

# EASY COMMIT: A NON-BLOCKING TWO-PHASE COMMIT PROTOCOL

**Suyash Gupta, Mohammad Sadoghi**

Dept. of Computer Science  
University of California Davis

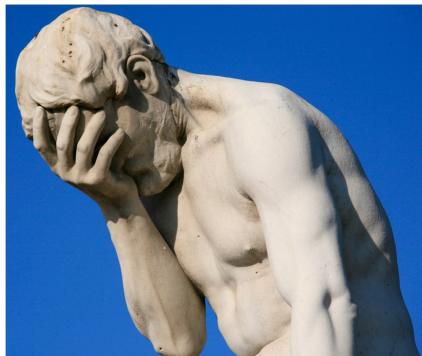
March 28, 2018

# WE ARE CHANGING THE WORLD!!!



copyright emojierra.com

# JUST KIDDING!



copyright – Alex E. Proimos (taken from Wikipedia)

# MOTIVATION

- Concerns of a Distributed Database Designer:
  - Availability
  - Consistency
- Lots of hardware, resources  $\Rightarrow$  Availability Solved.
- If **Not Consistent**  $\Rightarrow$  Correctness Compromised!

# MOTIVATION

- Concerns of a Distributed Database Designer:
  - Availability
  - Consistency
- Lots of hardware, resources  $\Rightarrow$  Availability Solved.
- If **Not Consistent**  $\Rightarrow$  Correctness Compromised!
- Need for **Agreement**!
  - Two Phase Commit Protocol (2PC)?  $\Rightarrow$  Blocked!
  - Three Phase Commit Protocol (3PC)?  $\Rightarrow$  Heavy!

# MOTIVATION

- Concerns of a Distributed Database Designer:
  - Availability
  - Consistency
- Lots of hardware, resources  $\Rightarrow$  Availability Solved.
- If **Not Consistent**  $\Rightarrow$  Correctness Compromised!
- Need for **Agreement!**
  - Two Phase Commit Protocol (2PC)?  $\Rightarrow$  Blocked!
  - Three Phase Commit Protocol (3PC)?  $\Rightarrow$  Heavy!
- **Easy Commit**  $\Rightarrow$  Two Phases! Non-Blocking!

# TWO PHASE COMMIT PROTOCOL

Coordinator

Start Commit



INITIAL

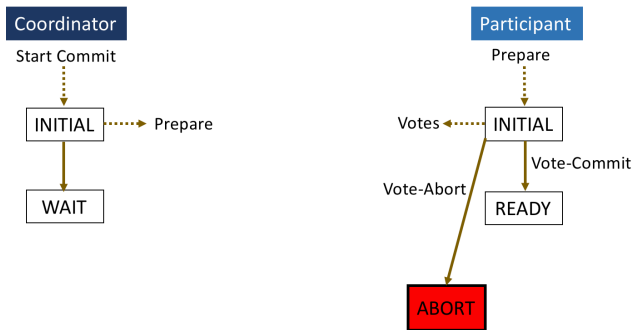
Participant

Prepare



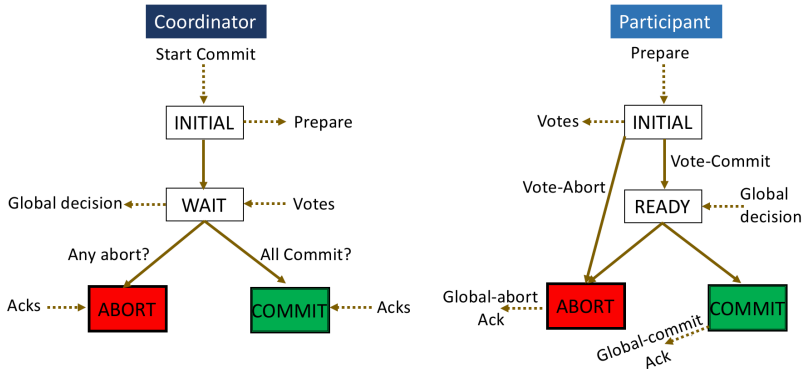
INITIAL

# TWO PHASE COMMIT PROTOCOL

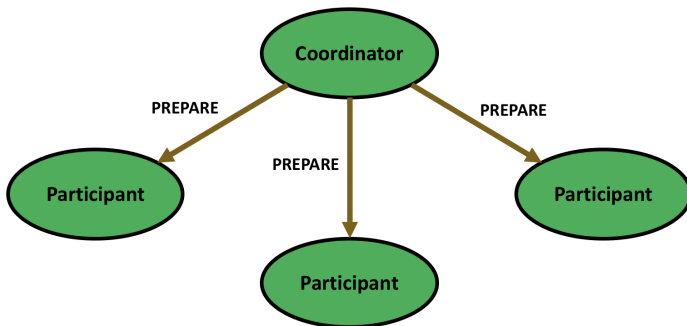




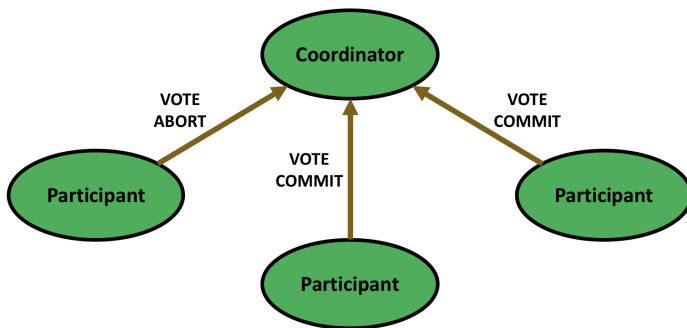
# TWO PHASE COMMIT PROTOCOL



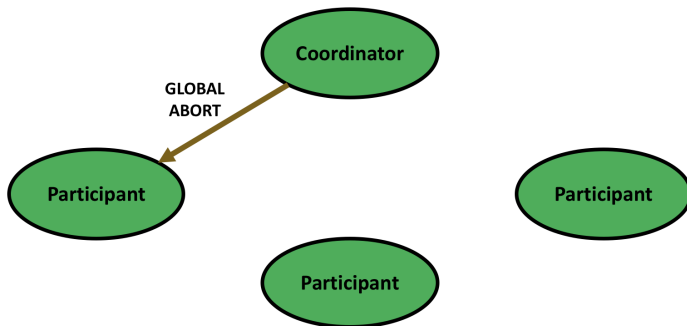
# 2PC IS BLOCKING



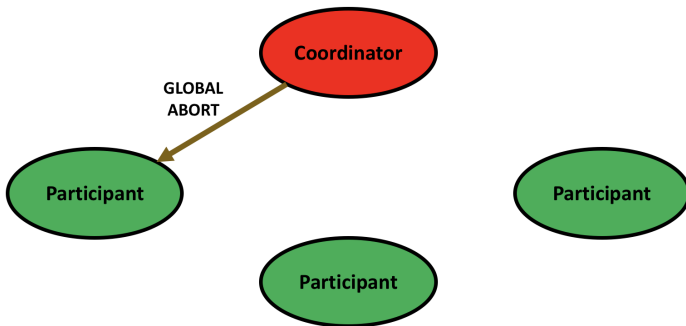
# 2PC IS BLOCKING



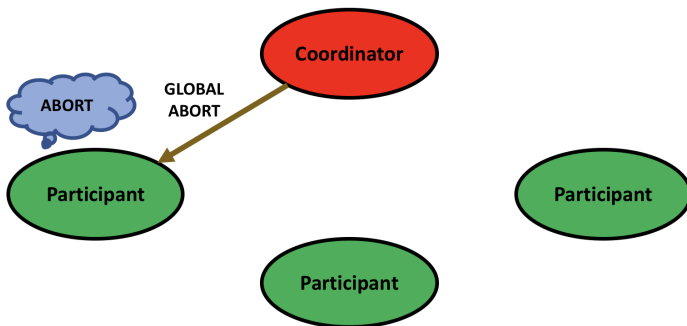
# 2PC IS BLOCKING



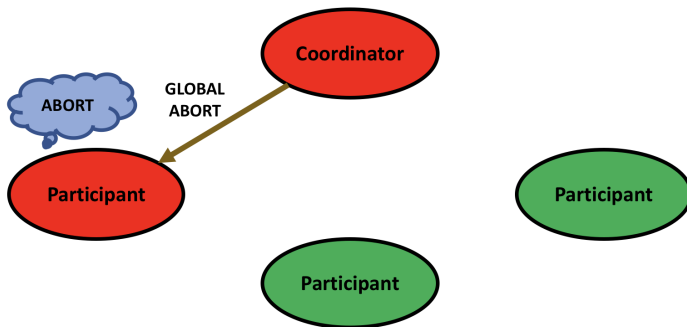
# 2PC IS BLOCKING



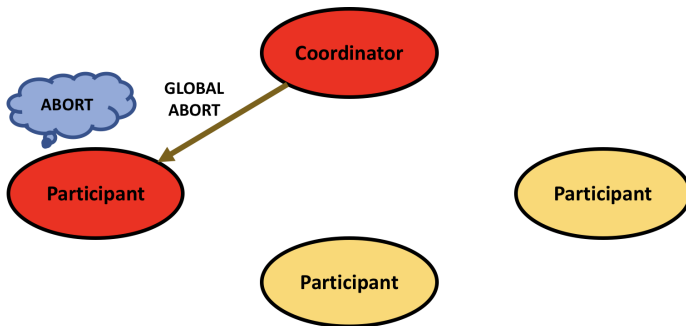
# 2PC IS BLOCKING



# 2PC IS BLOCKING

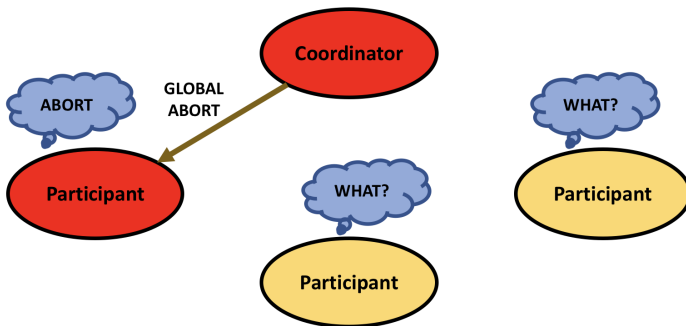


# 2PC IS BLOCKING





# 2PC IS BLOCKING



**System is in an Unstable State.**

# THREE PHASE COMMIT PROTOCOL

Coordinator

Start Commit



INITIAL

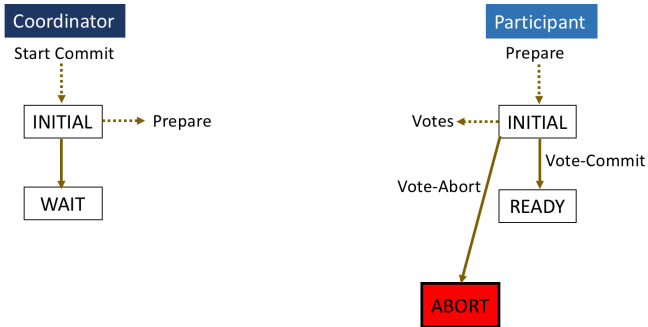
Participant

Prepare

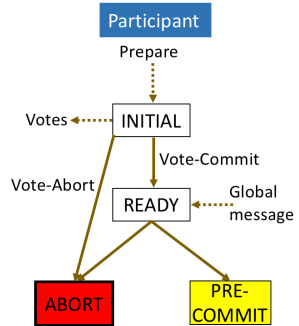
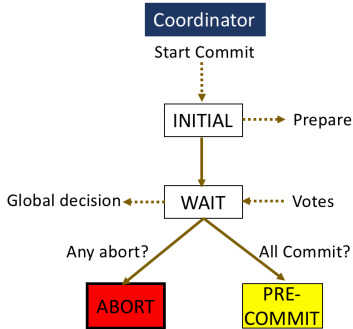


INITIAL

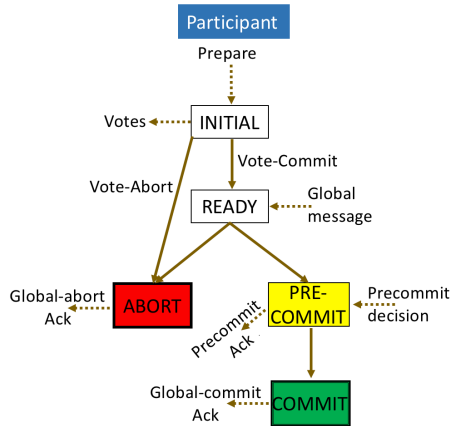
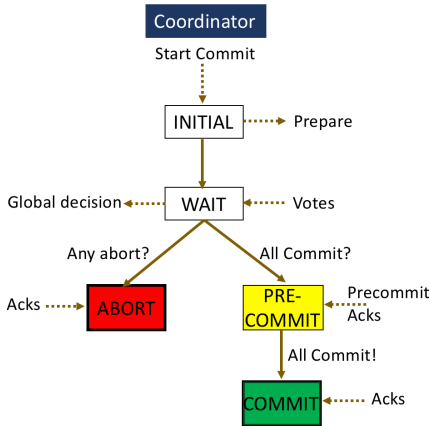
# THREE PHASE COMMIT PROTOCOL



# THREE PHASE COMMIT PROTOCOL



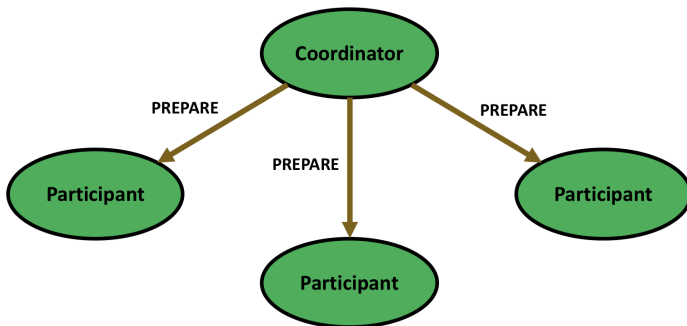
# THREE PHASE COMMIT PROTOCOL



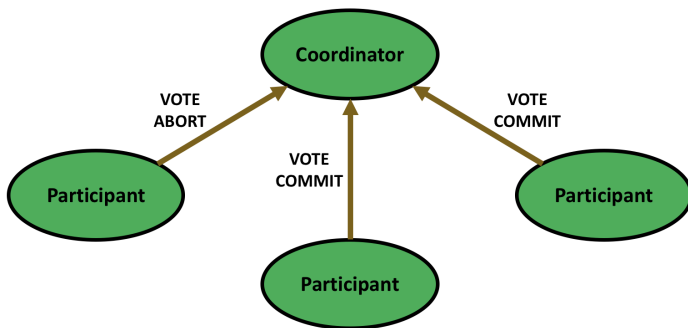
# 3PC IS NON-BLOCKING

- **Dale Skeen's** Requirements [SIGMOD 1981]:
- No state adjacent to both **Abort** and **Commit** states.
- No non-committable state adjacent to the **Commit** state.

# 3PC WORKS

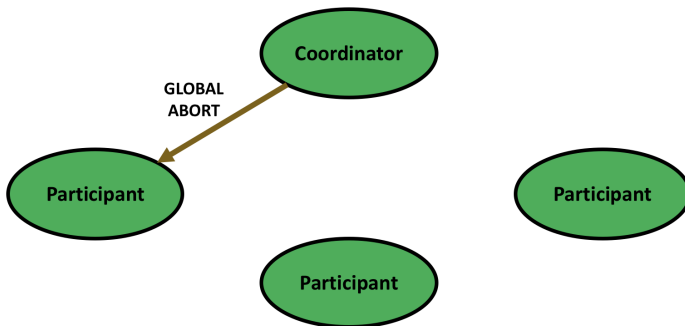


# 3PC WORKS

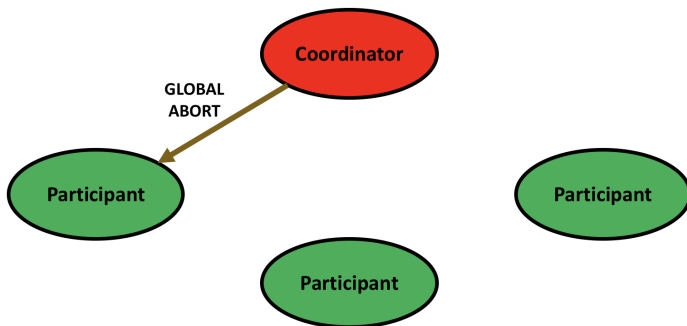




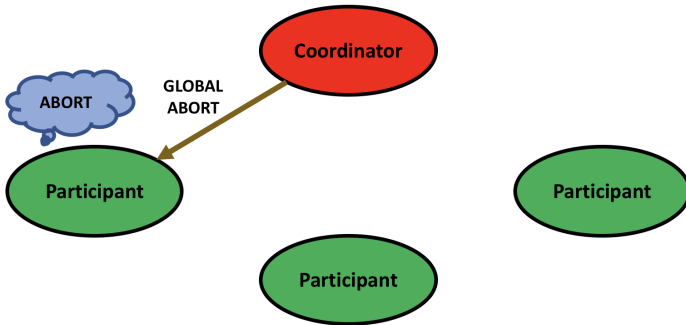
# 3PC WORKS



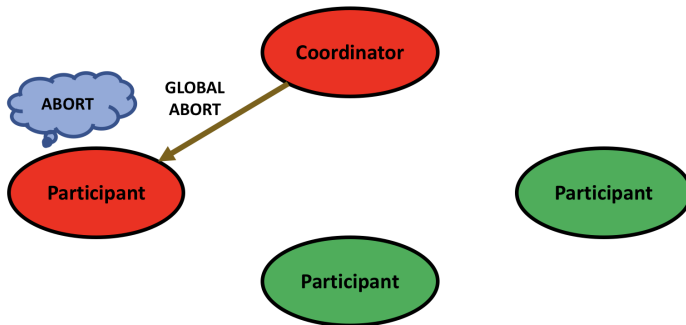
# 3PC WORKS



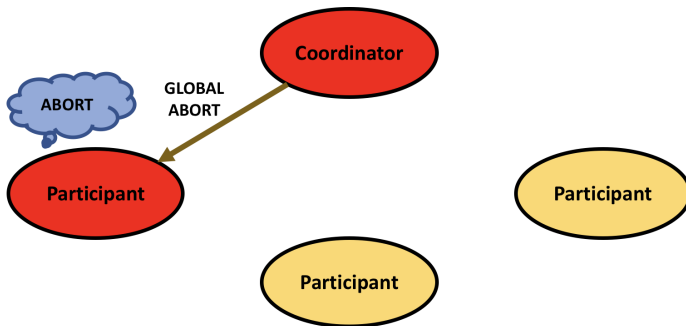
# 3PC WORKS



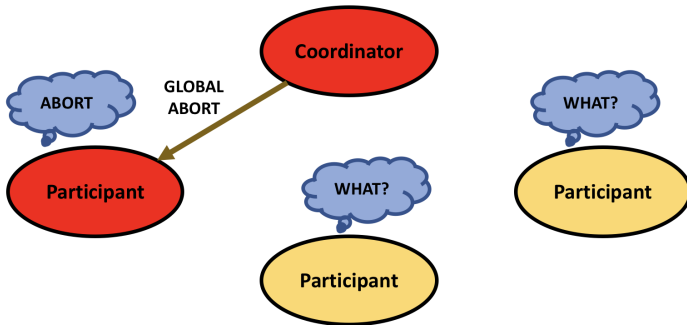
# 3PC WORKS



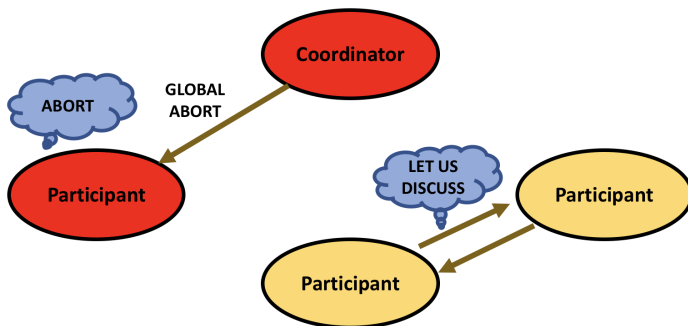
# 3PC WORKS



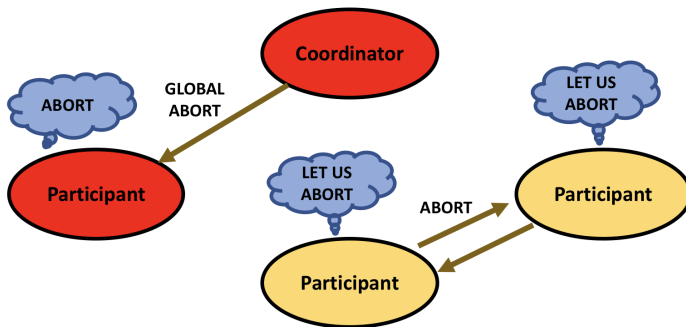
# 3PC WORKS



# 3PC WORKS

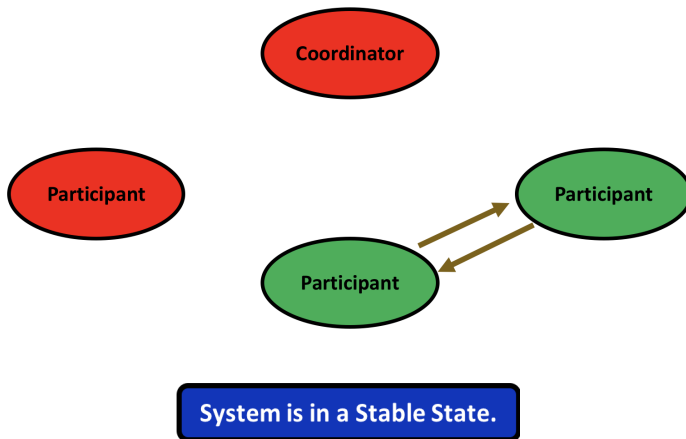


# 3PC WORKS





# 3PC WORKS



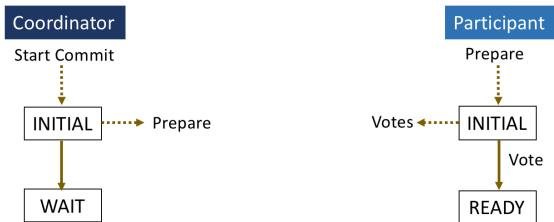
# EASY COMMIT PRINCIPLE

- First Transmit and then Commit
- Message Redundancy

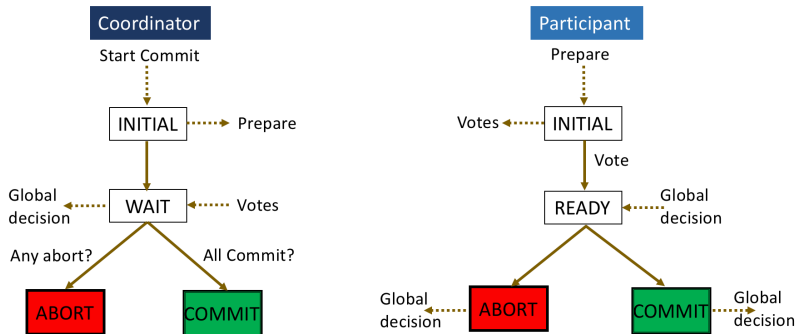
# EASY COMMIT PROTOCOL



# EASY COMMIT PROTOCOL



# EASY COMMIT PROTOCOL



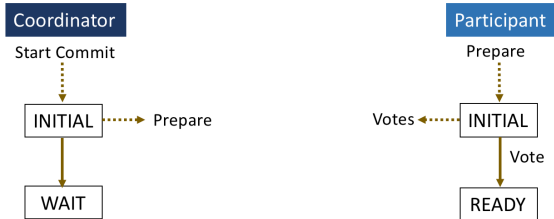
# EASY COMMIT OBSERVATIONS

- Participant cannot directly transition from **INITIAL** to **ABORT**.
- Each participant forwards the global decision to every other node.
- Participant need not wait for global decision from coordinator.
- Existence of **hidden** states: **TRANSMIT-A** and **TRANSMIT-C**.

# EASY COMMIT PROTOCOL – LOGICAL EXPANSION

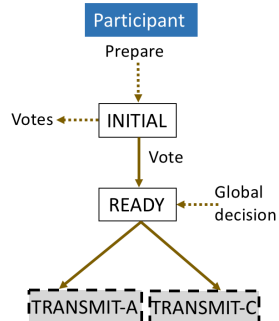
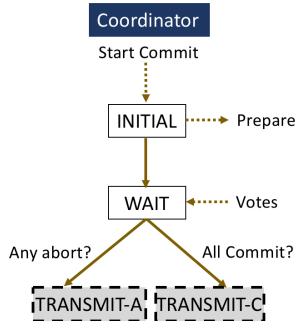


# EASY COMMIT PROTOCOL – LOGICAL EXPANSION

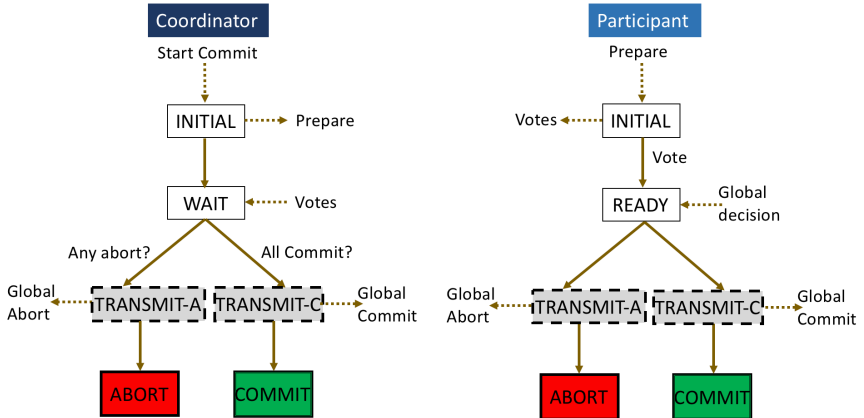




# EASY COMMIT PROTOCOL – LOGICAL EXPANSION



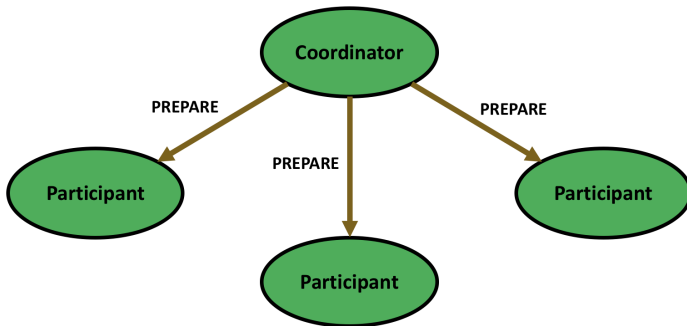
# EASY COMMIT PROTOCOL – LOGICAL EXPANSION



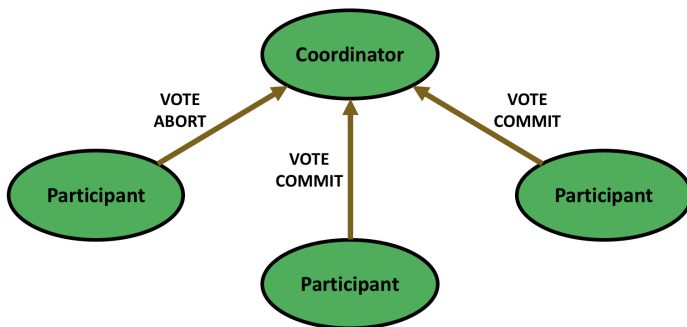
# EASY COMMIT TERMINATION PROTOCOL

- **Coordinator Timeout in WAIT state:**
  - Coordinator didn't receive **Votes**.
  - Adds a log entry, Transmits **Global-Abort**, Aborts Transaction.
- **Participant Timeout in INITIAL state:**
  - Participant didn't receive **PREPARE** message.
  - Communicates with other participants.
- **Participant Timeout in READY state:**
  - Participant didn't receive **Global Decision**.
  - Communicates with other participants.

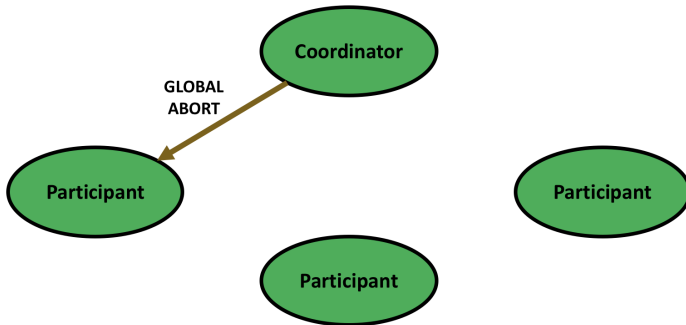
# EASY COMMIT WORKS



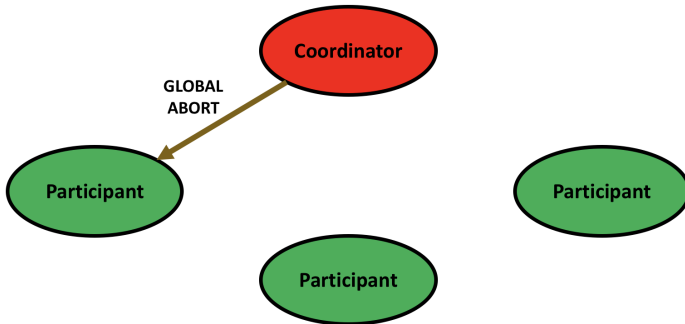
# EASY COMMIT WORKS



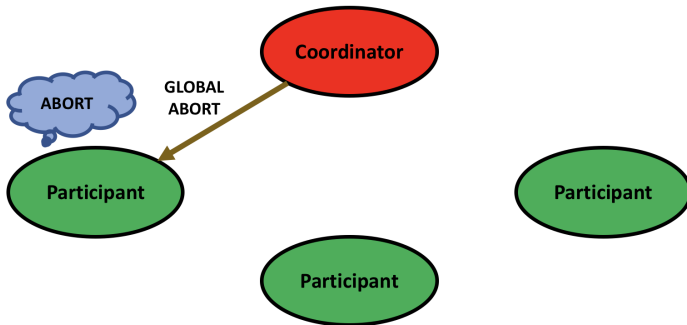
# EASY COMMIT WORKS



# EASY COMMIT WORKS

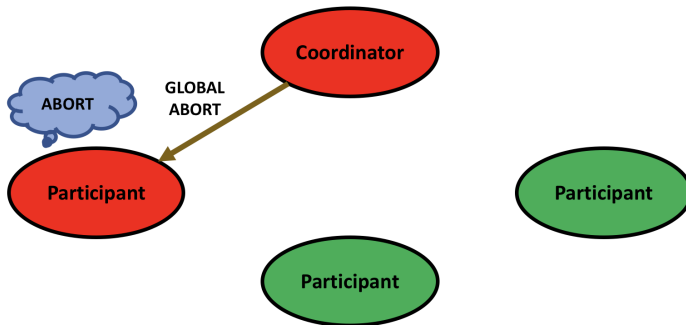


# EASY COMMIT WORKS

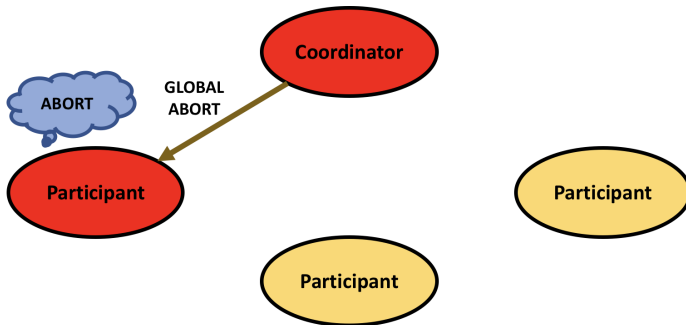




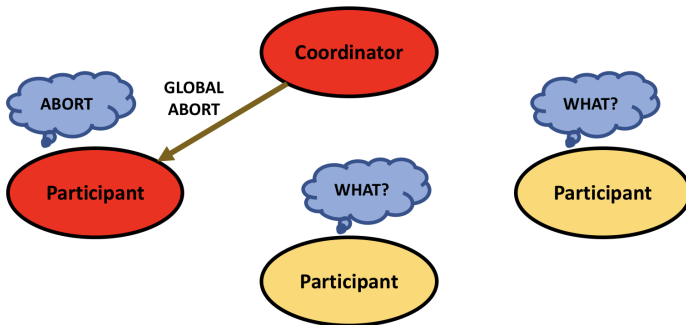
# EASY COMMIT WORKS



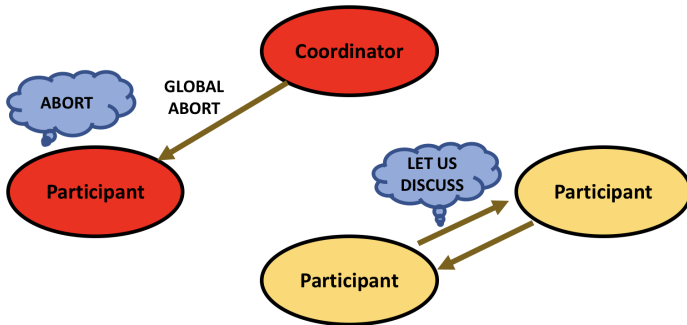
# EASY COMMIT WORKS



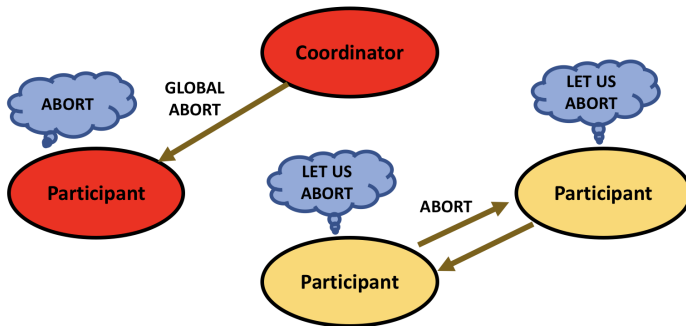
# EASY COMMIT WORKS



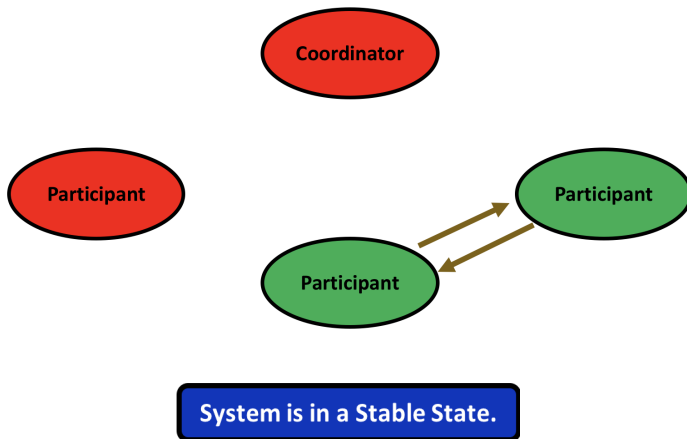
# EASY COMMIT WORKS



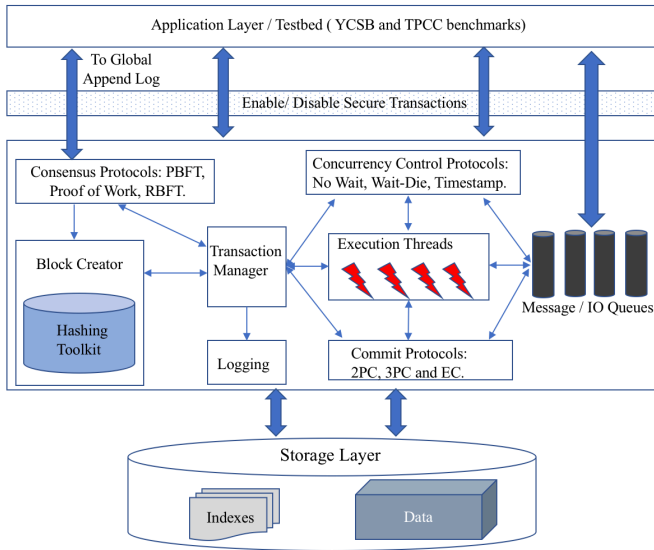
# EASY COMMIT WORKS



# EASY COMMIT WORKS



# IMPLEMENTATION: ExpODB

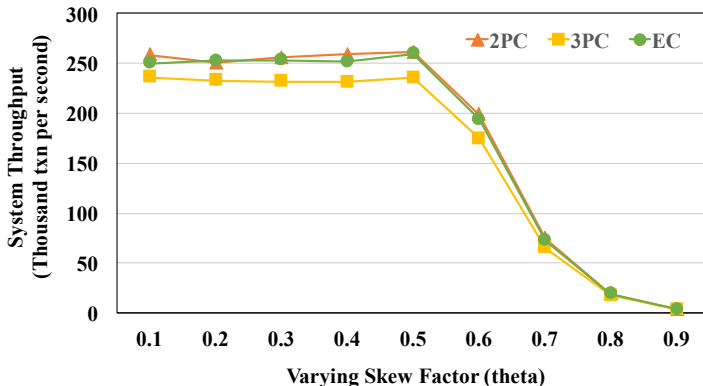


# EVALUATION

- 64 Standard\_D8S\_V3 Azure instances, deployed in US East.
  - Each machine has 8 cores and 32GB memory.
- 4 Worker threads attached to a dedicated core.
- Load of 10000 open client connections per node.
- First 60s warmup and next 60s execution.
- Results averaged over three runs.
- NO-WAIT concurrency control algorithm used.
- Two benchmark suites: YCSB and TPC-C.



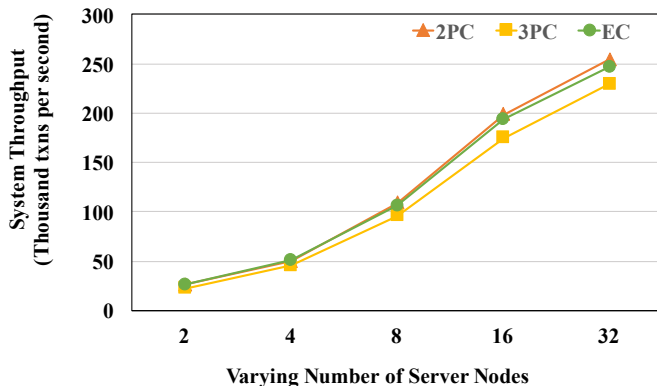
# VARYING SKEW FACTOR – YCSB ZIPFIAN THETA



Number of server nodes = 16 and partitions per transaction = 2.

**On varying the YCSB skew factor,  
EasyCommit throughput is equivalent to 2PC.**

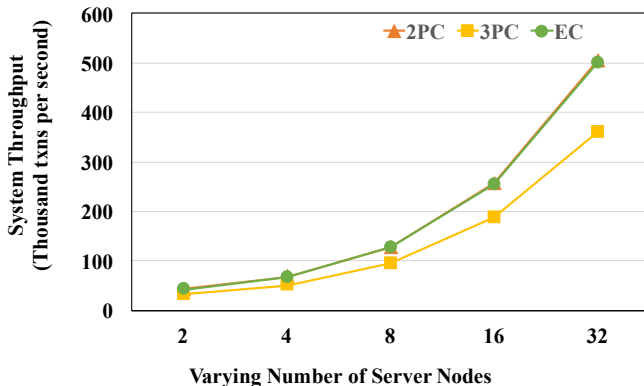
# VARYING SERVER NODES



Partitions per transaction = 2 and Skew factor = 0.6.

**EasyCommit is scalable as 2PC,  
on increasing the number of server nodes.**

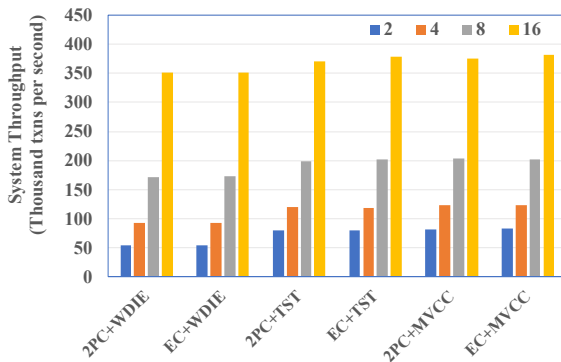
# TPC-C PAYMENT TRANSACTIONS



Number of warehouses per server = 128.

**EasyCommit is scalable as 2PC,  
on increasing the number of server nodes.**

# VARYING CONCURRENCY CONTROL ALGORITHMS



TPC-C benchmarking for 16 servers, 128 warehouses per server and 3 CC algorithms: WDIE (WAIT-DIE) and TST (TIMESTAMP).

**EasyCommit design is orthogonal to underlying Concurrency Control.**

# CONCLUSIONS

- We present novel commit protocol – **Easy Commit**.
- Leverages best of twin worlds (2PC and 3PC).
- Two key observations:
  - First transmit and then commit.
  - Message Redundancy.
- Easy Commit guarantees both **safety** and **liveness**.

# CONCLUSIONS

- We present novel commit protocol – **Easy Commit**.
- Leverages best of twin worlds (2PC and 3PC).
- Two key observations:
  - First transmit and then commit.
  - Message Redundancy.
- Easy Commit guarantees both **safety** and **liveness**.

