



2110625 - Data Science Architecture

Foundation of Scalable Architecture

Asst.Prof. Natawut Nupairoj, Ph.D.

Department of Computer Engineering
Chulalongkorn University
natawut.n@chula.ac.th

Scale Up vs Scale Out



មួយគ្រប់ងតែនឹងអាយុវត្ថុខ្លោយខ្សោយ

ទំនើស - ម៉ាស៊ីម

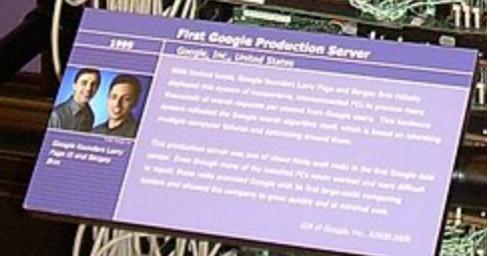
នៅលើ - network ការពេញចិត្ត សំណង់ថ្មី សេវានិភ័យរៀងរាល់
r/o bandwidth

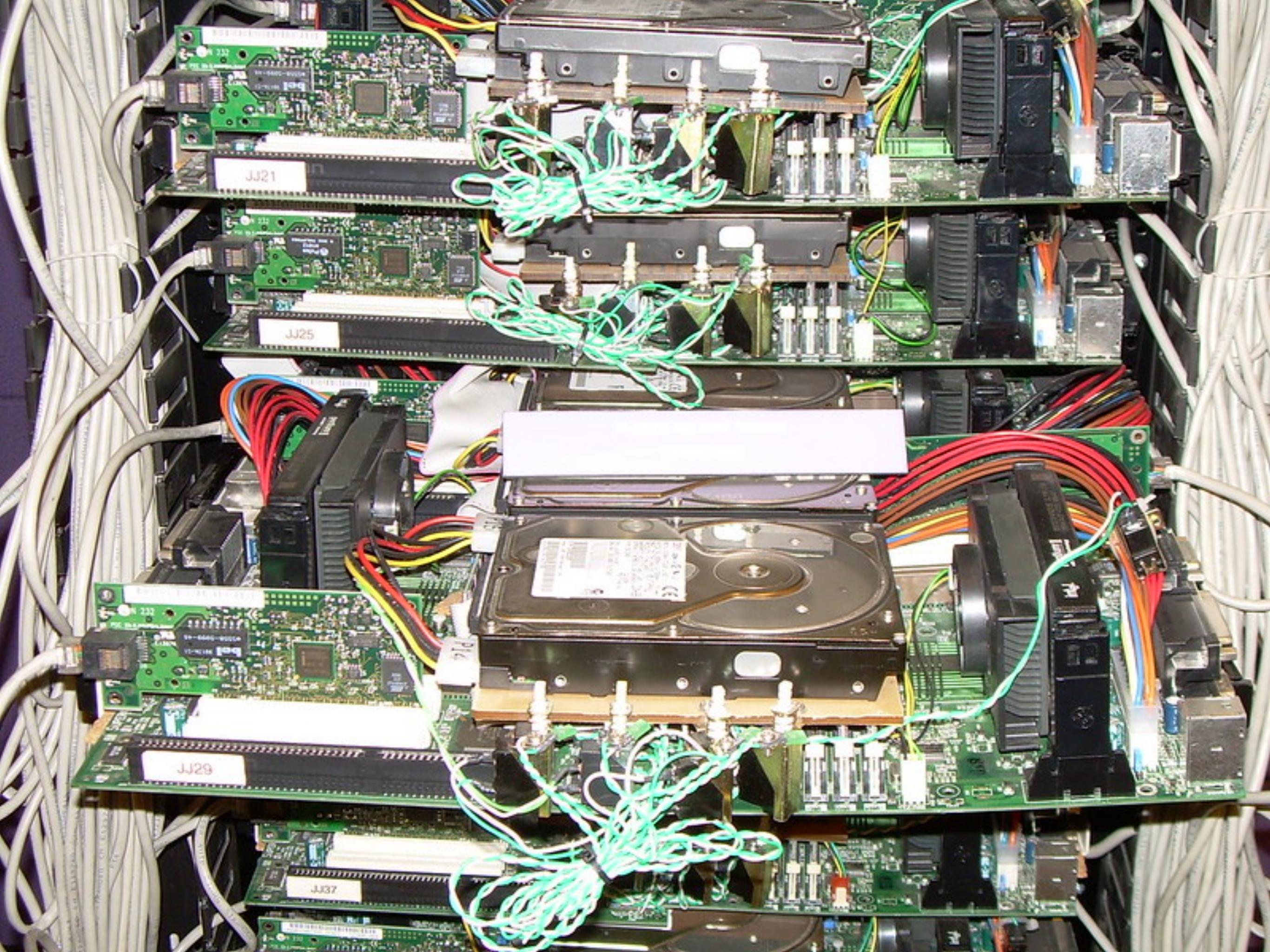


ធ្វើឱ្យជូនស្ថុ

Google
Cluster v.1

Rollable สูง
ต้านไฟฟ้าสถิต
กันฟลัก HW ใหม่
แล้ว auto sync
60V





Nature of Big Data Architecture

- Base on distributed and parallel processing paradigm
- Hardware or software components located at networked computers communicate and coordinate their actions only by passing messages
 - fast date នៃ ចំណុលកក censos
 - គីរារិន big data
- This leads to
 - Concurrency
 - distributed និង function ប្រមូលធម៌នការព័ត៌មាន
 - parallel និង performance
 - No global clock
 - Independent failure

Concurrency big data និង scale out គឺជាគារងាយការ

- Large number of resources (or nodes) are working together
 - Tencent - from 8,800 servers in 2014 to > 30,000 servers in 2017
 - Data is usually replicated into multiple copies to improve performance and availability **replication** គេងក្រម duplicate data ដើម្បីលើ server
 - Coordination and agreement are required
 - ensure pace of execution - everyone must reach at some execution point before continue **synchronization** នាំឱ្យតុលិភាគរបស់ពួកគេទៅការងារ អវត្សមក្តុំ / ការងារលេចក្បាតក្នុង
 - handle shared resources - data, devices
 - maintain **data consistency**

បែងការពីរវោសការលេខរបៀប cluster No Global Clock *** នឹងមានការចាប់ផ្តើម clock ក្នុងក្រឡុងវិវត្ថាកសនីភោគ (នានា nano sec)

- Coordination of programs' actions depends on the shared idea of time, but computers in a network cannot accurately synchronize their clocks សម្រាប់ clock មួយតួច precise ដូចជា ក្នុងក្រឡុង
- Each computer has its own internal clock; there is no single global clock for every computer
- Thus, no two computers running independently can have the same clock values
- How can we know the **order** of executions?
 - This is important for guaranteeing correctness e.g transaction តាមការណែនាំ ក្នុង global clock ដោយនឹងការទិន្នន័យ

Independent Failure

ຂະໜາດ ໂຕສັກະນະ

- Each component of the system (e.g. network, computer, program) can fail independently, leaving the others still running
- Google's studies showed 3% of hard disks will be failed after 3 months and up to 8% after 2 years
 - for 30,000 servers, 900 servers will experience hard disk failures within 3 months
- How can we **mask** (to hide or make it less severe) these failures?
 - ກຳປັບຄິດກວບນັ້ນໄປຮູ້ສຶກສາ failed
 - System can continue executing the jobs **correctly** even when failures occur

Ordering in Distributed Systems

តីមកាត់បែង
កំណើនពេលក្រែង

- Execution of distributed systems usually bases on series of events
- An event is the occurrence of a single action that a process carries out as it executes
- Since we cannot synchronize clocks perfectly, we cannot in general use physical time to find out the ordering of events that occur at different computers
- We usually rely on **ordering of events**, rather than exact timestamp of occurrence e.g. we are certain that A is executed before B even if A and B are on different nodes

A នៅលើ node
B នៅផ្លូវការ node

Ordering Models

- Total ordering *เรียงลำดับทุก event ได้แน่นอน*
 - Each event can be placed in a specific order, which is based on what time they occurred
- Partial ordering *มีนาฬิกาโลกที่ No global clock*
 - We do not know the exact order of every event, just the order of certain events that rely on one another
 - Causal ordering is a special case of partial ordering based on Happened-Before relation

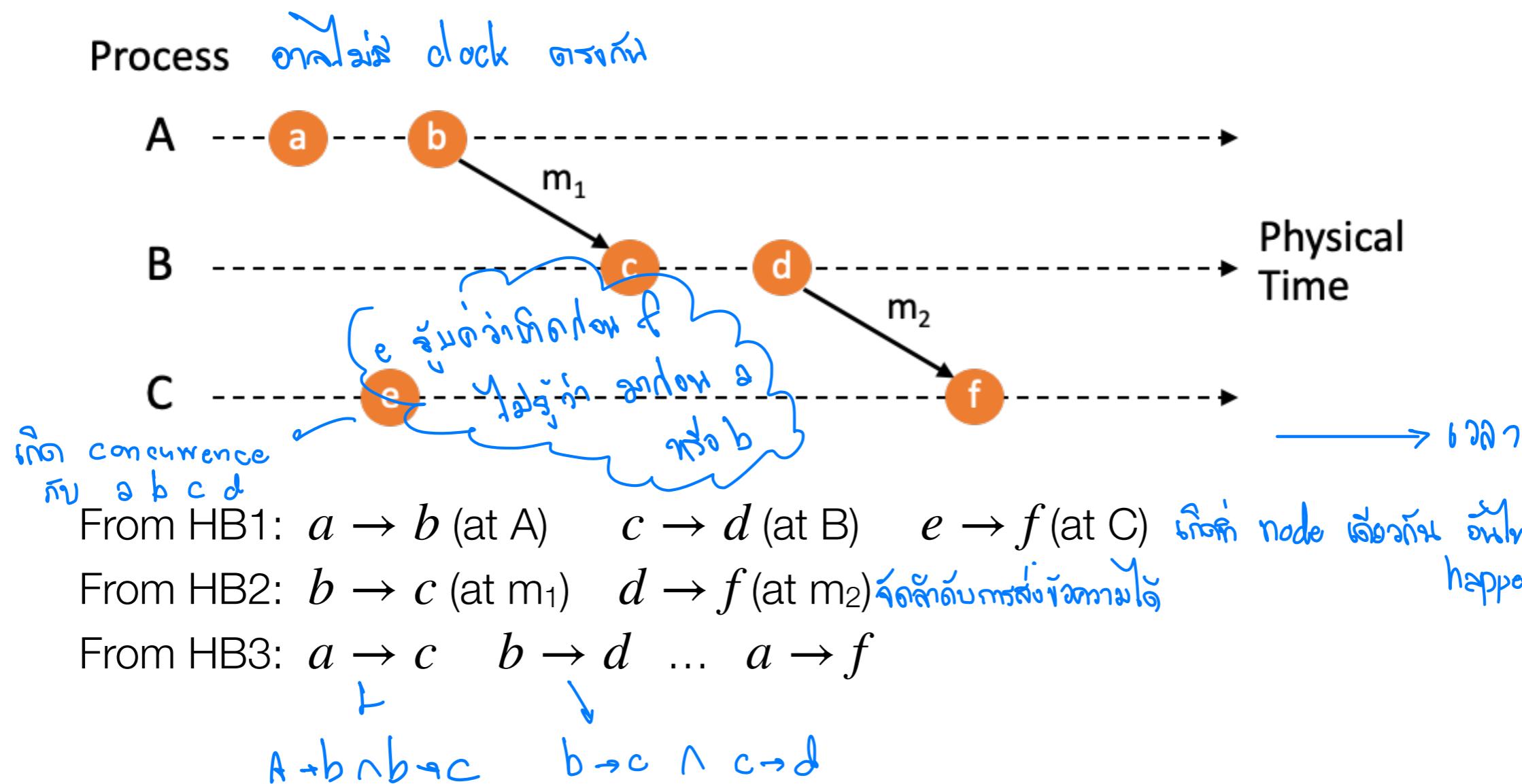
Happened-Before Relation

- A partial order of events that reflects a flow of information
 - HB1: for any events in process A, if e is executed before e' then $e \rightarrow e'$ (e is happened before e') e เกิดก่อน e'
 $\xrightarrow{\text{ถ้าฝ่ายนับข้อความ } m \text{ event } m \text{ แล้ว}}$ $e \rightarrow e'$ ต้องเกิดก่อนท่านรู้จัก m
 - HB2: for any message m, $send(m) \rightarrow recv(m)$ $send(m)$ เกิดก่อน $recv(m)$
 $\xrightarrow{\text{ถ้าฝ่ายนับข้อความ } m \text{ event } m \text{ แล้ว}}$ $send(m) \rightarrow recv(m)$ (คู่มีสมบัติความสอดคล้อง transitive)
 - HB3: if $e \rightarrow e'$ and $e' \rightarrow e''$ then $e \rightarrow e''$ e เกิดก่อน e'
 $\xrightarrow{e' \text{ เกิดก่อน } e''}$ e เกิดก่อน e''
 $\xrightarrow{\text{คู่มีสมบัติความสอดคล้อง transitive}}$

subset
 FIFO ⊂ Causal ordering
 ถ้าเกิดมีเส้นที่ยาวกว่า

Causal Ordering

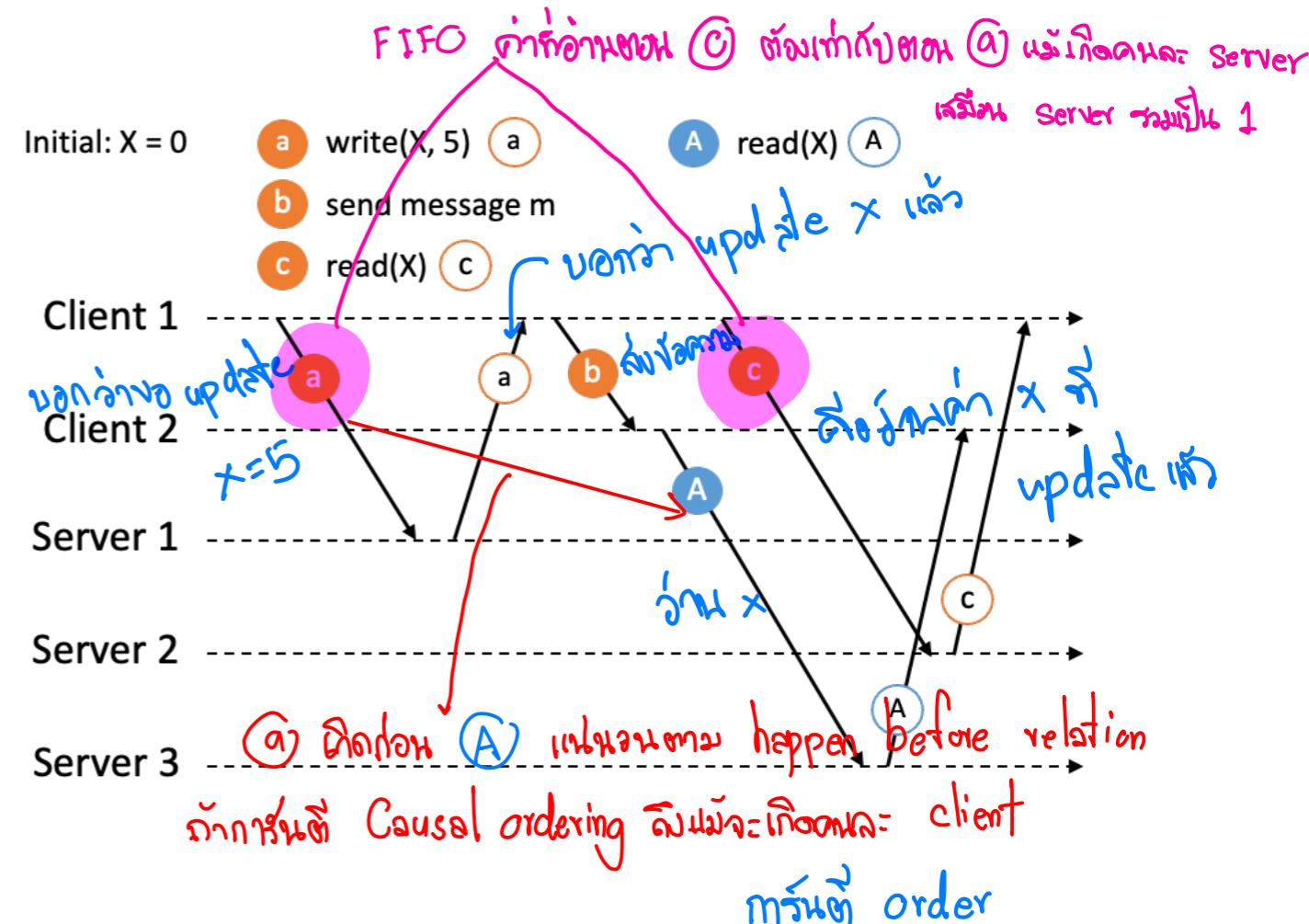
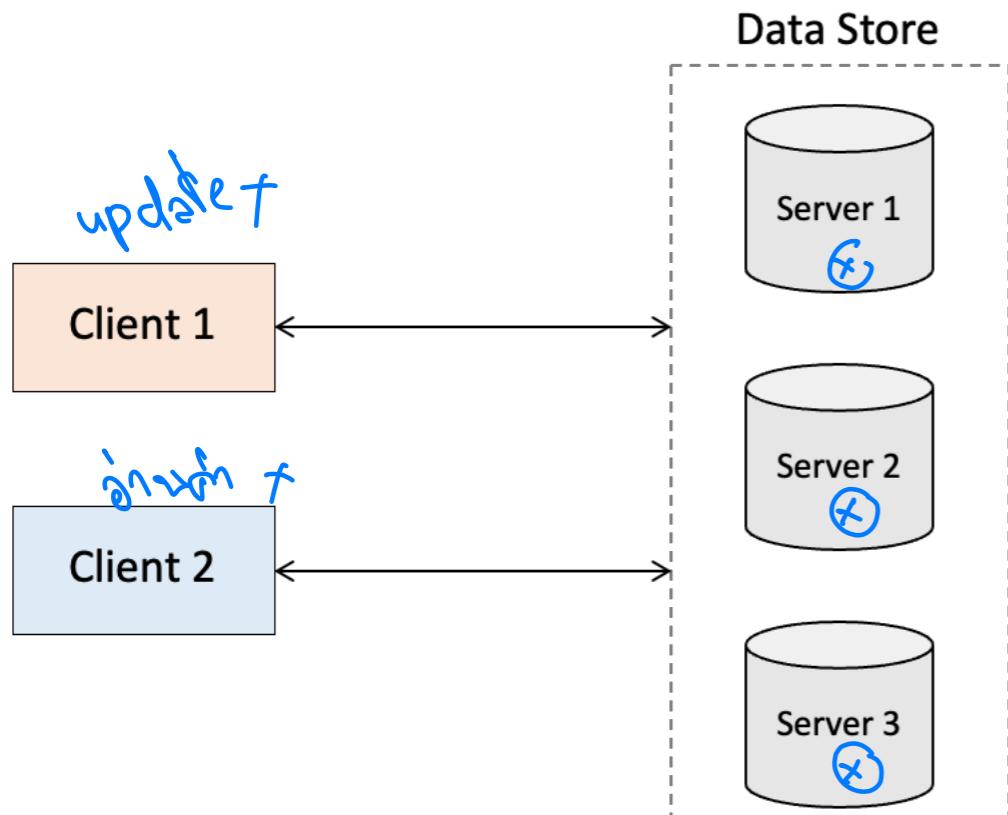
จัดตามลำดับที่เกิดก่อน happen before relation



Not all events are related by Happened-Before relation

Consider a and e (different processes and no chain of messages to relate them), they are not related by \rightarrow ; they are said to be concurrent; write as $a \parallel e$

Ordering and Data Store



- Suppose a client can send a request to any server and the system guarantees ordering
รูปแบบนี้ > 2 ภาระ node เต็มที่กัน ผลลัพธ์จะเป็น execute ตามลำดับที่เกิดขึ้นจริง
- For FIFO ordering: as long as execute (c) after (a)
คำศัพท์ว่า (C) ต้องเกิดก่อนหากไม่เกิดก่อนแล้วค่า ก ที่ update
- Causal ordering - as long as execute (A) after (a)
ผลลัพธ์จะเป็น client 2 จะเท็จก่อน client 1

Total Ordering

ที่ server ทุกตัว จะ执行 คำสั่ง
ที่ถูกส่งมาโดยลักษณะเดียวกัน

$$X = 5 - 10 - 15, 15 - 10 - 5 \text{ ได้ผล}$$

ต่อไปนี้จะ update ตามลำดับ ไม่สำคัญ

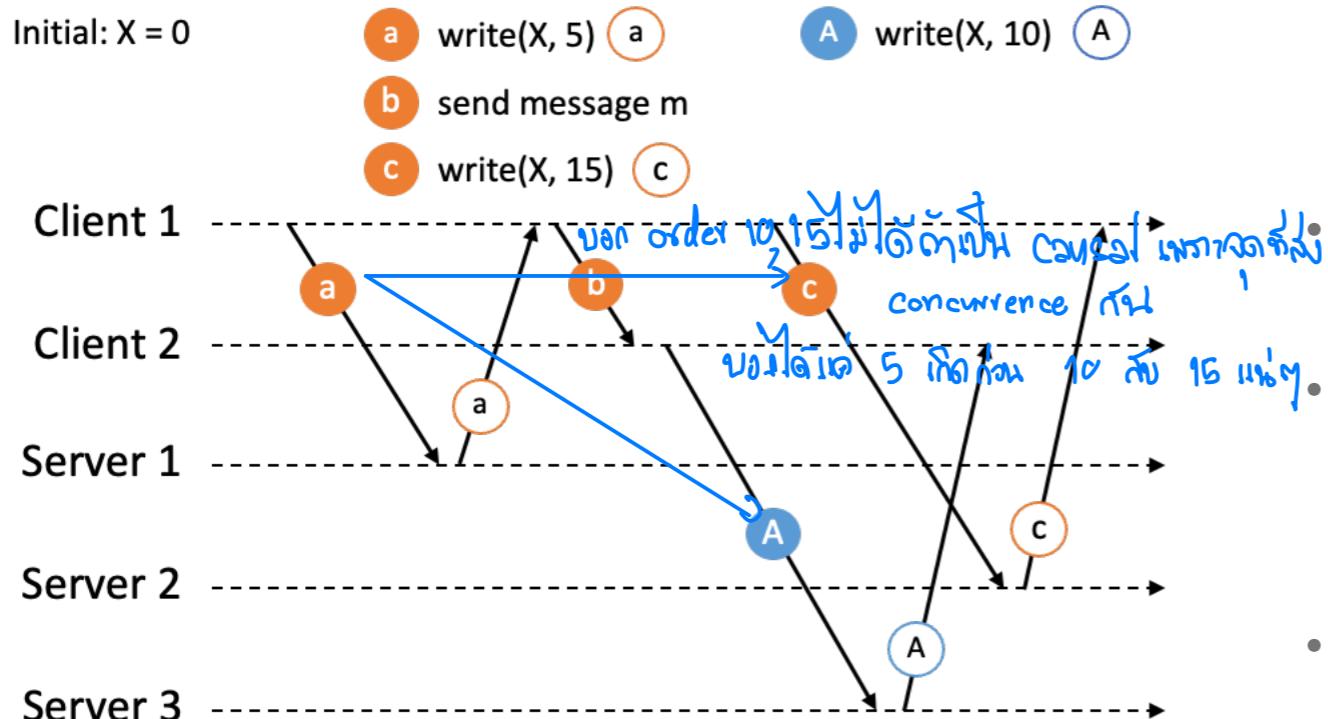
server ทุกตัวจะ execute ด้วยลักษณะเดียวกัน

- All servers will execute requests in the same order

Arbitrary-total - in any order
ซึ่งแต่ละคนจะจัดลำดับอย่างไร

FIFO-total - as long as
execute (c) after (a) คำสั่งที่มีเวลาเดียวกัน FIFO
ให้คำสั่ง update 5 ก่อน update 15 และ update 10 อยู่ทางขวา

- Causal-total - as long as
execute (A) after (a)



คำสั่งที่มีเวลาเดียวกัน FIFO

ให้คำสั่ง update 5 ก่อน update 15 และ update 10 อยู่ทางขวา

ไม่ต้องคำนึงถึง FIFO

$5 - 10 - 15, 15 - 10 - 5$ ได้ผล

ยกเว้น $15 - 5 - 10$

คืน $15 - 10 - 5$

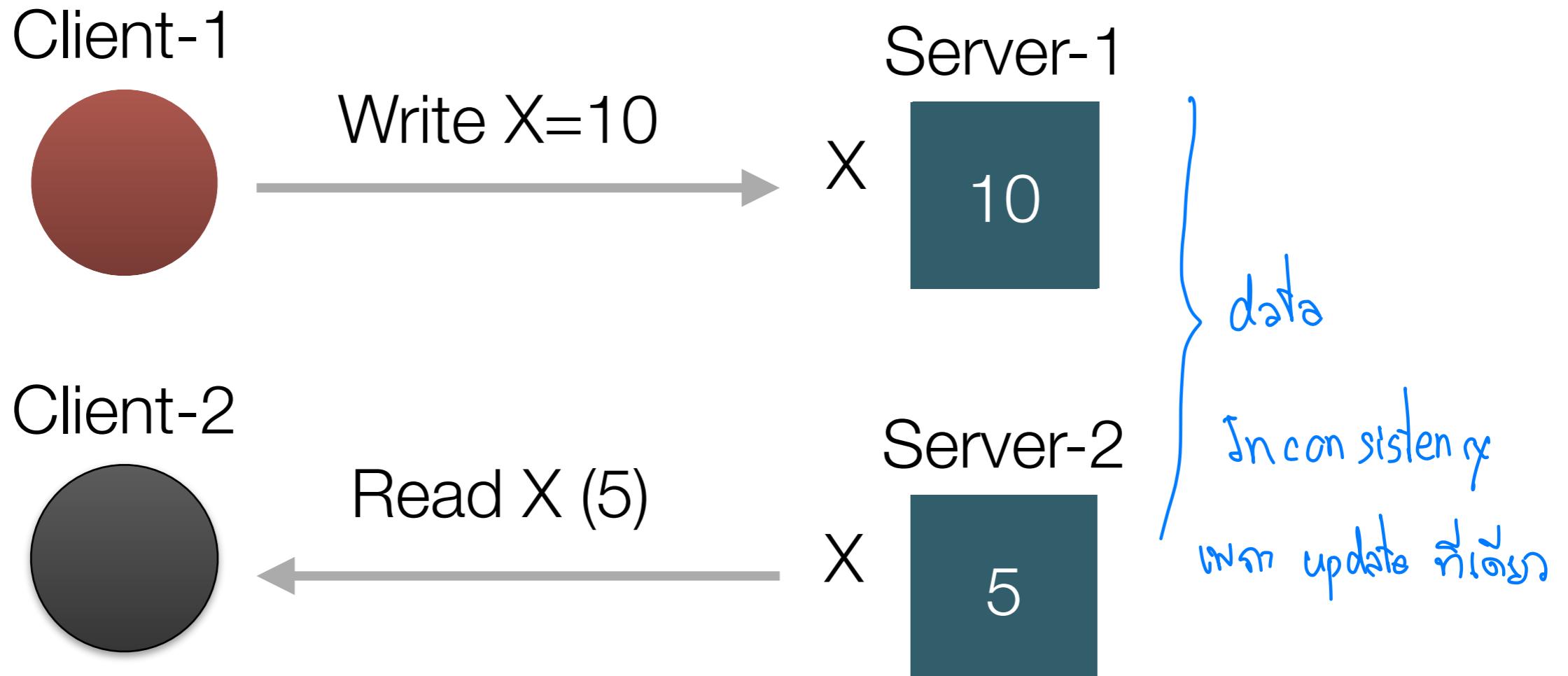
Scalability and Strong Consistency Model

- For scalability, the same data is stored on multiple servers (replica)
- Lead to difficulties to keep all copies in-sync
• Concurrent write
• Read and write on different replica
- Even more problems if
 - Data is replicated to lots of servers (hundreds or thousands)
 - Servers are distributed around the world

ក្នុង clone នៅលើ server
មិនបានបាត់បង្ហី

R/W នៅ server

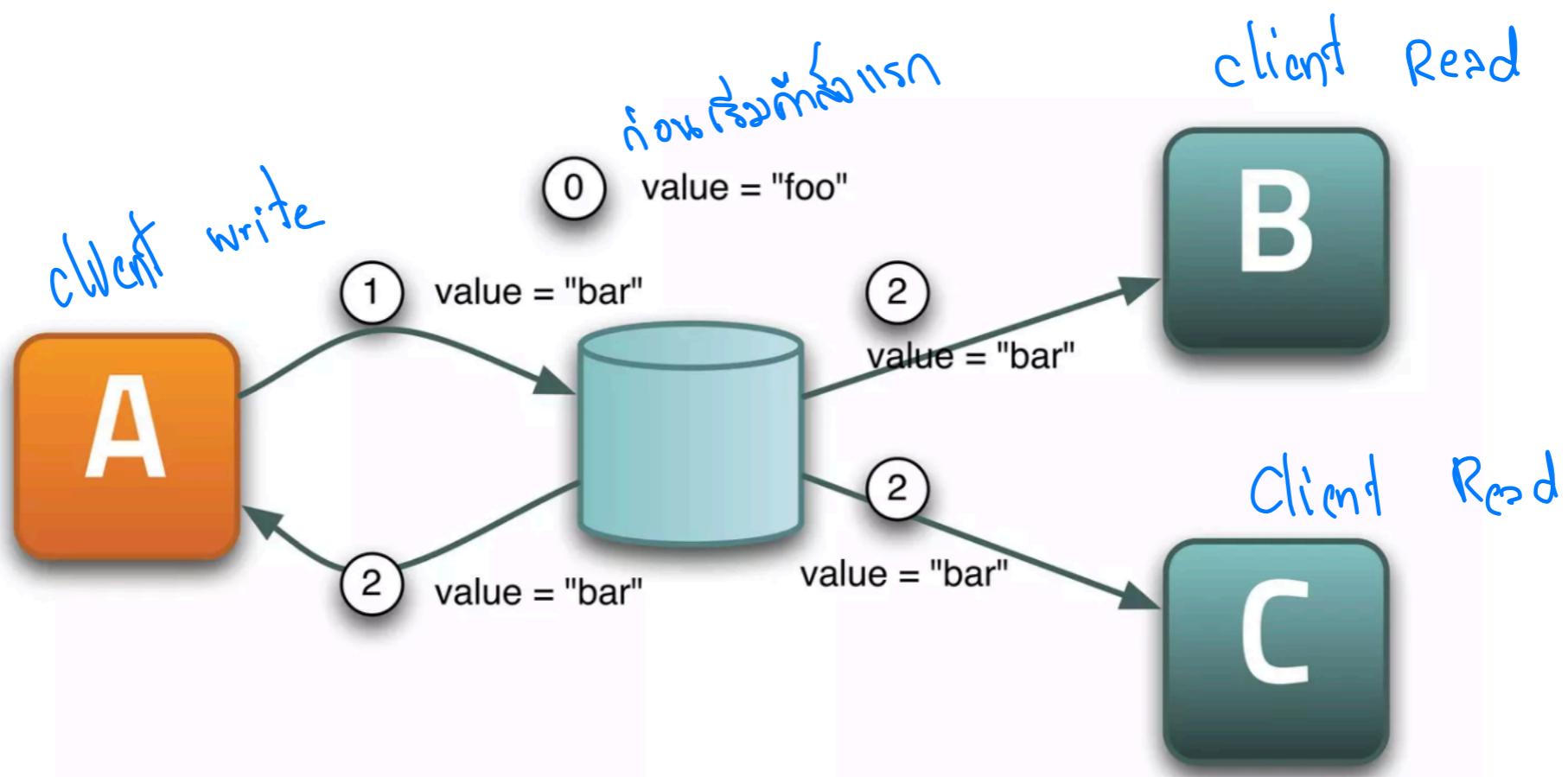
Replication and Data Consistency Problem



Consistency Models (based on Vogels' Paper)

- Strong consistency ເສົ້າມືອນດຳທຽບກັນ ແລ້ວອະນຸຍາ
 - After update completes, any access will return the updated value
 - High costs for large scale system
- Weak consistency ຢາງໄຫຼຍ້ວຍເລກຂຶ້ນທີ່ data ທີ່ມີກາດດັລວ
 - Certain conditions must be met before the consistency is guaranteed
 - Time between update completion and guaranteed consistency is called **inconsistency window** ເພີ້ນ ການ update calendar
ບໍ່ມາລາຍກຽມ ຢັ້ງມີ update ຈົ່ງເນື້ອ

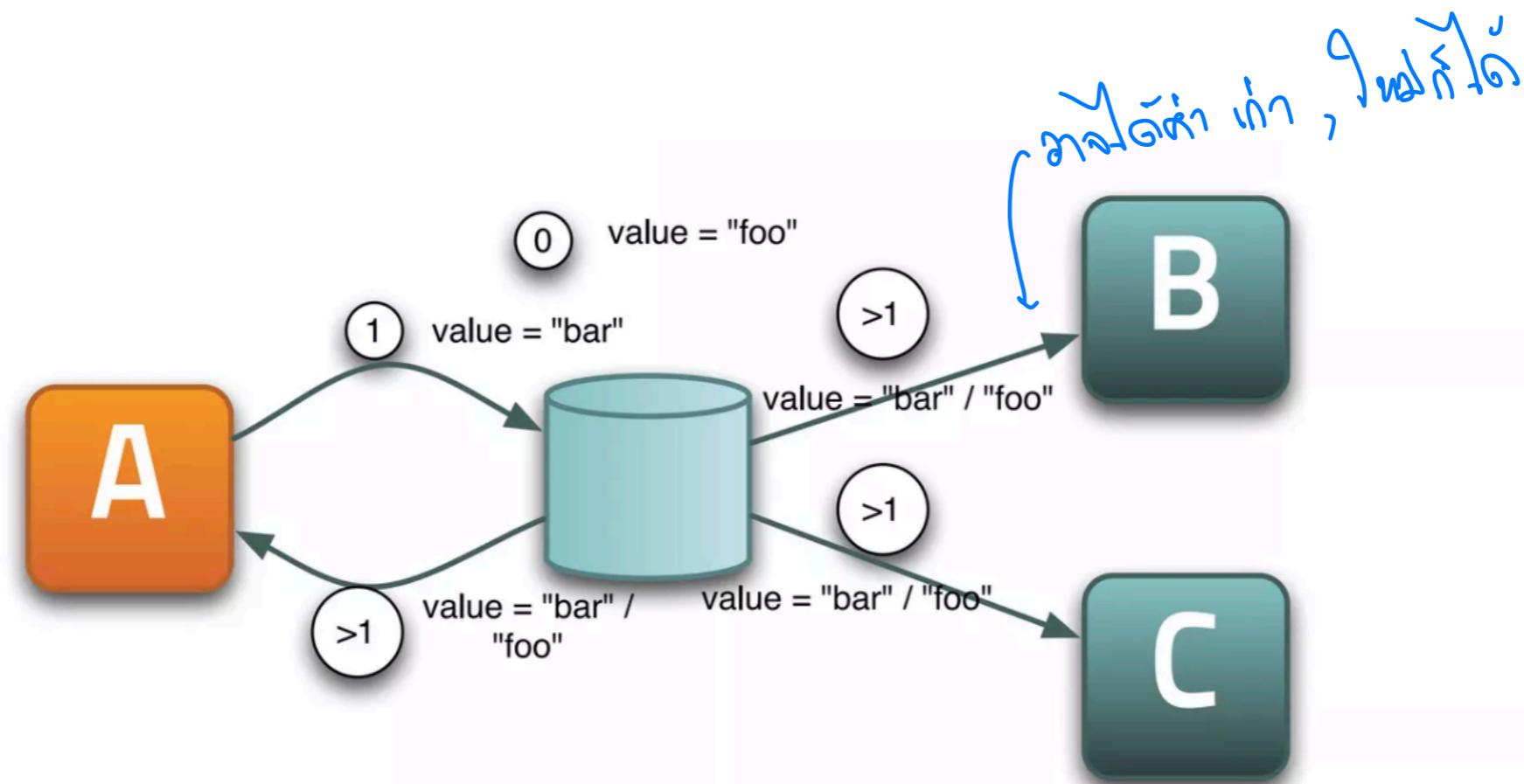
Strong Consistency



After the update, any subsequent access will return the updated value.

Weak Consistency

ເມື່ອມີການເງື່ອງຈາຕຽດ
ມາດີມວ່ານຫຼາຍດຽວ

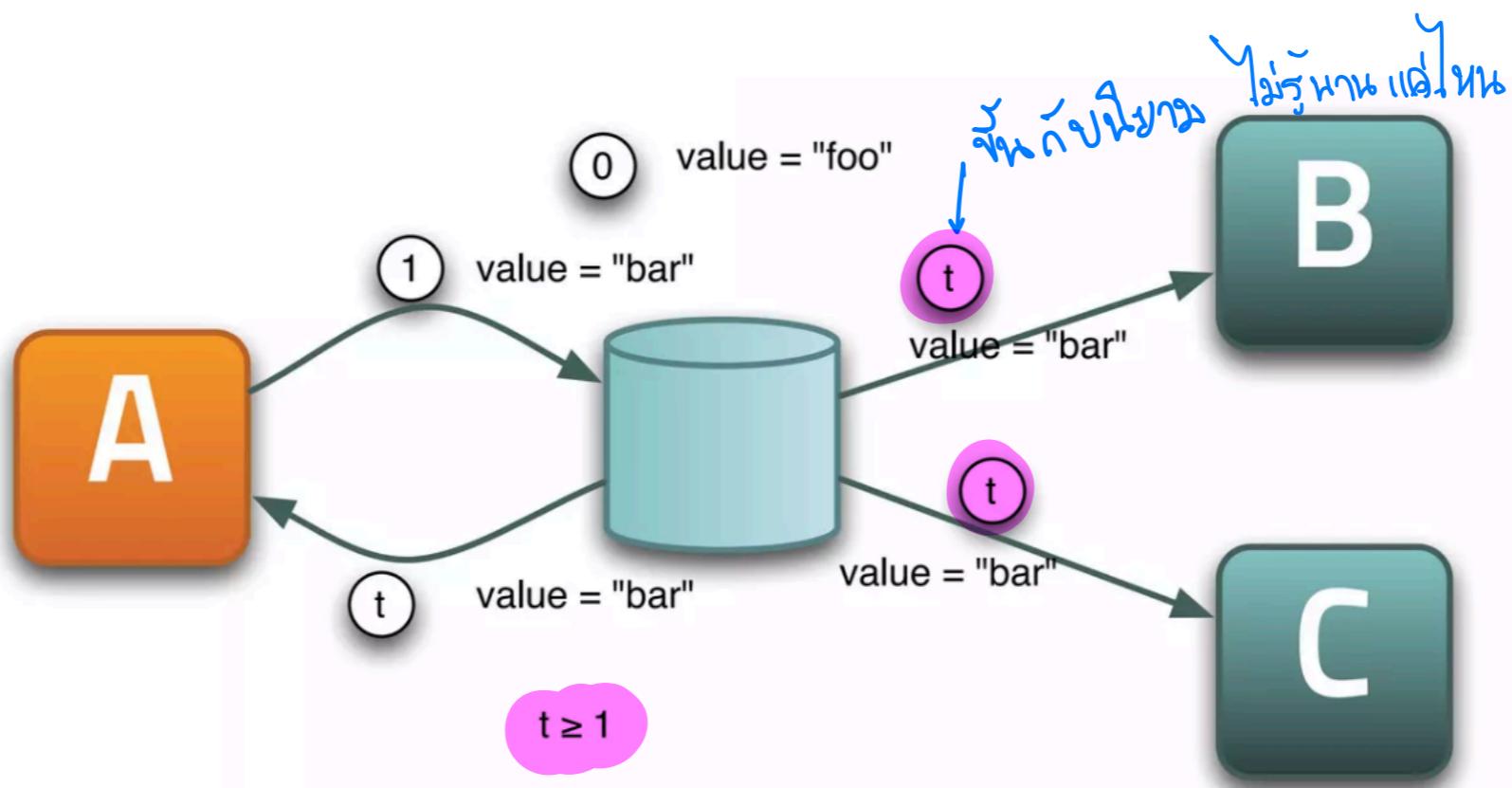


The system does not guarantee that at any given point in the future subsequent access will return the updated value

Eventual Consistency

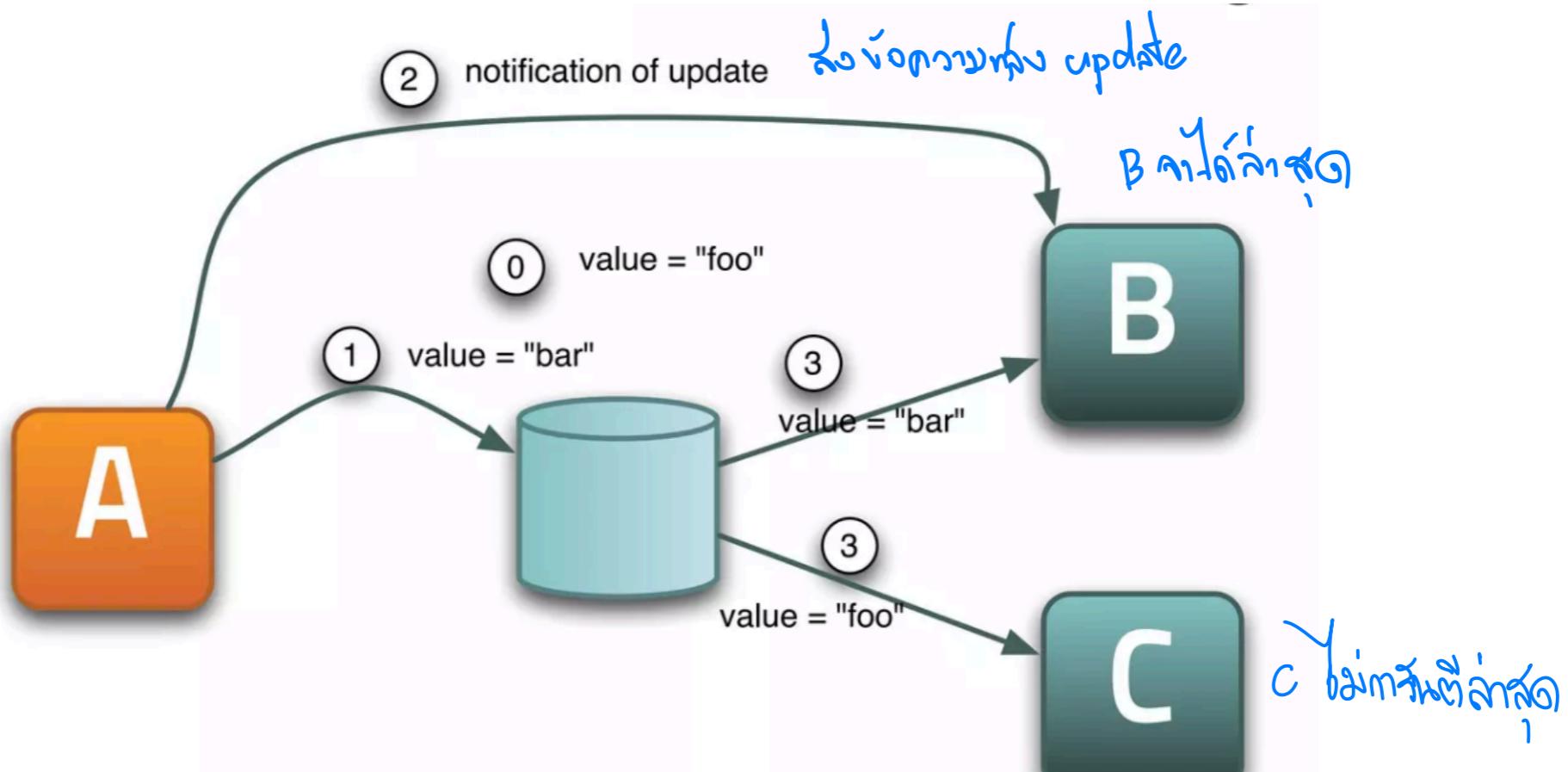
- Specific forms of weak consistency - if no new update are made to the object, **eventually** all accesses will return the last updated value *ក្នុង weak មានតម្លៃក្នុងក្រោមពេលពារសំនួល*
A ពេល B ដែលបានរៀបចំឡើង
- Causal consistency (A updates data, A informs B, B will always get updated value after that, also true for C if B informs C after B being informed)
- Read-your-write consistency (A updates data, A will always get updated value after that — special case of causal consistency) *A update នៅ A នៅពេល A បានរៀបចំឡើង ទៅវិញ ពាក្យនៅពេលវាស់គឺត្រូវបានរៀបចំឡើង*
- Others (reading assignments)

Eventual Consistency



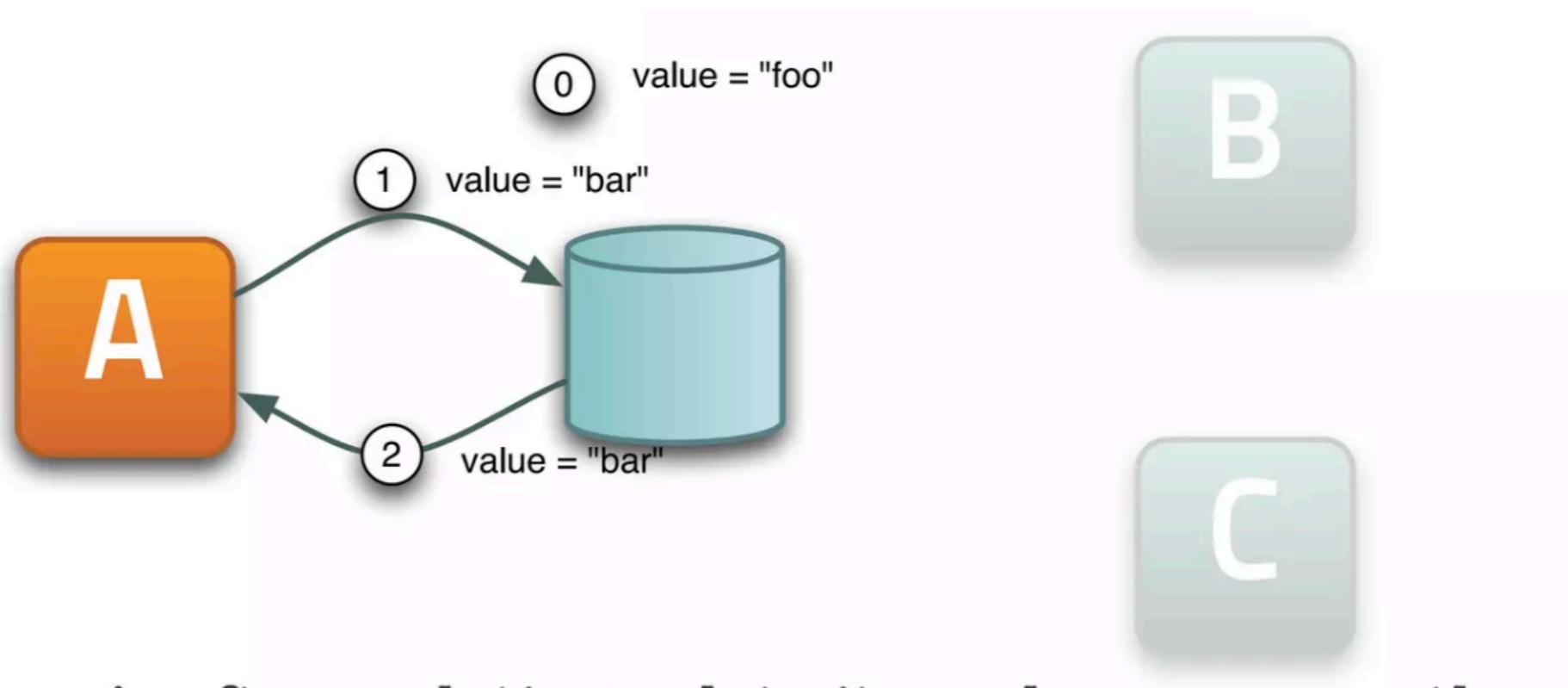
If no updates are made to the object, eventually all accesses will return the last updated value.

Causal Consistency



Subsequent access by process B will return the updated value, and a write is guaranteed to supersede the earlier write.

Read-Your-Write Consistency ការសន្ល់ពេលវាមួយ



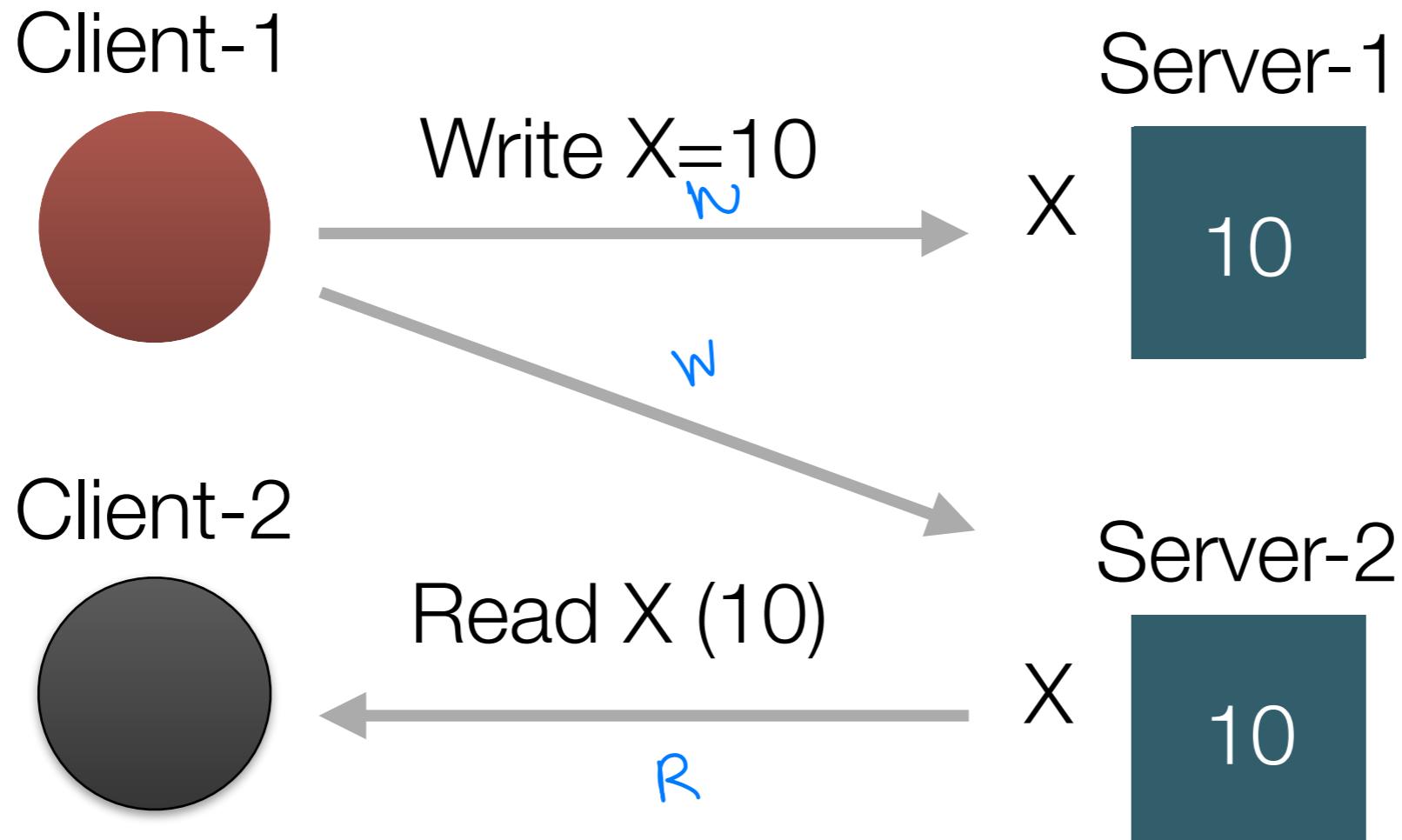
Process A, after updating a data item always access the updated value and never sees an older value

Providing Server-Side Consistency

- Read-only = no synchronization (HDFS)
Hadoop ແກ້ໄຂມາດ
ດ້ານທາເກົ່າຕໍ່ລະສົງ file ອຸນ
 - Handle consistency
 - Read and Write from multiple servers
 - Strong consistency ($R + W > N$)
 - Limit number of replications, limit synchronization time
 - Eventual data consistency ($R + W < N$)
 - Let it inconsistent and resolve it later
- ມີ R ຕັ້ງ (ແຕ່ລະຫວ່າງ
ຕ່ອງການໃໝ່ ທົ່ວ update W ຕັ້ງ
ຈຳນວນ server

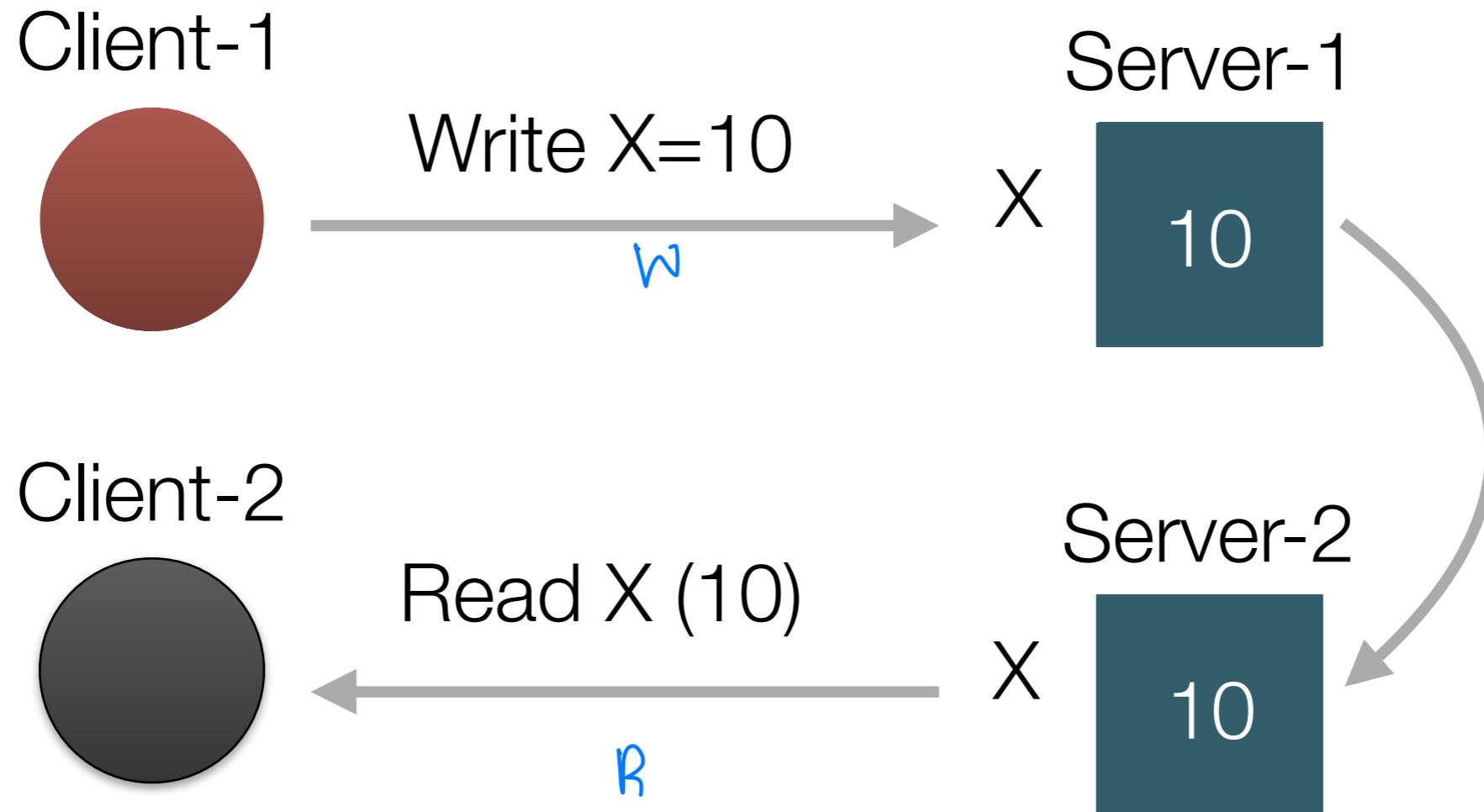
Providing Strong Consistency - Active Replication

($W = N$, $R = 1$, $R+W > N$) $\underline{N=2}$



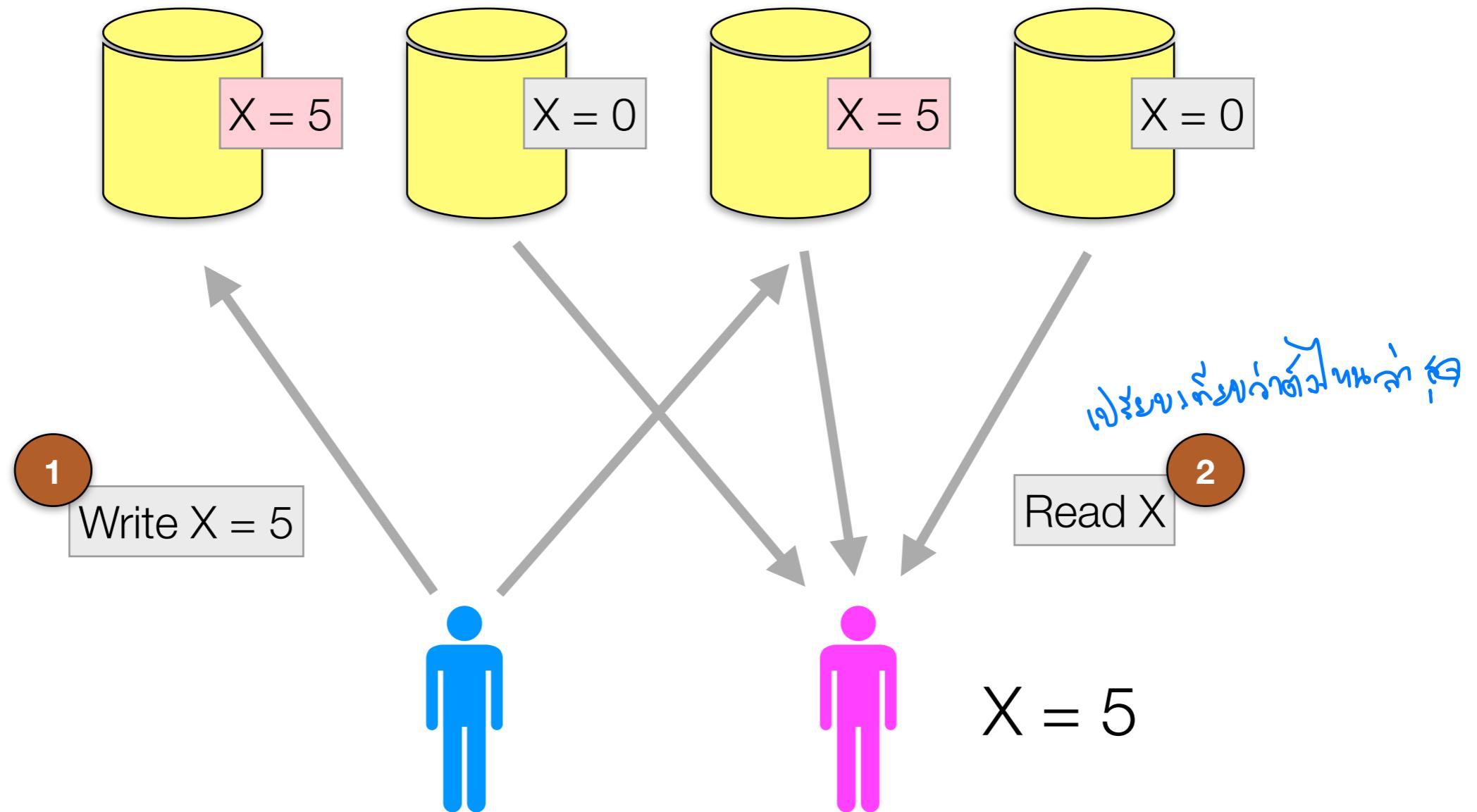
Providing Weak Consistency - Passive Replication

($W = 1$, $R = 1$, $R+W \leq N$)



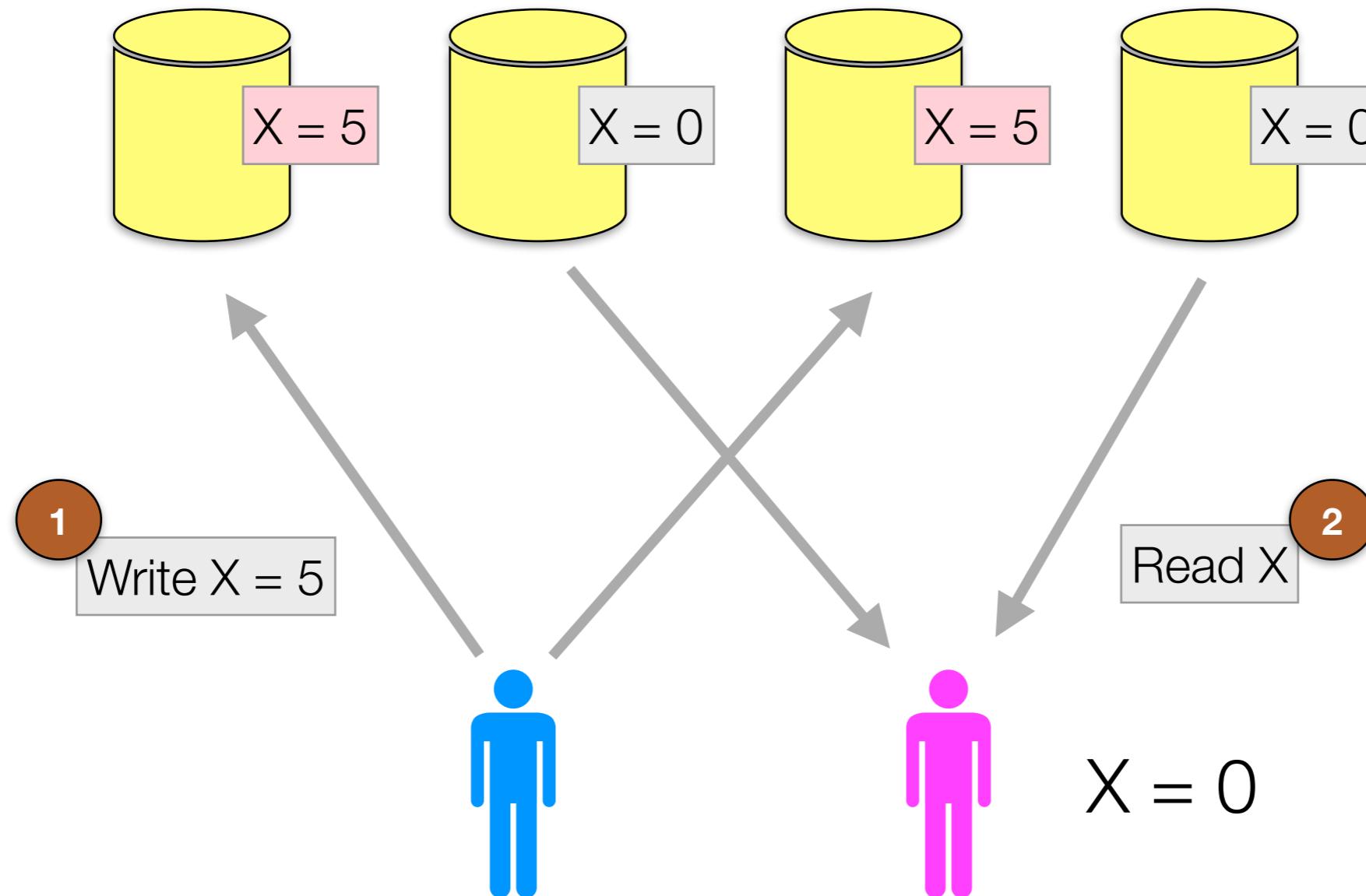
Consistency and Latency Trade-off

$N = 4, R = 3, W = 2$ (Consistency Guaranteed)



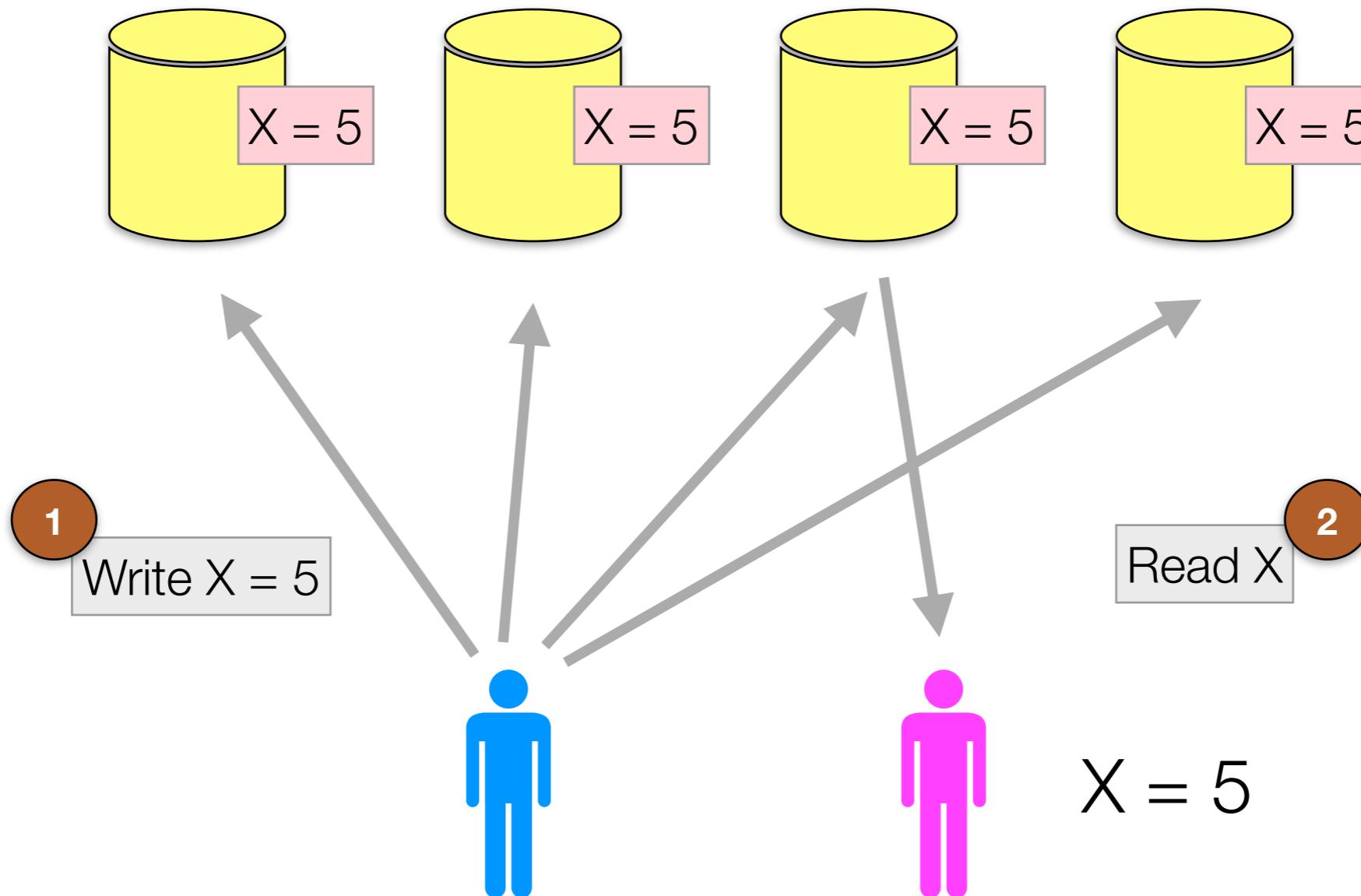
Consistency and Latency Trade-off

$N = 4, R = 2, W = 2$ (Consistency not Guaranteed)



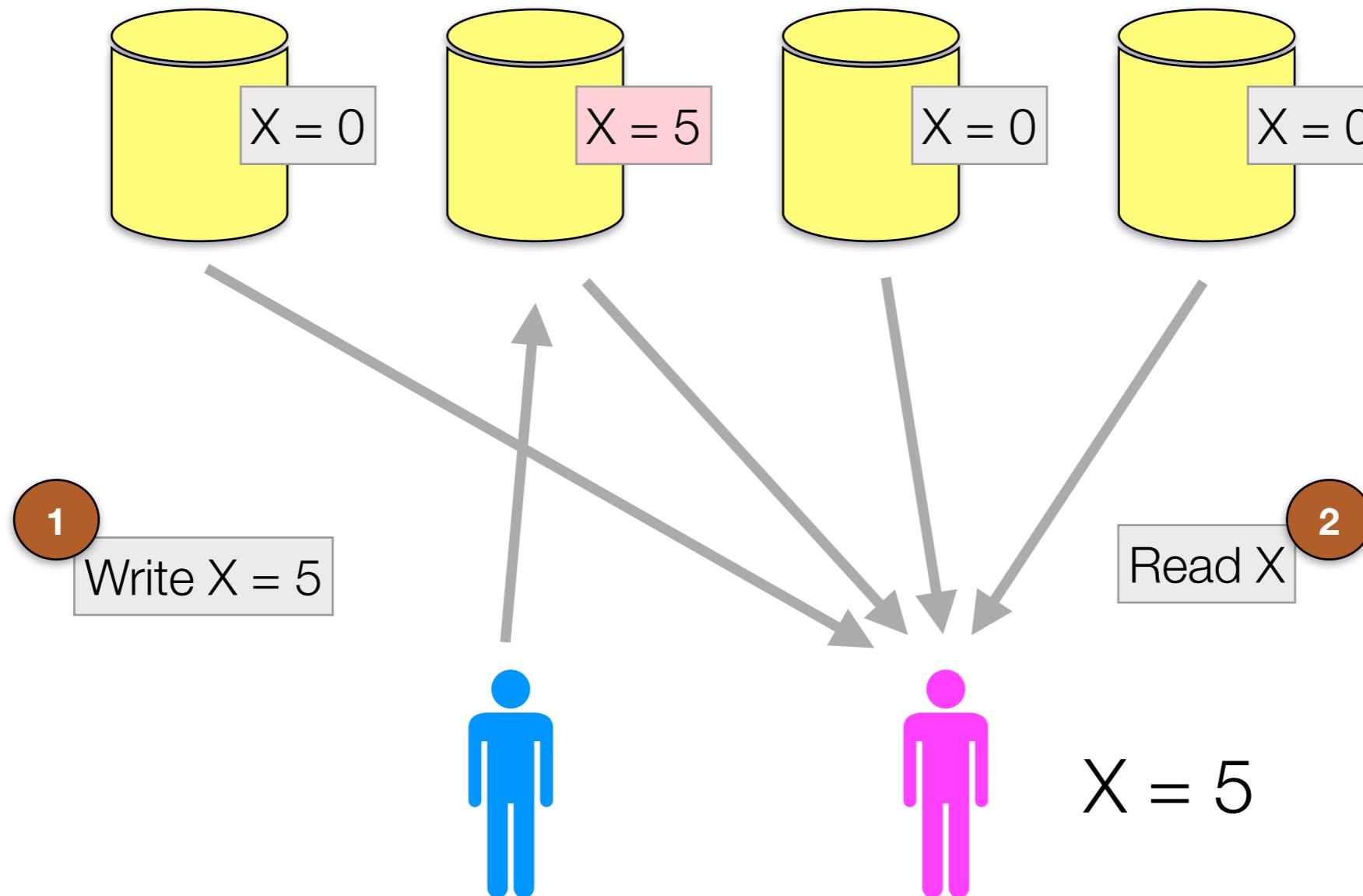
Consistency and Latency Trade-off

$N = 4, R = 1, W = 4$ (Fast Read, Slow Write)



Consistency and Latency Trade-off

$N = 4, R = 4, W = 1$ (Fast Write, Slow Read)



References

- V. Joshi, "Ordering Distributed Events", <https://medium.com/baseds/ordering-distributed-events-29c1dd9d1eff>
- W. Vogels, "Eventually consistent", Communications of the ACM, Vol. 52, pp. 40–44, 2009.
- W. Springer, "Eventually consistent", <https://www.slideshare.net/springerw/eventually-consistent>