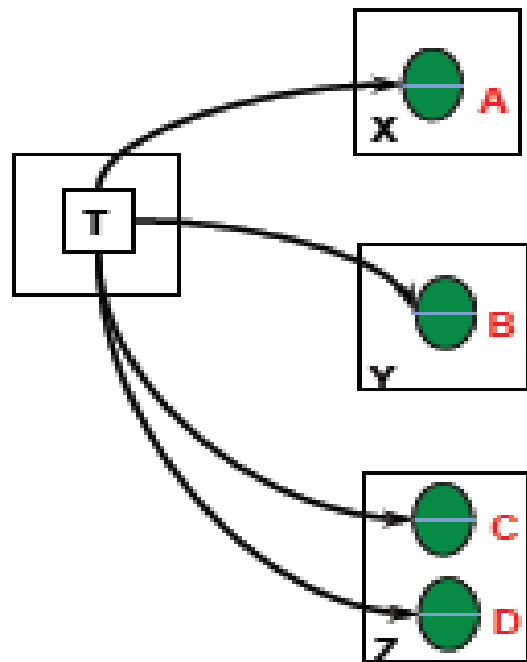
# Nested transactions

- When a **subtransaction aborts**, the **parent** can **decide** whether to **abort** or **not**.
- If the **top-level transaction commits**, then all of the **subtransactions** that have provisionally **committed** can **commit** too, provided that **none** of their **ancestors** has **aborted**.

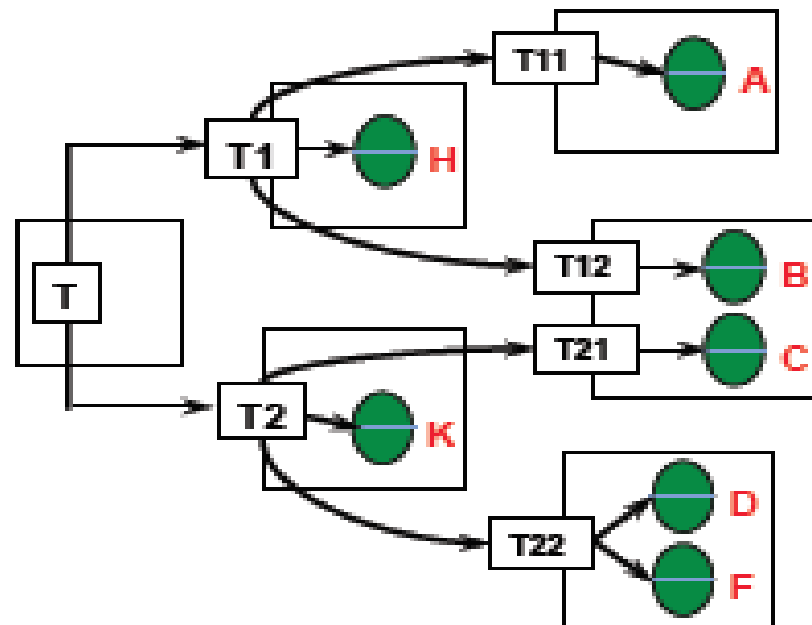


# Distributed Transactions

- Transactions that invoke operations at multiple servers



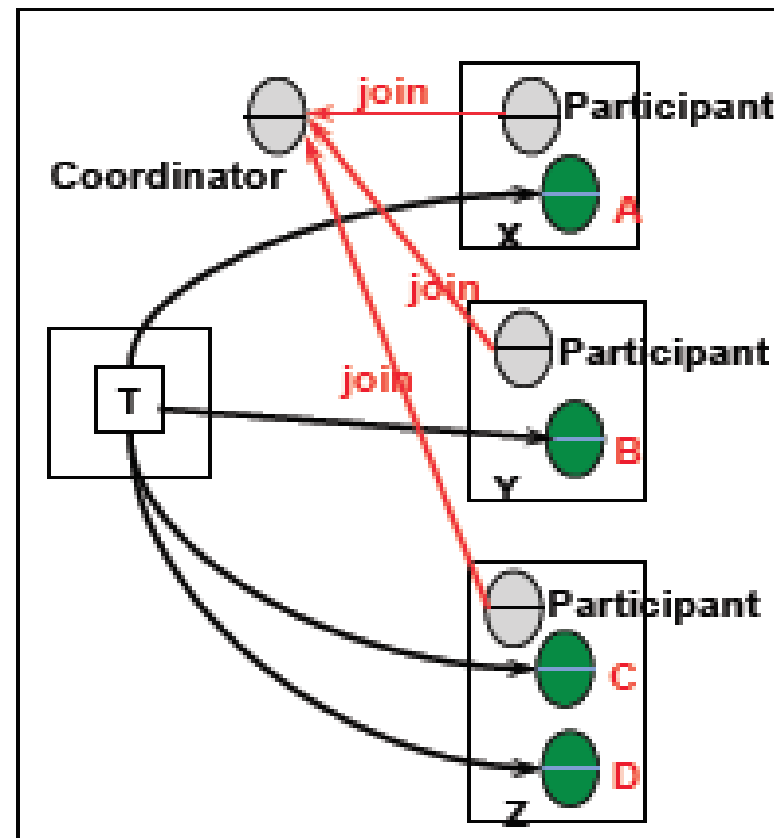
Flat Distributed Transaction



Nested Distributed Transaction

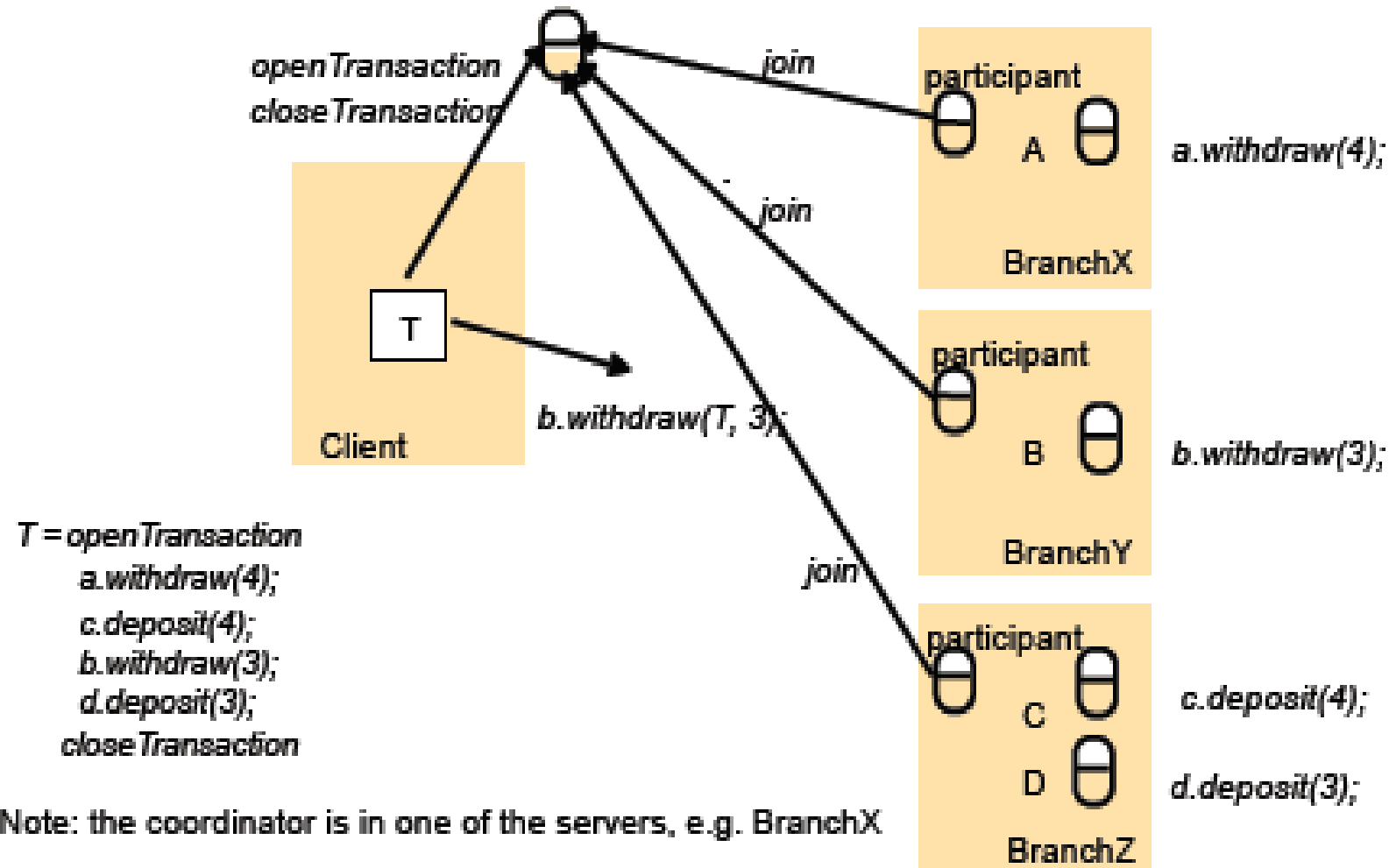
# Distributed Transactions

- **Coordinator**
  - In charge of begin, commit, and abort
- **Participants**
  - Server processes that handle local operations



Coordinator & Participants

# Distributed Transactions - Example



Thank You

