



THE UNIVERSITY OF  
MELBOURNE

# COMP90050: Advanced Database Systems

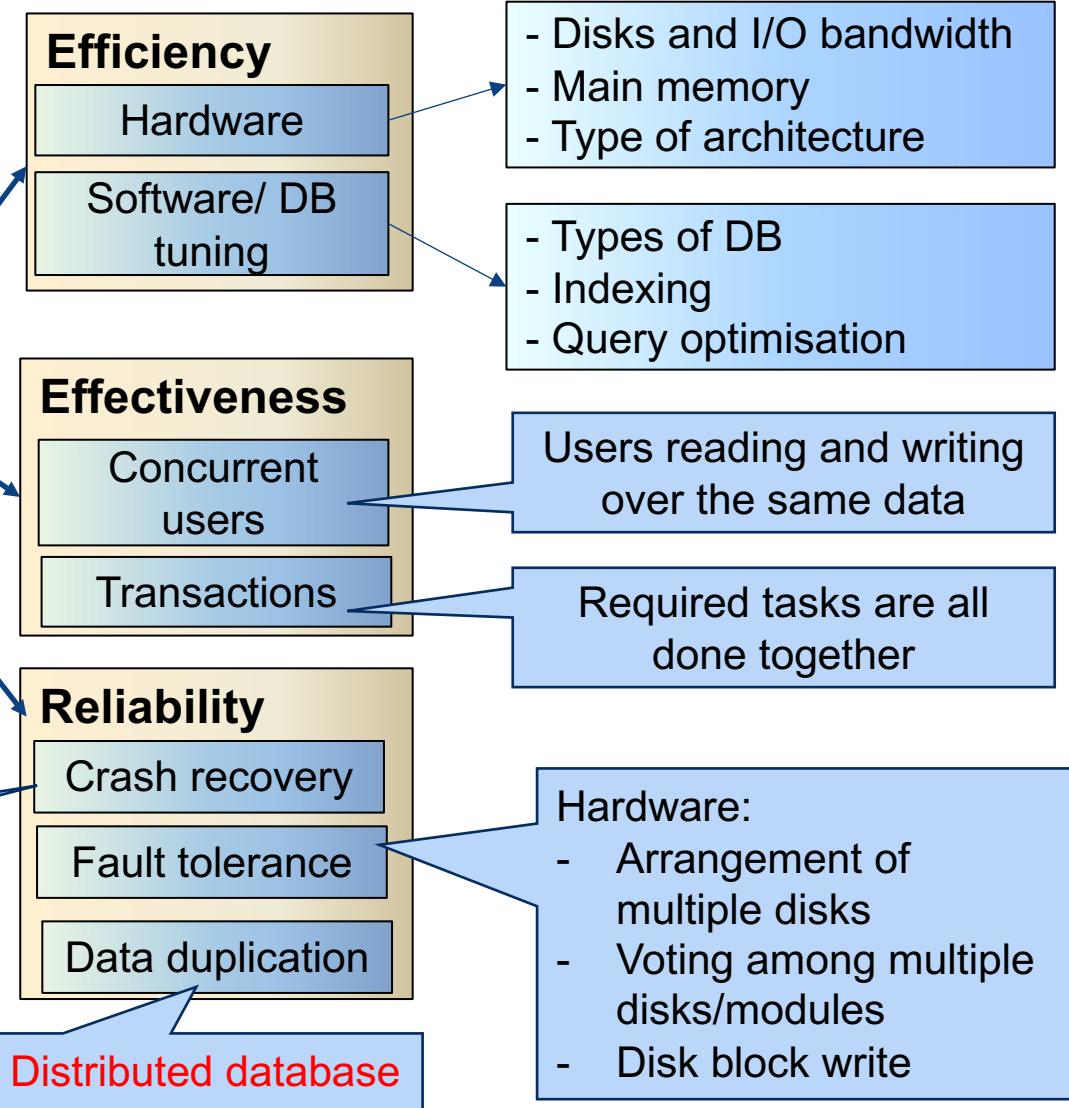
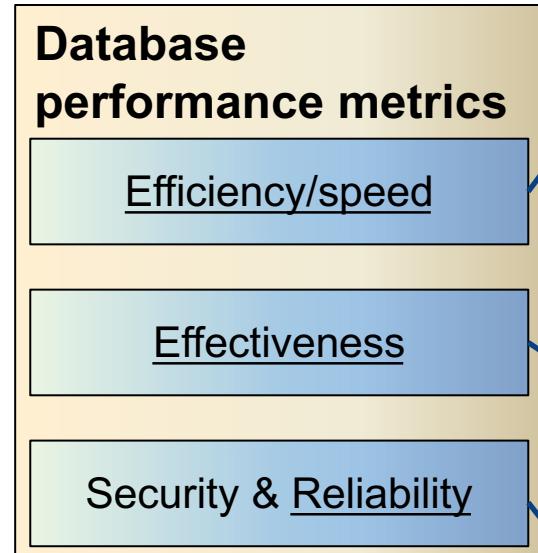


**Lecturer: Farhana Choudhury (PhD)**

**Distributed databases**

**Week 11 part 1**

# Core Concepts of Database management system



- Logging
- ARIES algorithm



# A related concept on distributed databases first - atomicity



# Atomicity in distributed transaction processing

The two-phase commit protocol (2PC) can help achieve atomicity in distributed transaction processing

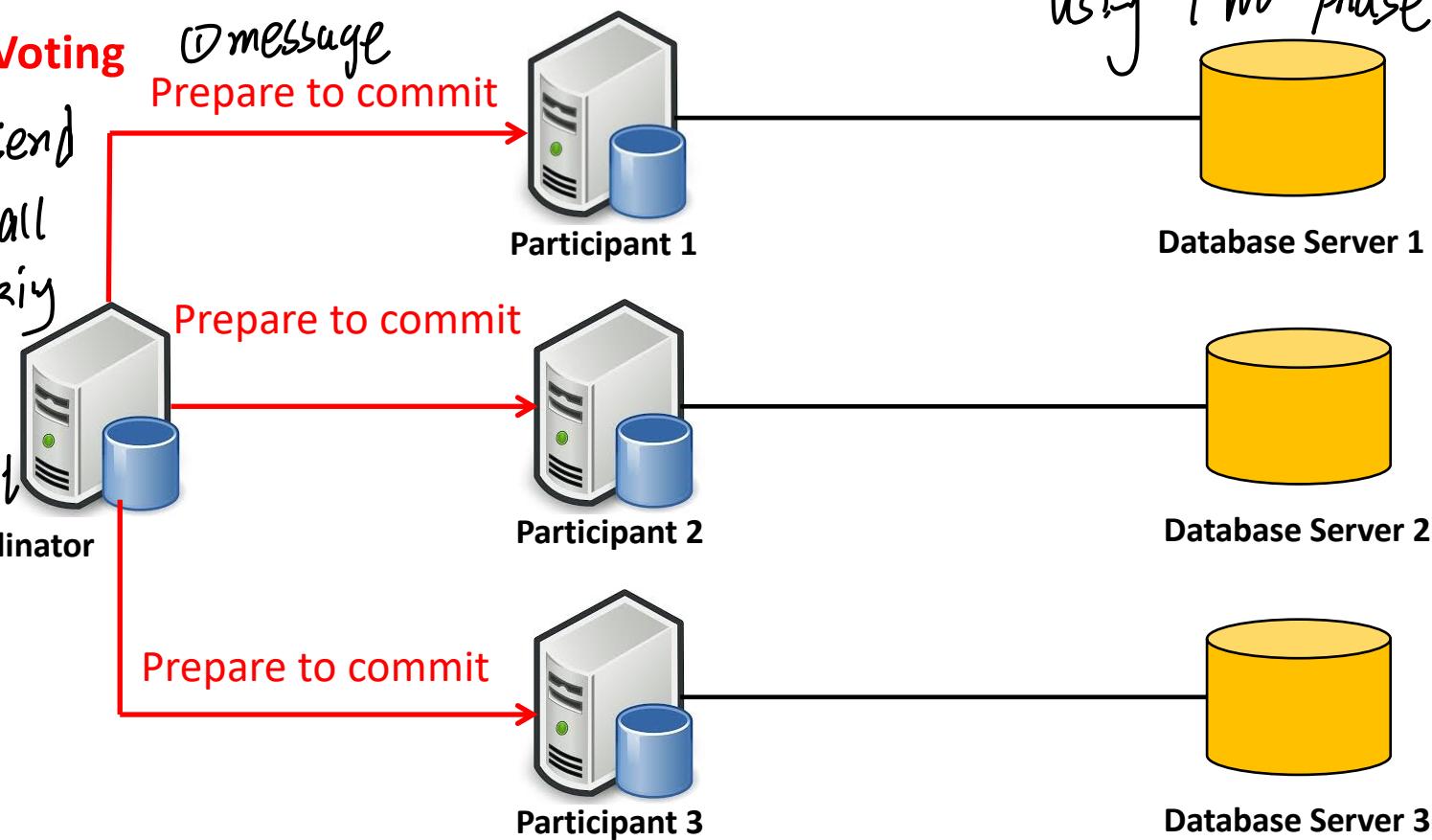
$$T = \begin{pmatrix} T_1 \\ T_2 \\ T_3 \end{pmatrix}$$

Sub transaction-

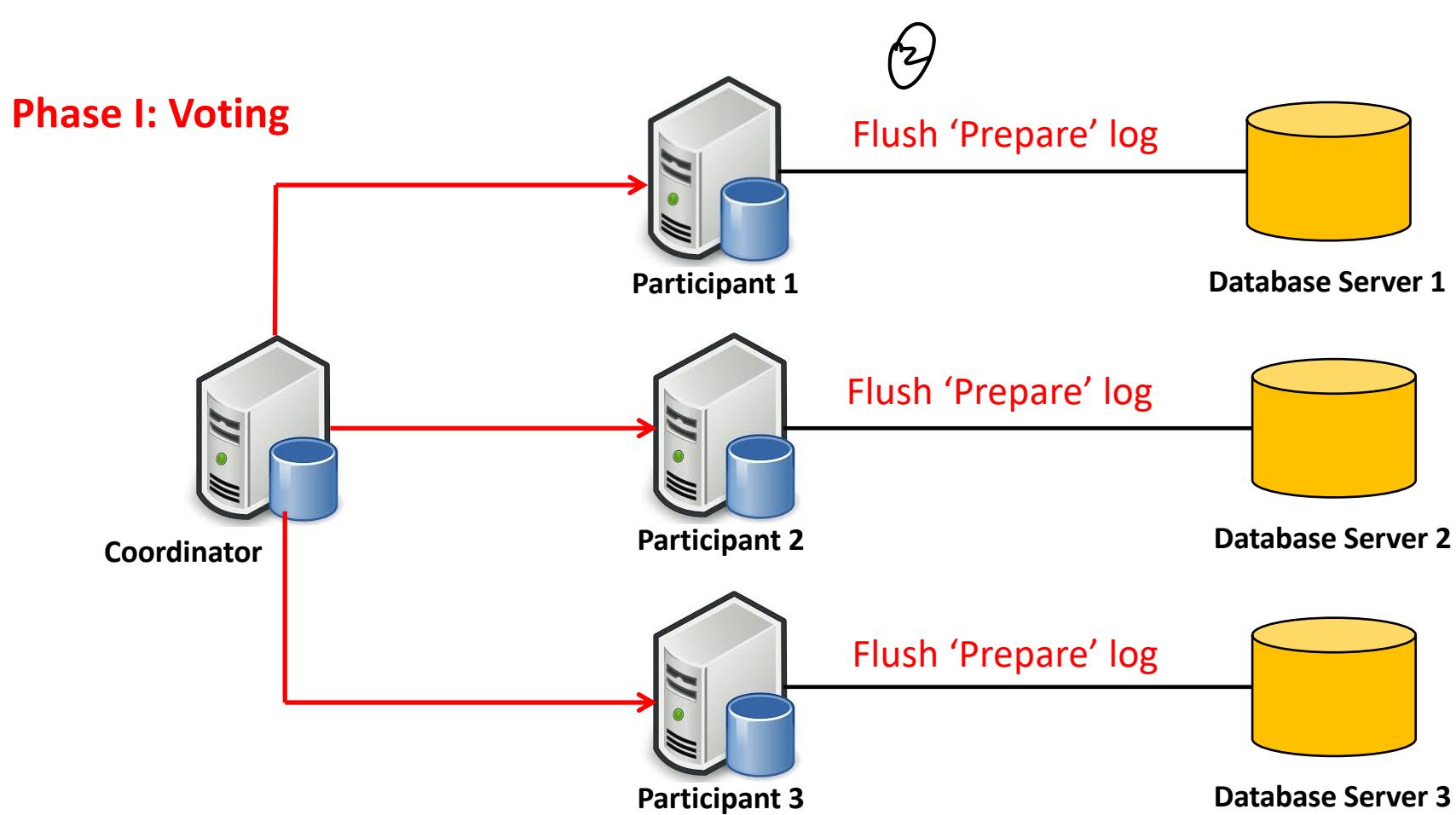
How to maintain atomicity for each of subtransaction?  
Using two phase commit.

## Phase I: Voting

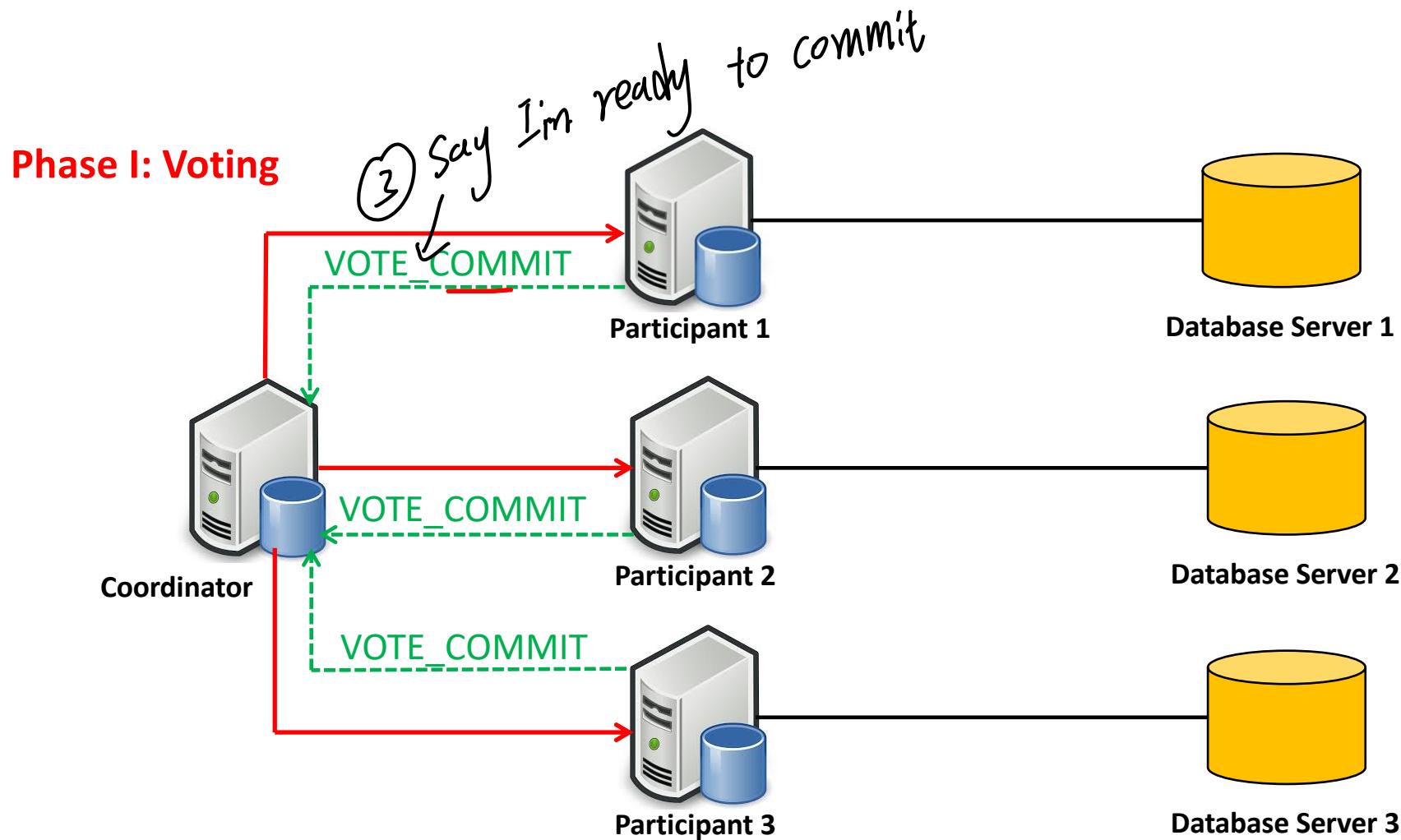
① Coordinator send message to all machines asking if they are prepared to commit



# Two phase commit protocol

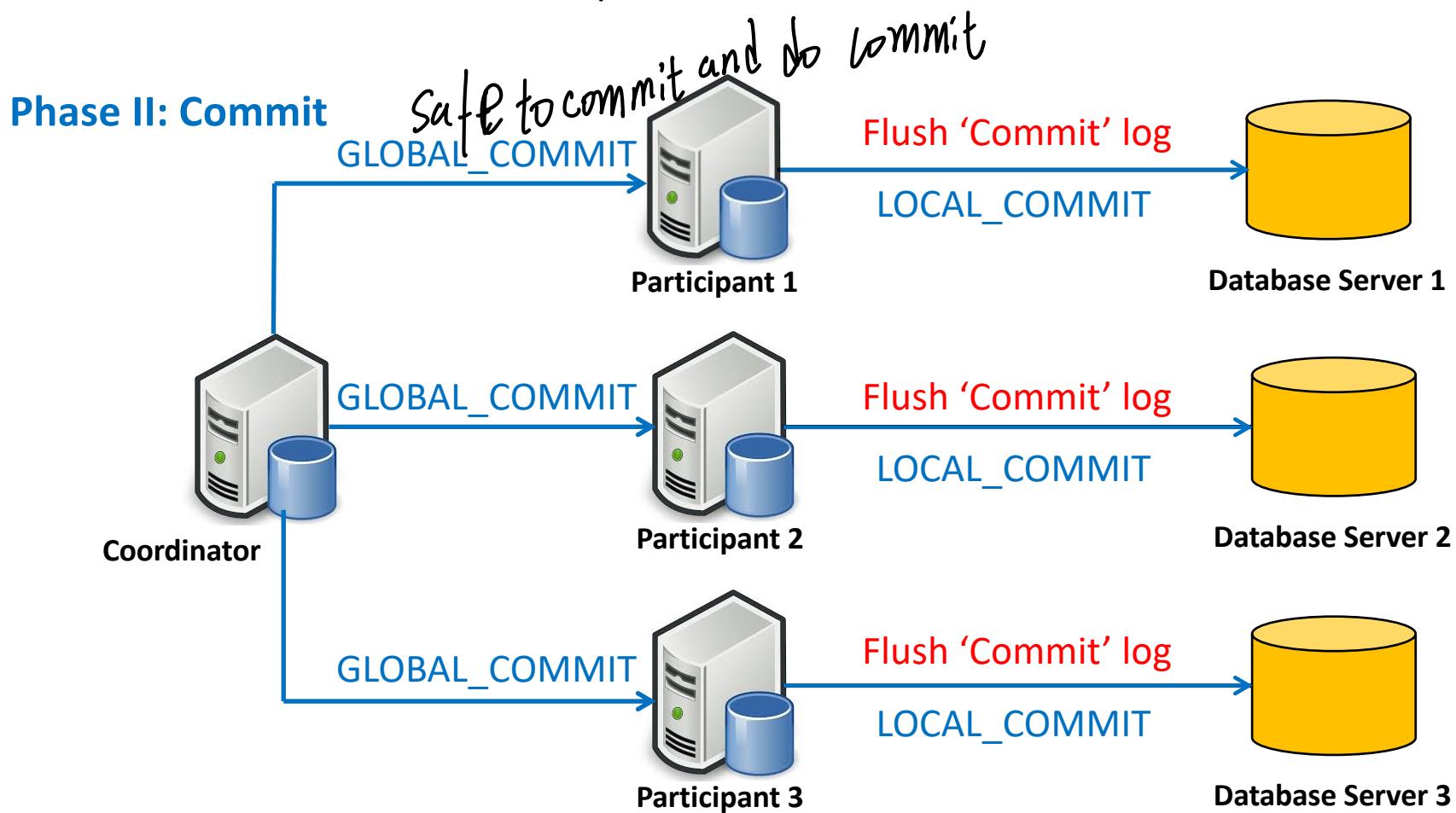


# Two phase commit protocol



# Two phase commit protocol

(Coordinator receives 'ready to commit' from all participants from that distributed transaction)





## Two phase commit protocol

What happens if some participates abort? (participate 4/5, abort 3/5 (coordinator))

- Coordinator or participant can abort transaction

- 通常情况
- If a participant abort, it must inform coordinator
- 若超时
- If a participant does not respond within a timeout period, coordinator will abort

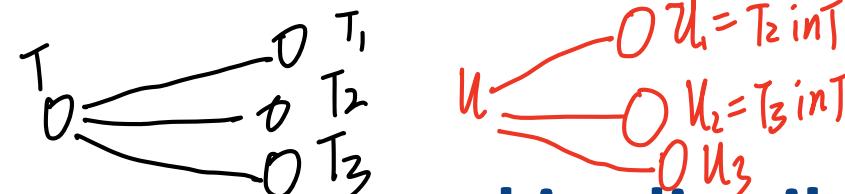
一投票便投票 one sub failed  $\Rightarrow$  all subs failed.

- If abort, coordinator asks all participates to rollback
- If abort, abort logs are forced to disk at coordinator and all participants

Voty and commit all abort to make atomicity.



# **Another related concept on distributed databases - concurrency control**



$U_1 = T_2 \text{ int}$   
 $U_2 = T_3 \text{ int}$   
 $U_3$

How to maintain concurrency  
Using  
① locks  
② timestamps.

## Concurrency control in distributed DBs

- Each server is responsible for applying concurrency control to its own objects
- The members of a collection of servers of distributed transactions are jointly responsible for ensuring that they are performed in a serially equivalent manner  
✓, 分步先准备  
3J
- BUT servers independently acting would not work
- If transaction **T** is before transaction **U** in their conflicting access to objects at one of the servers then:
  - They must be in that order at all of the servers whose objects are accessed in a conflicting manner by both **T** and **U**
- The central **Coordinator** should assure this



# Other Considerations for Locking-based systems

- A local lock manager cannot release any locks until it knows that the transaction has been committed or aborted at all the servers involved in the transaction.
- The objects remain locked and are unavailable for other transactions during the commit protocol.
  - An aborted transaction releases its locks after phase 1 of the protocol.



# Timestamp ordering concurrency control revisited

- The coordinator accessed by a transaction issues a globally unique timestamp
- The timestamp is passed with each object access
- The servers are jointly responsible for ensuring serial equivalence:
  - that is **if T access an object before U, then T is before U at all objects**



# Optimistic concurrency control revisited

For distributed transactions to work:

- 1) Validation takes place in phase 1 of 2PC protocol at each server
- 2) Transactions use a globally unique order for validation



## What if objects in different servers are replicas for increased availability

inconsistency: 若 Client1 在 B 上修改  $\rightarrow$  未广播到 A  $\rightarrow$  Client2 读到不将  
propagate 多值.

Client 1:	Client 2:	Server
$setBalance_B(x, 1)$		
$setBalance_A(y, 2)$		
	$getBalance_A(y) \rightarrow 2$	
	$getBalance_A(x) \rightarrow 0$	

Initial balance of x and y is \$0

The behaviour above cannot occur if A and B did not exist (that is, if we had only one server)



## Transactions with replicated data

- the effect of transactions on replicated objects should be the **same as if they had been performed one at a time on a single set of objects**
- this property is called ***one-copy serializability***
- **If all servers are available then no issue – but what if some servers are not available?**

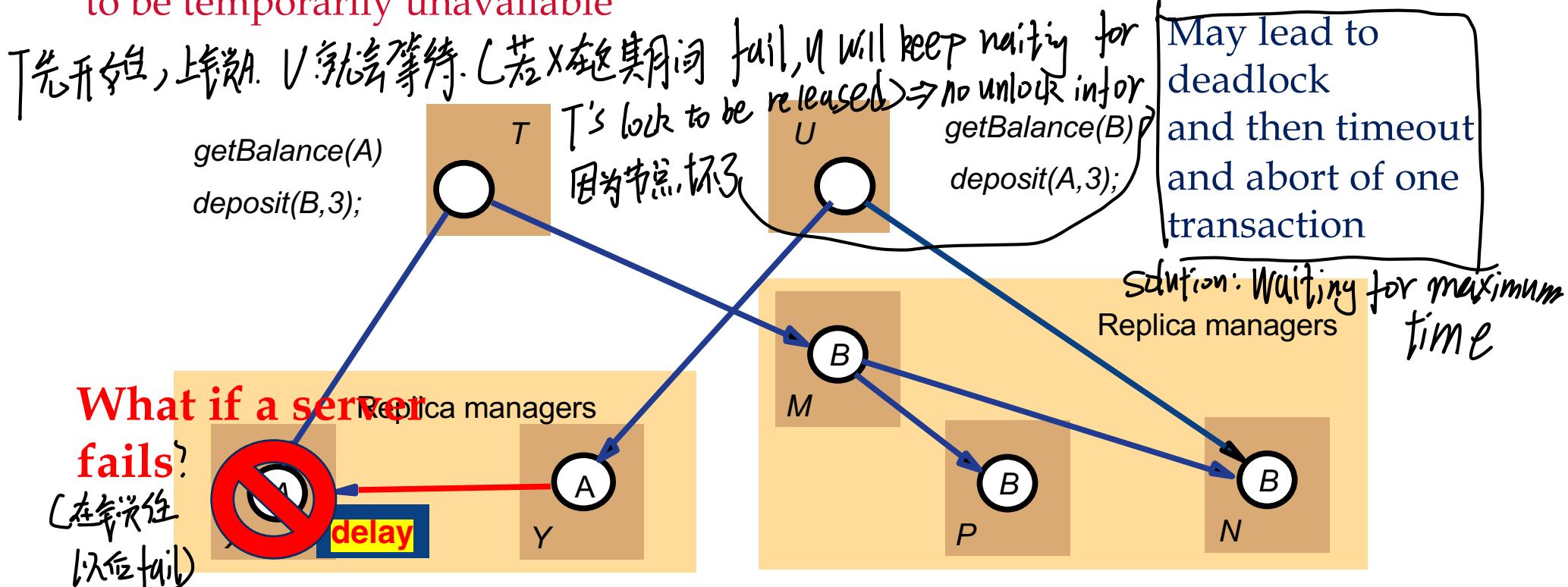


Object A (存在 X, 有一个 replica 在 Y)

## Lets build the solution step by step

Object B: 3个 replica.

The available copies replication scheme is designed to allow some servers to be temporarily unavailable



At X, T has read A and has locked it. Therefore U's deposit is delayed until T finishes.

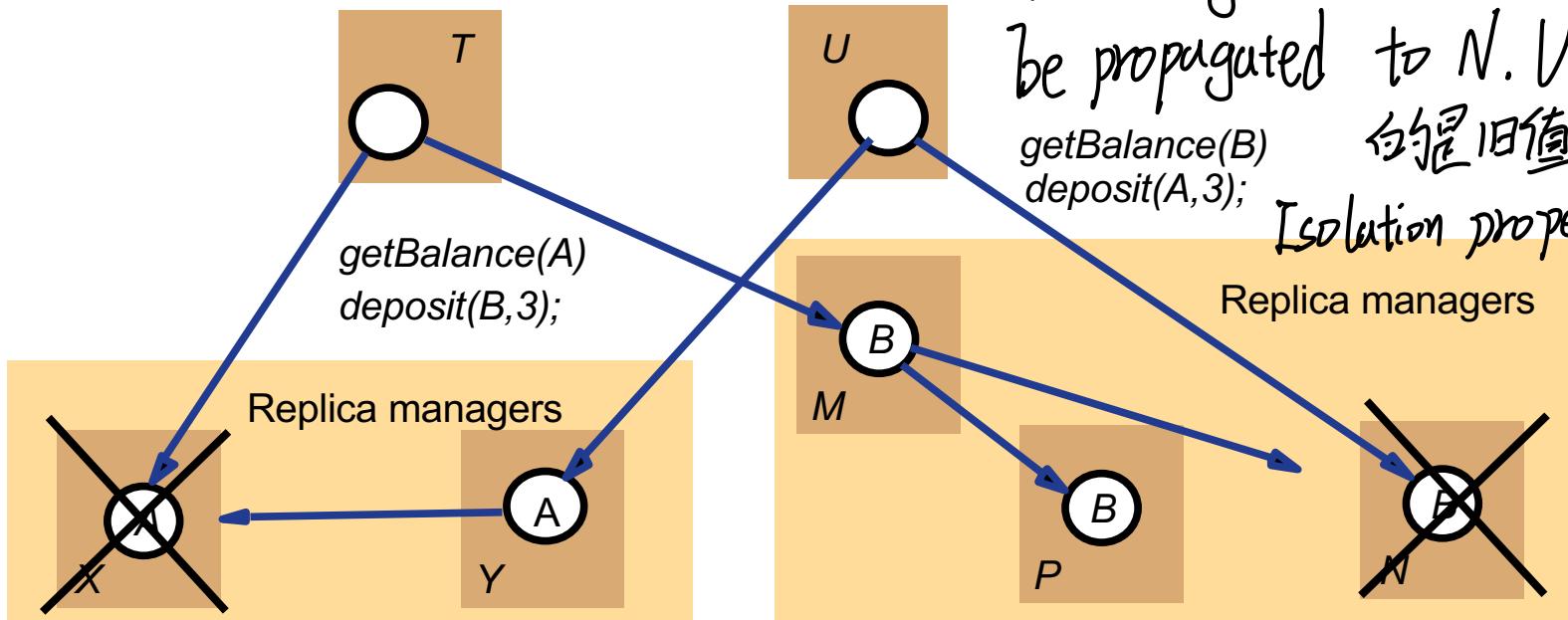
Normally, this leads to good concurrency control only if the servers do not fail....



X坏了，为什么不能轻易用其他好的？ $\Rightarrow$ Y?

## Cannot the other servers simply be used?

assume that X fails just after T has performed *getBalance* and N failing just after U has performed *getBalance*



therefore T's deposit will be performed at M and P (all available)  
and U's deposit will be performed at Y (all available)

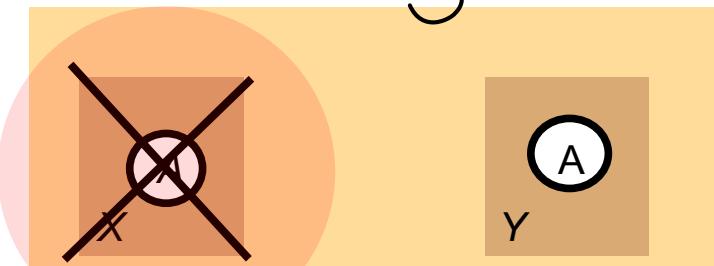
**NOT GOOD!!**

# Available copies replication rule

**Before a transaction commits, it checks for failed and available servers it has contacted, the set should not change during execution:**

- E.g.,  $T$  would check if  $X$  is still available among others.
- We said  $X$  fails before  $T$ 's deposit, in which case,  $T$  would have to abort.
- Thus no harm can come from this execution now.

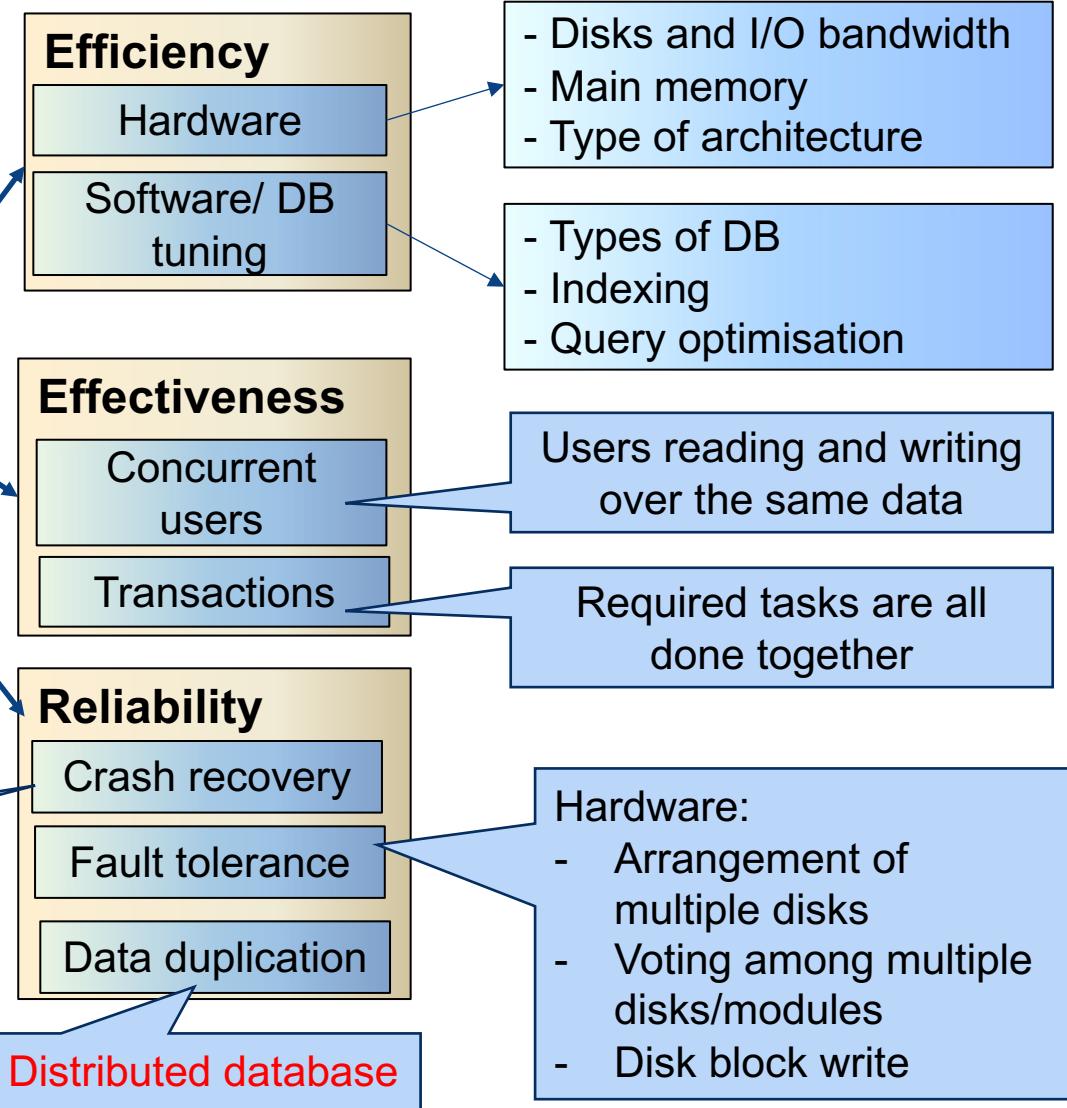
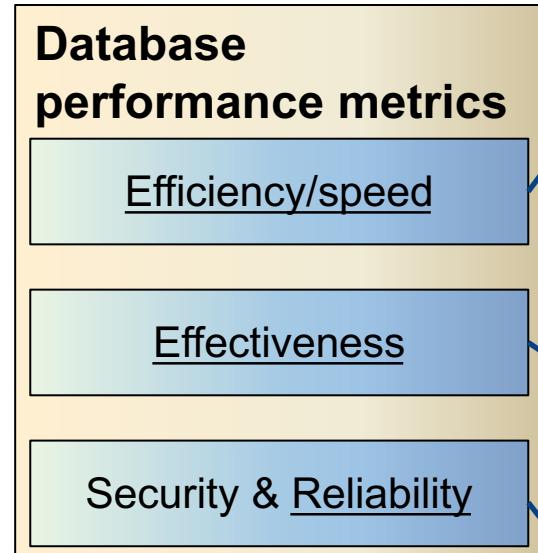
(1)  $T$  has  $(X, Y)$  during execution for  $A \rightarrow$  if one node xor  $y$  failed then  $T$  should abort.



(2)  $T(X, Y, Z) \rightarrow A$

node  $Z$  is recovering during execution of  $T$ ,  $T$  abort as well

# Core Concepts of Database management system



- Logging
- ARIES algorithm

# CAP Theorem

Some slides from Mohammad Hammoud,  
Dong Wang

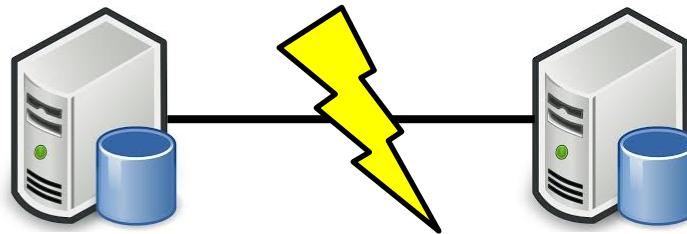
# The CAP Theorem

- The limitations of distributed databases can be described in the so called the **CAP theorem**
  - **Consistency**: every node always sees the same data at any given instance (i.e., strict consistency)
  - **Availability**: the system continues to operate, even if nodes crash, or some hardware or software parts are down due to upgrades(*still get data from read/write operations for other nodes*)
  - **Partition Tolerance**: the system continues to operate in the presence of network partitions

CAP theorem: any distributed database with shared data, can have at most two of the three desirable properties, C, A or P

# The CAP Theorem (Cont'd)

- Let us assume two nodes on opposite sides of a network partition:



- Availability + Partition Tolerance forfeit Consistency as changes in place cannot be propagated when the system is partitioned.  
~~~~~
- Consistency + Partition Tolerance entails that one side of the partition must act as if it is unavailable, thus forfeiting Availability  
↓ 该分区连一个 node .
- Consistency + Availability is only possible if there is no network partition, thereby forfeiting Partition Tolerance

# Large-Scale Databases

- When companies such as Google and Amazon were designing large-scale databases, 24/7 Availability was a key
  - A few minutes of downtime means lost revenue
- With databases in 1000s of machines, the likelihood of a node or a network failure increases tremendously
- Therefore, in order to have strong guarantees on Availability and Partition Tolerance, they had to sacrifice “strict” Consistency (*implied by the CAP theorem*)

企业：网速 $\downarrow$  + 网络故障  $\rightarrow$  重新启动机器，提供服务  $\Rightarrow$  Consistency 没了。  
Availability + partition tolerance

# Types of Consistency

- Strong Consistency
  - After the update completes, **any subsequent access** will return the **same** updated value.
- Weak Consistency
  - It is **not guaranteed** that subsequent accesses will return the updated value.
- Eventual Consistency  $O_{1-100} \rightarrow$  有的一返回旧值 | ,有的返回新值 100
  - Specific form of weak consistency  $\rightarrow$  某些没有更新发生且这些node都未被
  - It is guaranteed that if **no new updates** are made to object, **eventually** all accesses will return the last updated value (e.g., *propagate updates to replicas in a lazy fashion*)

# Eventual Consistency Variations

- Causal consistency  $A \rightarrow B$  *由引起修改的 user 先看到 A  
再次查询看到 B*
  - Processes that have causal relationship will see consistent data
- Read-your-write consistency *(自己改的数据总能看到自己的)*
  - A process always accesses the data item after its update operation and never sees an older value
- Session consistency *log in - log out* *算一个 session 保证一致性*
  - As long as session exists, system guarantees read-your-write consistency
  - Guarantees do not overlap sessions

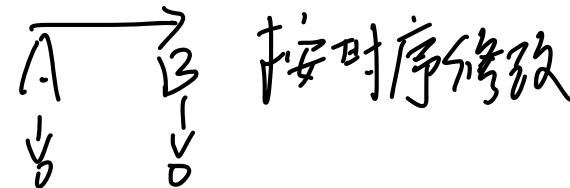
# Eventual Consistency Variations

- Monotonic read consistency *之后的不能看到A*

— If a process has seen a particular value of data item, any subsequent processes will never return any previous values

- Monotonic write consistency

— The system guarantees to serialize the writes by the *same process*



- In practice

— A number of these properties can be combined

→ Monotonic reads and read-your-writes are most

desirable

修改之后朋友看到，  
再次登陆查看还能是之前的  
facebook修改之后自己看到了，朋友看不到还是  
Cuz they need some time for all nodes get updated.

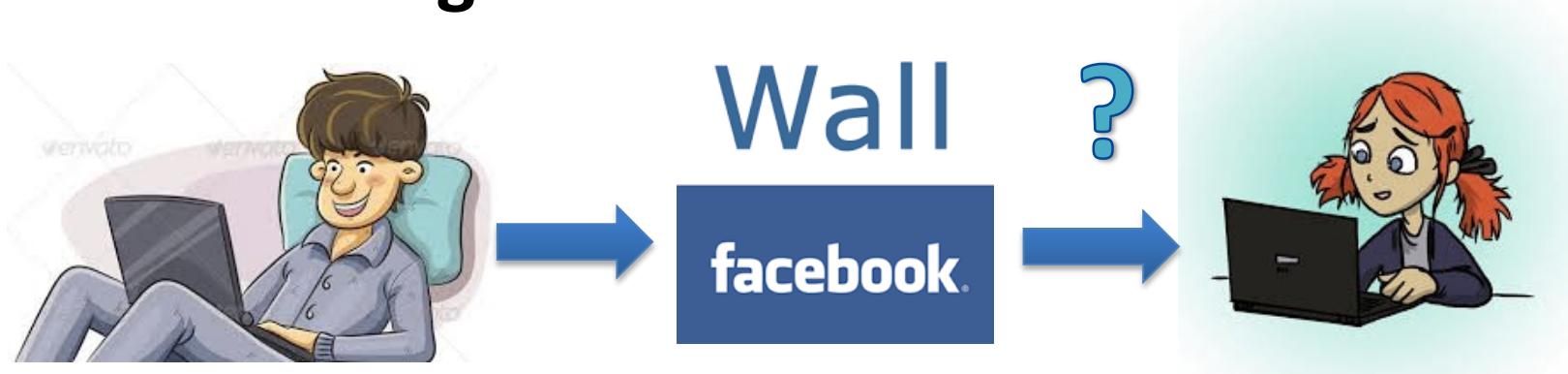
# Eventual Consistency

## - A Facebook Example

- Bob finds an interesting story and shares with Alice by posting on her Facebook wall
- Bob asks Alice to check it out
- Alice logs in her account, checks her Facebook wall but finds:
  - Nothing is there!



别的 post 及时看不见 (更新其他节点)  
- 需要时间



# Eventual Consistency

## - A Facebook Example

- Bob tells Alice to wait a bit and check out later
- Alice waits for a minute or so and checks back:
  - She finds the story Bob shared with her!

⇒过了会看到] 3. Face book上, 我朋友看到的新闻我看不到, 过了会儿才看到  
⇒eventual consistency. (the node I access from has not updated yet)



# Eventual Consistency

## - A Facebook Example

- Reason: it is possible because Facebook uses an **eventual consistent model**
- Why Facebook chooses eventual consistent model over the strong consistent one?
  - Facebook has more than 1 billion active users
  - It is non-trivial to efficiently and reliably store the huge amount of data generated at any given time

~~Strong consistency  
not always needed.~~ Eventual consistent model offers the option to **reduce the load and improve availability**

preference of the availability over strong consistency  $\Rightarrow$  無人感知使用 [即時],  
若实现 strong consistency, 有用户同时等待更新  $\Rightarrow$  annoying not desirable (银行等)  
for maintenance. (unlike bank, for social media,  
a bit less recent data is acceptable)

*Stock market . bank ⇒ Strong consistency is more desirable*

## Eventual Consistency

### - A Dropbox Example

- Dropbox enabled immediate consistency via synchronization in many cases.
- However, what happens in case of a network partition?



# Eventual Consistency

## - A Dropbox Example

- Let's do a simple experiment here:
  - Open a file in your drop box
  - Disable your network connection (e.g., WiFi, 4G)
  - Try to edit the file in the drop box: can you do that?
  - Re-enable your network connection: what happens to your dropbox folder?

# Eventual Consistency

## - A Dropbox Example

- Dropbox embraces eventual consistency:
  - Immediate consistency is impossible in case of a network partition
  - Users will feel bad if their word documents freeze each time they hit Ctrl+S , simply due to the large latency to update all devices across WAN
  - Dropbox is oriented to **personal syncing**, not on collaboration, so it is not a real limitation.

# Eventual Consistency

## - An ATM Example

- In design of automated teller machine (ATM):
  - Strong consistency appear to be a nature choice
  - However, in practice, **A beats C**
  - Higher availability means **higher revenue**

尽管银行追求  
强一致性，ATM will allow you to withdraw money even if the  
strong consistency machine is partitioned from the network

'however,  
However, it puts a **limit** on the amount of withdraw  
(e.g., \$200)

- The bank might also charge you a fee when a

In bank personal account, overdraft happens

Strong consistency can be relaxed,

在网络分区发生时，银行也可插入让客户取款（不强保证



Small risk -  
(limit amount)  
weak consistency  
+ partition tolerance  
⇒ although not consistent

文件修改的例子: (CAP)

Availability + Partition tolerance: 当 network partition 发生  $\Rightarrow$  保证 availability (可操作文件修改) 但仅对你可见  $\Rightarrow$  目标不变  $\Rightarrow$  磁盘再对齐 (但不同分区的值不同)  $\oplus$  available.

但要保证 consistency  $\Rightarrow$  当 network partition 发生时  $\Rightarrow$  就不能更新  
就不能操作

从 personal user 的角度: more desirable 是向用户提供服务.

$\Rightarrow$  still be able to update and access file when partition happen; (no freeze on my file).

---

Airline reservation: 预订机票  $\rightarrow$  强调 availability + consistency

(航班)取消: 表示余票量多少  
(E7A)  $\Rightarrow$  很多座位余量: (CAP) Weak consistency  $\Rightarrow$  close matter when we see 401 tickets left, or 399 left.  $\Rightarrow$  out of date data is acceptable.

not strong consistency,  $\Rightarrow$  keep availability during network partition.  
即使网络故障, 那也可下单.  $\Rightarrow$  无需等待系统 maintain (in strong consistency)

$\Rightarrow$  余量少  $\Rightarrow$  strong consistency. (ACID)  $\Rightarrow$  more desirable to see accurate data  $\rightarrow$  ensuring the plane is not overbooked.

# Dynamic Tradeoff between C and A

- An airline reservation system:
  - When most of seats are available: it is ok to rely on somewhat out-of-date data, availability is more critical
  - When the plane is close to be filled: it needs more accurate data to ensure the plane is not overbooked, consistency is more critical

AC → 红绿灯问题的变种 (因为在每个node  
要全部更新)

# Heterogeneity: Segmenting C and A

- No single uniform requirement *(用途 . 服务提供)*
  - Some aspects require strong consistency
  - Others require high availability
- Segment the system into different components
  - Each provides different types of guarantees
- Overall guarantees neither consistency nor availability
  - Each part of the service gets exactly what it needs
- Can be partitioned along different dimensions

Partitioning  $\Rightarrow$  按偏好分区，有的EAL, 有的AP, 有的CP.

不同区有不同程度的 consistency, availability or tolerance

shopping cart high availability, responsive

- ① Web上 products description, images : should be available, partition tolerance (无须等) 当 admin 修改时  
用户看最新的 (consistency) 但是可以接受 (很长时间)
- ② 对 shopping cart: add products to cart. (high availability)  
不想等 (no consistency)

加购物车时有的 item not available  $\rightarrow$  OK, just remove from cart.  
卖完

(inconsistency in terms of available products)  $\Leftarrow$  trade off between AVSL  
在 locks 中: 商品数减少 (有用冲突).  $\Rightarrow$  下单不了, 卖完了, 库存已扣过.  
 $\Rightarrow$  是一个厚边 (只从不同角度) { locks: consistency of data.  
(AP: shopping cart focus on more availability, not alway the updated info  
on the # of items available.)

checkout-  
shipping record

- ③ for checkout: financial bill involved: paying money: more updated info on: Whether the item is available for user  $\Rightarrow$  so that I can make payment.  
with OR for us to wait a little bit  $\Rightarrow$  maintain consistency and safety

locks中: exclusive lock for checkout.

Strong consistency over A  
when network partition happens

# Partitioning Examples

## Data Partitioning

- Different data may require different consistency and availability
- Example:
  - Shopping cart: high availability, responsive, can sometimes suffer anomalies
  - Product information need to be available, slight variation in inventory is sufferable
  - Checkout, billing, shipping records must be consistent

Availability: 网络分区后用户仍然可以访问 node 只不过不同分区数据可能不同，有的不是最新的。

# What if there are no partitions?

no network partition

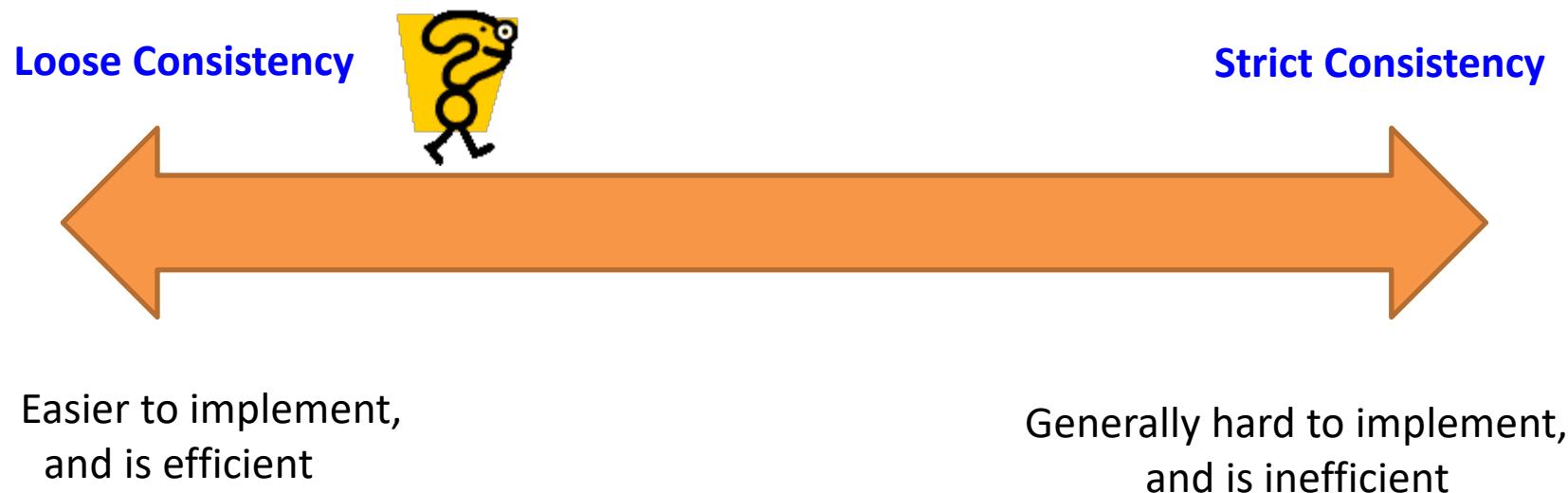
- Tradeoff between **Consistency** and **Latency**:
- Caused by the **possibility of failure** in distributed systems
  - High availability -> replicate data -> consistency problem
- Basic idea:
  - Availability and latency are arguably **the same thing**:  
unavailable -> extreme high latency [即使 system 是 available, 某些节点 down]  
– Achieving different levels of consistency/availability 也要等待时间,  
takes different amount of time  
  
latency.

# Trading-Off Consistency

- Maintaining consistency should balance between the strictness of consistency versus availability/scalability
  - Good-enough consistency *depends on your application*

# Trading-Off Consistency

- Maintaining consistency should balance between the strictness of consistency versus availability/scalability
  - Good-enough consistency depends on your application



# The BASE Properties

- The CAP theorem proves that it is impossible to guarantee strict Consistency and Availability while being able to tolerate network partitions
  - This resulted in databases with relaxed ACID guarantees  
*分区与公司订单的隔离 => relaxed ACID* *consistency is somewhat compromised.*
  - In particular, such databases apply the BASE properties:  
*(现代数据库都是这个)*
    - Basically Available: the system guarantees Availability
    - Soft-State: the state of the system may change over time
    - Eventual Consistency: the system will *eventually* become consistent
- partition tolerance.*  
*即分区发生=>不同分区能提供服务.*

why modern database rely on BASE rather than ACID:

- ① EC: ∵ large scale database are distributed  $\Rightarrow$  not possible to guarantee network tolerance  $\Rightarrow$  weaker consistency is applied to database. When system becomes eventually consistent, the all nodes are updated to the recent data, but it takes some time to propagate that info. State of system may change over time (whether focus more on consistency, or more constraint on data) State of system depends on application preference.

# CAP -> PACELC

- A more complete description of the space of potential tradeoffs for distributed system:
  - If there is a **partition (P)**, how does the system trade off **availability and consistency (A and C)**; **else (E)**, when the system is running normally in the absence of partitions, how does the system trade off **latency (L) and consistency (C)**?

Abadi, Daniel J. "Consistency tradeoffs in modern distributed database system design." Computer-IEEE Computer Magazine 45.2 (2012): 37.

consistency ↑ → latency ↑

## Examples

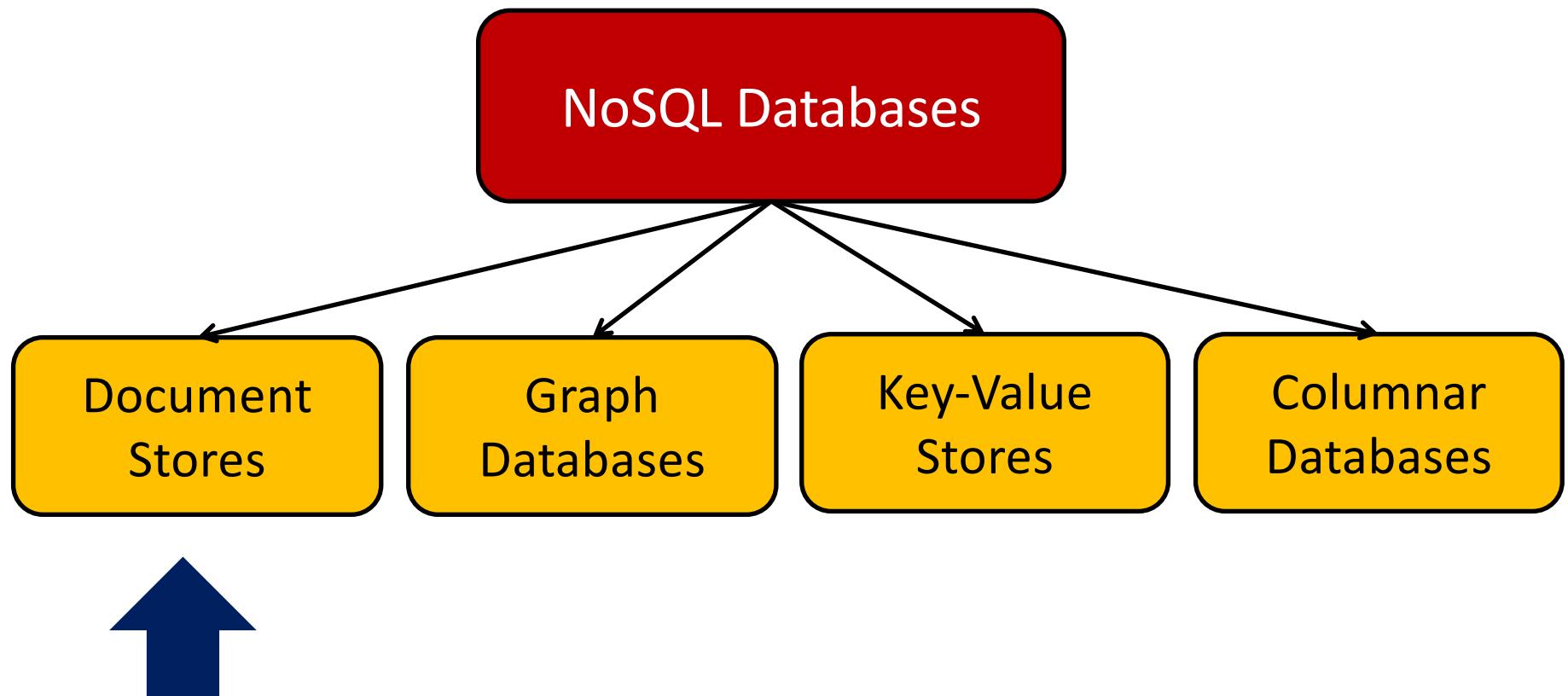
Partition tolerance always exists  $\Rightarrow$  network partition 不可避兔

- **PA/EL Systems:** Give up both Cs for availability and lower latency : system prefer to answer of accessing data with a shorter time by compromise consistency.
  - Dynamo, Cassandra, Riak
- **PC/EC Systems:** Refuse to give up consistency and pay the cost of availability and latency
  - BigTable, Hbase, VoltDB/H-Store
- **PA/EC Systems:** Give up consistency when a partition happens and keep consistency in normal operations
  - MongoDB
- **PC/EL System:** Keep consistency if a partition occurs but gives up consistency for latency in normal operations
  - Yahoo! PNUTS

# Types of NoSQL Databases

*BASE instead of ACID*

- Here is a limited taxonomy of NoSQL databases:

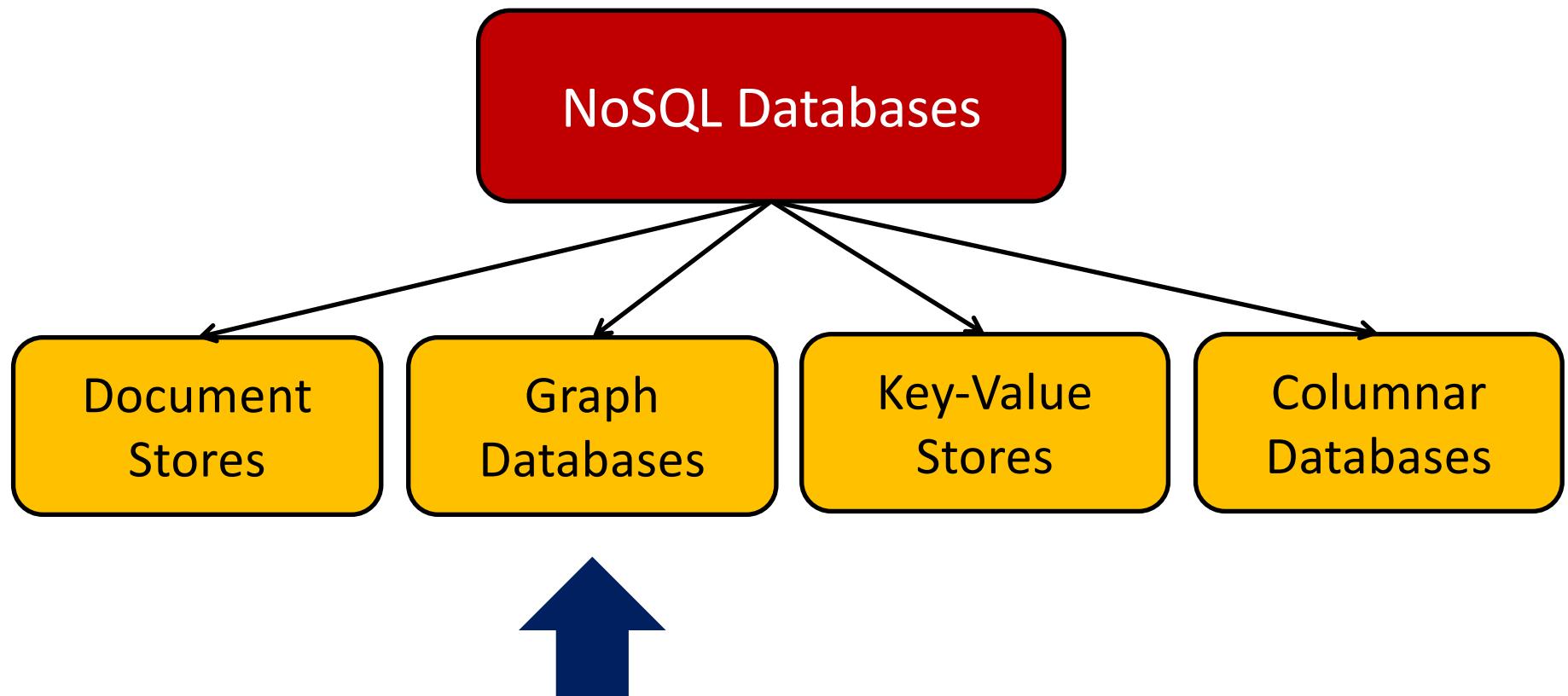


# Document Stores

- Documents are stored in some standard format or encoding (e.g., XML, JSON, PDF or Office Documents)
  - These are typically referred to as Binary Large Objects (BLOBs)
- Documents can be indexed
  - This allows document stores to outperform traditional file systems
- E.g., MongoDB and CouchDB (both can be queried using MapReduce)

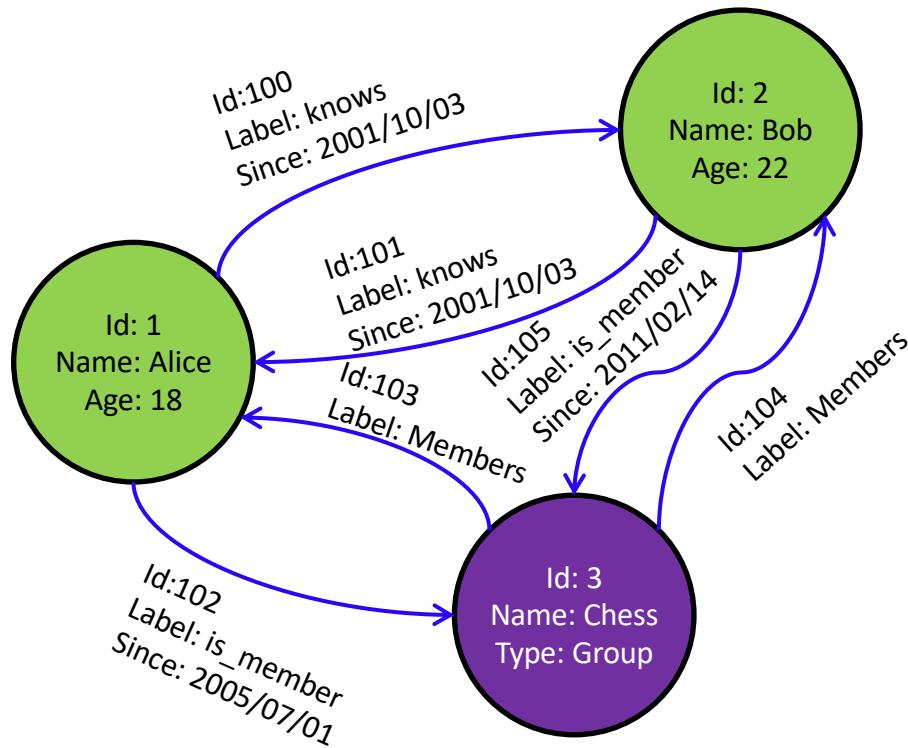
# Types of NoSQL Databases

- Here is a limited taxonomy of NoSQL databases:



# Graph Databases

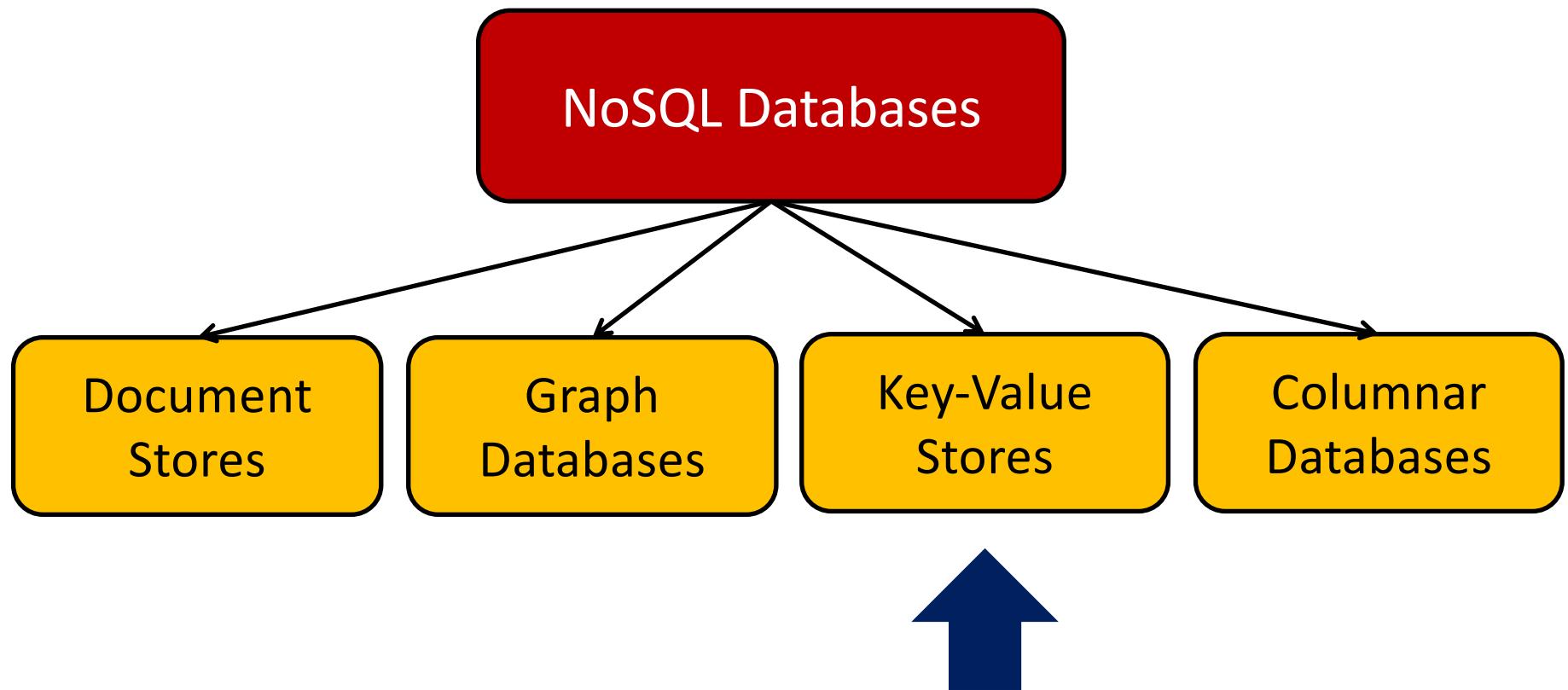
- Data are represented as vertices and edges



- Graph databases are powerful for graph-like queries (e.g., find the shortest path between two elements)
- E.g., Neo4j and VertexDB

# Types of NoSQL Databases

- Here is a limited taxonomy of NoSQL databases:

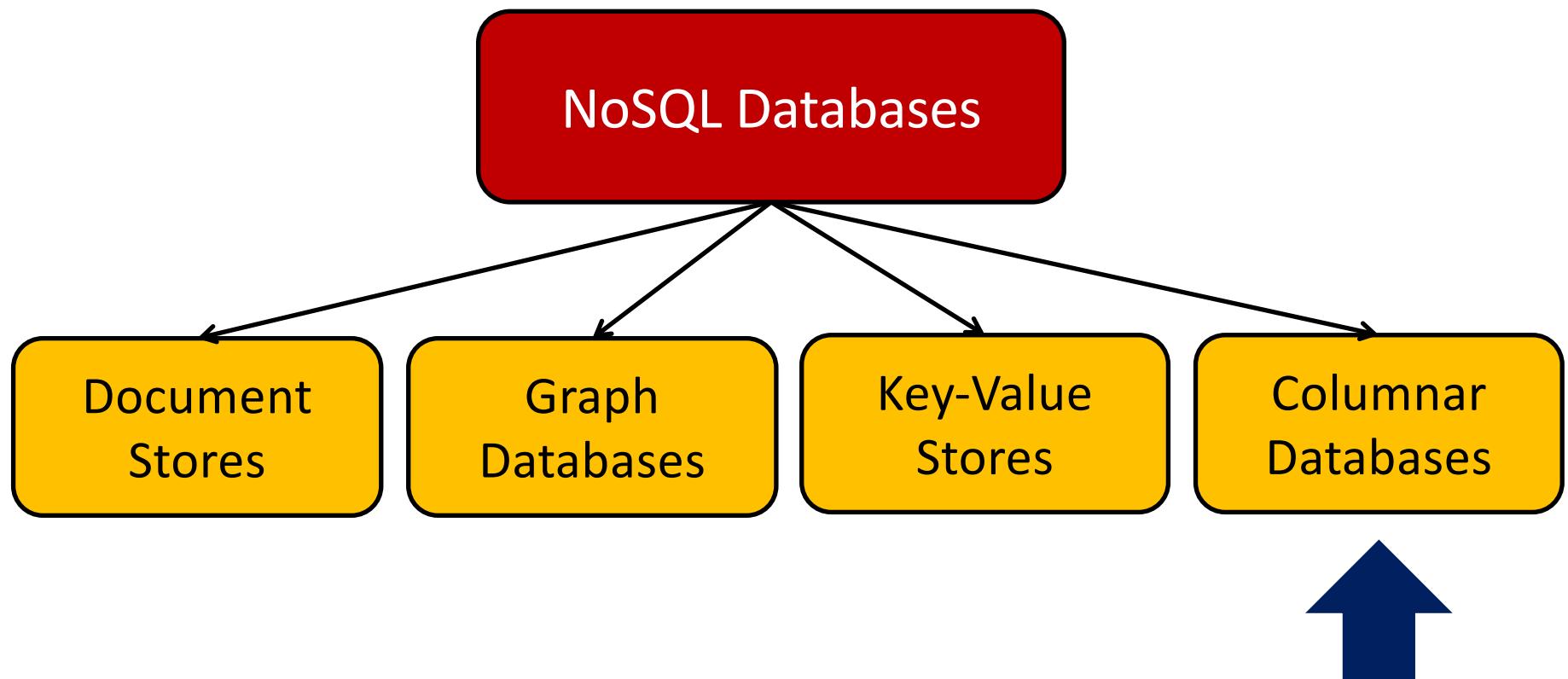


# Key-Value Stores

- Keys are mapped to (possibly) more complex value (e.g., lists)
- Keys can be stored in a hash table and can be distributed easily
- Such stores typically support regular CRUD (create, read, update, and delete) operations
  - That is, no joins and aggregate functions
- E.g., Amazon DynamoDB and Apache Cassandra

# Types of NoSQL Databases

- Here is a limited taxonomy of NoSQL databases:



# Columnar Databases

- Columnar databases are a hybrid of RDBMSs and Key-Value stores
  - Values are stored in groups of zero or more columns, but in Column-Order (as opposed to Row-Order)

Record 1

|       |    |       |     |
|-------|----|-------|-----|
| Alice | 3  | 25    | Bob |
| 4     | 19 | Carol | 0   |
| 45    |    |       |     |

Column A

|       |     |       |
|-------|-----|-------|
| Alice | Bob | Carol |
| 3     | 4   | 0     |
| 19    | 45  |       |

Column A = Group A

|       |     |       |
|-------|-----|-------|
| Alice | Bob | Carol |
| 3     | 25  | 4     |
| 0     | 45  | 19    |

Row-Order

Columnar (or Column-Order)

Columnar with Locality Groups

- Values are queried by matching keys
- E.g., HBase and Vertica

# Summary

- The *CAP theorem* states that any distributed database with shared data can have at most two of the three desirable properties:
  - Consistency
  - Availability
  - Partition Tolerance
- The CAP theorem leads to various designs of databases with *relaxed* ACID guarantees

# Summary (*Cont'd*)

- *NoSQL* (or *Not-Only-SQL*) databases follow the *BASE properties*:
  - Basically Available
  - Soft-State
  - Eventual Consistency

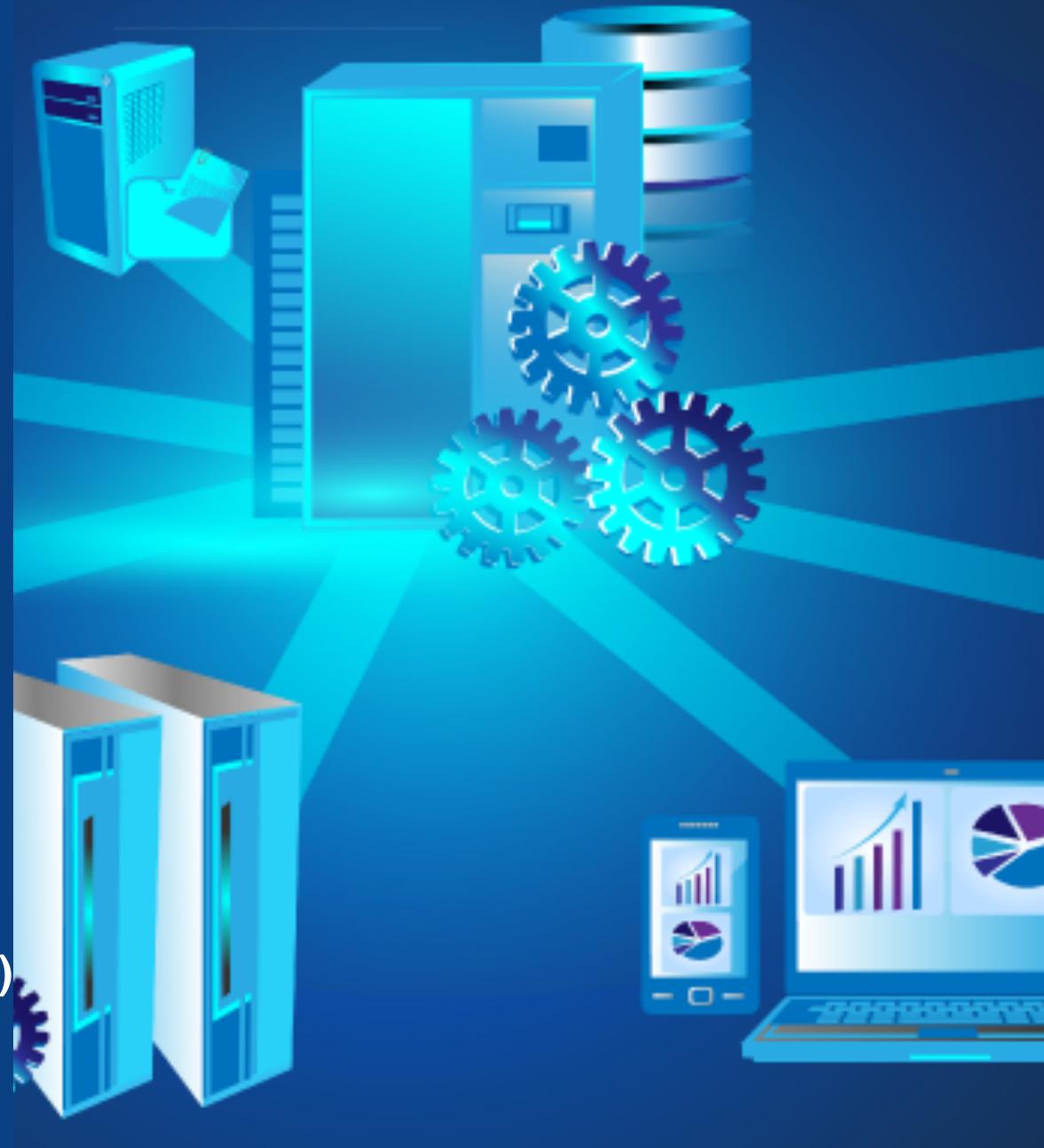


# COMP90050 Advanced Database Systems

## Semester 2, 2024

Lecturer: Farhana Choudhury (PhD)

Live lecture – Week 11





# Today's topics

- Activities on distributed databases
- Q/A
  - Project, final exam, quiz questions



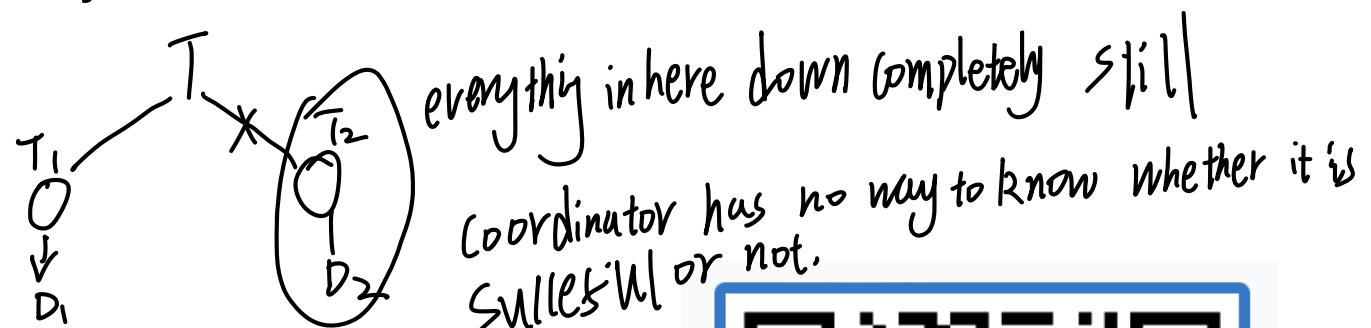
# DBMSs in the era of Networking

For distributed transaction processing:

- Atomicity: two-phase commit protocol

# Atomicity in distributed DB

In a distributed DB, a transaction T is being executed as T<sub>1</sub> in node D<sub>1</sub> and as T<sub>2</sub> in node D<sub>2</sub> (T<sub>1</sub> and T<sub>2</sub> are parts of T). Both T<sub>1</sub> and T<sub>2</sub> successfully finished their operations in their own nodes, but D<sub>2</sub> lost connection with the coordinator, and couldn't send the 'yes' vote in phase 1. Can T successfully commit?



**Time for a poll –**

[Pollev.com/farhanachoud585](https://Pollev.com/farhanachoud585)

To maintain Atomicity,  
both T<sub>1</sub> and T<sub>2</sub> need to be  
completed because overall  
it's one transaction.

So even if T<sub>1</sub> committed, both of T<sub>1</sub> and T<sub>2</sub>



will be aborted to maintain atomicity.

信息传递不过去他们也无法 Commit.





# DBMSs in the era of Networking

For distributed transaction processing:

- Atomicity: two-phase commit protocol
- Concurrency:
  - Lock release
  - Timestamping
  - Optimistic locking



# Concurrency in distributed DB

A local lock manager cannot release any locks until it knows that the transaction has been committed or aborted at all the servers involved in the transaction.

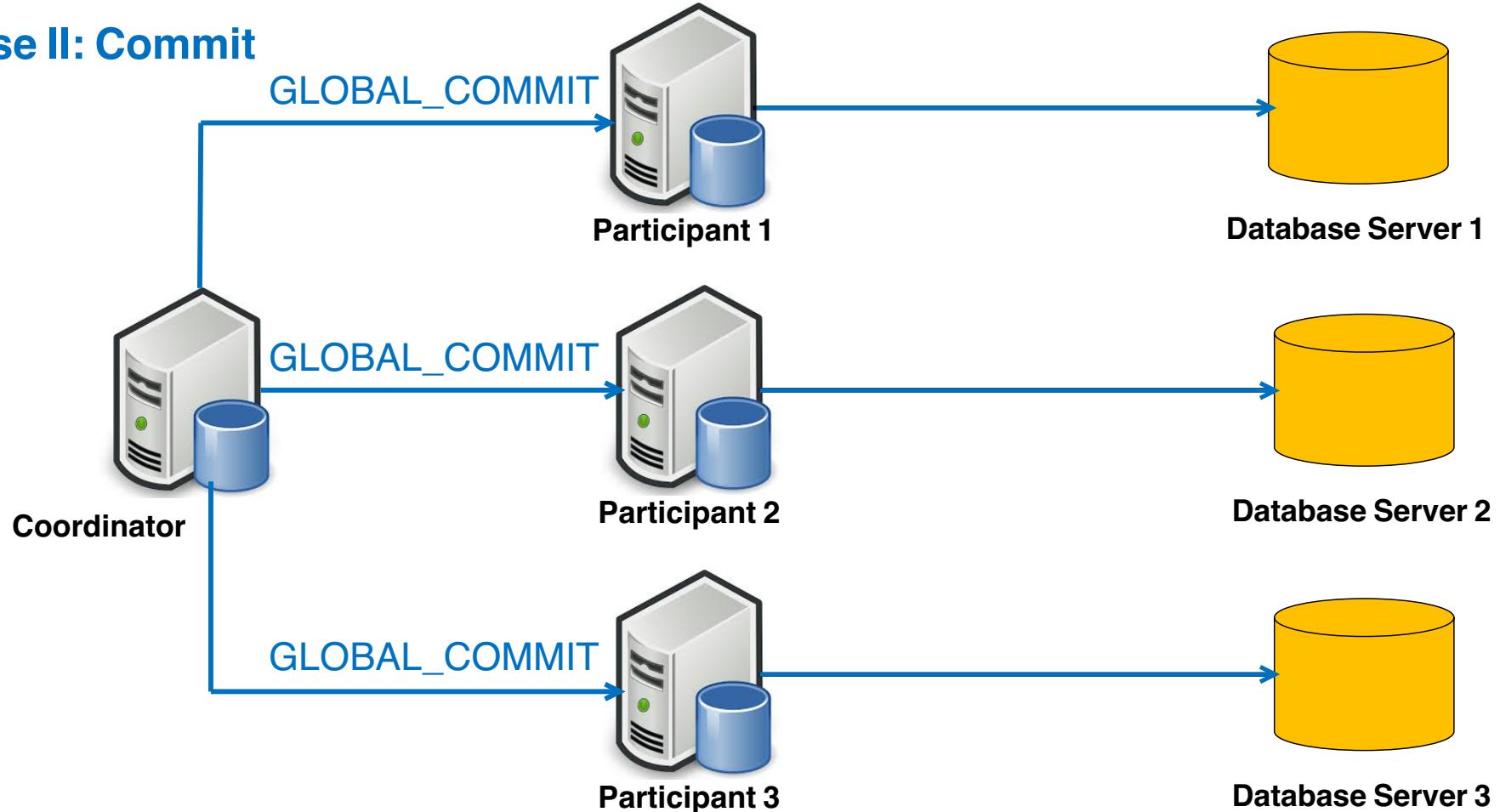
When/how will a node know when to release the locks of a successfully transaction?

When global\_commit sent, it means that coordinator tells <sup>the</sup> participants : others have finished their part of the transaction, if they receive that, then each of them can release their own locks for that transaction.

当一个node down 了, 这个问题还涉及到 [2]: When any of sub transactions abort, others will also abort as well. So In phase I, the down node ( $T_i$ ) sending 'no' message at that time,  $T_i$  releases its lock in phase I, others will know this in phase II, and then they will release their locks of  $T_i$ .

# When/how will a node know when to release the locks of a transaction?

## Phase II: Commit



What if a sub-transaction aborts in a node?



# DBMSs in the era of Networking

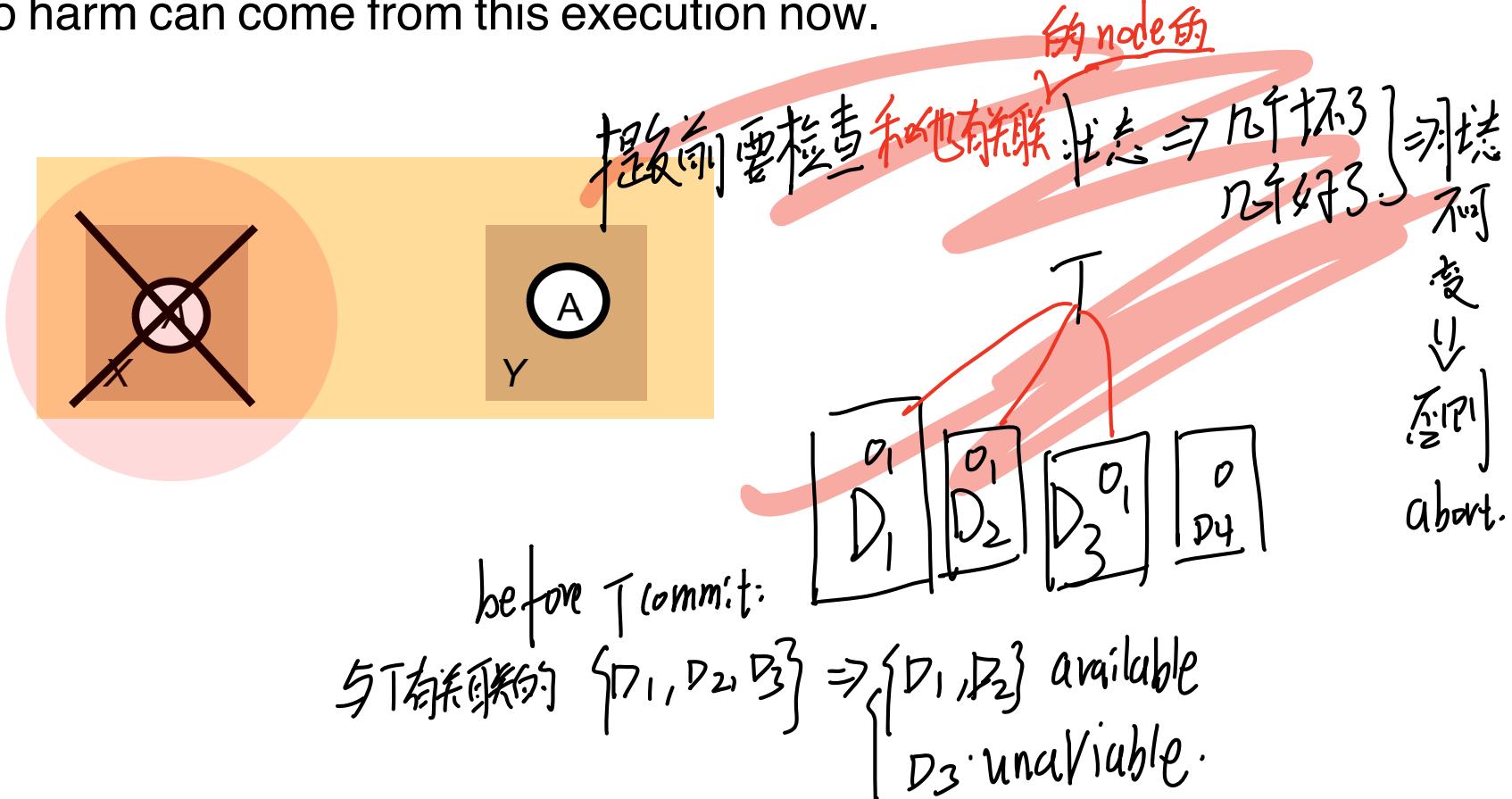
For distributed transaction processing:

- Atomicity: two-phase commit protocol
- Concurrency:
  - Lock release
  - Timestamping
  - Optimistic locking
- Data replication
  - Rule for one-copy serializability
  - CAP However, not always desirable because data replication come with some benefits.

# Available copies replication rule

Before a transaction commits, it checks for failed and available servers it has contacted, the set should not change during execution : make sure consistency.

- E.g.,  $T$  would check if  $X$  is still available among others.
- We said  $X$  fails before  $T$ 's deposit, in which case,  $T$  would have to abort.
- Thus no harm can come from this execution now.





☆☆☆ After committed  $\Rightarrow \{D_1, D_2\}$ , available

~~D<sub>3</sub> unavailable~~  $\rightarrow$  available, D<sub>3</sub> unavailable

## Available copies replication rule

但 D<sub>4</sub> 与 T 无 contact  $\Rightarrow$  不影响执行.

still the same set  
Compare with starting time -

In a distributed DB, a transaction T only needs to read an object o1. o1 is replicated in the nodes D1, D2, and D3 of this distributed system. There is another node D4, but D4 does not store any copies of o1.

When T started its execution, the nodes D1, D2 are available, and the other nodes are unavailable. When T commits, the nodes D1, D2, D4 are available and D3 is unavailable. Can T successfully commit?

Time for a poll –

[Pollev.com/farhanachoud585](https://pollev.com/farhanachoud585)



# The CAP Theorem

- ? The limitations of distributed databases can be described in the so called the CAP theorem

CAP theorem: any distributed database with shared data, can have at most two of the three desirable properties, Consistency, Availability, or Partition tolerance

# Different design choices based on CAP theorem: Examples

- Different data may require different consistency and availability
- Example:
  - Shopping cart: high availability, responsive, can sometimes suffer anomalies/inconsistencies
  - Product information need to be available, slight variation in inventory is sufferable
  - Checkout, billing, shipping records must be consistent

different part of application have different level of consistency -

For Browsing the product : inconsistency / weak consistency on # of items available ; description on items are acceptable and tolerable .

For checking out, billing : Strong consistency, otherwise we'll lose money

⇒ locking mechanism  
more strict locking mechanism VS relaxed locking mechanism .

Strict: even finish transaction don't release locks (writing, making changes).  
relaxed locking: reading (甚至读到不同的值) (for reading more relaxed locks used).

for adding items into shopping cart : sometimes we'll find there's items are unavailable,

so we just remove it from shopping cart → available by the time of check out ,

Someone has bought . (inconsistency)

readily lock & relaxed lock: 所有人都想读这个产品 and all need to add to the shopping cart, (inconsistency)

在 distributed system 中：每个 node 存着 item count, 如果 item 被用户购买了 # of item has not updated immediately by coordinator ⇒ We'll also meet the same situation

with readily lock .

# Different design choices based on CAP theorem: Examples

More real-life applications/examples where high availability is important, can sometimes tolerate some inconsistency?

For social media Apps, usually there are millions of users, their data are stored across many nodes. The nodes you are reading from maybe not the same node that your friends are reading from. And the new comment hasn't been propagated to all those nodes (so in reality, if we change or edit some post, those changes may not get immediately propagated to everyone  $\Rightarrow$  inconsistency but API).



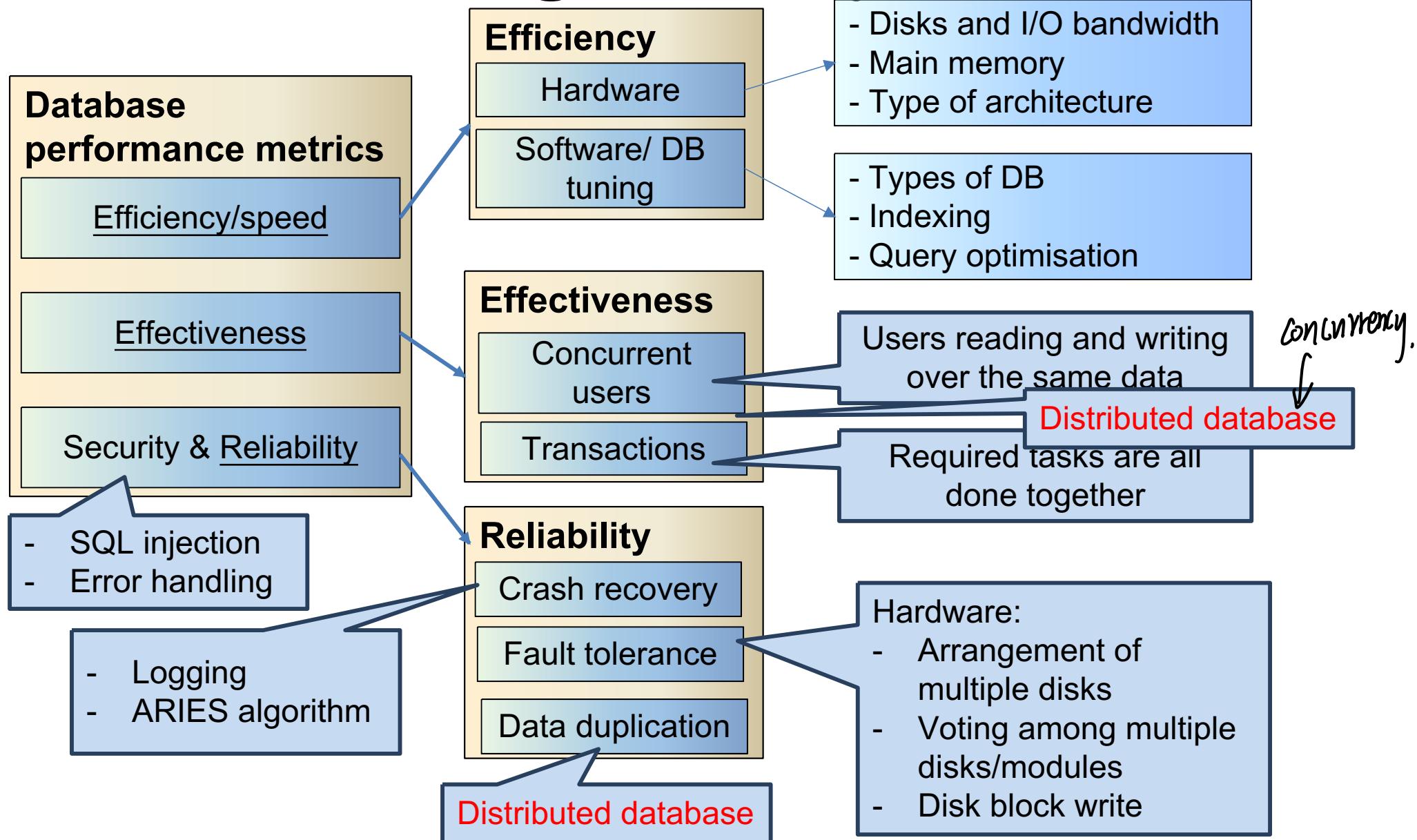
flight prices, hotel prices 数据库的机制可以随应用状态、数据切換, (emphasize).  
inconsistency → 是回多 or → 但看到的还是 old price  $\Rightarrow$  new price hasn't  
been propagated to it.  $\Rightarrow$  等一會再看就更新了, 但不影响当前使用 (waited consistency)  
although the updated info would be propagated to all nodes eventually,  
but give some higher priority for something.

When room is empty  $\Rightarrow$  It will keep consistency eventually, but it goes down  
in the priority list.

Inconsistency: 我看见、我朋友看还没 了。 (没被 propagated).

If rooms have been almost fully booked  $\Rightarrow$  emphasize on consistency

# Core Concepts of Database management system



RAID 0 with 3 disks no replica , no additional info.

flush LSN  $\geq$  page LSN

DPT records the  $\bigtriangleup$  next LSN for  $\bigtriangleup$  each page .

Q/A