

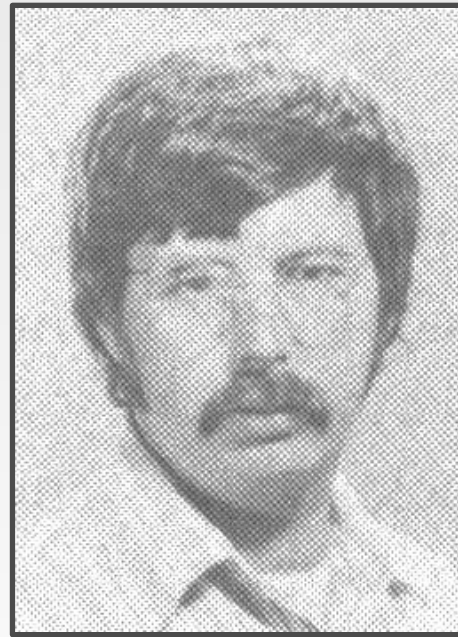
**MY DATABASE SYSTEM IS
THE ONLY THING I CAN**



TRUST

@ANDY_PAVLO

Thirty Years Ago...





INTERACTIVE TRANSACTIONS



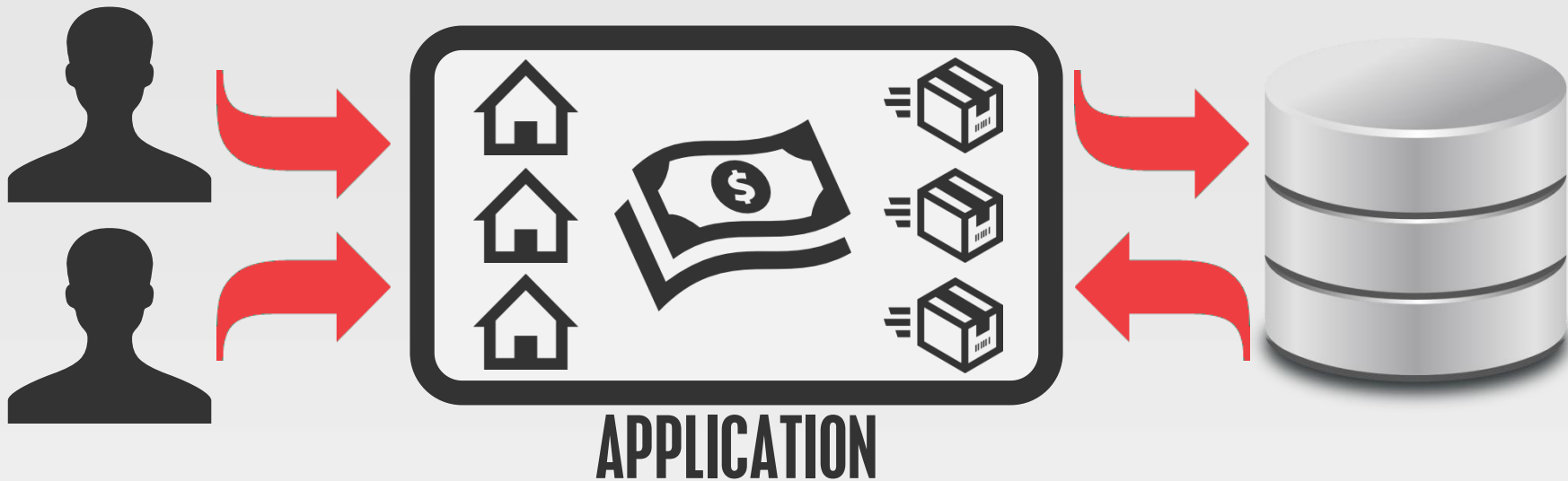
SMALL # OF CPU CORES



SMALL MEMORY SIZES

TPC-C BENCHMARK

Warehouse Order Processing

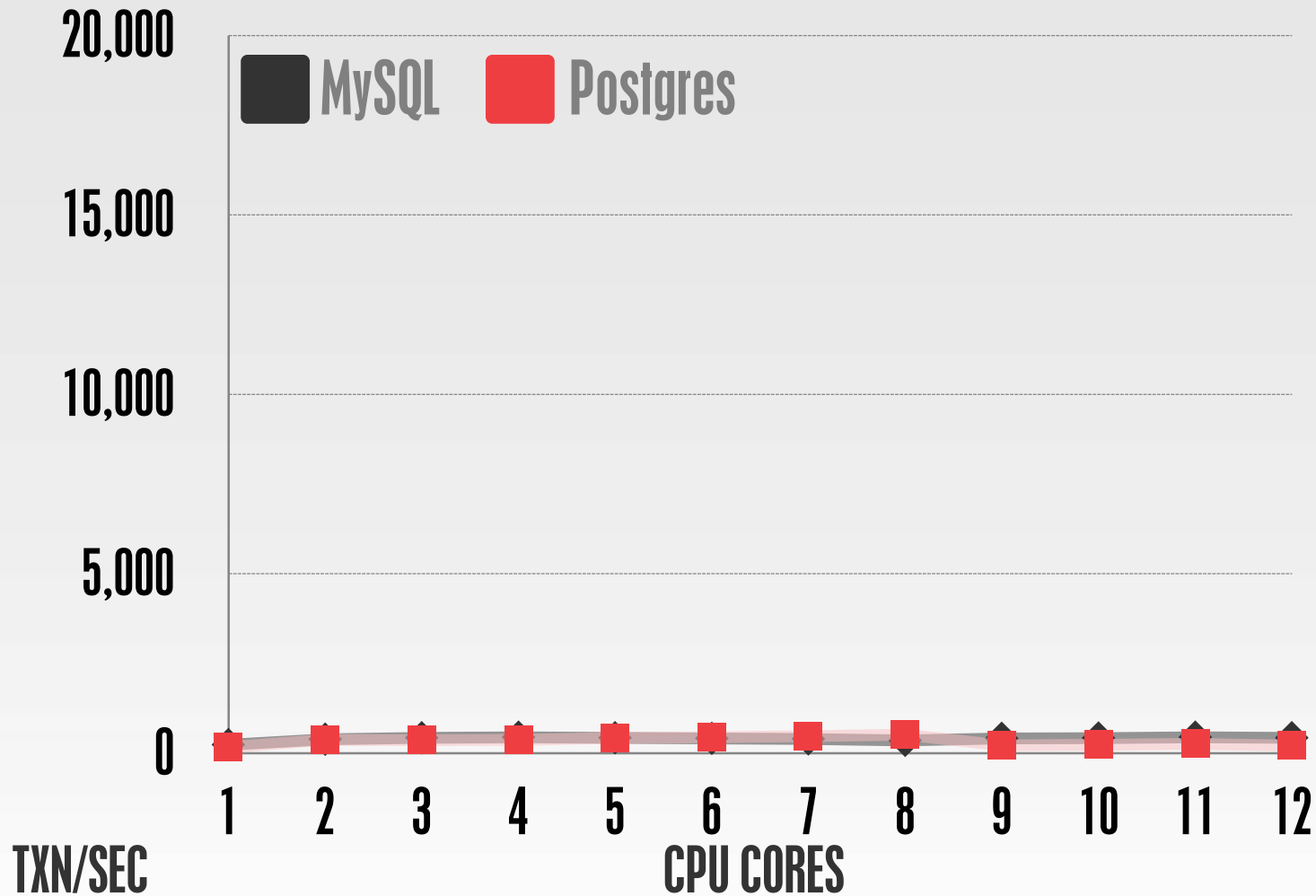


NewOrder Transaction

1. Check item stock level.
2. Create new order information.
3. Update item stock levels.

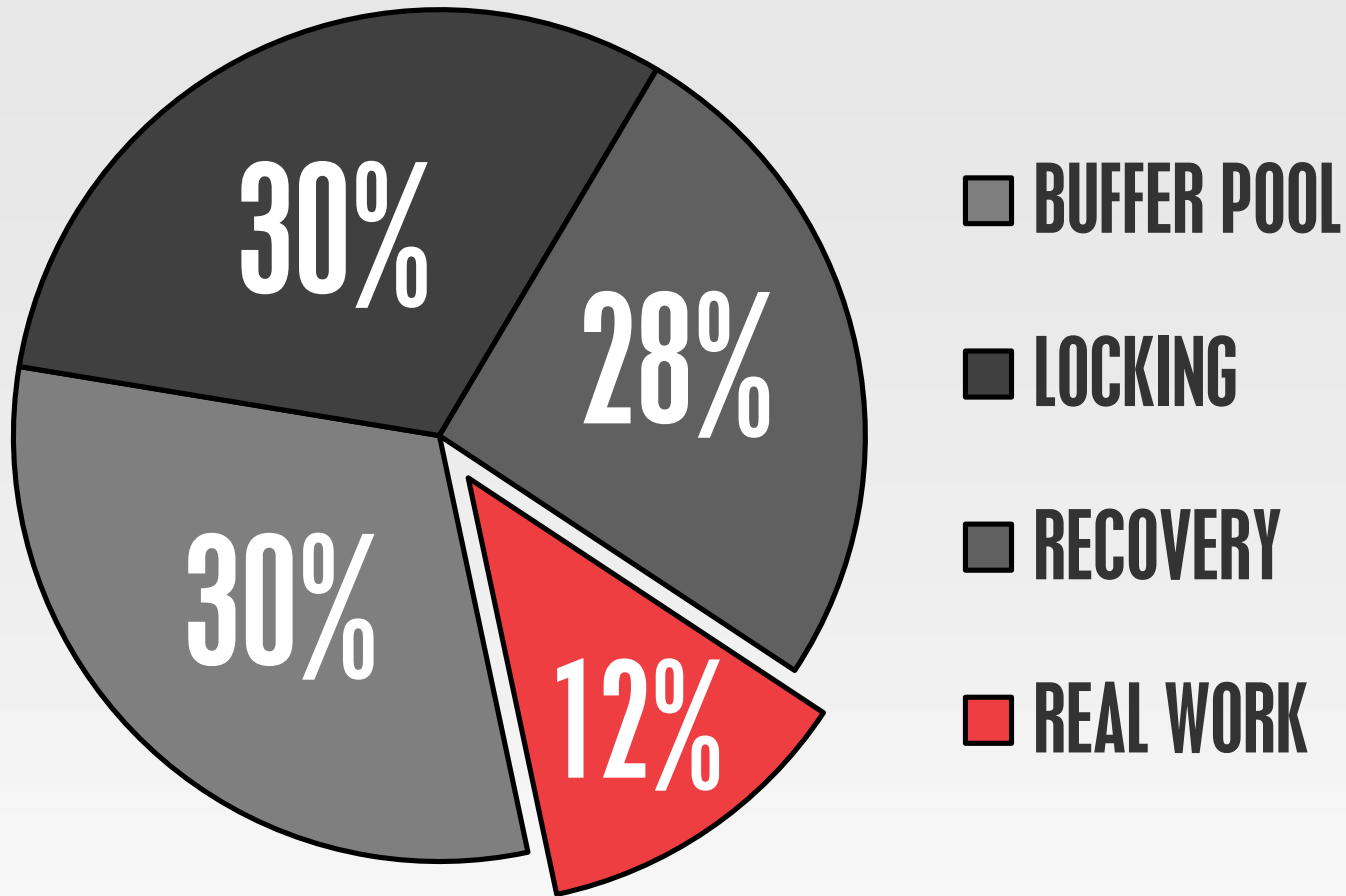
TPC-C BENCHMARK

Warehouse Order Processing



TRADITIONAL DBMS

Measured CPU Cycles



OLTP THROUGH THE LOOKING GLASS,
AND WHAT WE FOUND THERE
SIGMOD, pp. 981-992, 2008.



HARDWARE UPGRADE



REPLICATION



DISTRIBUTED CACHE



SHARDING MIDDLEWARE



NOSQL

**HOW TO SCALE UP
WITHOUT GIVING UP
TRANSACTIONS?**

-Store

Distributed Main Memory
Transaction Processing System



H-STORE: A HIGH-PERFORMANCE, DISTRIBUTED
MAIN MEMORY TRANSACTION PROCESSING SYSTEM
Proc. VLDB Endow., vol. 1, iss. 2, pp. 1496-1499, 2008.



**DISK ORIENTED
MAIN MEMORY STORAGE**

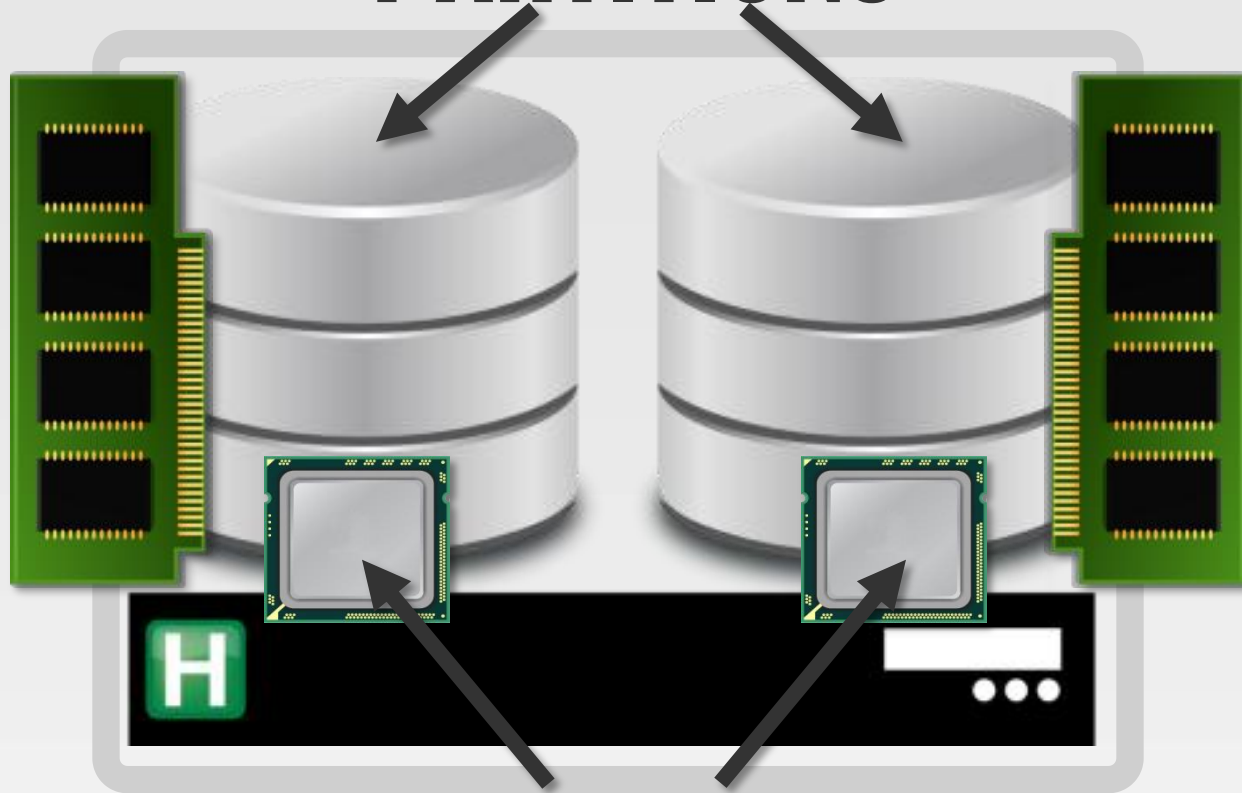


**CONCURRENT EXECUTION
SERIAL EXECUTION**

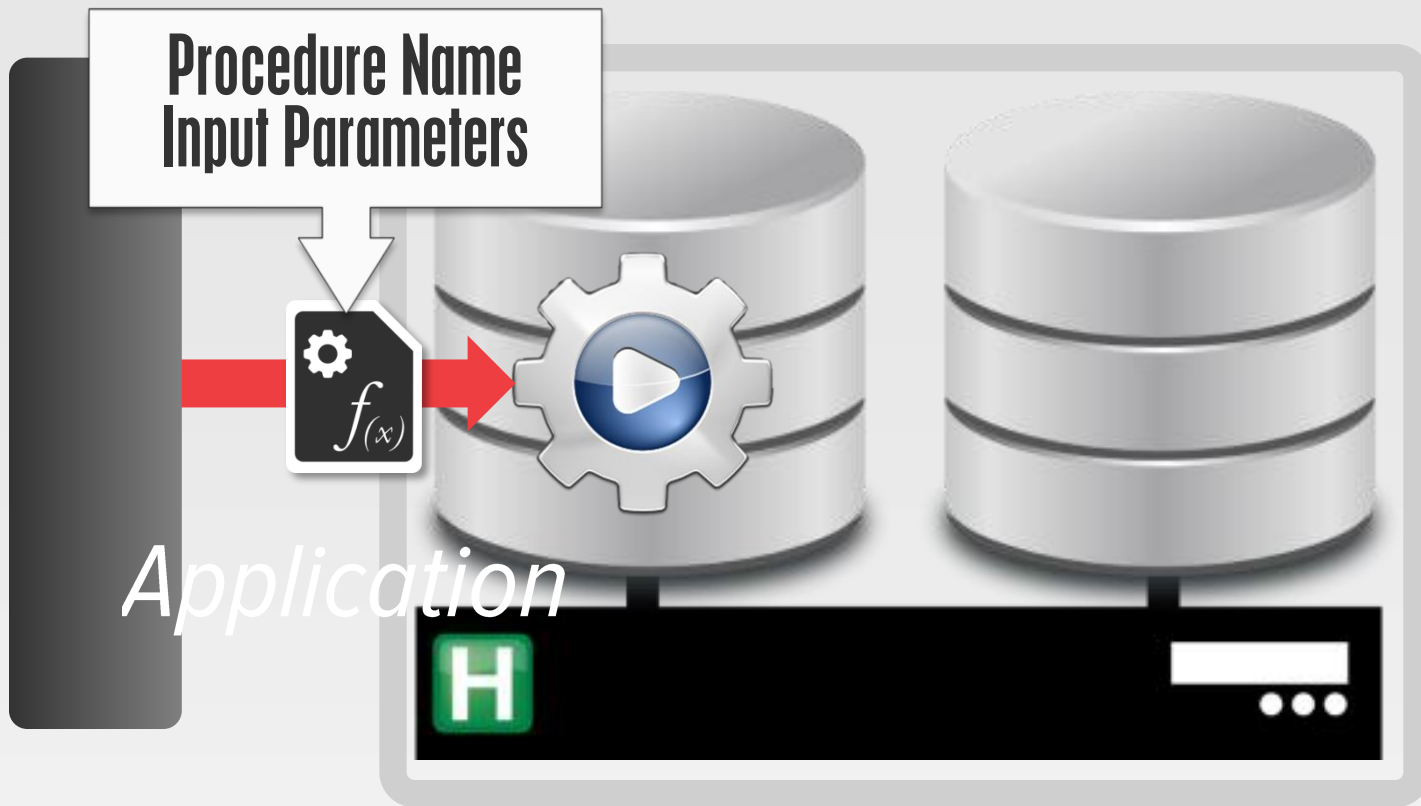


**HEAVYWEIGHT RECOVERY
COMPACT LOGGING**

PARTITIONS



SINGLE-THREADED EXECUTION ENGINES



STORED PROCEDURE

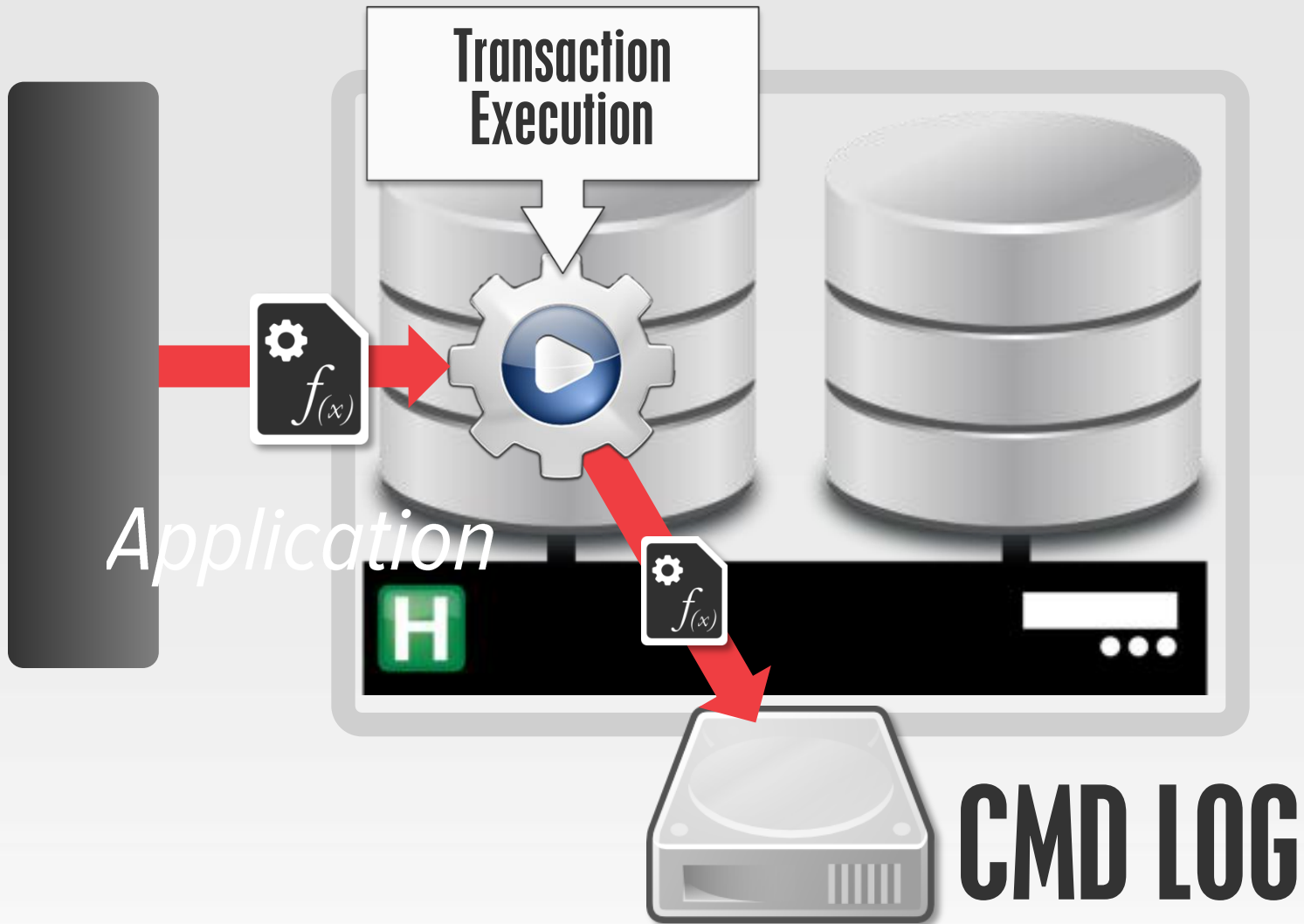
VoteCount:

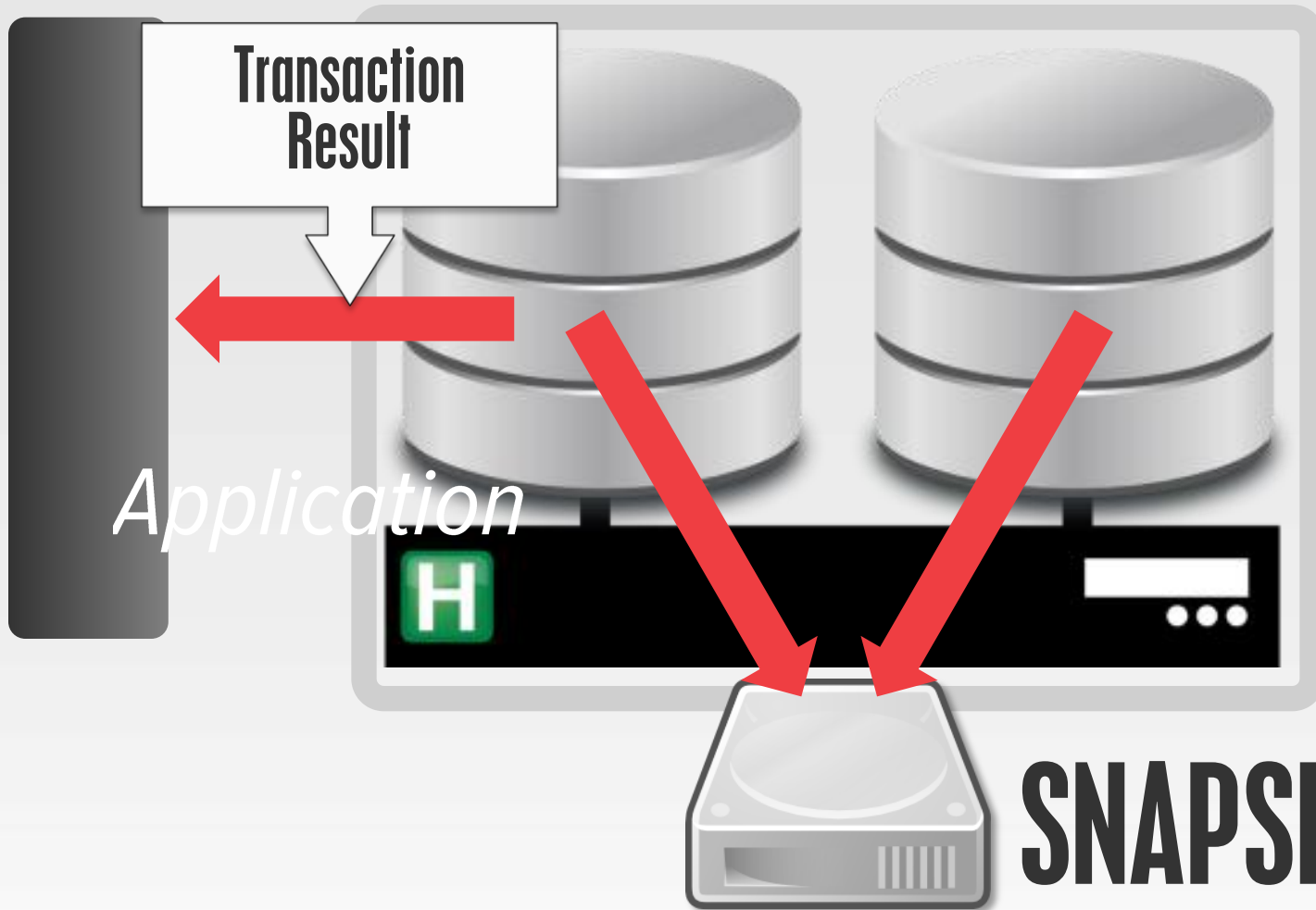
```
SELECT COUNT(*)  
  FROM votes  
 WHERE phone_num = ?;
```

InsertVote:

```
INSERT INTO votes  
VALUES (?, ?, ?);
```

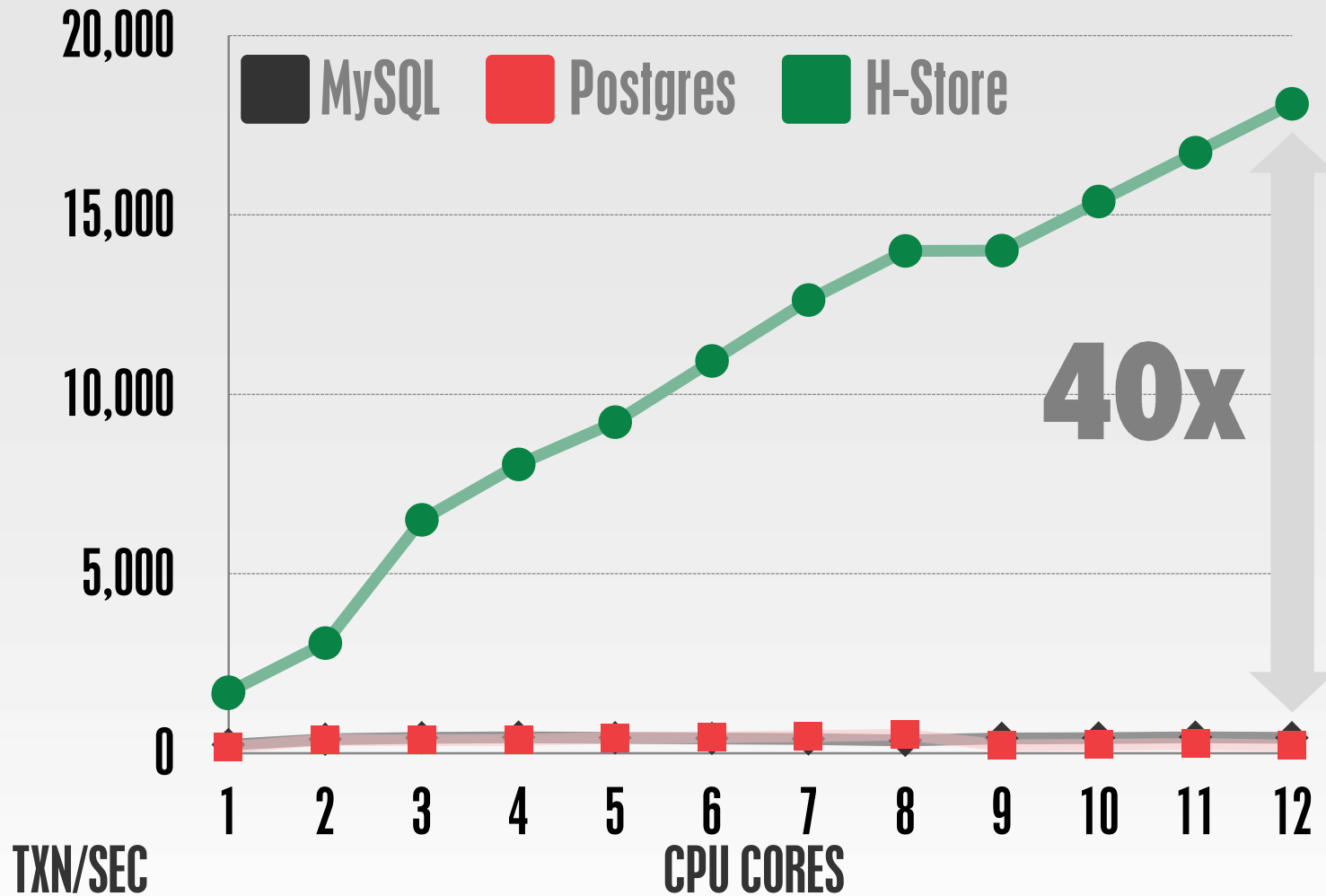
```
run(phoneNum, contestantId, currentTime) {  
    result = execute(VoteCount, phoneNum);  
    if (result > MAX_VOTES) {  
        return (ERROR);  
    }  
    execute(InsertVote, phoneNum,  
           contestantId,  
           currentTime);  
    return (SUCCESS);  
}
```





TPC-C BENCHMARK

Warehouse Order Processing



DISTRIBUTED TRANSACTIONS

TPC-C BENCHMARK

8 Cores per Node
10% Distributed Transactions



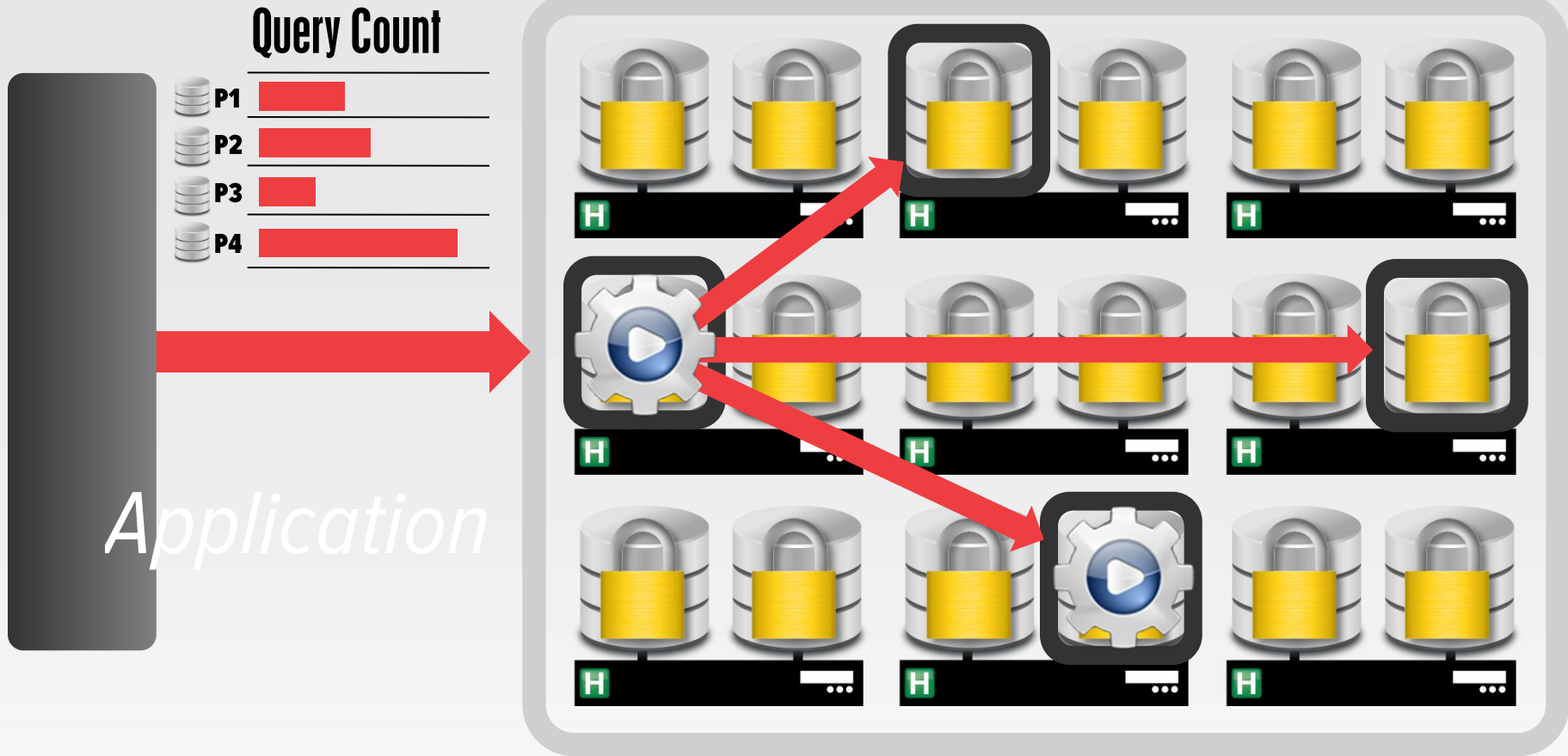
DISTRIBUTED TRANSACTIONS



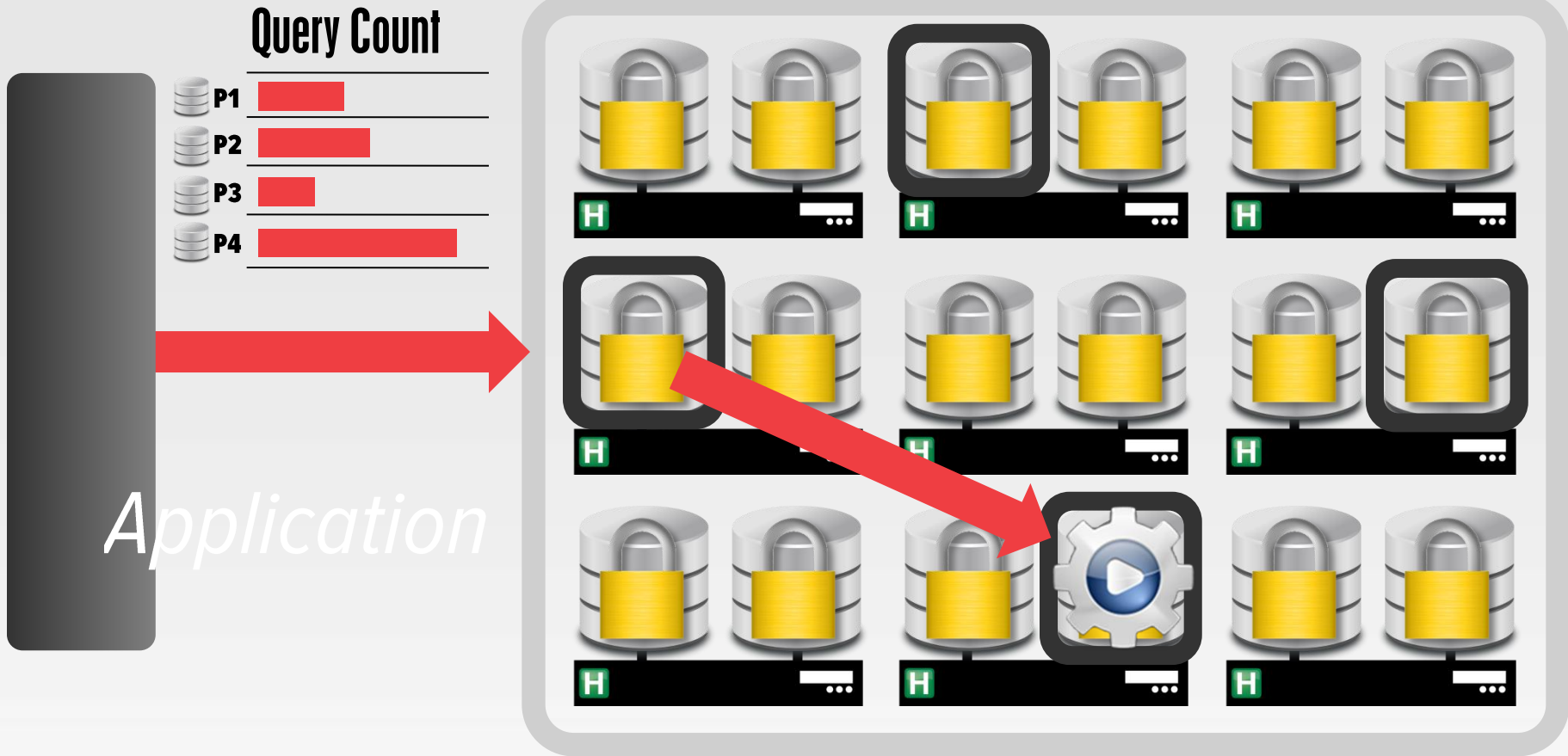
DISTRIBUTED TRANSACTIONS



DISTRIBUTED TRANSACTIONS



DISTRIBUTED TRANSACTIONS



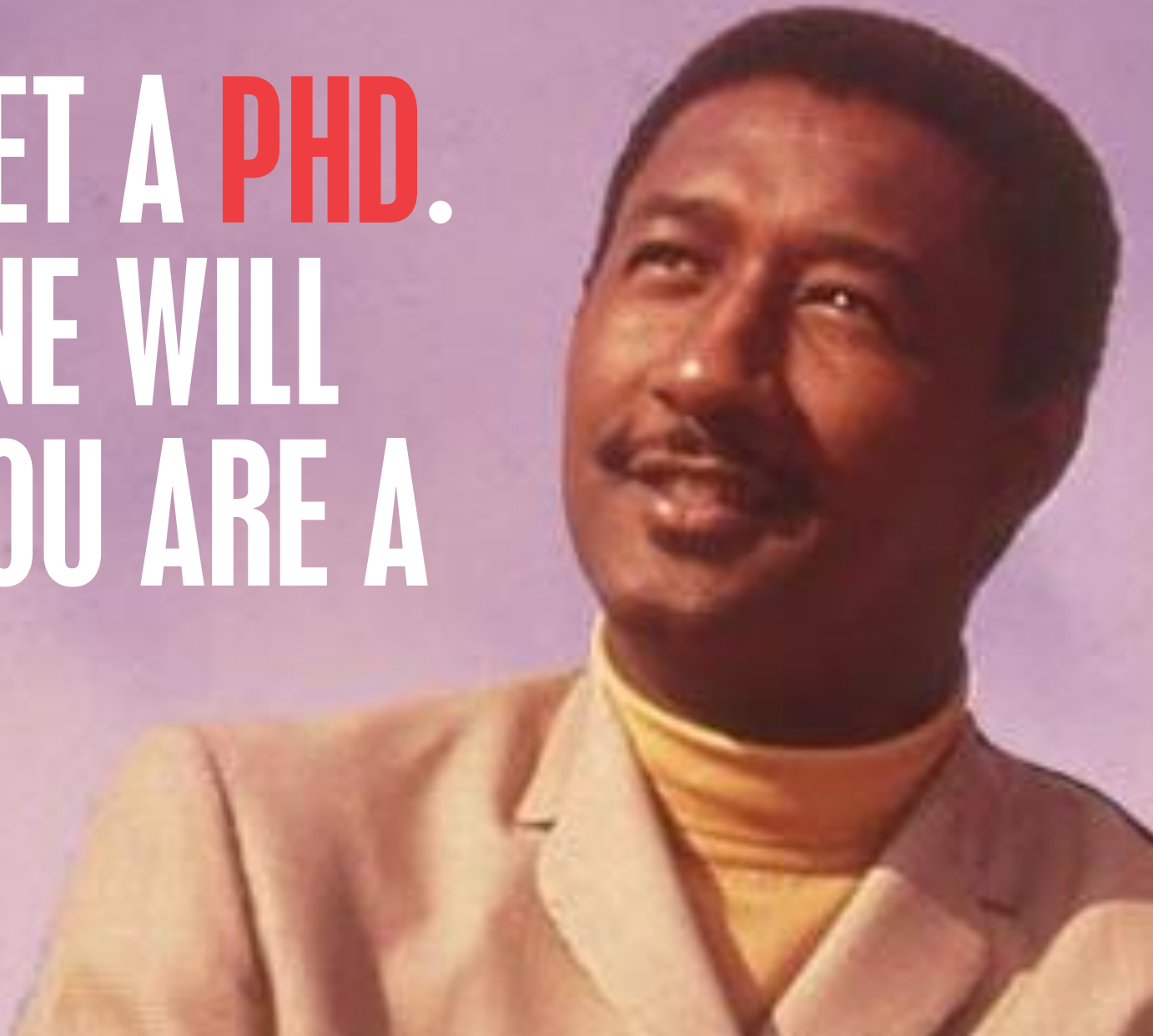
**KNOW WHAT TRANSACTIONS
WILL DO BEFORE THEY START**

**BUT PEOPLE ALWAYS
GIVE ME BAD ADVICE**

**DON'T GET
INVOLVED WITH
COMPUTERS.
YOU'LL NEVER
MAKE ANY
MONEY.**



**DON'T GET A PHD.
EVERYONE WILL
THINK YOU ARE A
JERK.**



**THE DATABASE SYSTEM
ALWAYS HAS MORE**



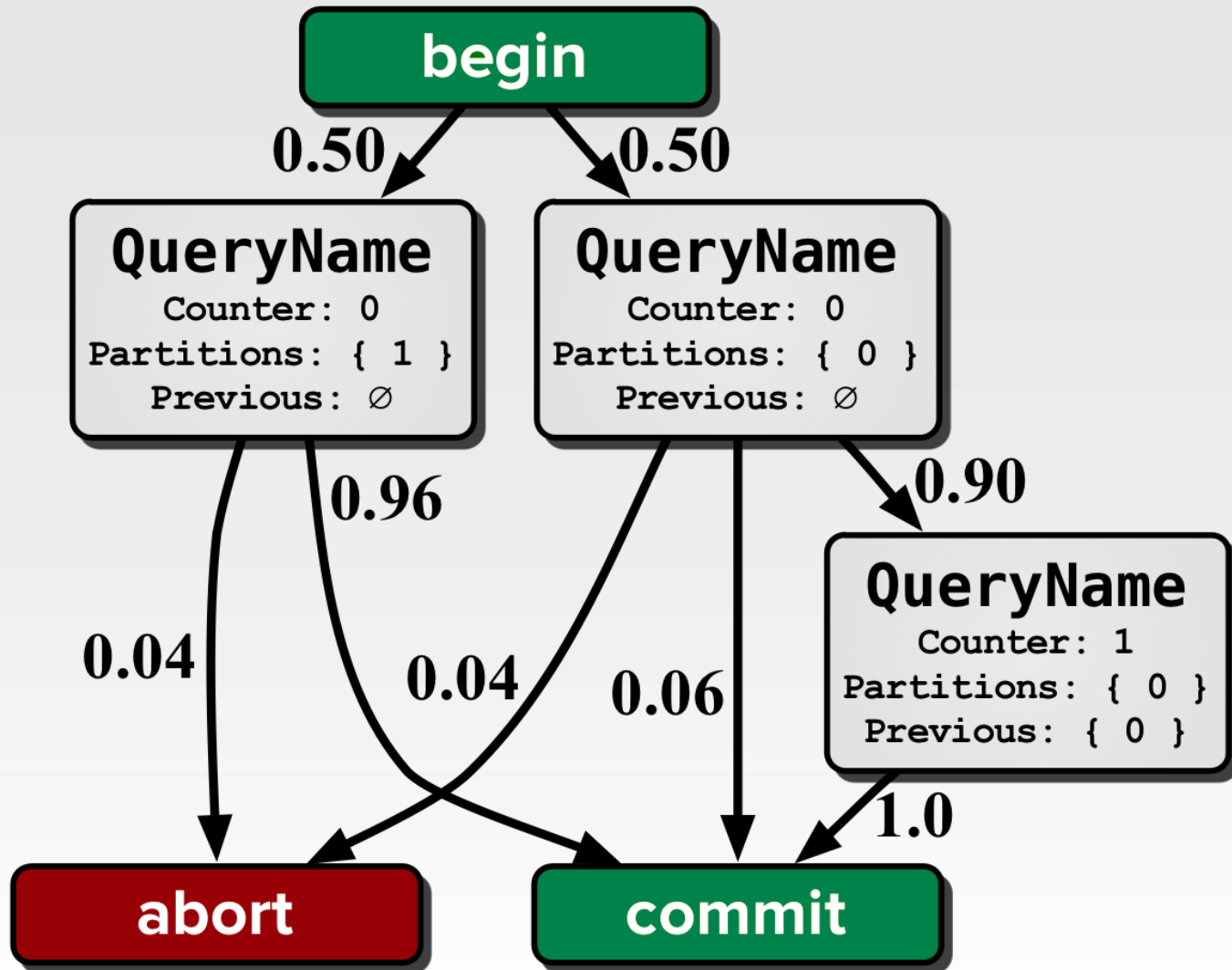
INFORMATION

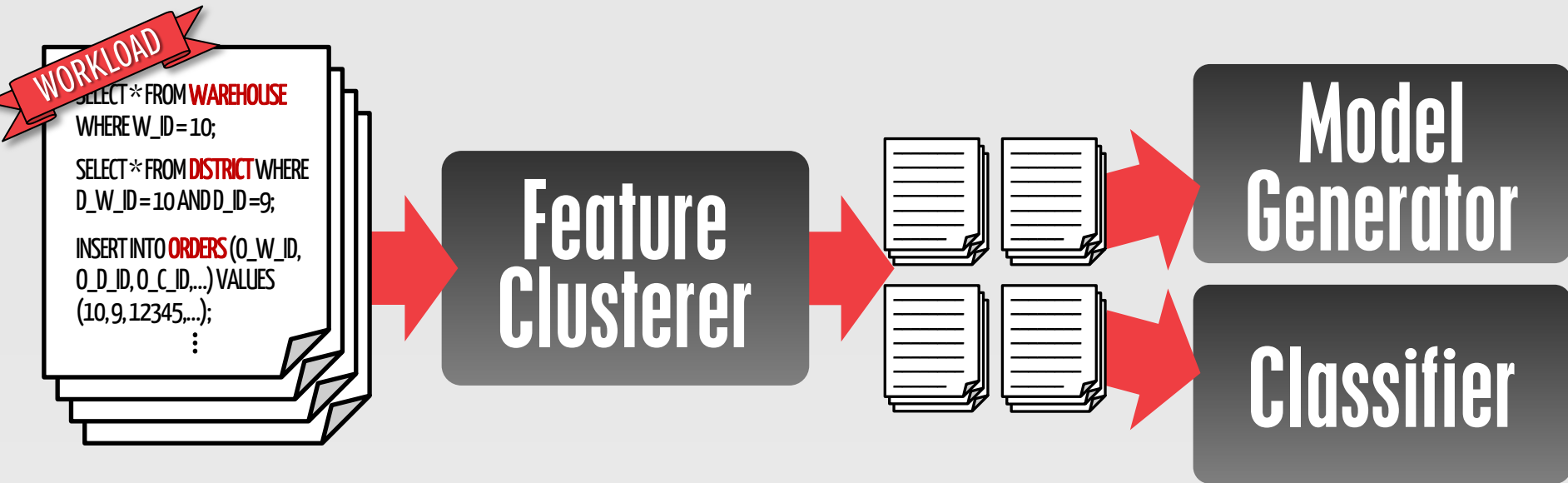
**DO USE MACHINE LEARNING
TO PREDICT
TRANSACTION BEHAVIOR.**



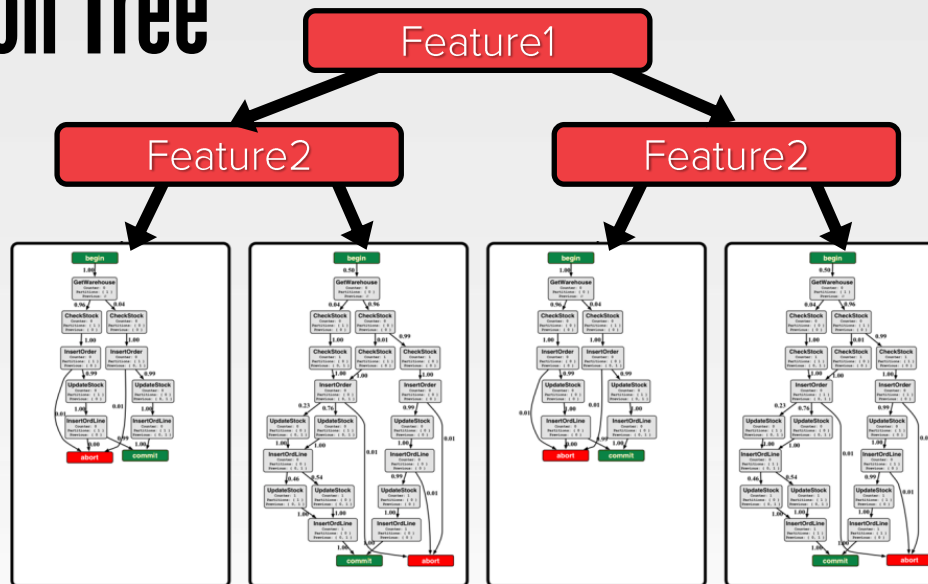
ON PREDICTIVE MODELING FOR OPTIMIZING
TRANSACTION EXECUTION IN PARALLEL OLTP SYSTEMS
Proc. VLDB Endow., Vol 5, Iss. 2, pp. 85-96, 2011

PREDICTIVE MODELS



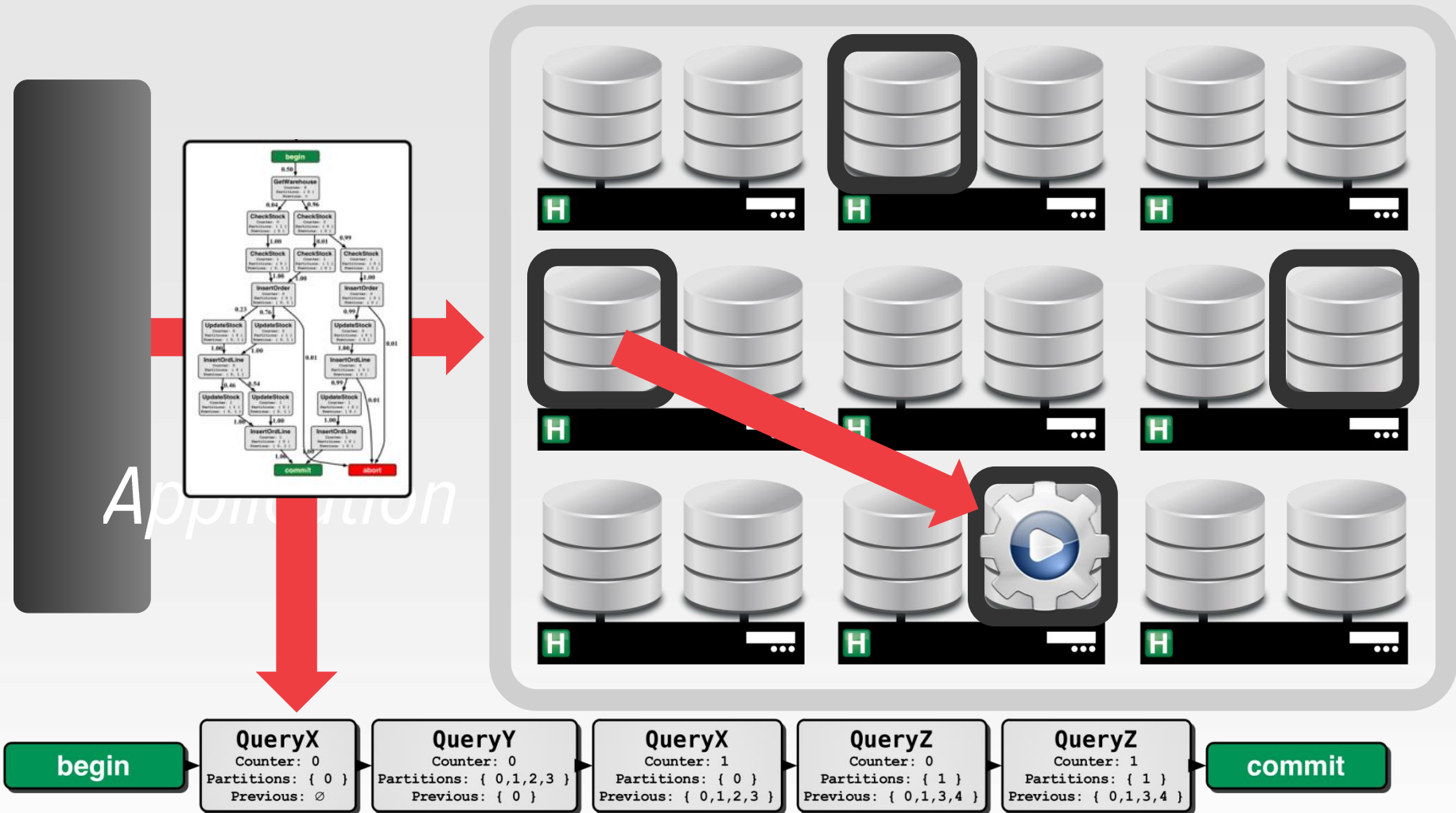


Decision Tree

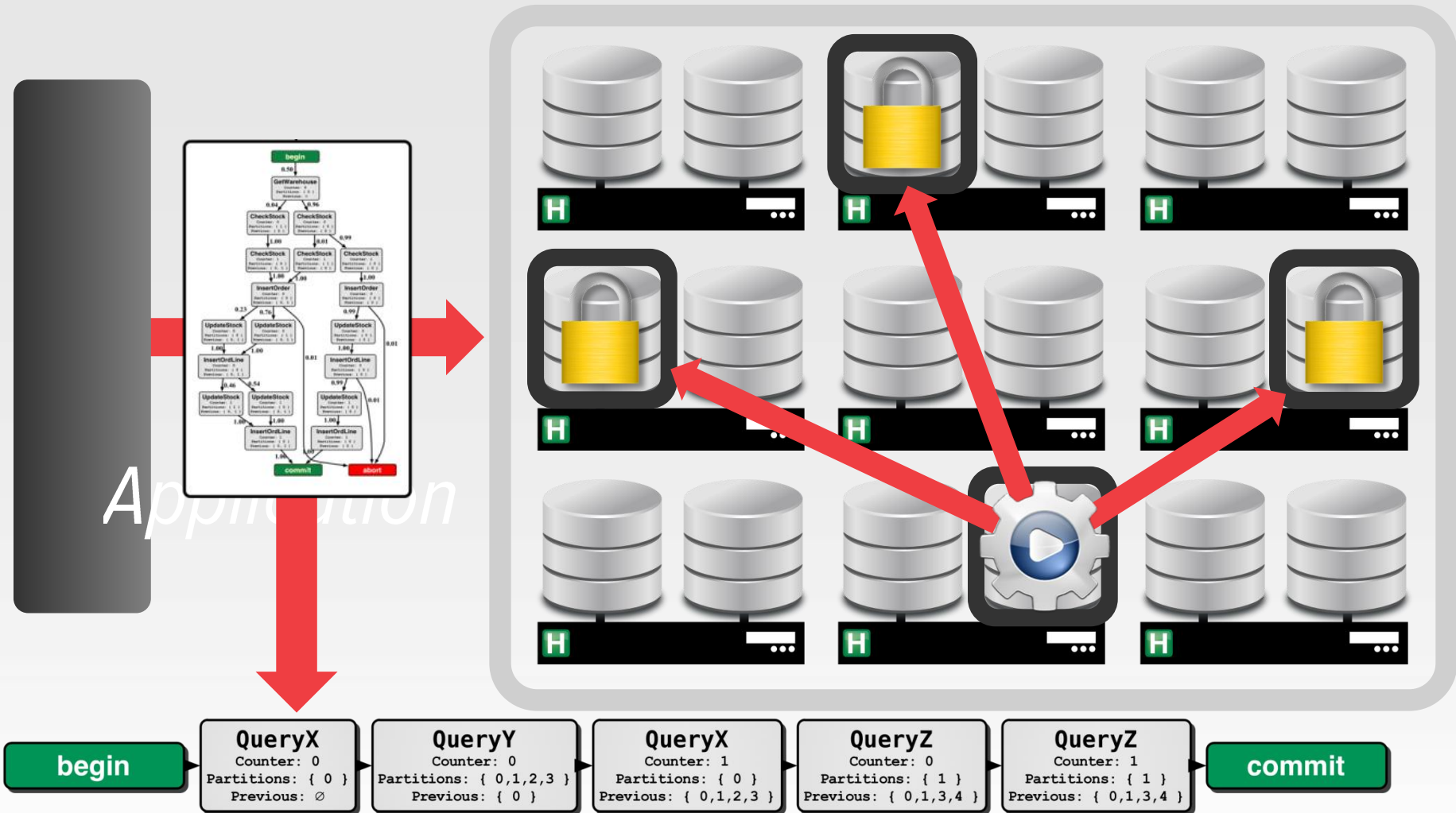


Markov Models

DISTRIBUTED TRANSACTIONS

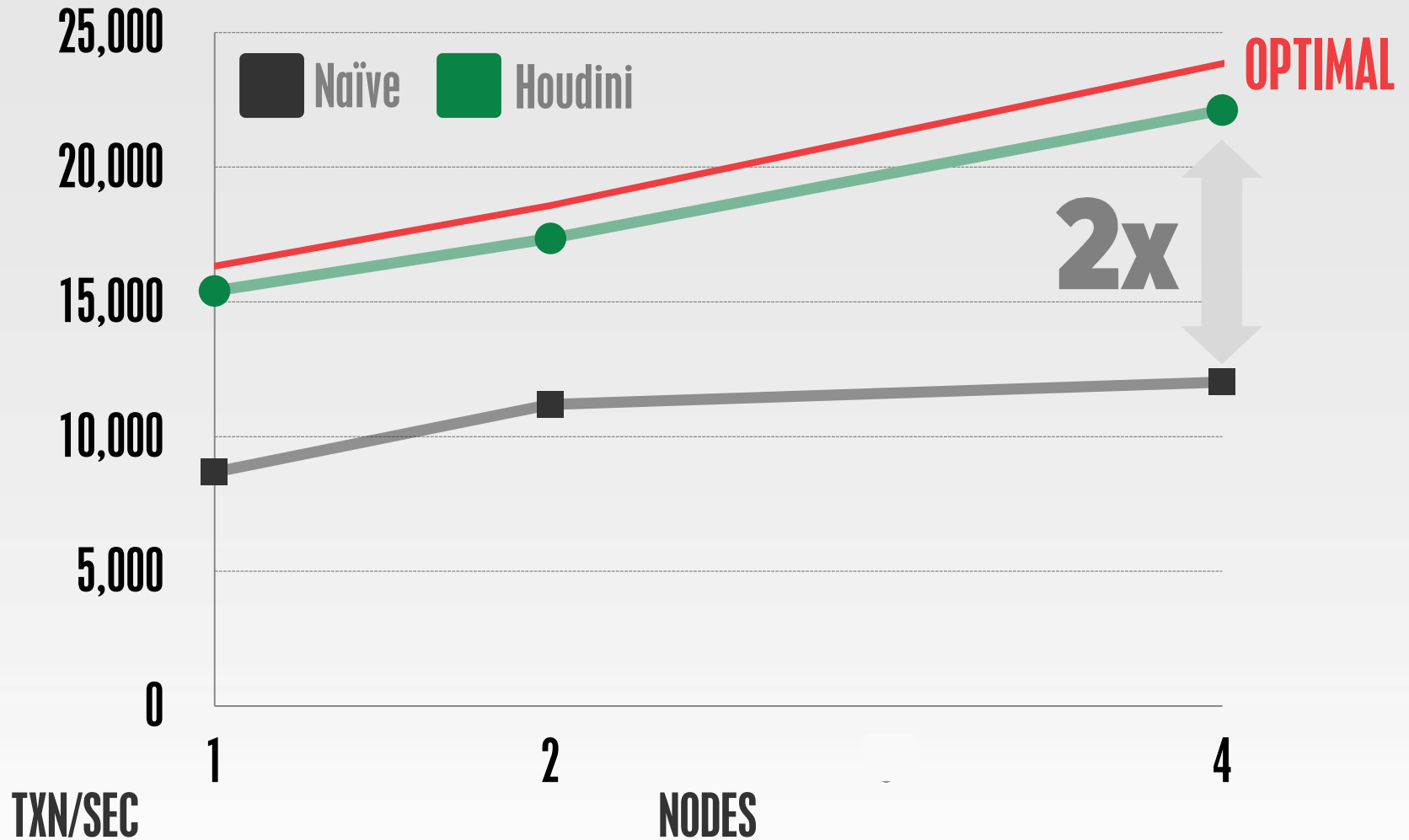


DISTRIBUTED TRANSACTIONS



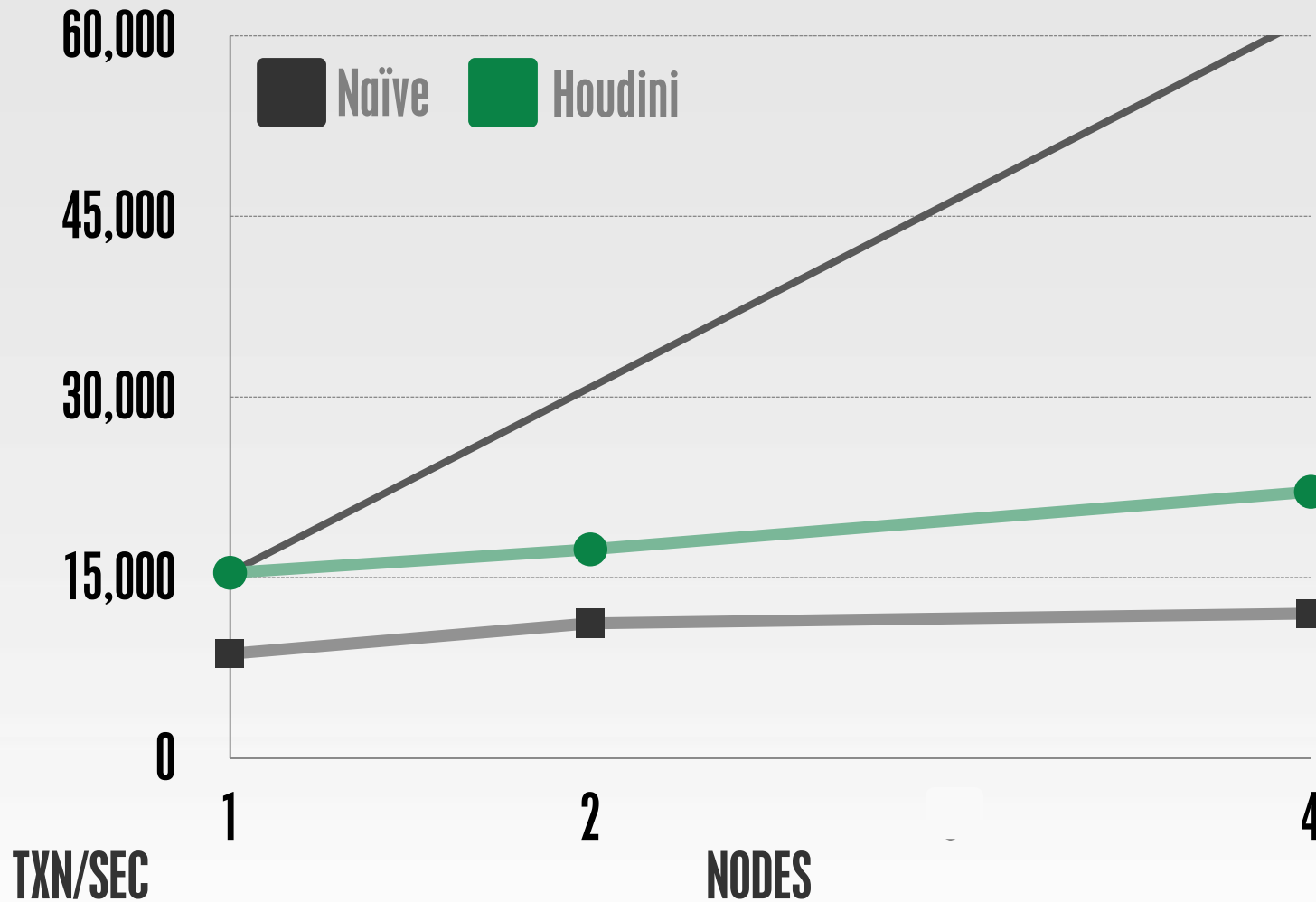
TPC-C BENCHMARK

8 Cores per Node
10% Distributed Transactions



TPC-C BENCHMARK

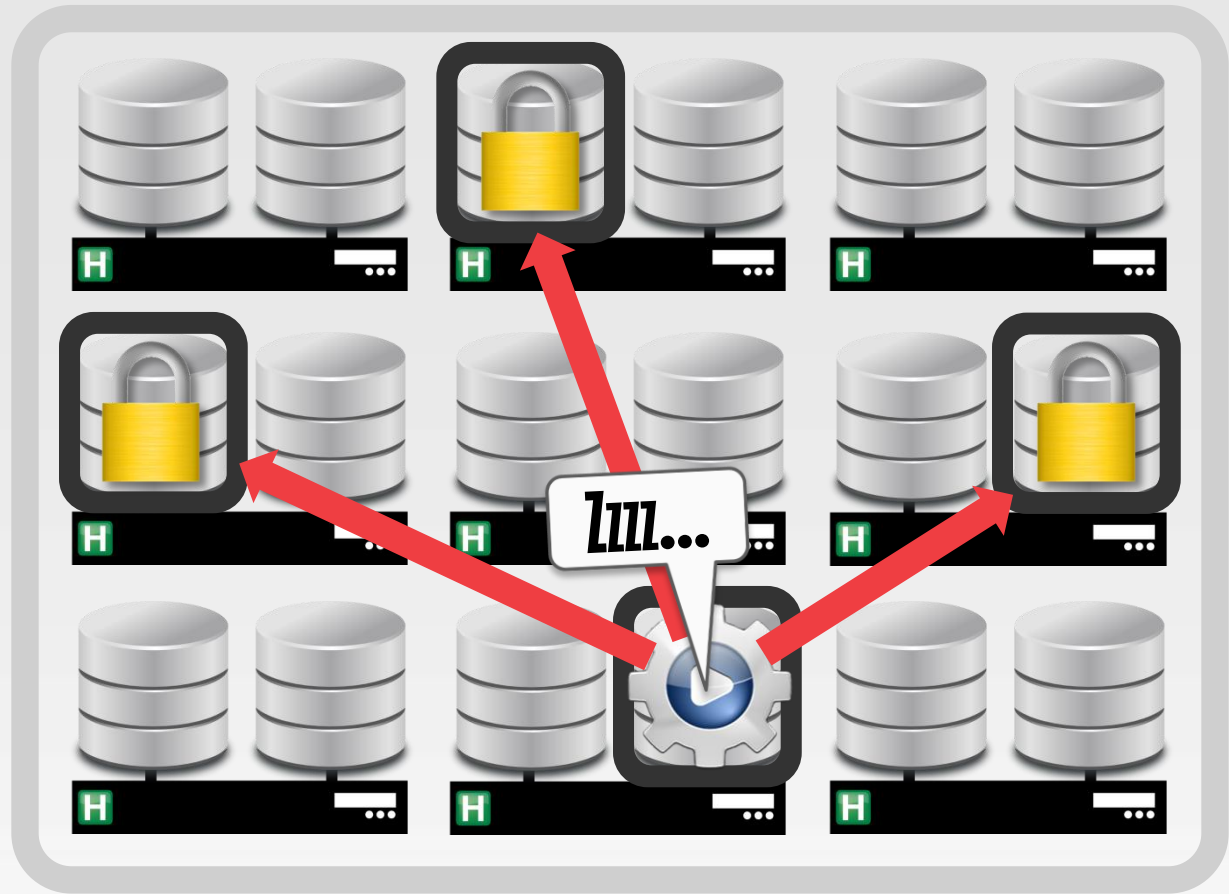
8 Cores per Node
10% Distributed Transactions



DISTRIBUTED TRANSACTIONS



Application



SP1 - Waiting for Query Result

DISTRIBUTED TRANSACTIONS



SP1 - Waiting for Query Result **SP2** - Waiting for Query Request

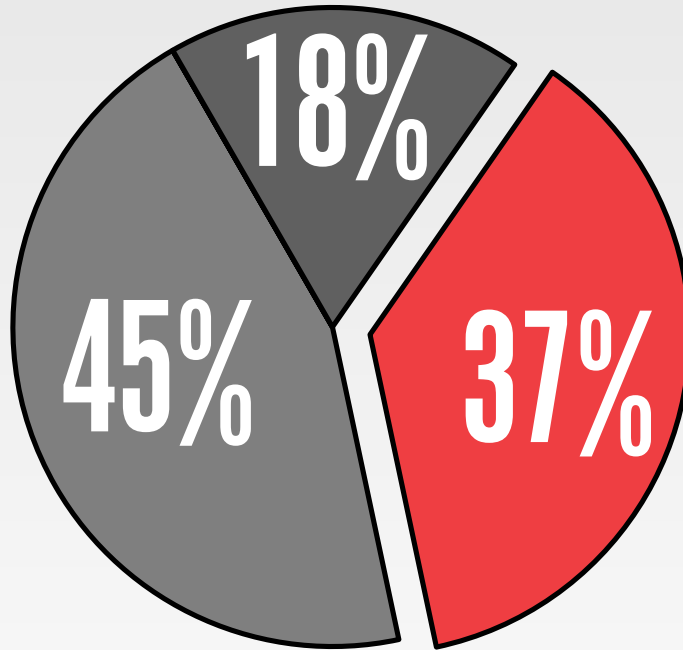
DISTRIBUTED TRANSACTIONS



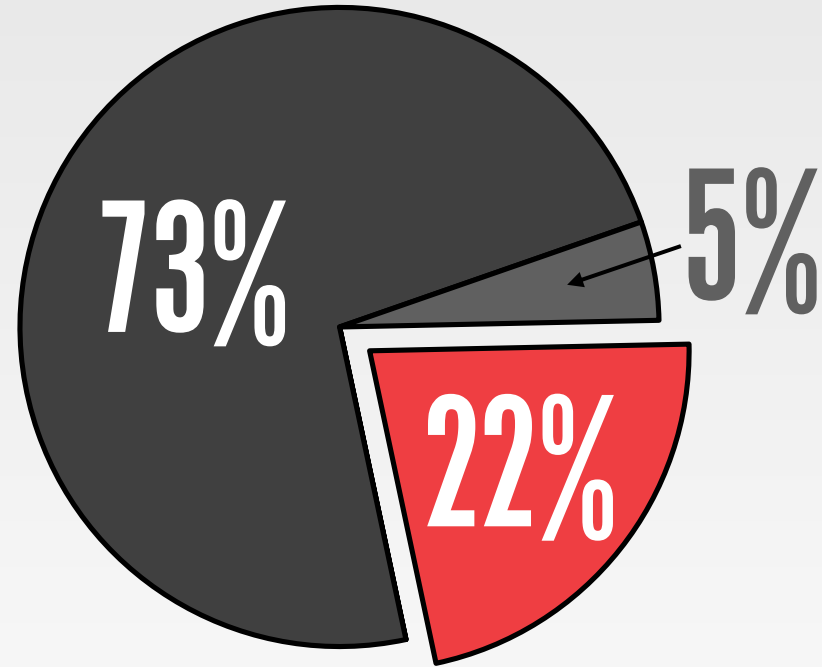
SP1 - Waiting for Query Result **SP2** - Waiting for Query Request
SP3 - Two-Phase Commit

TRANSACTION STALL POINTS

BASE PARTITION



REMOTE PARTITION



■ **SP1** - Waiting for Query Result
■ **SP3** - Two-Phase Commit

■ **SP2** - Waiting for Query Request
■ **Real Work**

**DO SOMETHING USEFUL
WHEN STALLED**

**DON'T BE SURPRISED
IF YOU & KB DON'T
LAST THROUGH
GRAD SCHOOL.**



**DON'T BE STAN'S
STUDENT
IF YOU GO TO
BROWN.**



**DO USE MACHINE LEARNING
TO SCHEDULE
SPECULATIVE TASKS.**



THE ART OF SPECULATIVE EXECUTION
In Progress (August 2013)

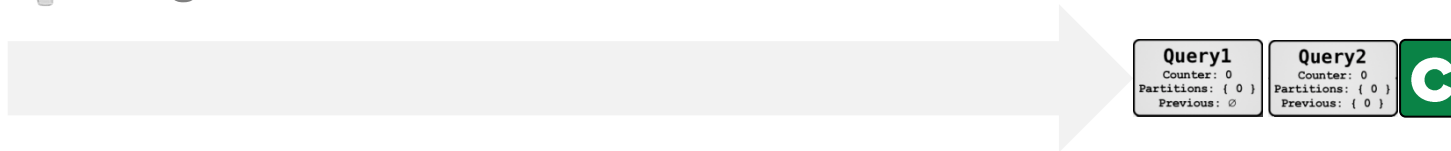
SERIALIZABLE SCHEDULE



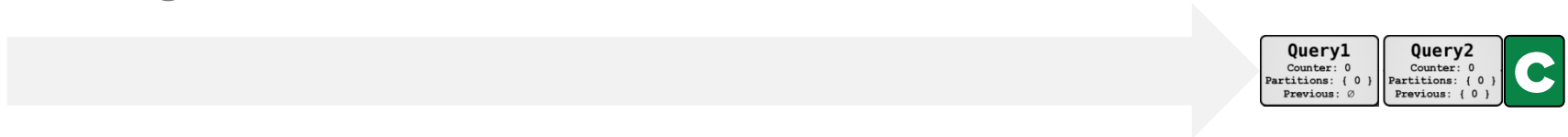
Distributed Transaction



Single-Partition Transaction



Single-Partition Transaction



SERIALIZABLE SCHEDULE



Distributed Transaction

QueryX
Counter: 0
Partitions: { 0 }
Previous: ∅

1m...

QueryZ
Counter: 1
Partitions: { 1 }
Previous: { 0,1,3,4 }

QueryZ
Counter: 0
Partitions: { 1 }
Previous: { 0,1,3,4 }



VERIFY



Speculative Transaction

Query1
Counter: 0
Partitions: { 0 }
Previous: ∅

Query2
Counter: 0
Partitions: { 0 }
Previous: { 0 }



Speculative Transaction

Query1
Counter: 0
Partitions: { 0 }
Previous: ∅

Query2
Counter: 0
Partitions: { 0 }
Previous: { 0 }



SPECULATIVE TRANSACTIONS

1111...



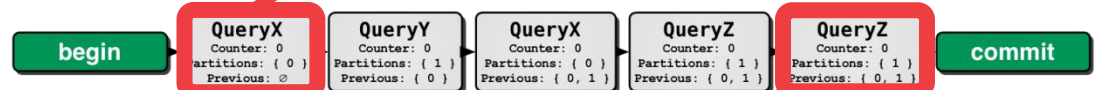
Transaction Queue

Speculation Candidate:



WRITE X

Distributed Transaction:

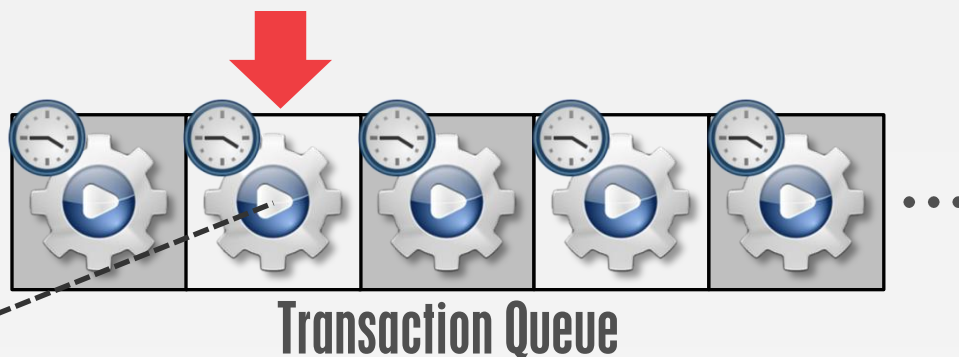


READ X

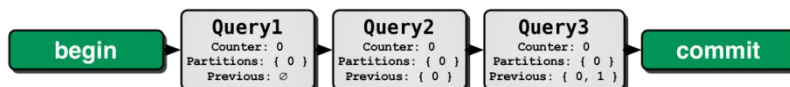
READ X

SPECULATIVE TRANSACTIONS

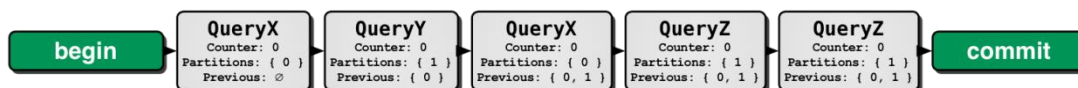
1111...



Speculation Candidate:



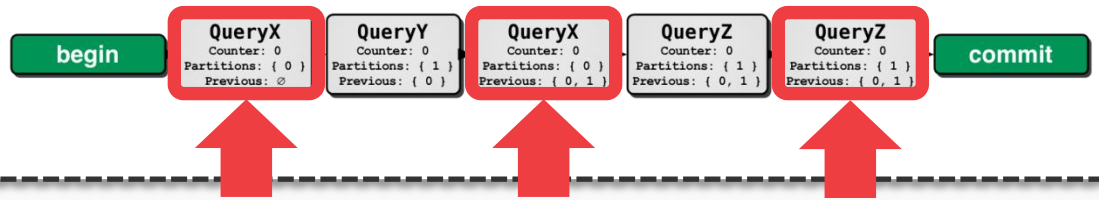
Distributed Transaction:



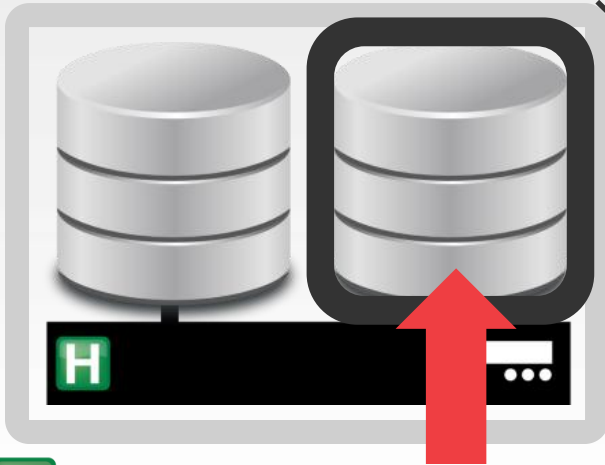
SPECULATIVE QUERIES



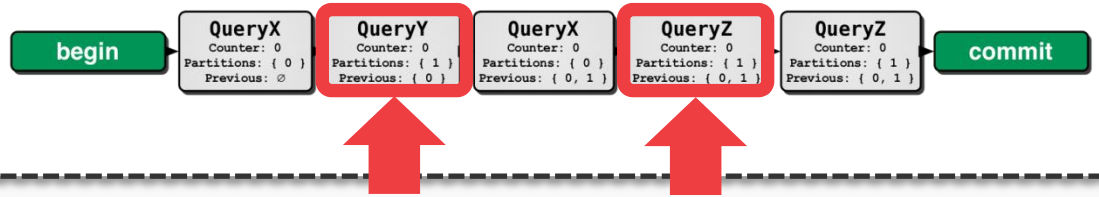
Distributed Transaction:



SPECULATIVE QUERIES



Distributed Transaction:



SPECULATIVE QUERIES

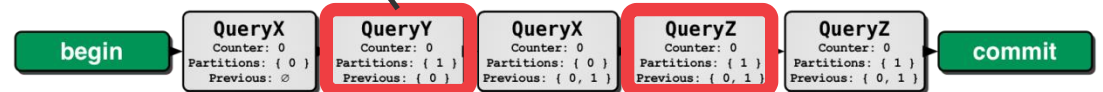


QueryY:

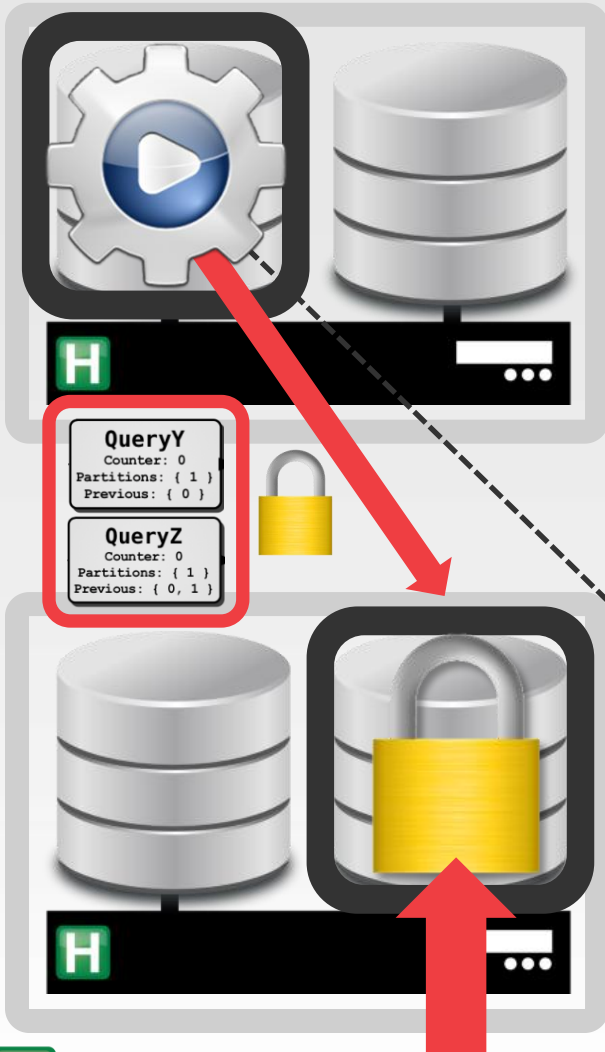
```
SELECT S_QTY FROM STOCK
WHERE S_W_ID = ?
AND S_I_ID = ?
```



Distributed Transaction:



SPECULATIVE QUERIES

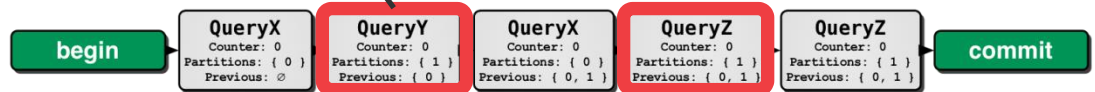


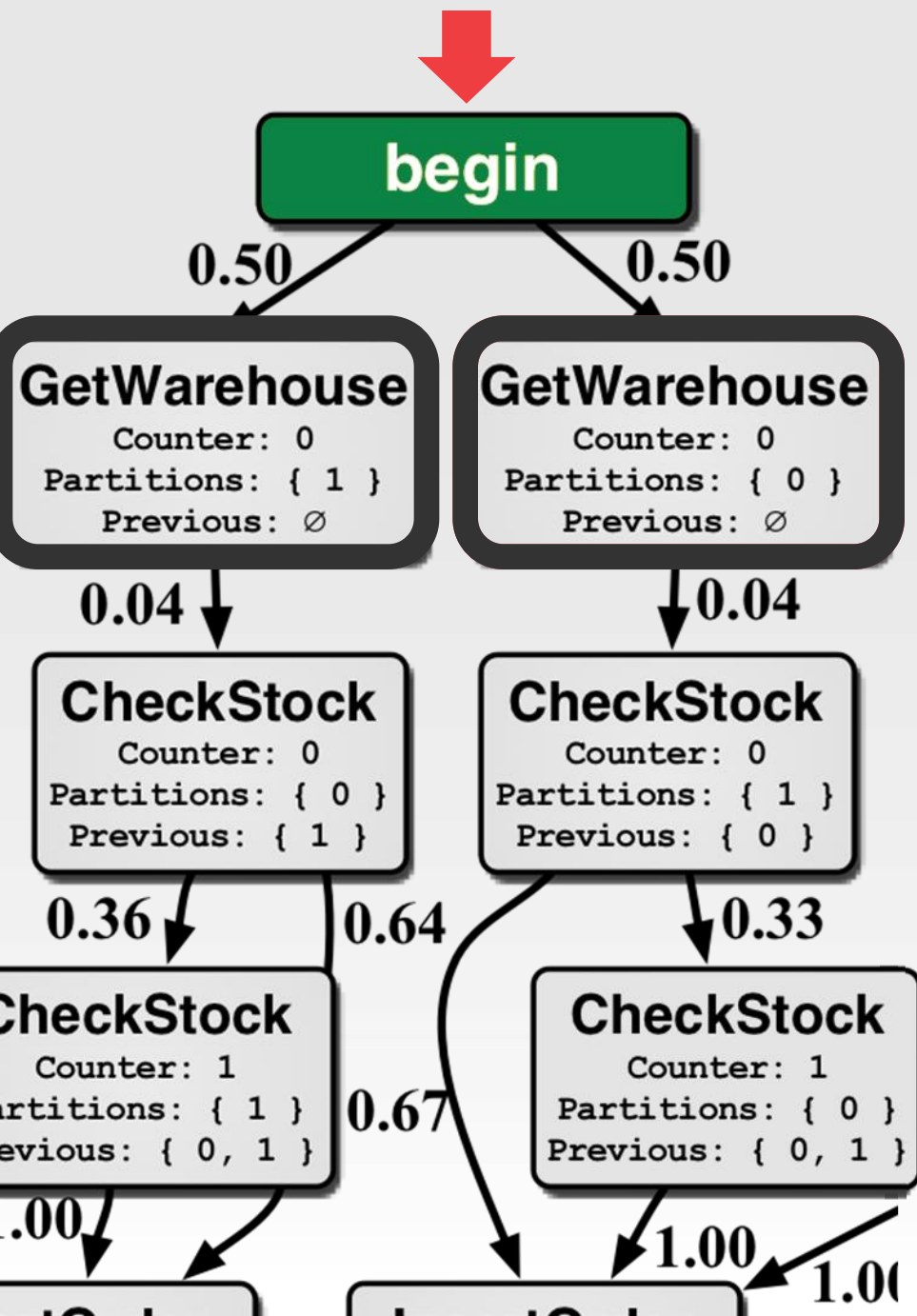
QueryY:

```
SELECT S_QTY FROM STOCK
WHERE S_W_ID = ?
AND S_I_ID = ?
```



Distributed Transaction:



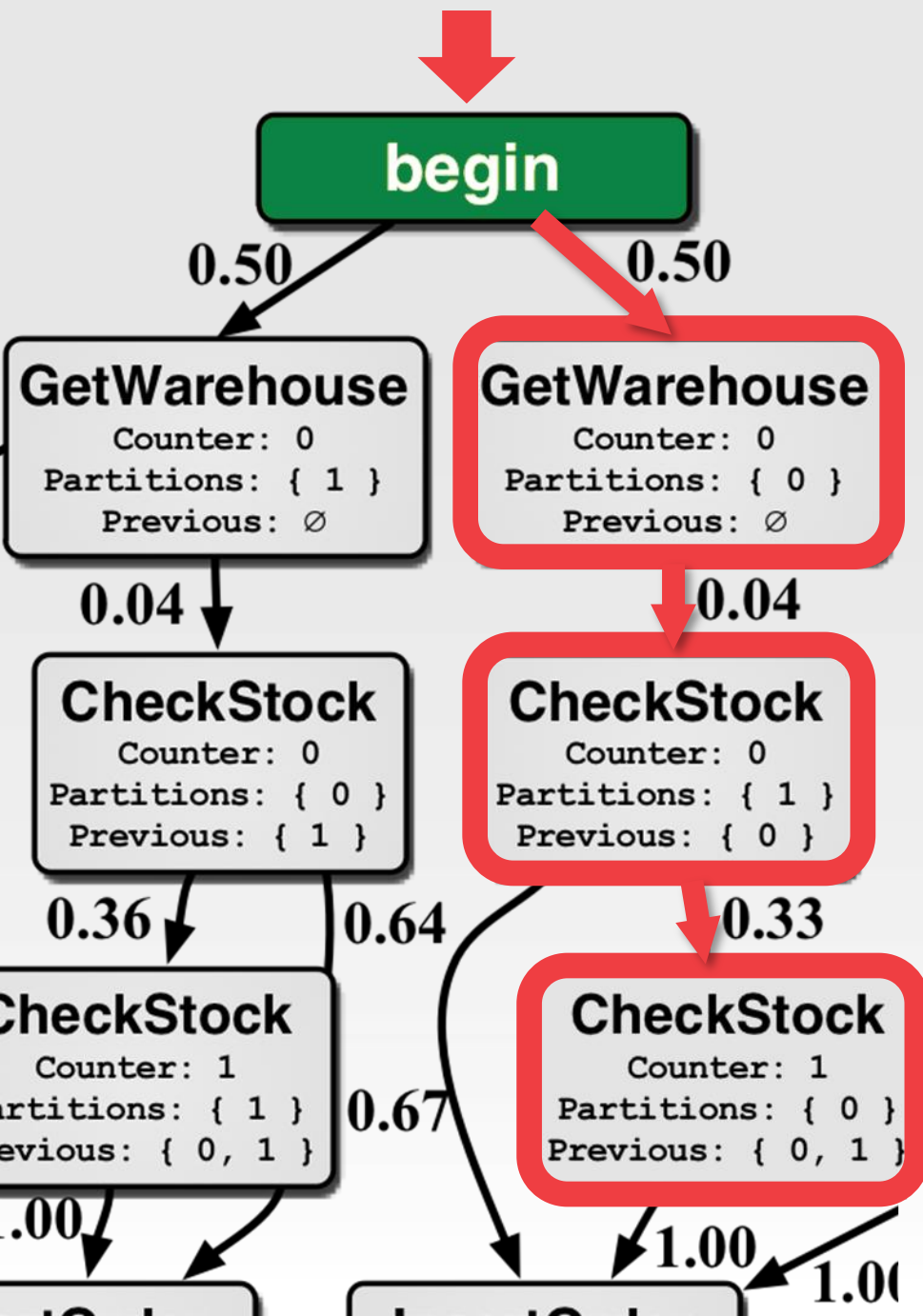


Transaction Parameters:

`w_id=0`
`i_w_ids=[1,0]`
`i_ids=[1001,1002]`

GetWarehouse:

`SELECT * FROM WAREHOUSE`
`WHERE W_ID = ?`



Transaction Parameters:

```

w_id=0
i_w_ids=[1 0]
i_ids=[1001,1002]
  
```

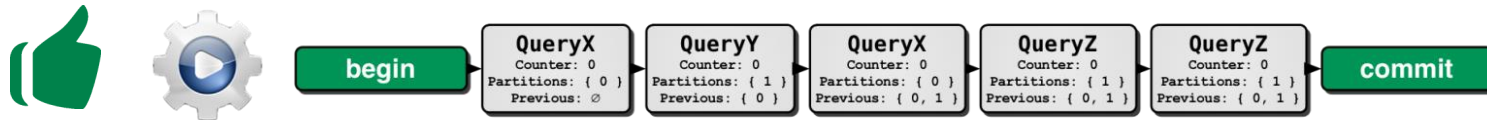
CheckStock:

```

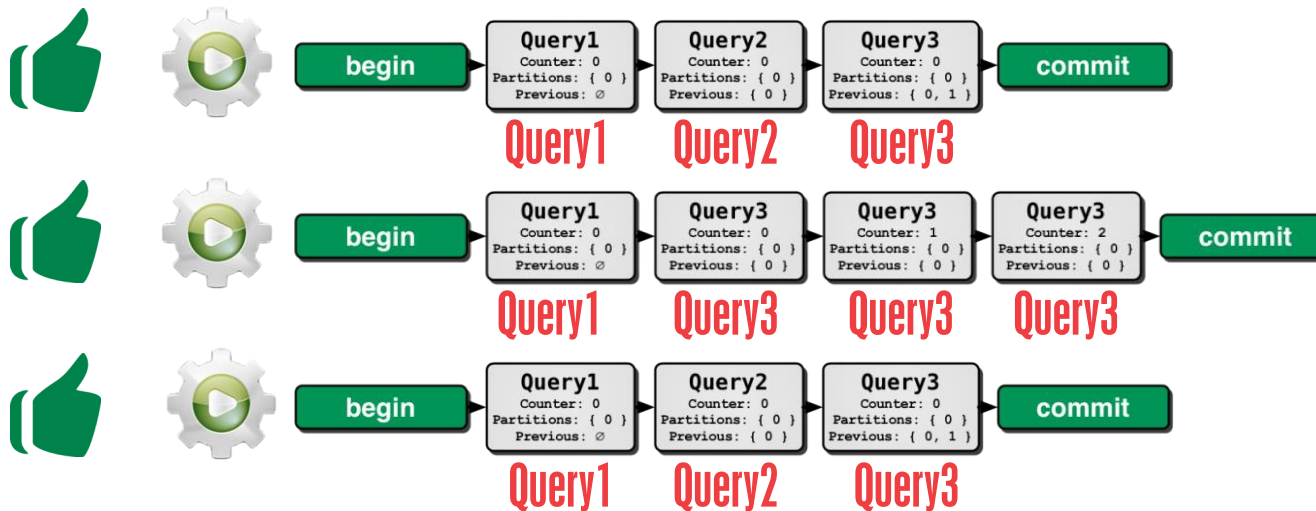
SELECT S_QTY FROM STOCK
WHERE S_W_ID = ?
AND S_I_ID = ?;
  
```

VERIFICATION

Distributed Transaction

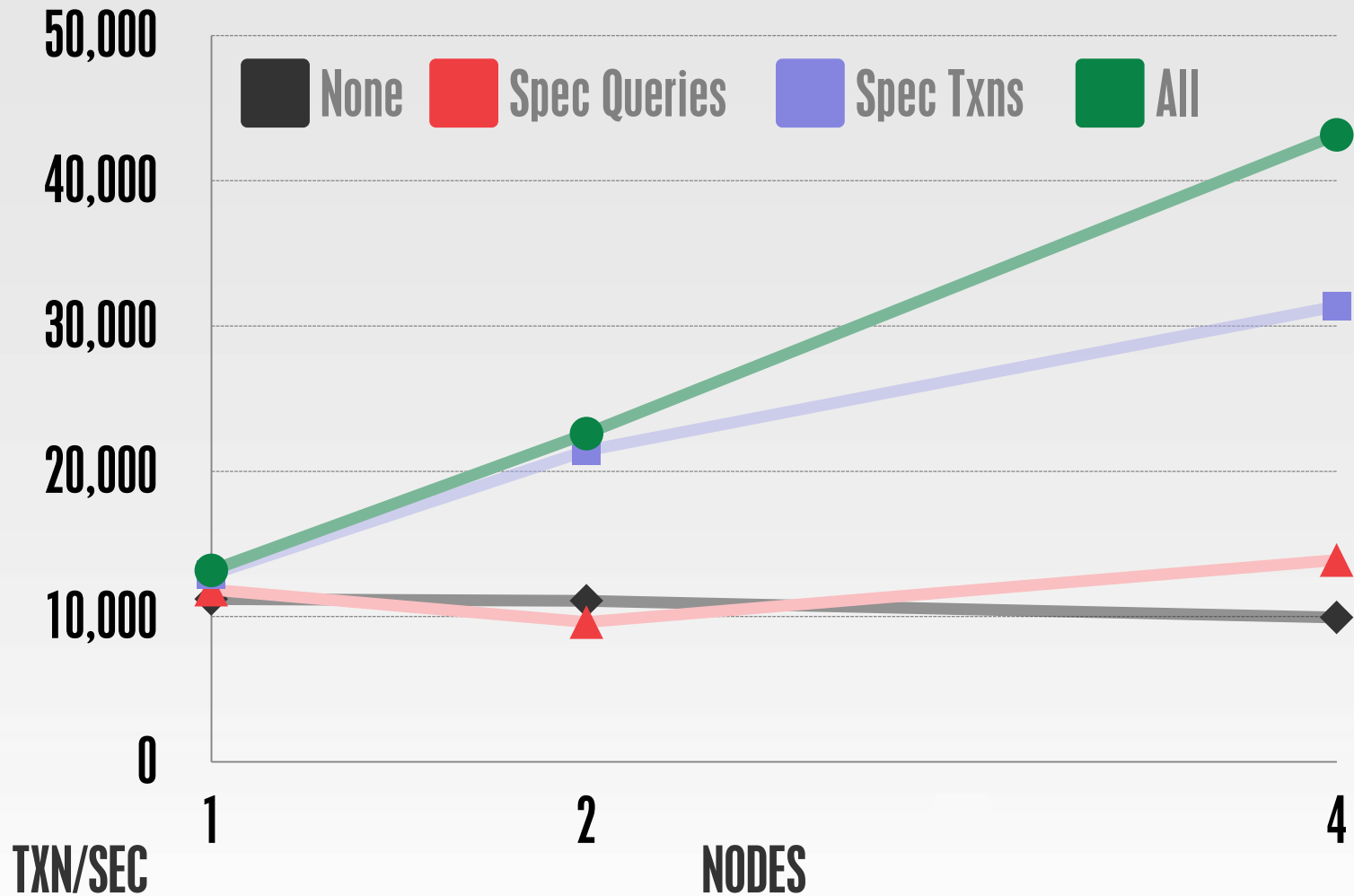


Speculative Transactions



TPC-C BENCHMARK

8 Cores per Node
10% Distributed Transactions





Optimize Single-Partition Execution

H-STORE: A HIGH-PERFORMANCE, DISTRIBUTED MAIN MEMORY TRANSACTION PROCESSING SYSTEM

Proc. VLDB Endow., vol. 1, iss. 2, pp. 1496-1499, 2008.



Minimize Distributed Transactions

SKEW-AWARE AUTOMATIC DATABASE PARTITIONING IN SHARED-NOTHING, PARALLEL OLTP SYSTEMS

Proceedings of SIGMOD, 2012.



Identify Distributed Transactions

ON PREDICTIVE MODELING FOR OPTIMIZING TRANSACTION EXECUTION IN PARALLEL OLTP SYSTEMS

Proc. VLDB Endow., vol. 5, pp. 85-96, 2011.



Utilize Transaction Stalls

THE ART OF SPECULATIVE EXECUTION

In Progress (August 2013)

FUTURE WORK

One Size
Almost
Fits AllTM



H-STORE



S-STORE



N-STORE



ISTC
BIG DATA



Escape From Planet Zdonik

(i.e., Andy Needs to Get Tenure)

Beyond the 'Stores

- Non-Partitionable Workloads.
- The Poor Man's Spanner.
- Scientific Databases.

DON'T MESS IT UP
WITH KB.





**Stan
Zdonik**



**"The Thrill"
Stonebraker**



**Sam
Madden**



**Ugur
Cetintemel**



**David
DeWitt**



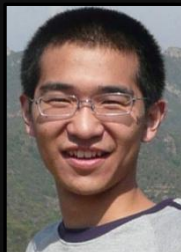
**Dan
Abadi**



**Evan
Jones**



**Saurya
Velagapudi**



**Xin
Jia**



**Carlo
Curino**



**Justin
DeBrabant**



**Yang
Zou**



**Visawee
Angkana.**



**Ning
Shi**



**John
Meehan**