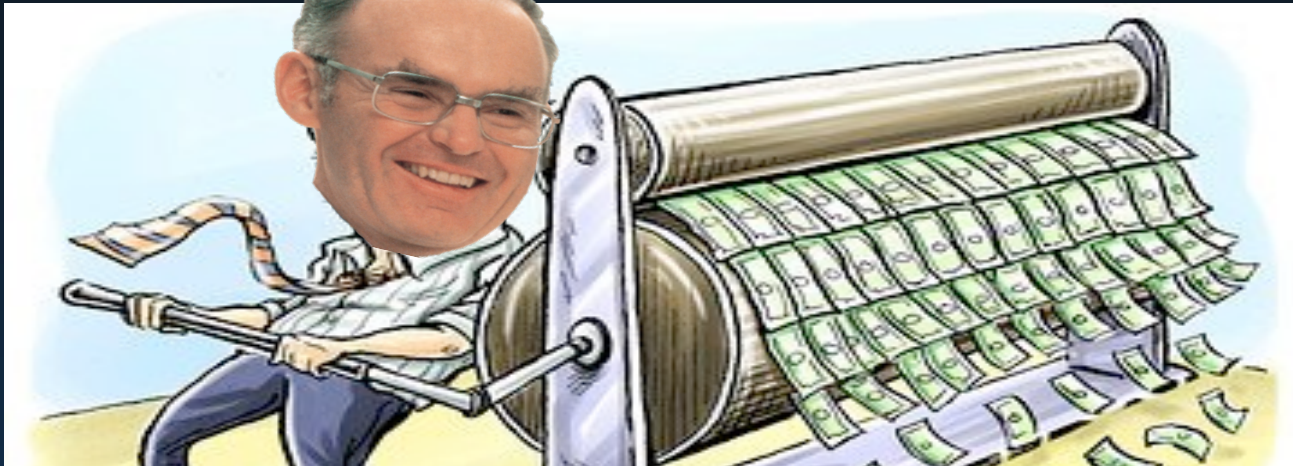


# TIME TRAVELING HARDWARE AND SOFTWARE SYSTEMS

Xiangyao Yu, Srini Devadas  
CSAIL, MIT

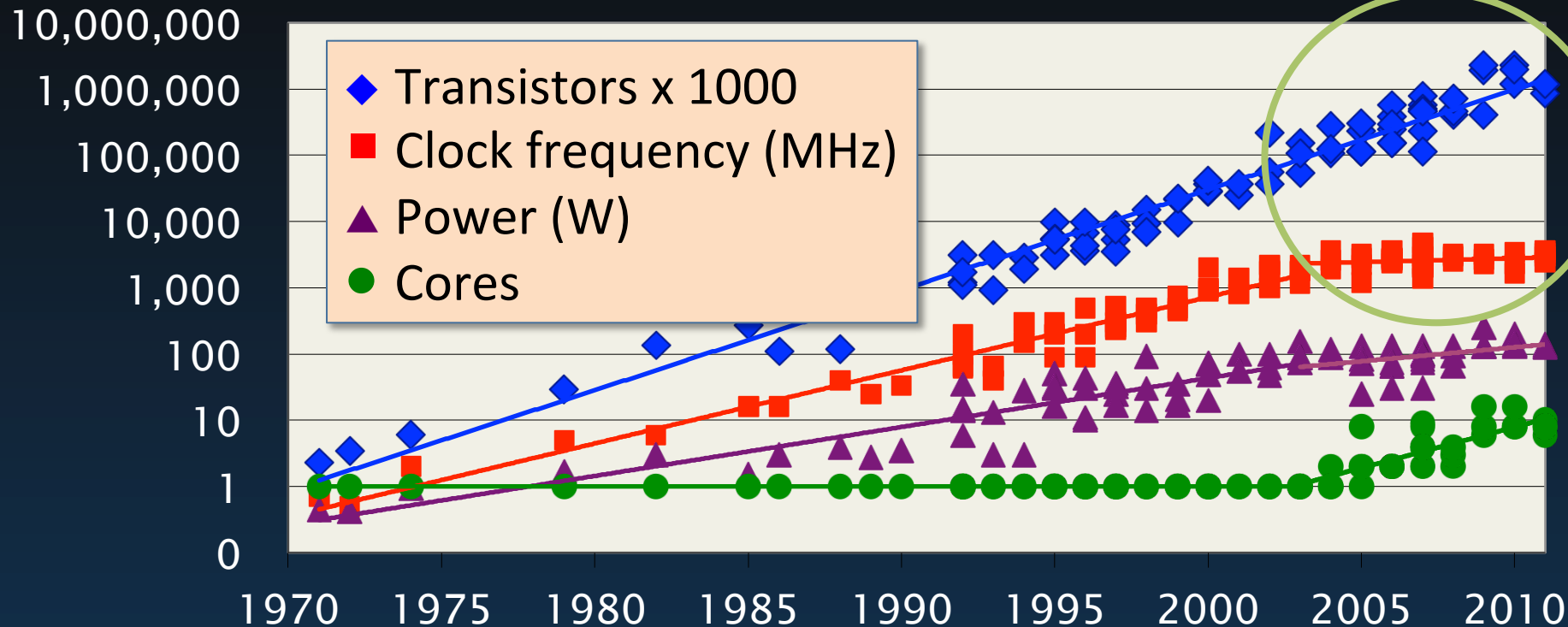


# FOR FIFTY YEARS, WE HAVE RIDDEN MOORE'S LAW



Moore's Law and the scaling of clock frequency  
= printing press for the currency of performance

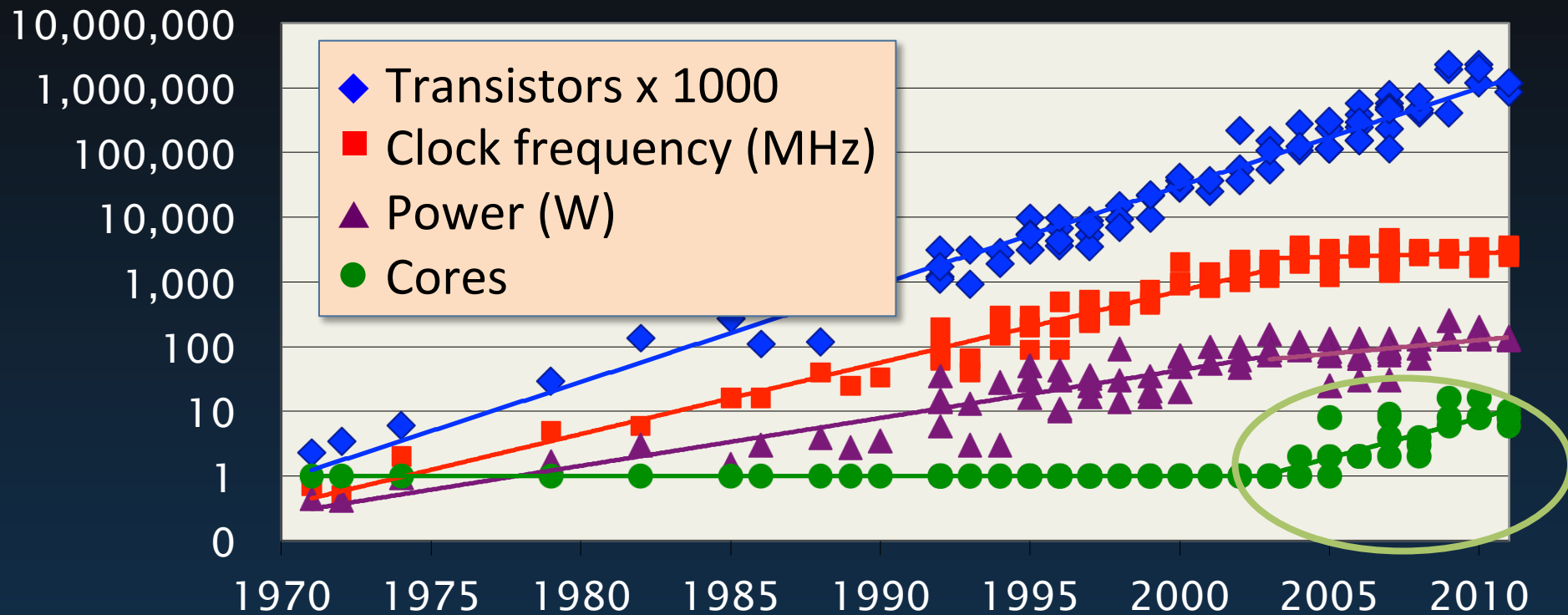
# TECHNOLOGY SCALING



Each generation of Moore's law doubles the number of transistors but clock frequency has stopped increasing.



# TECHNOLOGY SCALING



To increase performance, need to exploit parallelism.

# DIFFERENT KINDS OF PARALLELISM - 1

## Instruction Level

$a = b + c$

$d = e + f$

$g = d + b$

## Transaction Level

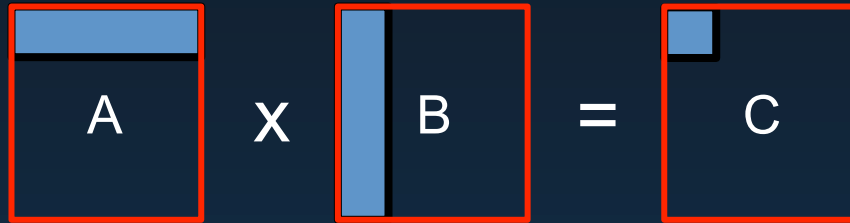
Read A  
Read B  
Compute C

Read A  
Read D  
Compute E

Read C  
Read E  
Compute F

# DIFFERENT KINDS OF PARALLELISM - 2

## Thread Level



Different thread computes each entry of product matrix C

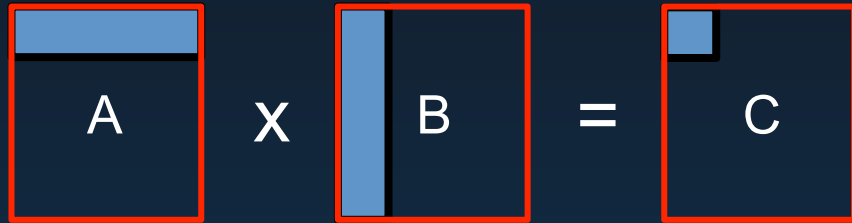
## Task Level

Search("image")



# DIFFERENT KINDS OF PARALLELISM - 3

## Thread Level



Different thread computes each entry of product matrix C

## User Level

Search("image")

Lookup("data")

Query("record")



# DEPENDENCY DESTROYS PARALLELISM

for  $i = 1$  to  $n$

$$a[b[i]] = (a[b[i - 1]] + b[i]) / c[i]$$

Need to compute  $i^{\text{th}}$  entry after  $i - 1^{\text{th}}$   
has been computed ☹️



# DIFFERENT KINDS OF DEPENDENCY

Read A

No

Write A

WAW:

Semantics

Read A

dependency!

Write A

decide  
order

Write A

RAW:

Read needs  
new value

Read A

Read A

Write A

WAR:

We have  
flexibility  
here!

# DEPENDENCE IS ACROSS TIME, BUT WHAT IS TIME?

- Time can be physical time
- Time could correspond to **logical timestamps** assigned to instructions
- Time could be a combination of the above

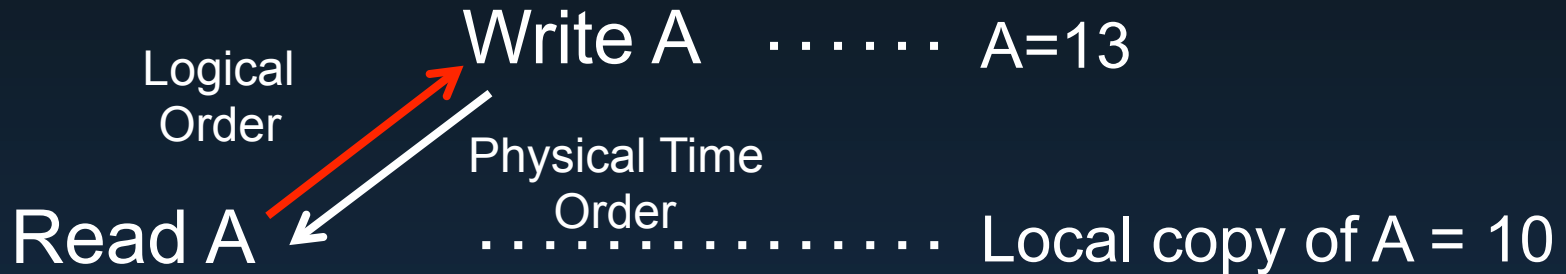
→ Time is a definition of ordering

# WAR DEPENDENCE

Initially A = 10

Thread 0

Thread 1



**Read happens later than Write in physical time but is before Write in logical time.**

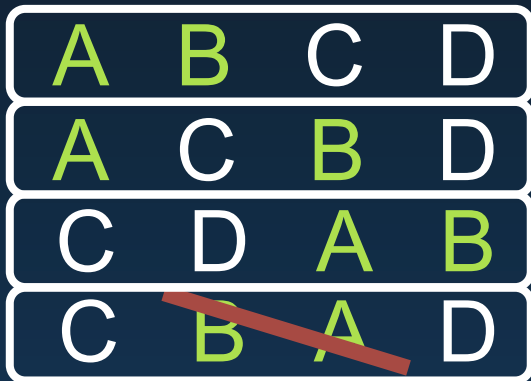
# WHAT IS CORRECTNESS?

- We define correctness of a parallel program based on its outputs in relation to the program run sequentially

# SEQUENTIAL CONSISTENCY



Global Memory Order



Can we exploit  
this freedom in  
correct  
execution to  
avoid  
dependency?

# AVOIDING DEPENDENCY ACROSS THE STACK



Circuit

Efficient atomic instructions



Multicore  
Processor

**Tardis coherence protocol**



Multicore  
Database

**TicToc concurrency control**  
*with Andy Pavlo and Daniel Sanchez*



Distributed  
Database

Distributed TicToc

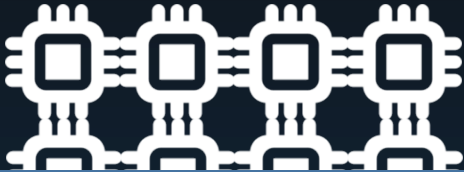


Distributed  
Shared Memory

Transaction processing with  
fault tolerance



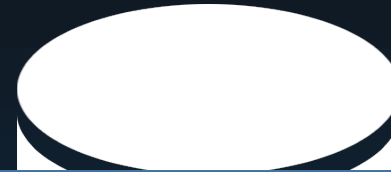
# SHARED MEMORY SYSTEMS



**Cache  
Coherence**



Multi-core  
Processor

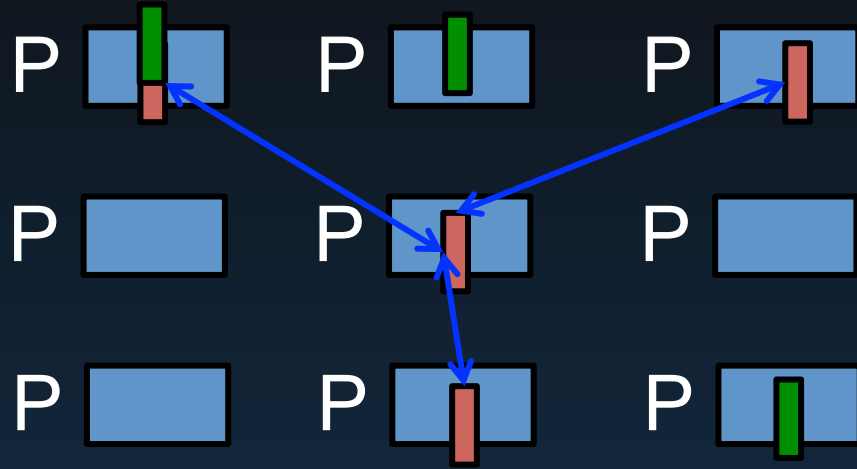


**Concurrency  
Control**



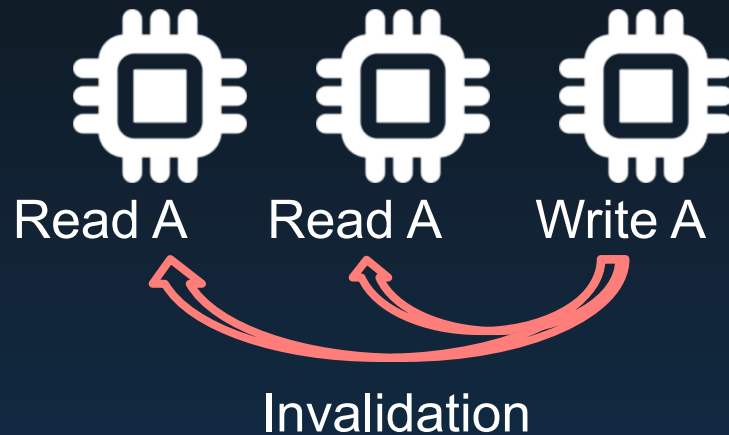
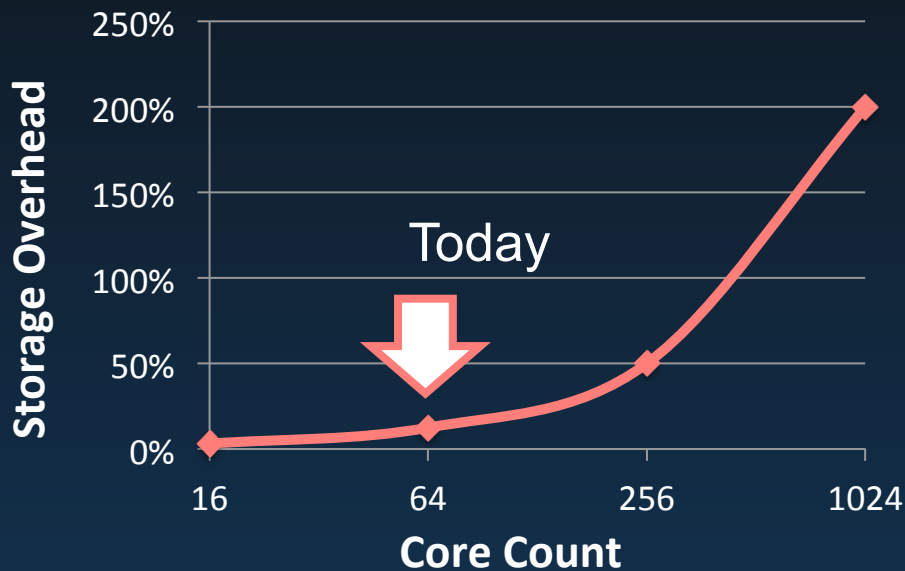
OLTP  
Database

# DIRECTORY-BASED COHERENCE



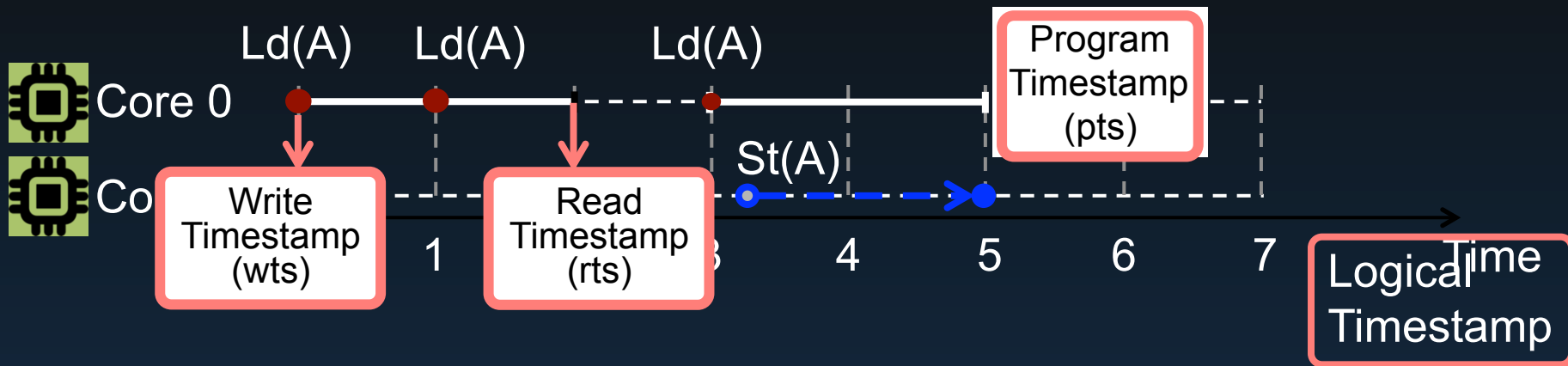
- Data replicated and cached locally for access
- Uncached data copied to local cache, writes invalidate data copies

# CACHE COHERENCE SCALABILITY



$O(N)$  Sharer List

# LEASE-BASED COHERENCE



- A read gets a lease on a cacheline
- Lease renewal after lease expires
- A store can only commit after leases expire
- Tardis: logical leases

# LOGICAL TIMESTAMP

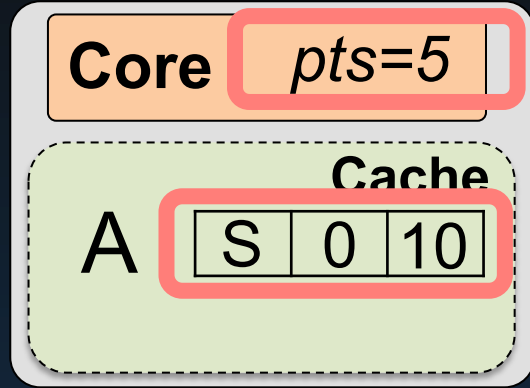
Invalidation ← Physical Time Order

**Tardis**  
(No Invalidation)

← Logical Time Order  
(concept borrowed from database)

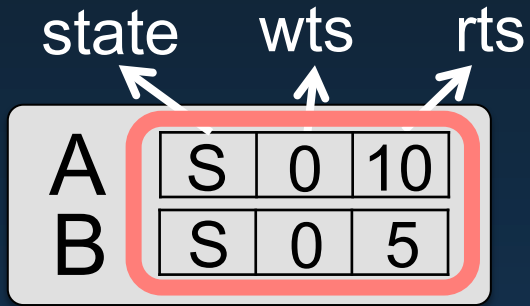


# TIMESTAMP MANAGEMENT



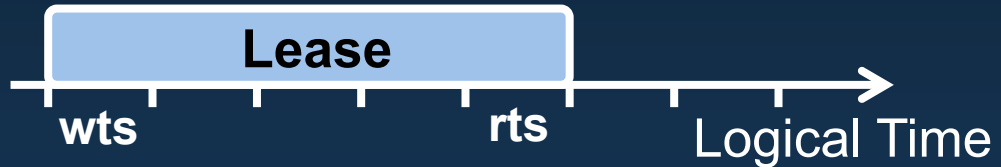
⇒ Program Timestamp (*pts*)  
Timestamp of last memory operation

⇒ Write Timestamp (*wts*)  
Data created at *wts*



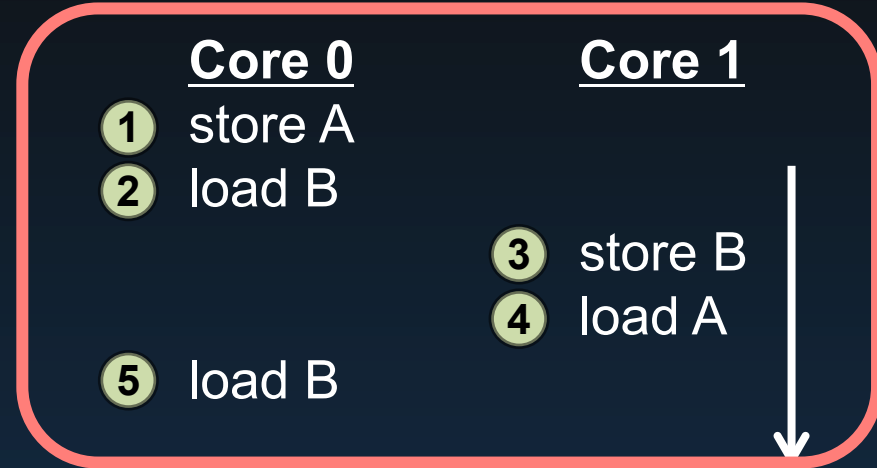
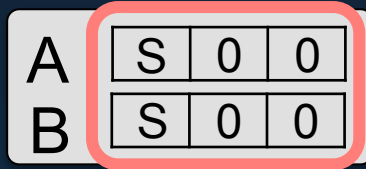
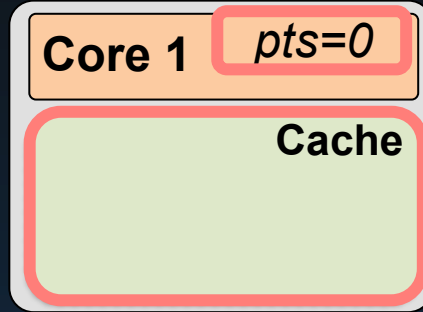
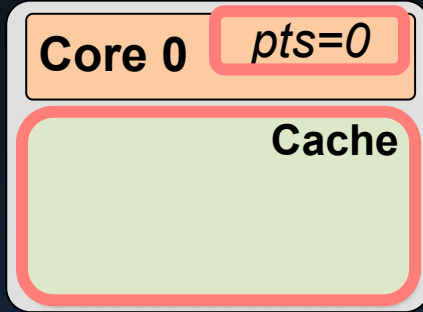
⇒ Read Timestamp (*rts*)  
Data valid from *wts* to *rts*

Shared LLC

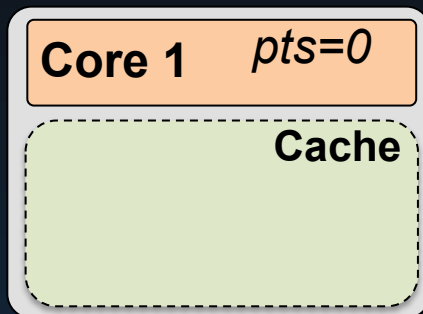
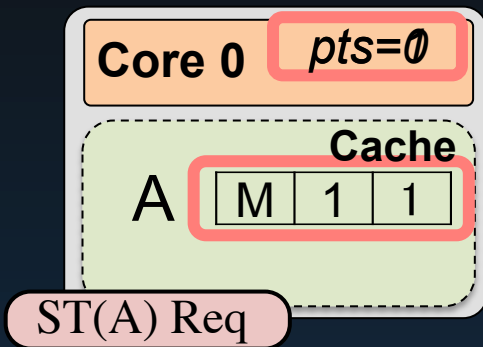




# TWO-CORE EXAMPLE



# STORE A @ CORE 0



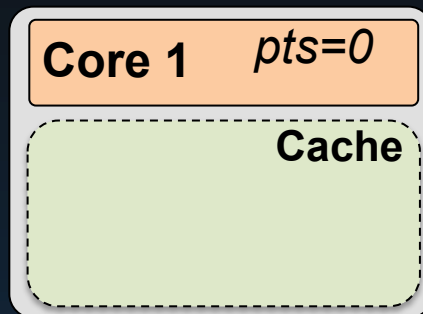
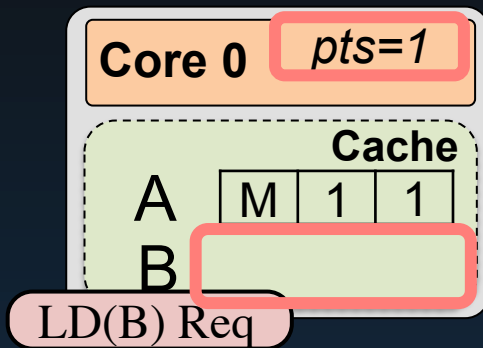
- Core 0
- 1 **store A**
  - 2 load B
  - 5 load B

- Core 1
- 3 store B
  - 4 load A

A	<del>M</del>	owner	0
B	S	0	0

Write at *pts = 1*

# LOAD B @ CORE 0



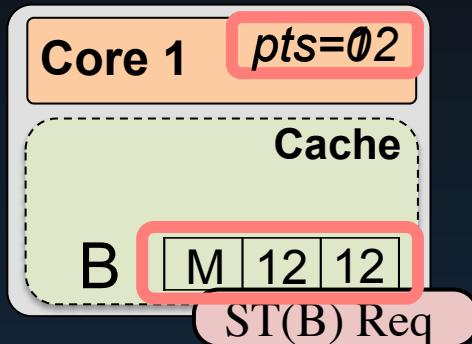
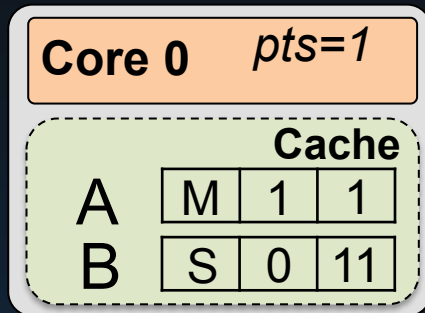
- Core 0
- 1 store A
  - 2 **load B**
  - 5 load B

- Core 1
- 3 store B
  - 4 load A

A	M	owner:0	
B	S	0	11

Reserve *rts* to *pts* + *lease* = 11

# STORE B @ CORE 1



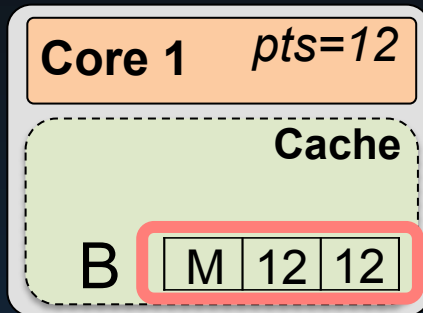
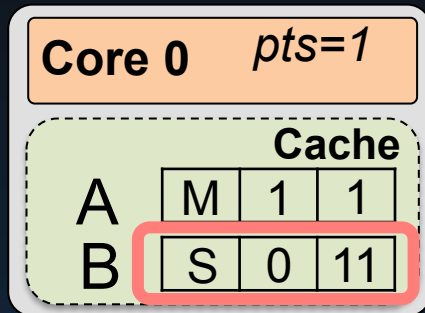
- Core 0**
- 1 store A
  - 2 load B
  - 5 load B

- Core 1**
- 3 **store B**
  - 4 load A



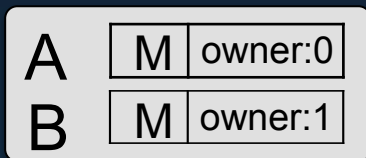
Exclusive ownership returned  
No invalidation

# Two VERSIONS COEXIST



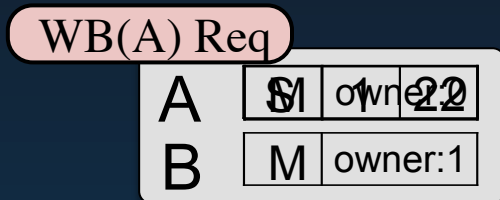
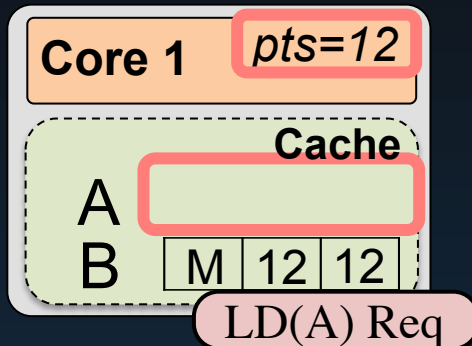
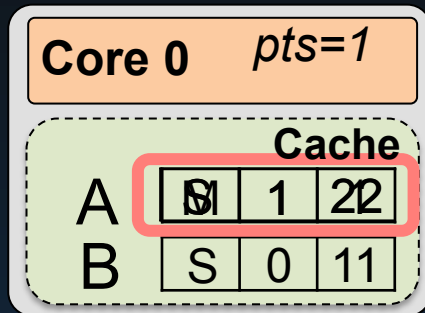
- Core 0**
- 1 store A
  - 2 load B
  - 5 load B

- Core 1**
- 3 **store B**
  - 4 load A



Core 1 traveled ahead in time  
Versions ordered in logical time

# LOAD A @ CORE 1



- Core 0**
- 1 store A
  - 2 load B
  - 5 load B

- Core 1**
- 3 store B
  - 4 **load A**

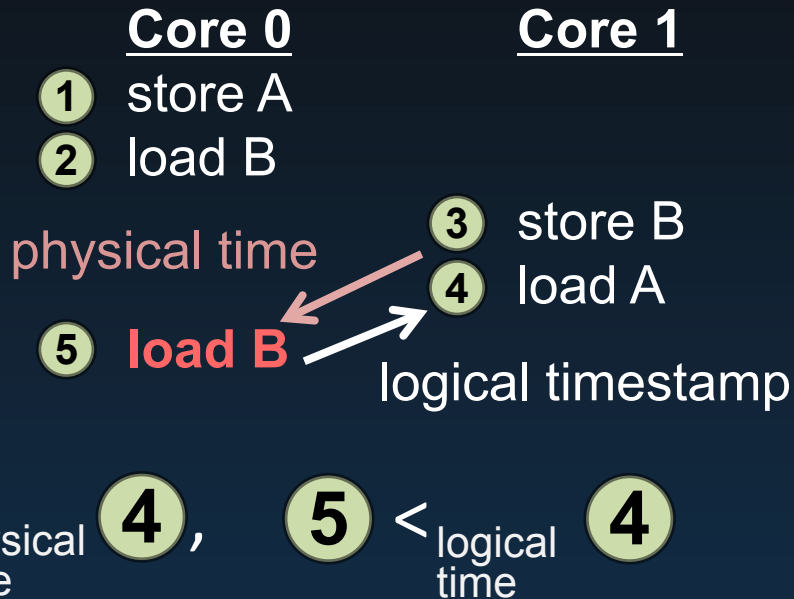
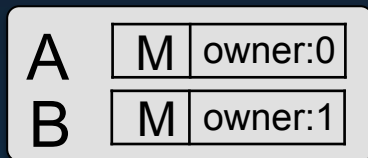
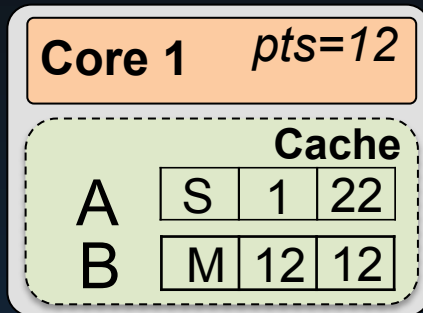
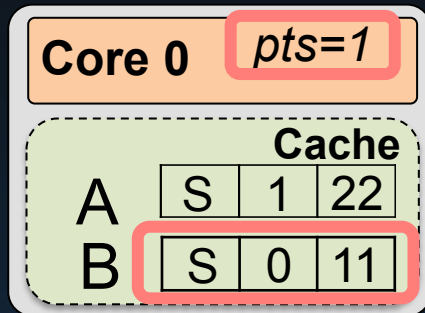
Write back request to Core 0

Downgrade from M to S

Reserve *rts* to  $pts + lease = 22$



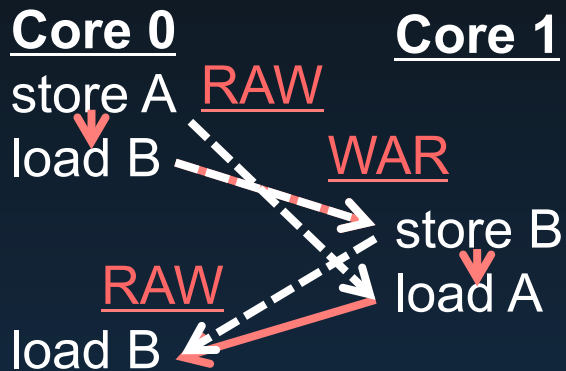
# LOAD B @ CORE 0



global memory order  $\neq$  physical time order

# SUMMARY OF EXAMPLE

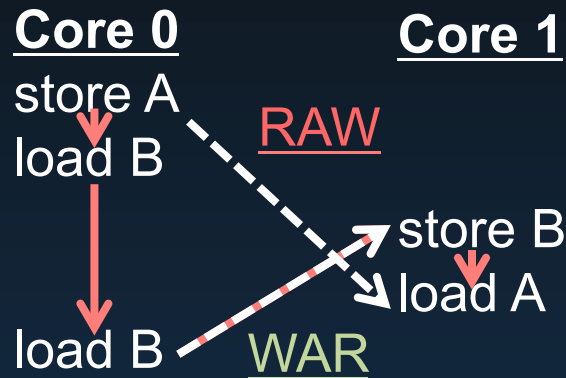
## Directory



physical time order

## Tardis

physical time



physical + logical time order

# PHYSIOLOGICAL TIME

Tardis

Global Memory Order

Core 0

store A (1)

load B (1)

load B (1)

Core 1

store B (12)

load A (12)



Physical Time

+



Logical Time

=



Physiological Time

$$X <_{PL} Y := X <_L Y \text{ or } (X =_L Y \text{ and } X <_P Y)$$

**Thm:** Tardis obeys Sequential Consistency

# TARDIS PROS AND CONS



Scalability

No Invalidation,  
Multicast or  
Broadcast



Lease Renew

Speculative Read



Timestamp Size

Timestamp Compression



Time Stands Still

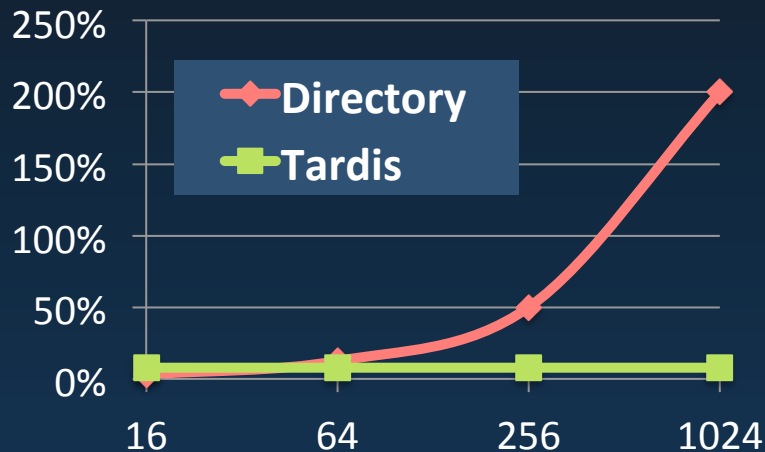
Livelock Avoidance

# EVALUATION

## Storage overhead per cacheline (N cores)

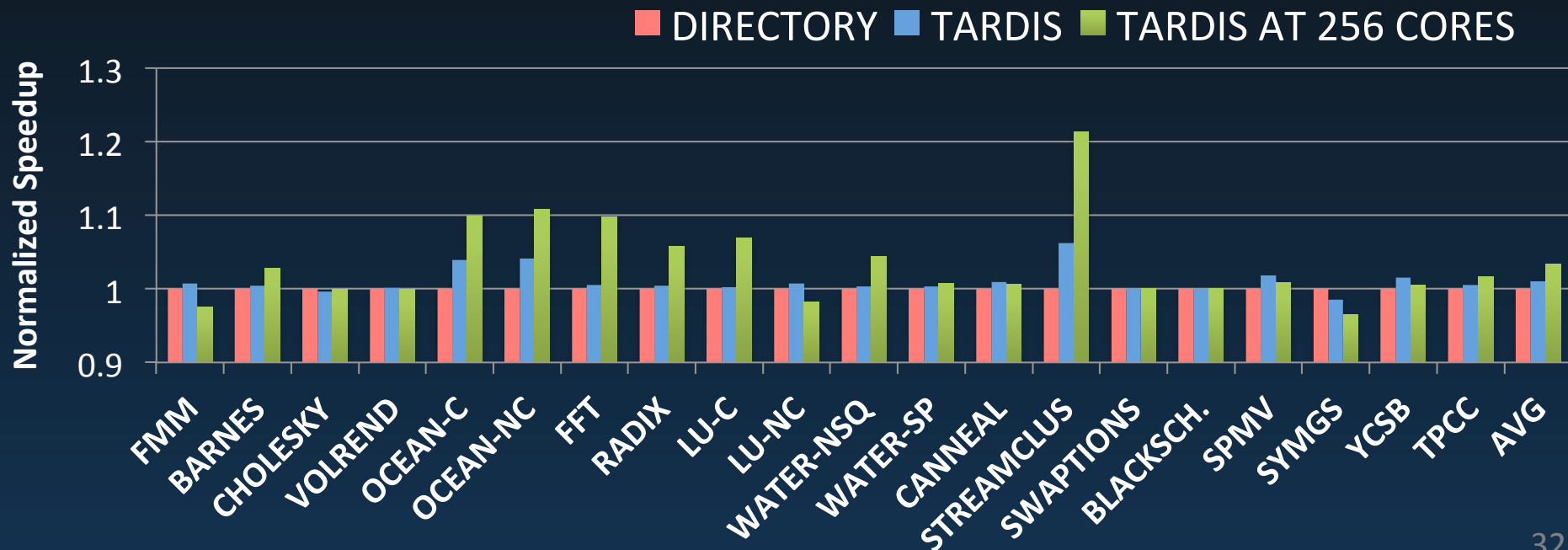
Directory:  $N$  bits per cacheline

Tardis:  $\text{Max}(\text{Const}, \log(N))$  bits per cacheline



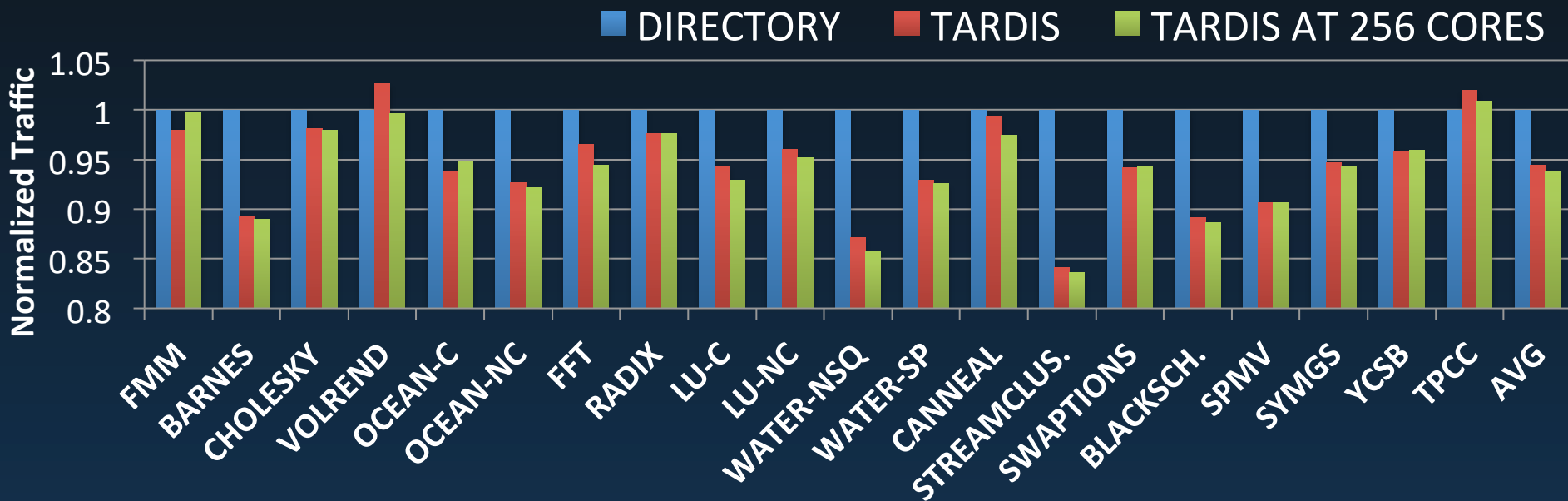
# SPEEDUP

## Graphite Multi-core Simulator (64 cores)

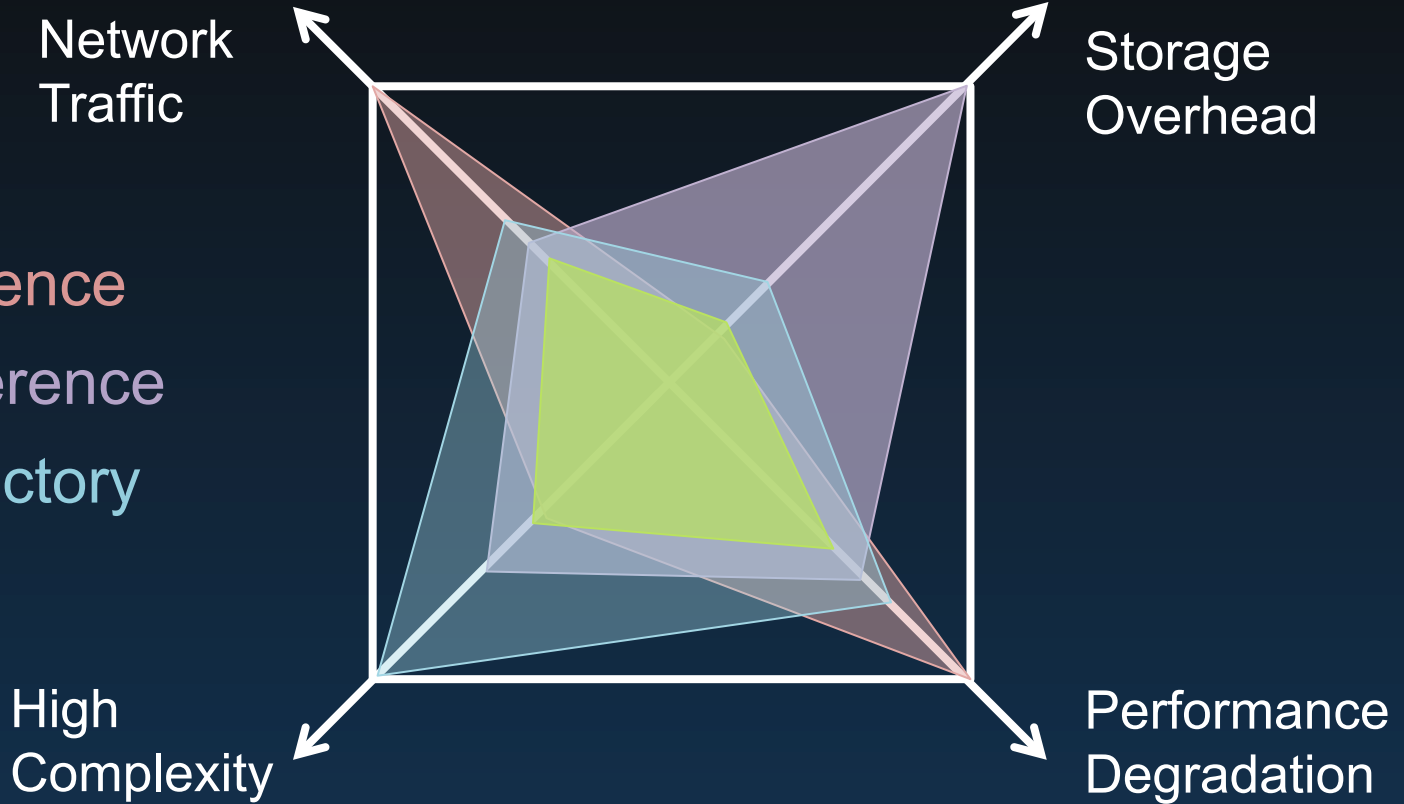




# NETWORK TRAFFIC



Snoopy Coherence  
Directory Coherence  
Optimized Directory  
**TARDIS**



# CONCURRENCY CONTROL

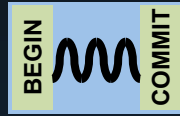
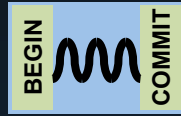


Serializability



Results should correspond to **some** serial order of **atomic** execution

# CONCURRENCY CONTROL



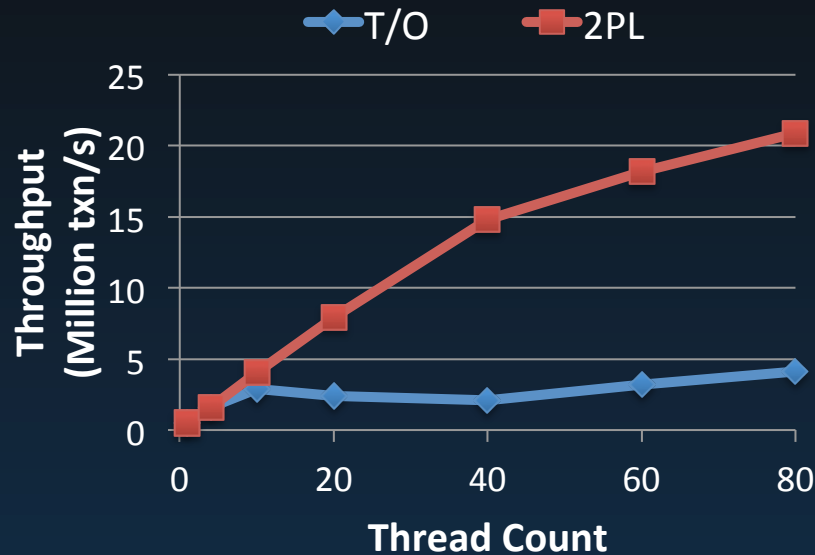
Can't Have This



Results should correspond to **some** serial order of **atomic** execution

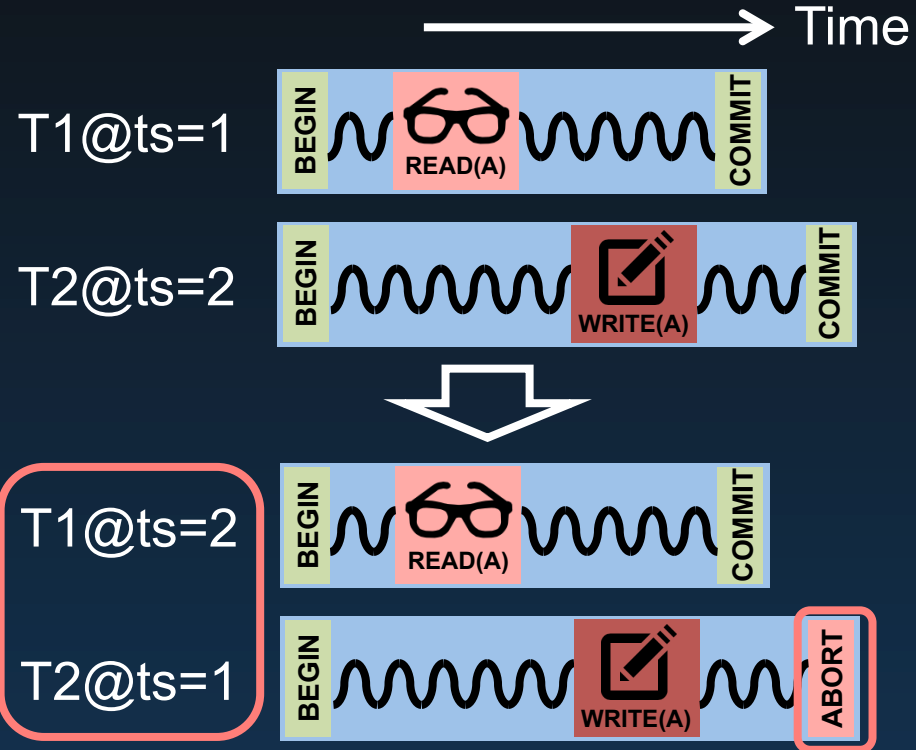
# BOTTLENECK 1: TIMESTAMP ALLOCATION

- Centralized Allocator
  - Timestamp allocation is a scalability bottleneck
- Synchronized Clock
  - Clock skew causes unnecessary aborts



# BOTTLENECK 2: STATIC ASSIGNMENT

- Timestamps assigned before a transaction starts
- Suboptimal assignment leads to unnecessary aborts.



# KEY IDEA: DATA DRIVEN TIMESTAMP MANAGEMENT

## Traditional T/O

1. Acquire timestamp (TS)
2. Determine tuple visibility using TS



Timestamp Allocation



Static Timestamp Assignment

## TicToc

1. Access tuples and remember their timestamp info.
2. Compute commit timestamp (CommitTS)

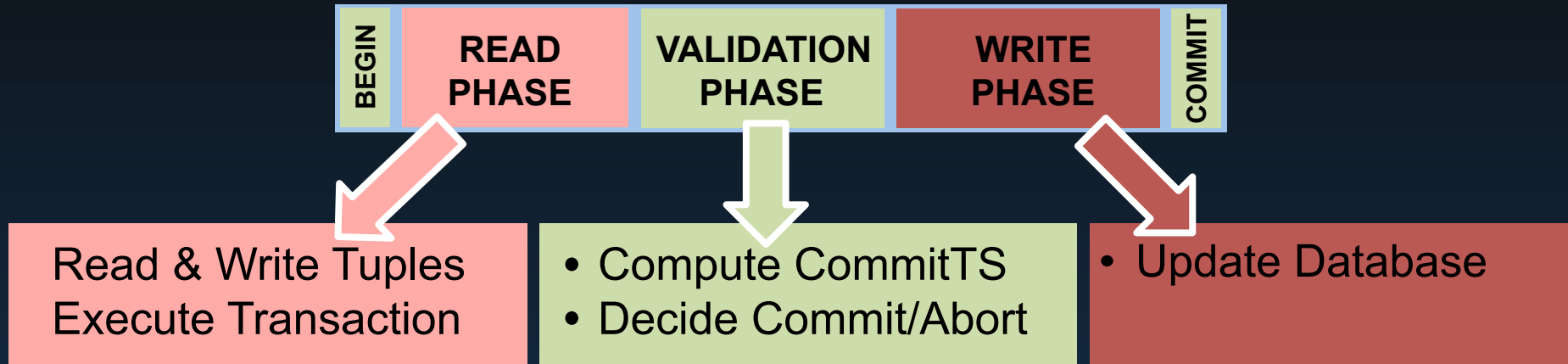


No Timestamp Allocation



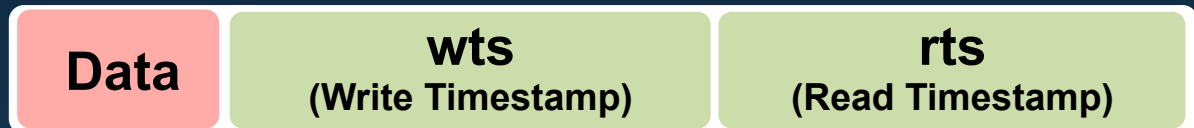
Dynamic Timestamp Assignment

# TICTOC TRANSACTION EXECUTION



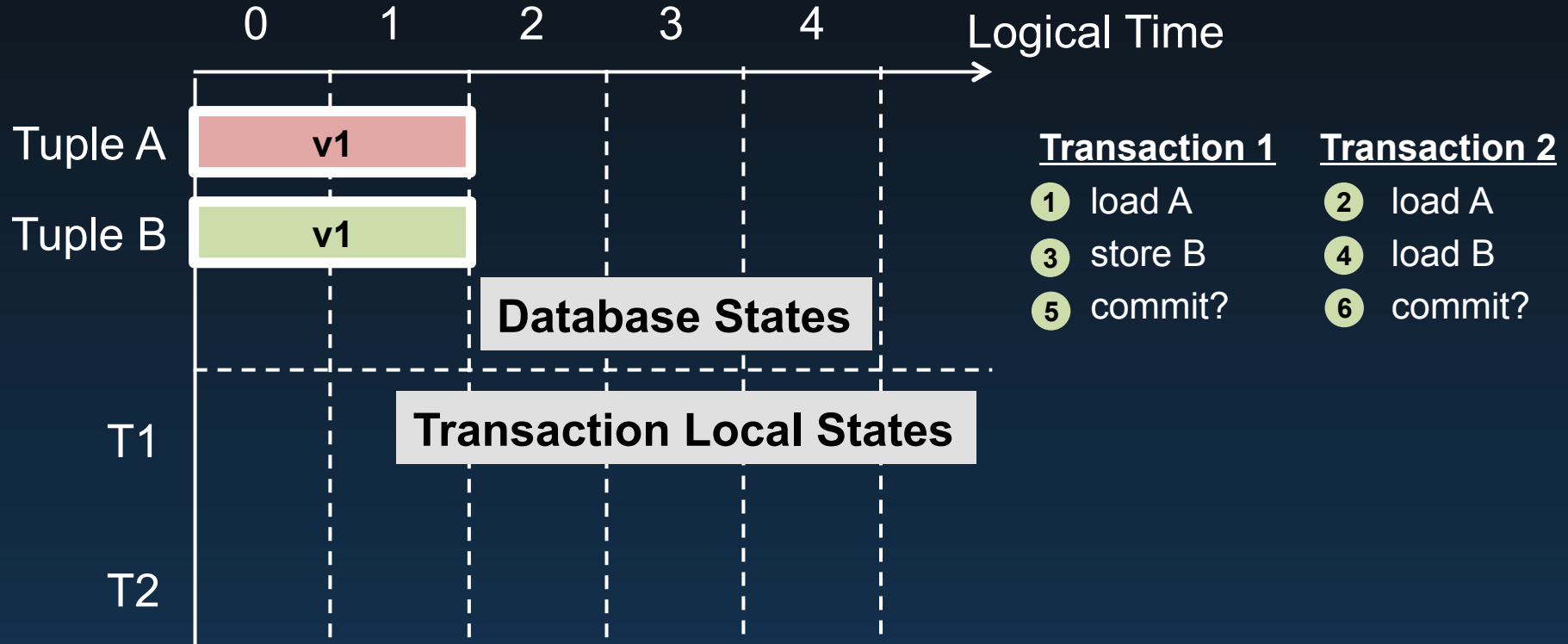
wt<sub>s</sub> : last data write @ wt<sub>s</sub>  
rt<sub>s</sub> : last data read @ rt<sub>s</sub> } data valid between wt<sub>s</sub> and rt<sub>s</sub>

Tuple Format

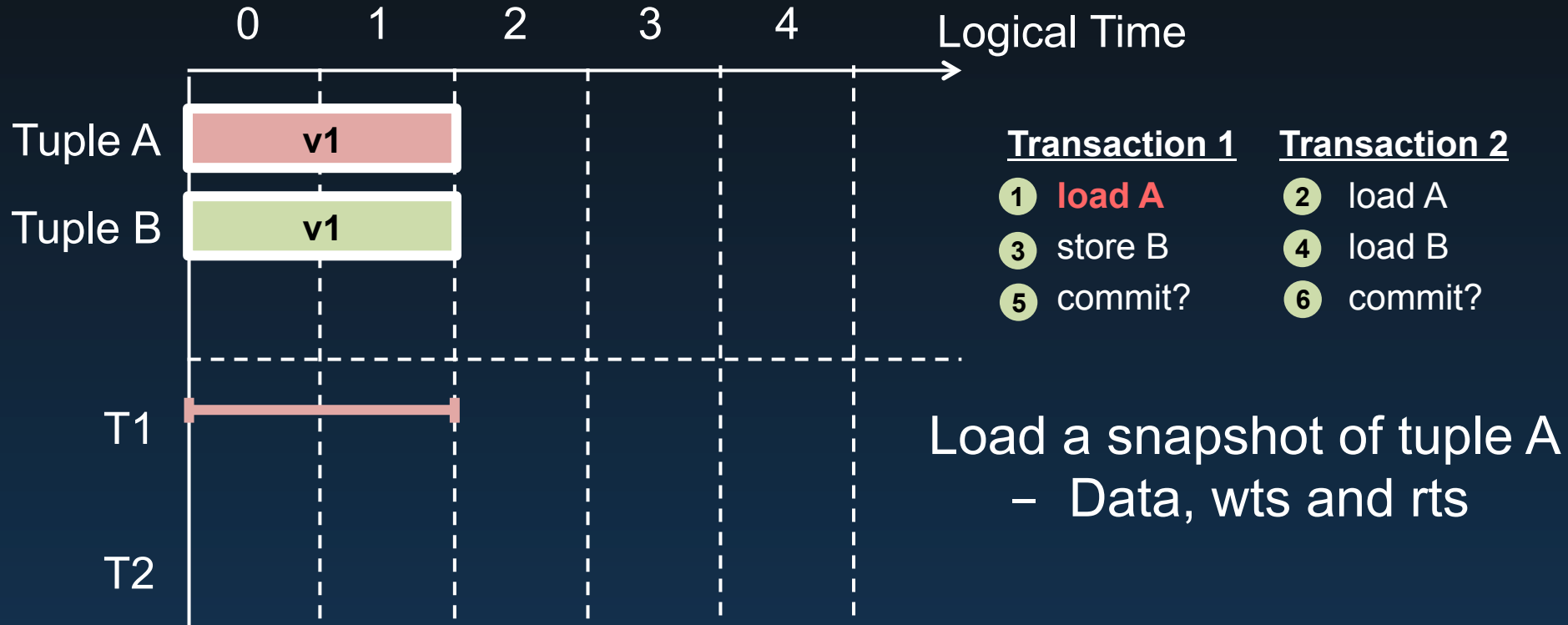




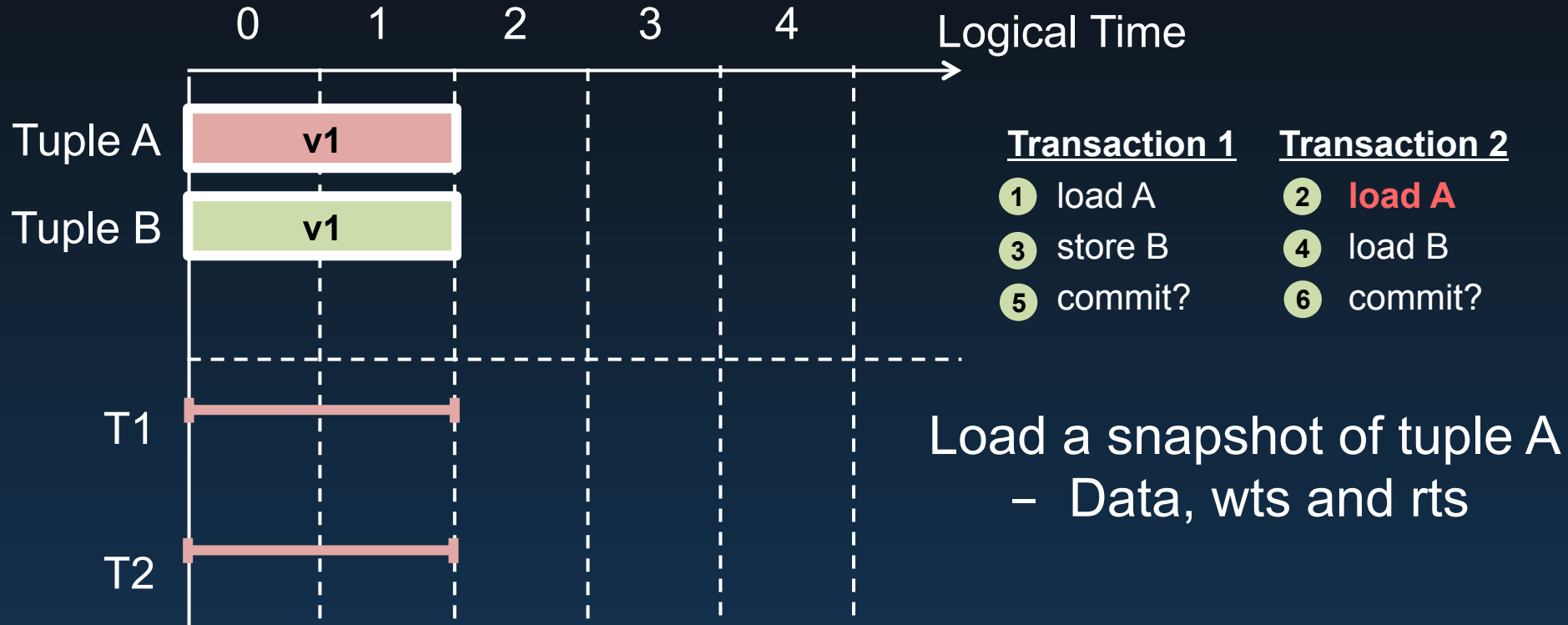
# TicTOC EXAMPLE



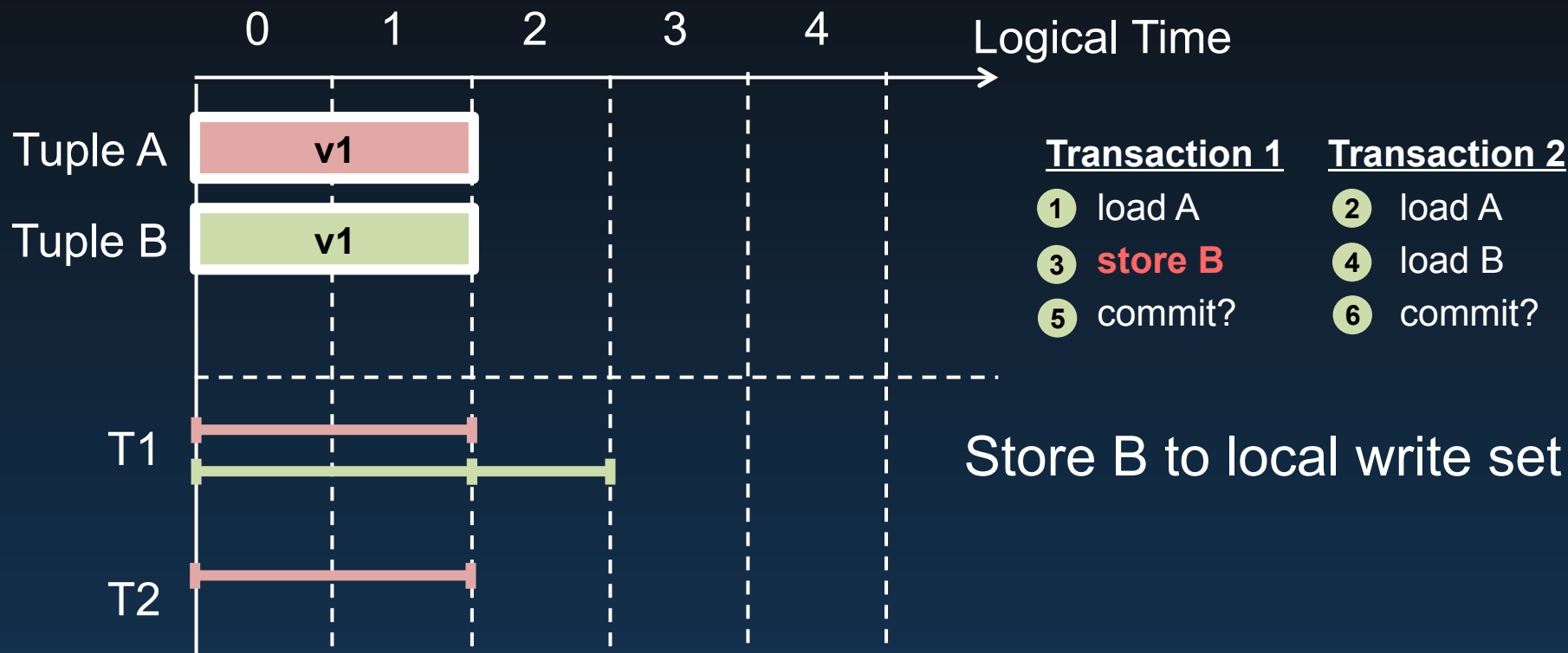
# LOAD A FROM T1



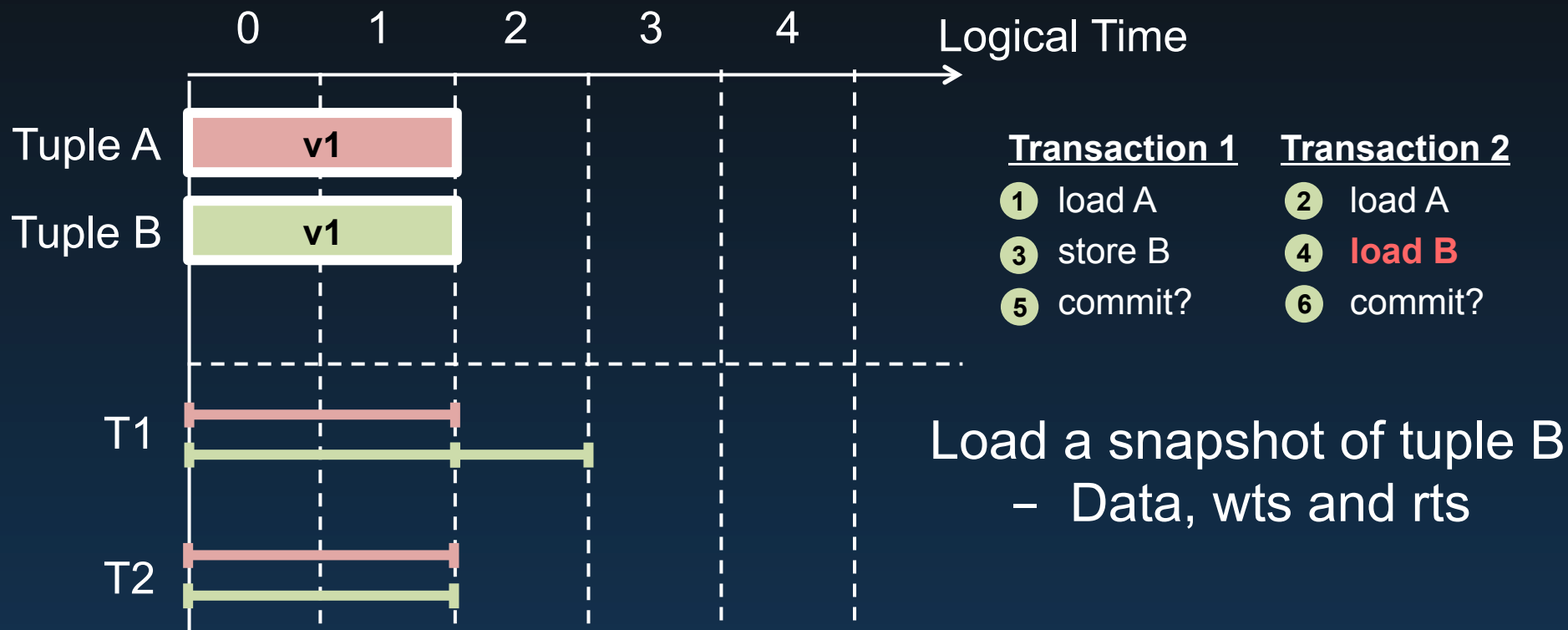
# LOAD A FROM T2



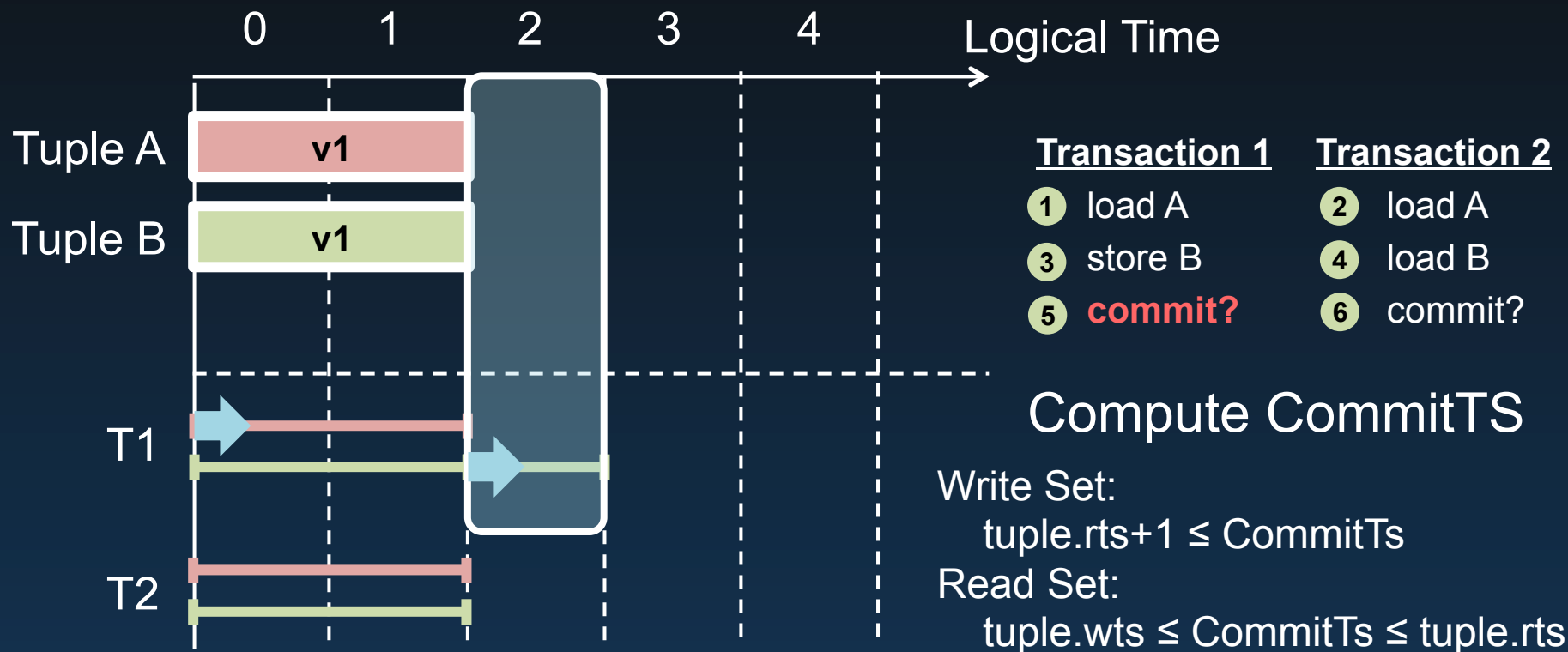
# STORE B FROM T1



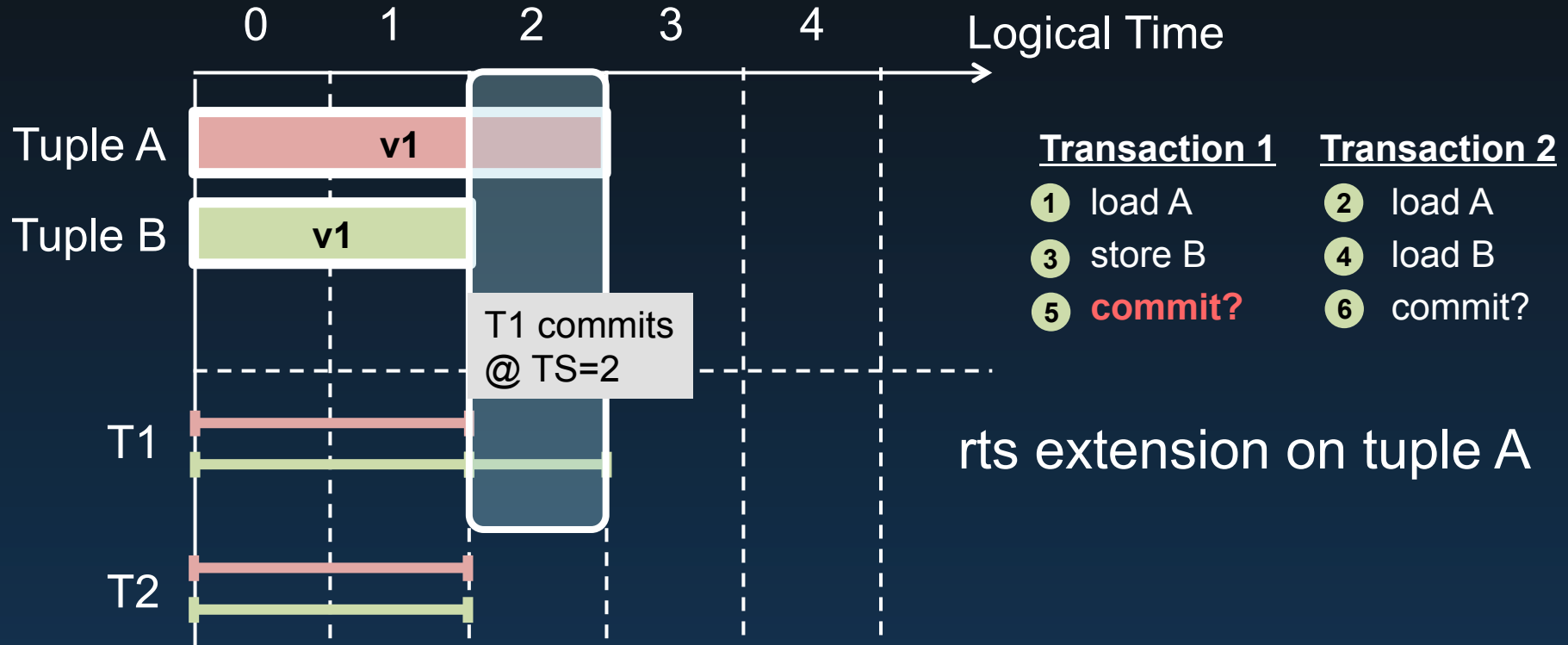
# LOAD B FROM T2



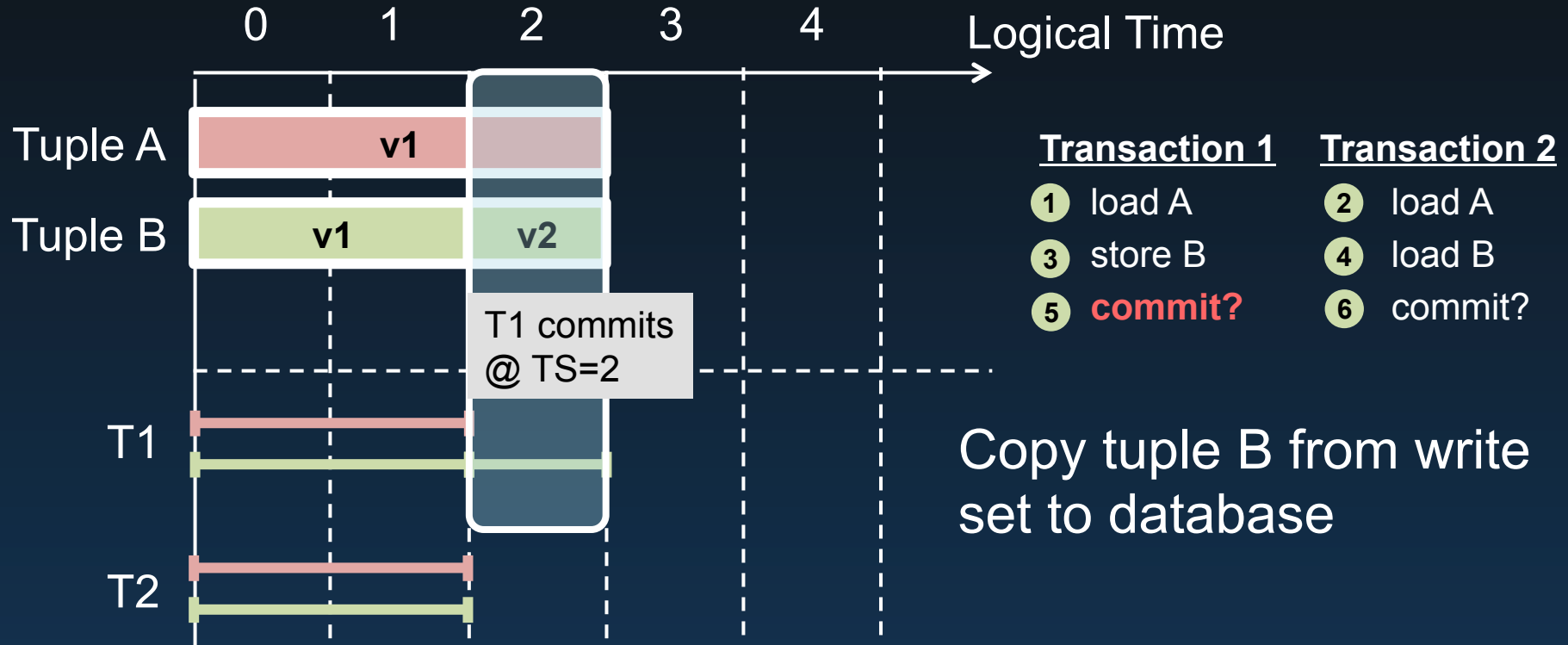
# COMMIT PHASE OF T1



# COMMIT PHASE OF T1

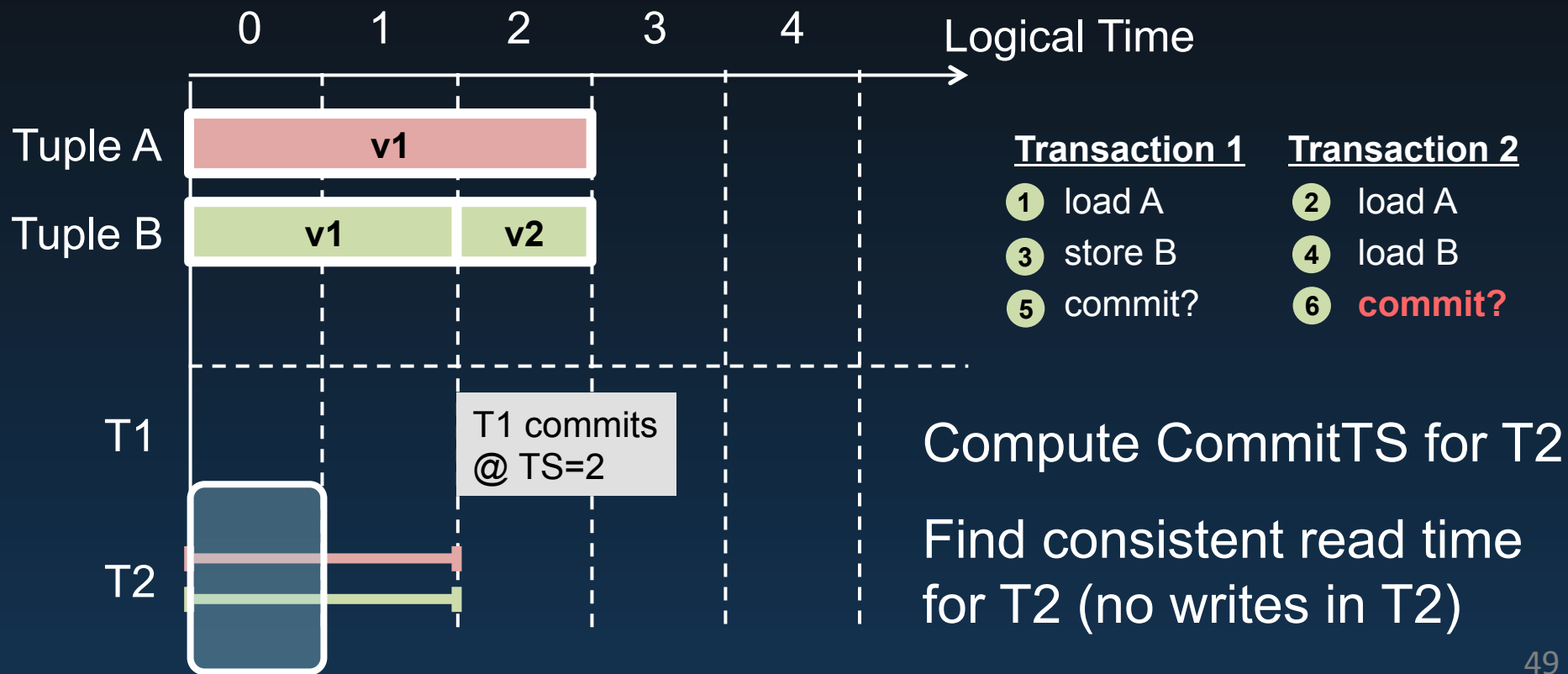


# COMMIT PHASE OF T1

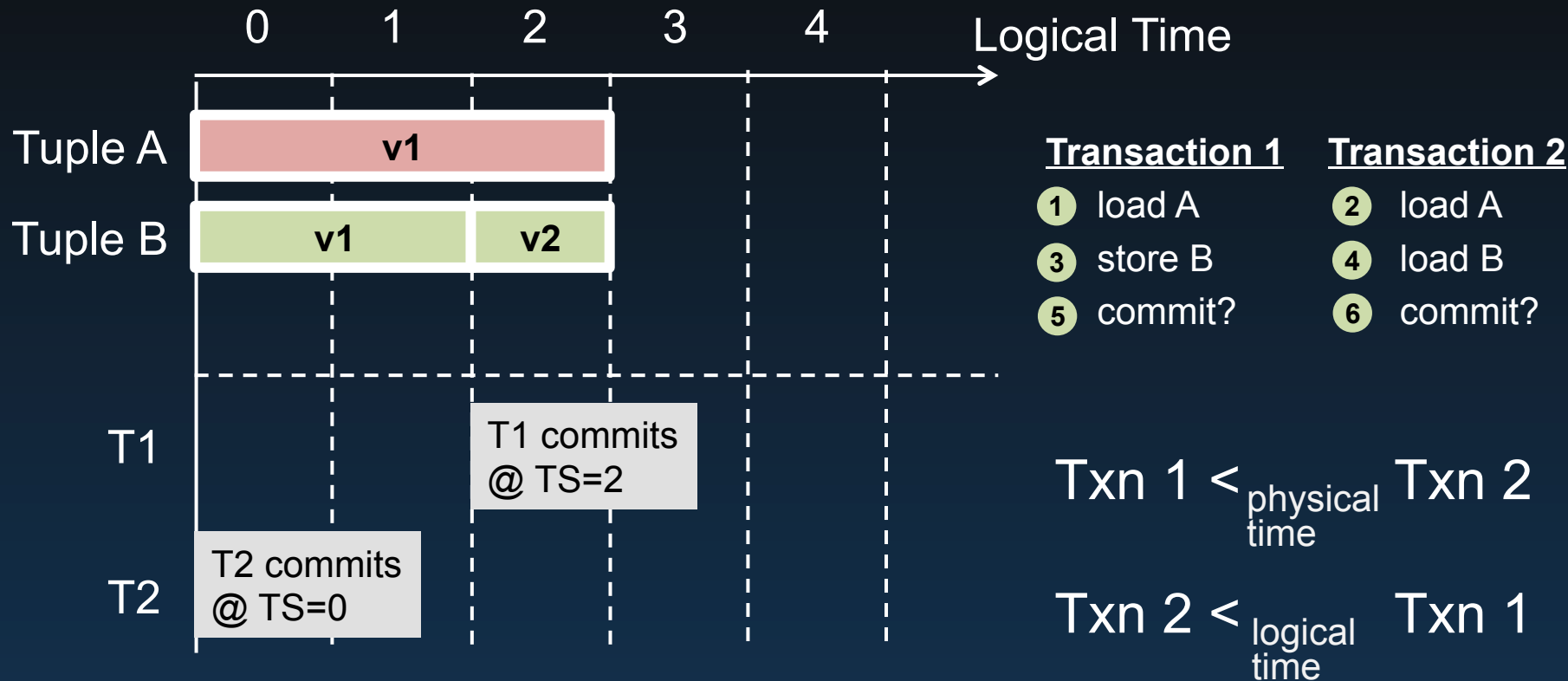




# COMMIT PHASE OF T2



# FINAL STATE



**Thm: Serializability = All operations valid at CommitTS**

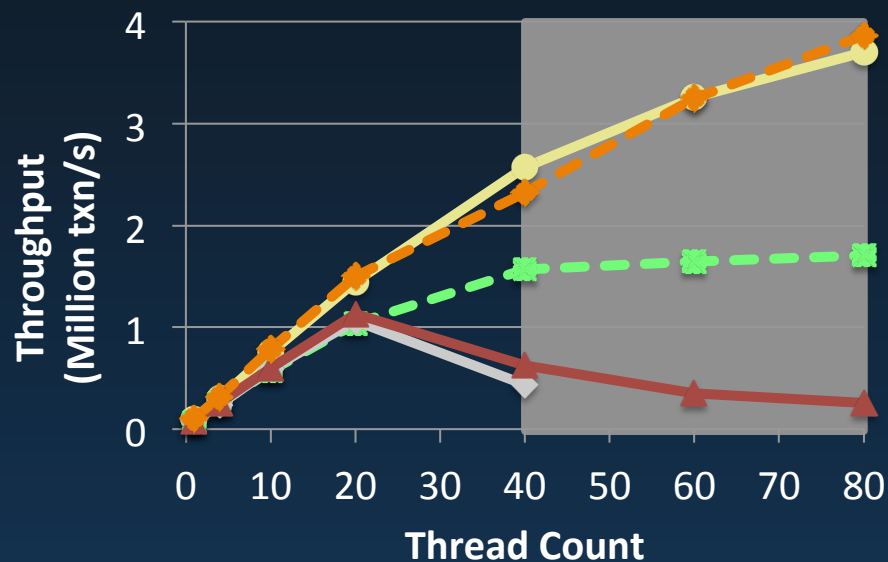
# EXPERIMENTAL SETUP

- DBx1000: Main Memory DBMS
  - No logging
  - No B-tree (hash indexing)
- Concurrency Control Algorithms
  - MVCC: HEKATON (Microsoft)
  - OCC: SILO (Harvard/MIT)
  - 2PL: DL\_DETECT, NO\_WAIT
- 10 GB YCSB Benchmark

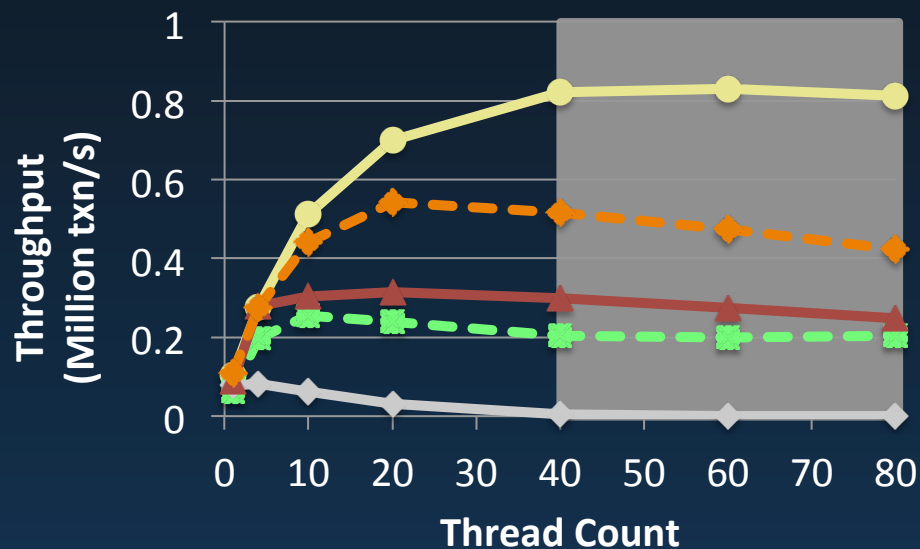
# EVALUATION

TICTOC HEKATON DL\_DETECT NO\_WAIT SILO

## YCSB -- Medium Contention



## YCSB -- High Contention

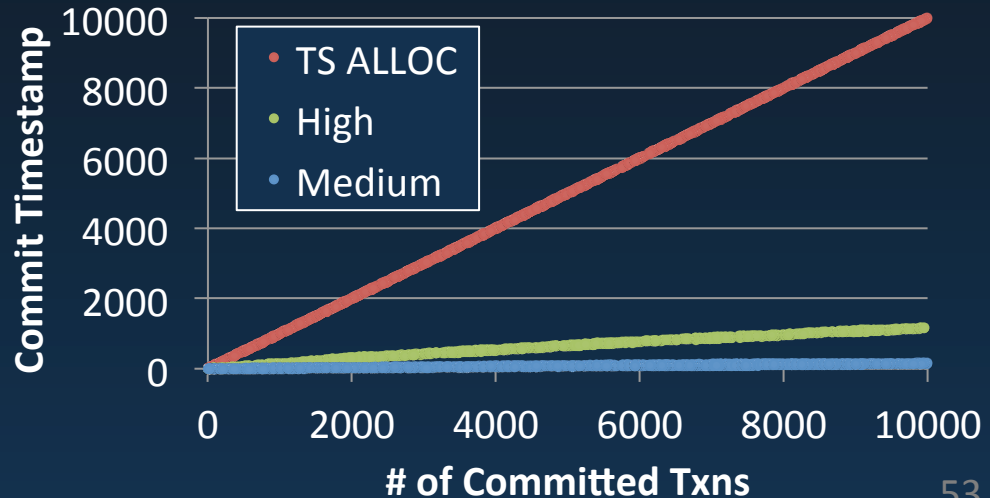


# TicToc DISCUSSION

**Thm:** Serializability = All ops valid at CommitTS

Transactions may  
have same CommitTS

Logical timestamp  
growing rate indicates  
inherent parallelism



# PHYSIOLOGICAL TIME ACROSS THE STACK



Circuit

Efficient atomic instructions



Multicore  
Processor

Tardis coherence protocol



Multicore  
Database

TicToc concurrency control



Distributed  
Database

Distributed TicToc

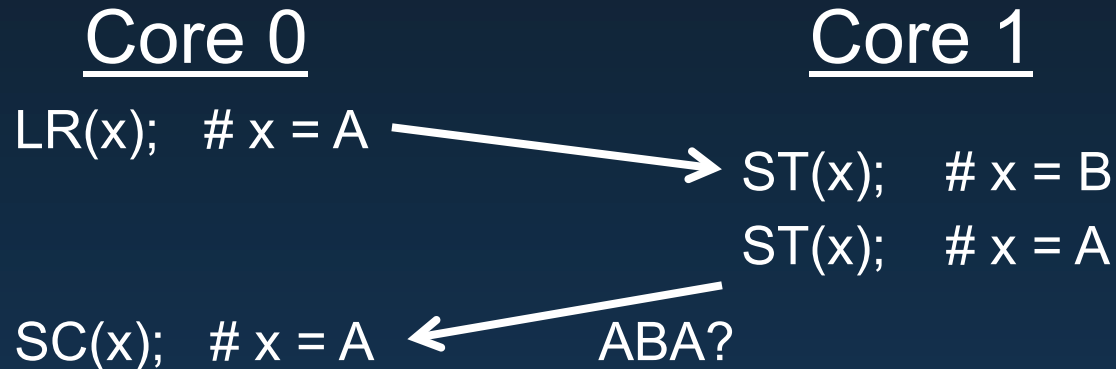


Distributed  
Shared Memory

Transaction processing with  
fault tolerance

# ATOMIC INSTRUCTION (LR/SC)

- ABA Problem
- Detect ABA using timestamp (wts)

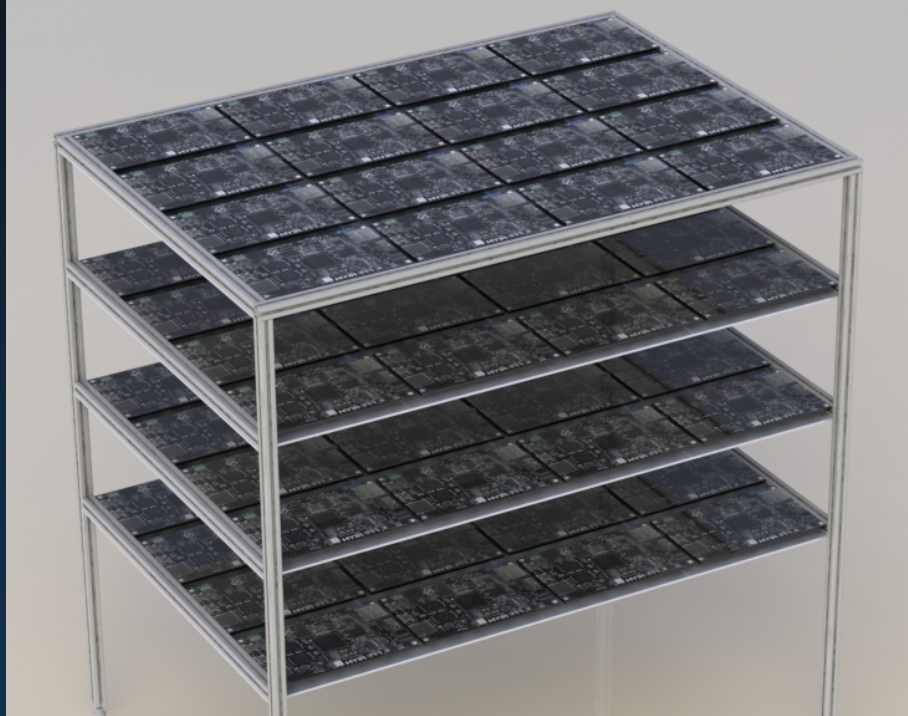


# TARDIS CACHE COHERENCE

- Simple: No Invalidation
- Scalable:
  - $O(\log N)$  storage
  - No Broadcast, No Multicast
  - No Clock Synchronization
- Support Relaxed Consistency Models



# T1000: PROPOSED 1000-CORE SHARED MEMORY PROCESSOR



# TicToc CONCURRENCY CONTROL

- Data Driven Timestamp Management
- No Central Timestamp Allocation
- Dynamic Timestamp Assignment

# DISTRIBUTED TicToc

- Data Driven Timestamp Management
- Efficient Two-Phase Commit Protocol
- Support Local Caching of Remote Data

# FAULT TOLERANT DISTRIBUTED SHARED MEMORY

- Transactional Programming Model
- Distributed Command Logging
- Dynamic Dependency Tracking Among Transactions (WAR dependency can be ignored)



# TIME TRAVELING TO ELIMINATE WAR

Xiangyao Yu, Srini Devadas  
CSAIL, MIT

