


Carnegie Mellon Univ.
Dept. of Computer Science
15-415/615 - DB Applications


C. Faloutsos – A. Pavlo
 Lecture#23: Distributed Database Systems
 (R&G ch. 22)



Administrivia – Final Exam

- **Who:** You
- **What:** R&G Chapters 15-22
- **When:** Tuesday May 6th 5:30pm- 8:30pm
- **Where:** WEH 7500
- **Why:** Databases will help your love life.

Faloutsos/Pavlo CMU SCS 15-415/615 2



Today’s Class

- High-level overview of distributed DBMSs.
- Not meant to be a detailed examination of all aspects of these systems.

Faloutsos/Pavlo CMU SCS 15-415/615 3

CMU SCS

Today's Class

- Overview & Background
- Design Issues
- Distributed OLTP
- Distributed OLAP
- Real-world Examples

Faloutsos/Pavlo CMU SCS 15-415/615 4

CMU SCS

Why Do We Need Parallel/Distributed DBMSs?

- PayPal in 2008...
- Single, monolithic Oracle installation.
- Had to manually move data every xmas.
- Legal restrictions.

Faloutsos/Pavlo CMU SCS 15-415/615 5

CMU SCS

Why Do We Need Parallel/Distributed DBMSs?

- Increased Performance.
- Increased Availability.
- Potentially Lower TCO.

Faloutsos/Pavlo CMU SCS 15-415/615 6

CMU SCS

Parallel/Distributed DBMS

- Database is spread out across multiple resources to improve parallelism.
- Appears as a single database instance to the application.
 - SQL query for a single-node DBMS should generate same result on a parallel or distributed DBMS.

Faloutsos/Pavlo CMU SCS 15-415/615 7

CMU SCS

Parallel vs. Distributed

- **Parallel DBMSs:**
 - Nodes are physically close to each other.
 - Nodes connected with high-speed LAN.
 - Communication cost is assumed to be small.
- **Distributed DBMSs:**
 - Nodes can be far from each other.
 - Nodes connected using public network.
 - Communication cost and problems cannot be ignored.

Faloutsos/Pavlo CMU SCS 15-415/615 8

CMU SCS

Database Architectures

- The goal is parallelize operations across multiple resources.
 - CPU
 - Memory
 - Network
 - Disk

Faloutsos/Pavlo CMU SCS 15-415/615 9

CMU SCS

Database Architectures

Shared Memory

Shared Disk

Shared Nothing

Faloutsos/Pavlo CMU SCS 15-415/615 10

CMU SCS

Shared Memory

- CPUs and disks have access to common memory via a fast interconnect.
 - Very efficient to send messages between processors.
 - Sometimes called “shared everything”
- Examples: All single-node DBMSs.

Faloutsos/Pavlo CMU SCS 15-415/615 11

CMU SCS

Shared Disk

- All CPUs can access all disks directly via an interconnect but each have their own private memories.
 - Easy fault tolerance.
 - Easy consistency since there is a single copy of DB.
- Examples: Oracle Exadata, ScaleDB.

Faloutsos/Pavlo CMU SCS 15-415/615 12

CMU SCS

Shared Nothing

- Each DBMS instance has its own CPU, memory, and disk.
- Nodes only communicate with each other via network.
 - Easy to increase capacity.
 - Hard to ensure consistency.
- Examples: Vertica, Parallel DB2, MongoDB.

Faloutsos/Pavlo CMU SCS 15-415/615 13

CMU SCS

Early Systems

- **MUFFIN** – UC Berkeley (1979)
- **SDD-1** – CCA (1980)
- **System R*** – IBM Research (1984)
- **Gamma** – Univ. of Wisconsin (1986)
- **NonStop SQL** – Tandem (1987)

Stonebraker
Bernstein
Mohan
DeWitt
Gray

Faloutsos/Pavlo CMU SCS 15-415/615 14

CMU SCS

Inter- vs. Intra-query Parallelism

- **Inter-Query:** Different queries or txns are executed concurrently.
 - Increases throughput but not latency.
 - Already discussed for shared-memory DBMSs.
- **Intra-Query:** Execute the operations of a single query in parallel.
 - Increases latency for long-running queries.

Faloutsos/Pavlo CMU SCS 15-415/615 15

CMU SCS

Parallel/Distributed DBMSs

- Advantages:
 - Data sharing.
 - Reliability and availability.
 - Speed up of query processing.
- Disadvantages:
 - May increase processing overhead.
 - More database design issues.
 - Harder to ensure ACID guarantees.

Faloutsos/Pavlo CMU SCS 15-415/615 16

CMU SCS

Today's Class

- Overview & Background
- ➔ • Design Issues
- Distributed OLTP
- Distributed OLAP
- Real-world Examples

Faloutsos/Pavlo CMU SCS 15-415/615 17

CMU SCS

Design Issues

- How do we store data across nodes?
- How does the application find data?
- How to execute queries on distributed data?
 - Push query to data.
 - Pull data to query.
- How does the DBMS ensure correctness?

Faloutsos/Pavlo CMU SCS 15-415/615 18

CMU SCS

Database Partitioning

- Split database across multiple resources:
 - Disks, nodes, processors.
 - Sometimes called “sharding”
- The DBMS executes query fragments on each partition and then combines the results to produce a single answer.

Faloutsos/Pavlo CMU SCS 15-415/615 19

CMU SCS

Naïve Table Partitioning

- Each node stores one and only table.
- Assumes that each node has enough storage space for a table.

Faloutsos/Pavlo CMU SCS 15-415/615 20

CMU SCS


Naïve Table Partitioning

Table1 Table2

Tuple1									
Tuple2									
Tuple3									
Tuple4									
Tuple5									

→

Partitions



Ideal Query:

`SELECT * FROM table`

Faloutsos/Pavlo CMU SCS 15-415/615 21

CMU SCS

Horizontal Partitioning

- Split a table's tuples into disjoint subsets.
 - Choose column(s) that divides the database equally in terms of size, load, or usage.
 - Each tuple contains all of its columns.
- Three main approaches:
 - Round-robin Partitioning.
 - Hash Partitioning.
 - Range Partitioning.

Faloutsos/Pavlo CMU SCS 15-415/615 22

CMU SCS

Horizontal Partitioning

Table

Tuple1				
Tuple2				
Tuple3				
Tuple4				
Tuple5				

Partitioning Key

➔

Partitions

Ideal Query:
SELECT * FROM table
WHERE partitionKey = ?

Faloutsos/Pavlo CMU SCS 15-415/615 23

CMU SCS

Vertical Partitioning

- Split the columns of tuples into fragments:
 - Each fragment contains all of the tuples' values for column(s).
- Need to include primary key or unique record id with each partition to ensure that the original tuple can be reconstructed.

Faloutsos/Pavlo CMU SCS 15-415/615 24

CMU SCS

Vertical Partitioning

Table

Tuple1				
Tuple2				
Tuple3				
Tuple4				
Tuple5				

➔

Partitions

Ideal Query:
`SELECT column FROM table`

Faloutsos/Pavlo CMU SCS 15-415/615 25

CMU SCS

Replication

- **Partition Replication:** Store a copy of an entire partition in multiple locations.
 - Master – Slave Replication
- **Table Replication:** Store an entire copy of a table in each partition.
 - Usually small, read-only tables.
- The DBMS ensures that updates are propagated to all replicas in either case.

Faloutsos/Pavlo CMU SCS 15-415/615 26

CMU SCS

Replication

Partition Replication

Table Replication

Faloutsos/Pavlo CMU SCS 15-415/615 27

CMU SCS

Data Transparency

- Users should not be required to know where data is physically located, how tables are partitioned or replicated.
- A SQL query that works on a single-node DBMS should work the same on a distributed DBMS.

Faloutsos/Pavlo CMU SCS 15-415/615 28

CMU SCS

OLTP vs. OLAP

- On-line Transaction Processing:
 - Short-lived txns.
 - Small footprint.
 - Repetitive operations.
- On-line Analytical Processing:
 - Long running queries.
 - Complex joins.
 - Exploratory queries.

Faloutsos/Pavlo CMU SCS 15-415/615 29

CMU SCS

Workload Characterization

Michael Stonebraker – "Ten Rules For Scalable Performance In Simple Operation" Databases"
<http://cmu.edu/~stonebr/10rules/10rules2011.06.108651>

CMU SCS

Today's Class

- Overview & Background
- Design Issues
- ➔ • Distributed OLTP
- Distributed OLAP
- Real-world Examples

Faloutsos/Pavlo CMU SCS 15-415/615 31

CMU SCS

Distributed OLTP

- Execute txns on a distributed DBMS.
- Used for user-facing applications:
 - Example: Credit card processing.
- Key Challenges:
 - Consistency
 - Availability

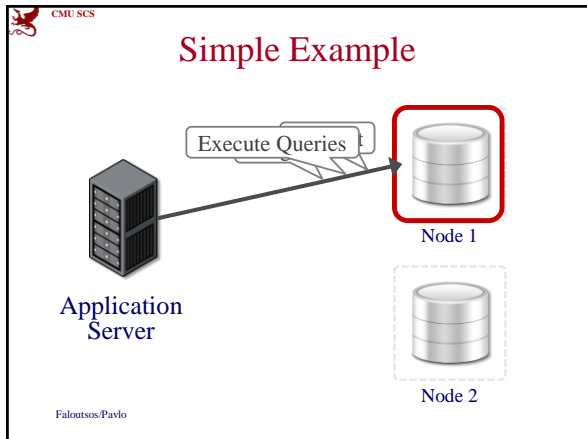
Faloutsos/Pavlo CMU SCS 15-415/615 32

CMU SCS

Single-Node vs. Distributed Transactions

- Single-node txns do not require the DBMS to coordinate behavior between nodes.
- Distributed txns are any txn that involves more than one node.
 - Requires expensive coordination.

Faloutsos/Pavlo CMU SCS 15-415/615 33



Transaction Coordination

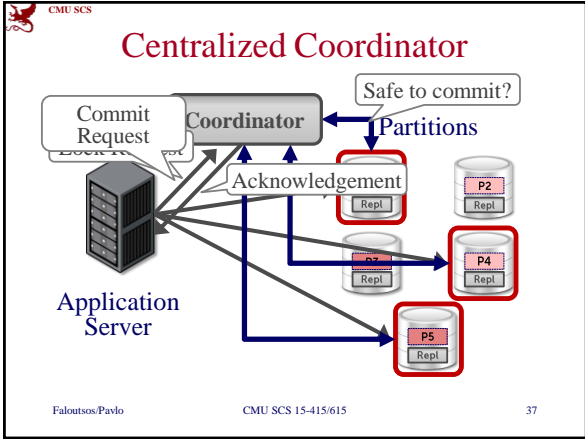
- Assuming that our DBMS supports multi-operation txns, we need some way to coordinate their execution in the system.
- Two different approaches:
 - **Centralized:** Global “traffic cop”.
 - **Decentralized:** Nodes organize themselves.

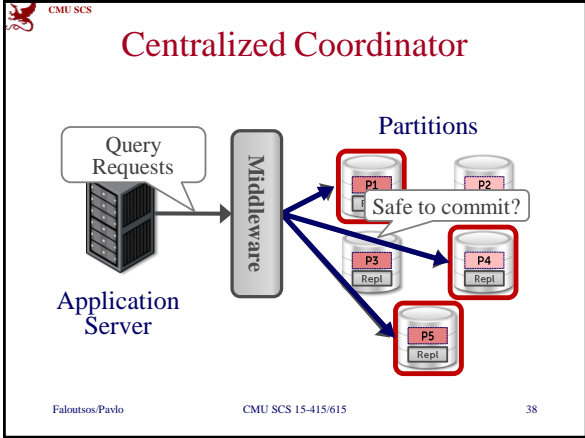
Faloutsos/Pavlo CMU SCS 15-415/615 35

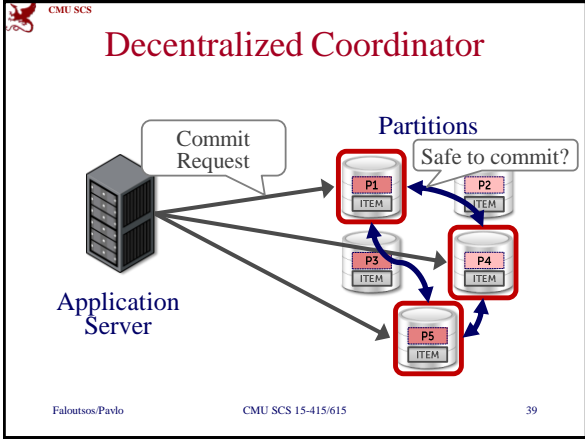
TP Monitors

- Example of a centralized coordinator.
- Originally developed in the 1970-80s to provide txns between terminals + mainframe databases.
 - Examples: ATMs, Airline Reservations.
- Many DBMSs now support the same functionality internally.

Faloutsos/Pavlo CMU SCS 15-415/615 36







CMU SCS

Observation

- **Q:** How do we ensure that all nodes agree to commit a txn?
 - What happens if a node fails?
 - What happens if our messages show up late?


Faloutsos/Pavlo CMU SCS 15-415/615 40

CMU SCS

CAP Theorem

- Proposed by Eric Brewer that it is impossible for a distributed system to always be:
 - Consistent
 - Always Available
 - Network Partition Tolerant
- Proved in 2002.

Pick Two!



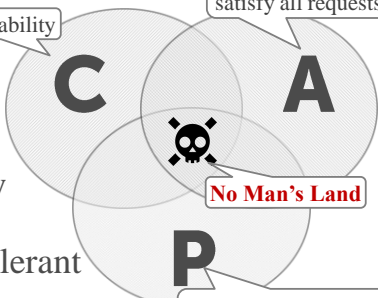
Brewer

Faloutsos/Pavlo CMU SCS 15-415/615 41

CMU SCS

CAP Theorem

Linearizability



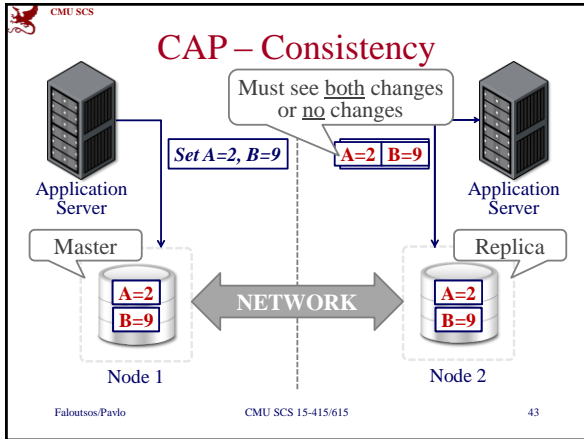
All up nodes can satisfy all requests.

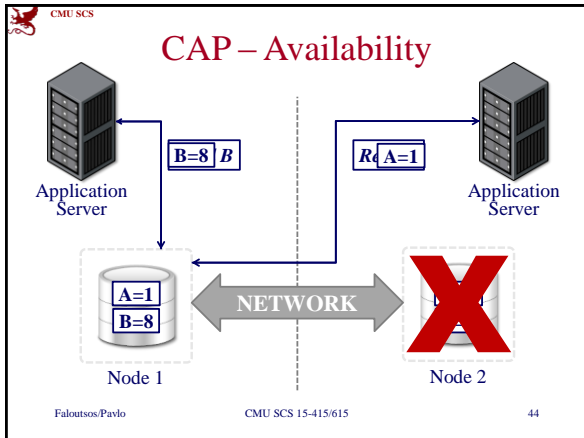
Consistency
Availability
Partition Tolerant

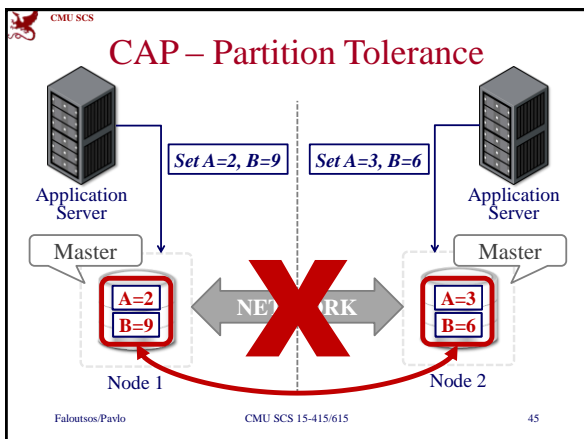
No Man's Land

Still operate correctly despite message loss.

Faloutsos/Pavlo CMU SCS 15-415/615







CMU SCS

CAP Theorem

These are essentially the same!

- **Relational DBMSs: CA/CP**
 - Examples: IBM DB2, MySQL Cluster, VoltDB
- **NoSQL DBMSs: AP**
 - Examples: Cassandra, Riak, DynamoDB

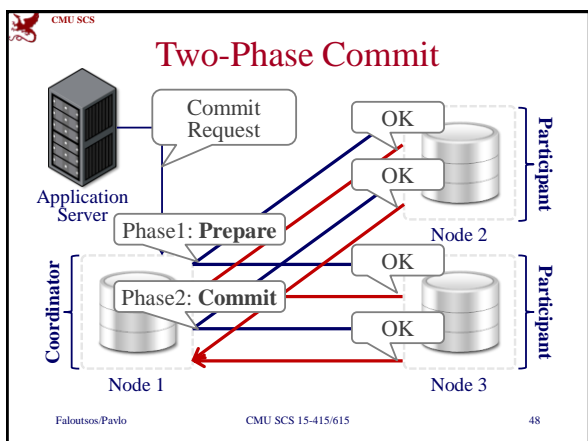
Faloutsos/Pavlo CMU SCS 15-415/615 46

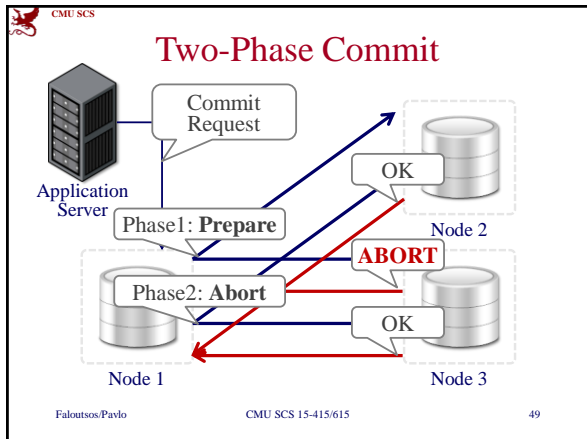
CMU SCS

Atomic Commit Protocol

- When a multi-node txn finishes, the DBMS needs to ask all of the nodes involved whether it is safe to commit.
 - All nodes must agree on the outcome
- Examples:
 - Two-Phase Commit
 - Three-Phase Commit
 - Paxos

Faloutsos/Pavlo CMU SCS 15-415/615 47





- ### Two-Phase Commit
- Each node has to record the outcome of each phase in a stable storage log.
 - **Q:** What happens if coordinator crashes?
 - Participants have to decide what to do.
 - **Q:** What happens if participant crashes?
 - Coordinator assumes that it responded with an abort if it hasn't sent an acknowledgement yet.
 - The nodes have to block until they can figure out the correct action to take.

- ### Three-Phase Commit
- The coordinator fails. Failure doesn't always mean a hard crash. That it intends to commit.
 - If the coordinator fails, then the participants elect a new coordinator and finish commit.
 - Nodes do not have to block if there are no network partitions.

CMU SCS

Paxos

- Consensus protocol where a coordinator proposes an outcome (e.g., commit or abort) and then the participants vote on whether that outcome should succeed.
- Does not block if a majority of participants are available and has provably minimal message delays in the best case.
 - First correct protocol that was provably resilient in the face asynchronous networks

Faloutsos/Pavlo CMU SCS 15-415/615 52

CMU SCS

2PC vs. Paxos

- **Two-Phase Commit:** blocks if coordinator fails after the prepare message is sent, until coordinator recovers.
- **Paxos:** non-blocking as long as a majority participants are alive, provided there is a sufficiently long period without further failures.

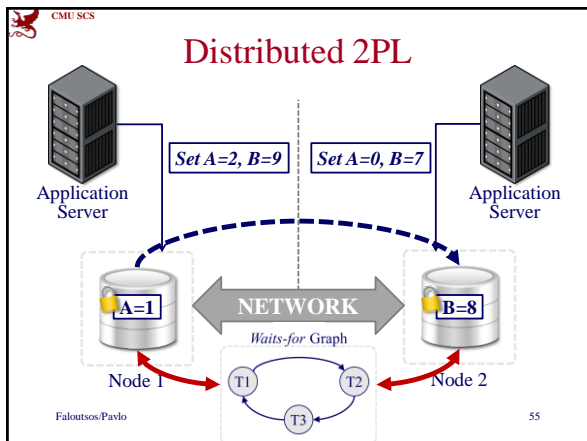
Faloutsos/Pavlo CMU SCS 15-415/615 53

CMU SCS

Distributed Concurrency Control

- Need to allow multiple txns to execute simultaneously across multiple nodes.
 - Many of the same protocols from single-node DBMSs can be adapted.
- This is harder because of:
 - Replication.
 - Network Communication Overhead.
 - Node Failures.

Faloutsos/Pavlo CMU SCS 15-415/615 54



Recovery

- **Q:** What do we do if a node crashes in CA/CP DBMS?
- If node is replicated, use Paxos to elect a new primary.
 - If node is last replica, halt the DBMS.
- Node can recover from checkpoints + logs and then catch up with primary.

Today's Class

- Overview & Background
- Design Issues
- Distributed OLTP
- ➔ • Distributed OLAP
- Real-world Examples

CMU SCS

Distributed OLAP

- Execute analytical queries that examine large portions of the database.
- Used for back-end data warehouses:
 - Example: Data mining
- Key Challenges:
 - Data movement.
 - Query planning.

Faloutsos/Pavlo CMU SCS 15-415/615 58

CMU SCS

Distributed OLAP

Application Server

Single Complex Query

Partitions

P1 ITEM

P2 ITEM

P3 ITEM

P4 ITEM

P5 ITEM

Faloutsos/Pavlo CMU SCS 15-415/615 59

CMU SCS

Distributed Joins Are Hard

```
SELECT * FROM table1, table2
WHERE table1.val = table2.val
```

- Assume tables are horizontally partitioned:
 - Table1 Partition Key → table1.key
 - Table2 Partition Key → table2.key
- **Q:** How to execute?
- Naïve solution is to send all partitions to a single node and compute join.

Faloutsos/Pavlo CMU SCS 15-415/615 60

CMU SCS

Semi-Joins

- Main Idea: First distribute the join attributes between nodes and then recreate the full tuples in the final output.
 - Send just enough data from each table to compute which rows to include in output.
- Lots of choices make this problem hard:
 - What to materialize?
 - Which table to send?

Faloutsos/Pavlo CMU SCS 15-415/615 61

CMU SCS

Today's Class

- Overview & Background
- Design Issues
- Distributed OLTP
- Distributed OLAP
- ➔ • Real-world Examples

Faloutsos/Pavlo CMU SCS 15-415/615 62

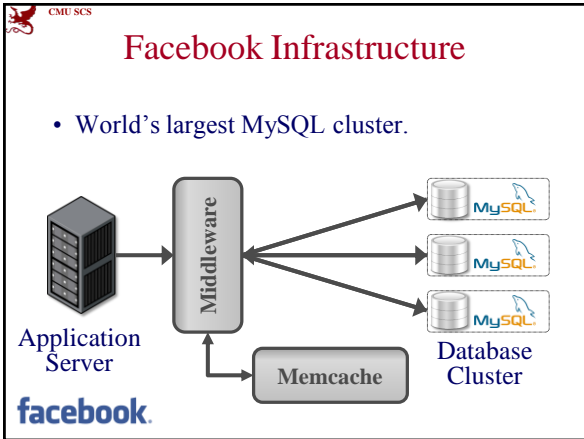
CMU SCS

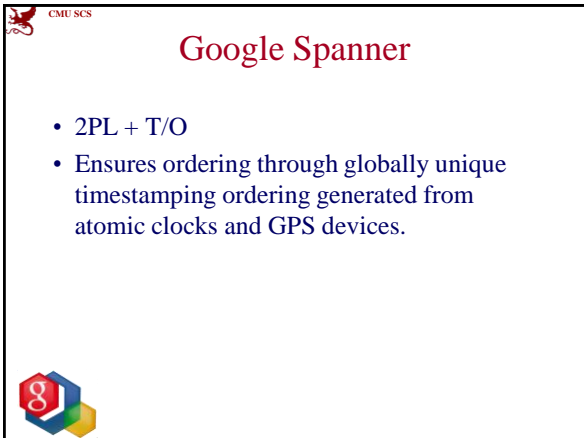
NuoDB

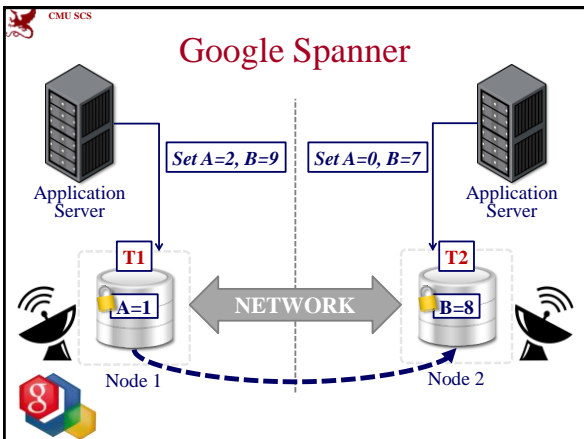
- Distributed MVCC+OCC.
- Split the database into “atoms” (i.e., blocks)
 - Nodes assigned as executor and storage nodes.
 - Move atoms to the node where a txn is executing, push writes at commit

The diagram shows a central 'Executor Node' (represented by a gear icon) receiving a 'Query Request' from the left. From the Executor Node, two arrows labeled 'Atom Request' point to two separate 'Storage Nodes' (represented by disk icons) on the right. The NuoDB logo is in the bottom left corner.

Faloutsos/Pavlo CMU SCS 15-415/615 63







CMU SCS

Summary

- Everything is harder in a distributed setting:
 - Concurrency Control
 - Query Execution
 - Recovery

Faloutsos/Pavlo CMU SCS 15-415/615 67

CMU SCS

Next Class

- Discuss distributed OLAP more.
 - You'll learn why MapReduce was a bad idea.
- Compare NoSQL vs. NewSQL
- Learn the answer to the #1 student question:
 - What DBMS should I use for my start-up?

Faloutsos/Pavlo CMU SCS 15-415/615 68
