

Carnegie Mellon University

Database Systems

Distributed Databases



15-445/645 SPRING 2025

» PROF. JIGNESH PATEL

ADMINISTRIVIA

Project #4 is due Sunday April 20th @ 11:59pm

→ Recitation: Friday, April 11th in GHC 4303 from 3:00 - 4:00 PM

HW6 is due Sunday, April 20, 2025 @ 11:59pm

Final Exam is on Monday, April 28, 2025, from
05:30pm - 08:30pm

→ Early exam will not be offered. Do not make travel plans.

This course is recruiting TAs for the next semester

→ Apply at: <https://www.ugrad.cs.cmu.edu/ta/F25/>

UPCOMING DATABASE TALKS

MariaDB (DB Seminar)

- Monday, April 14 @ 4:30pm
- MariaDB's New Query Optimizer
- Speaker: Michael Widenius
- <https://cmu.zoom.us/j/93441451665>



Gel (DB Seminar)

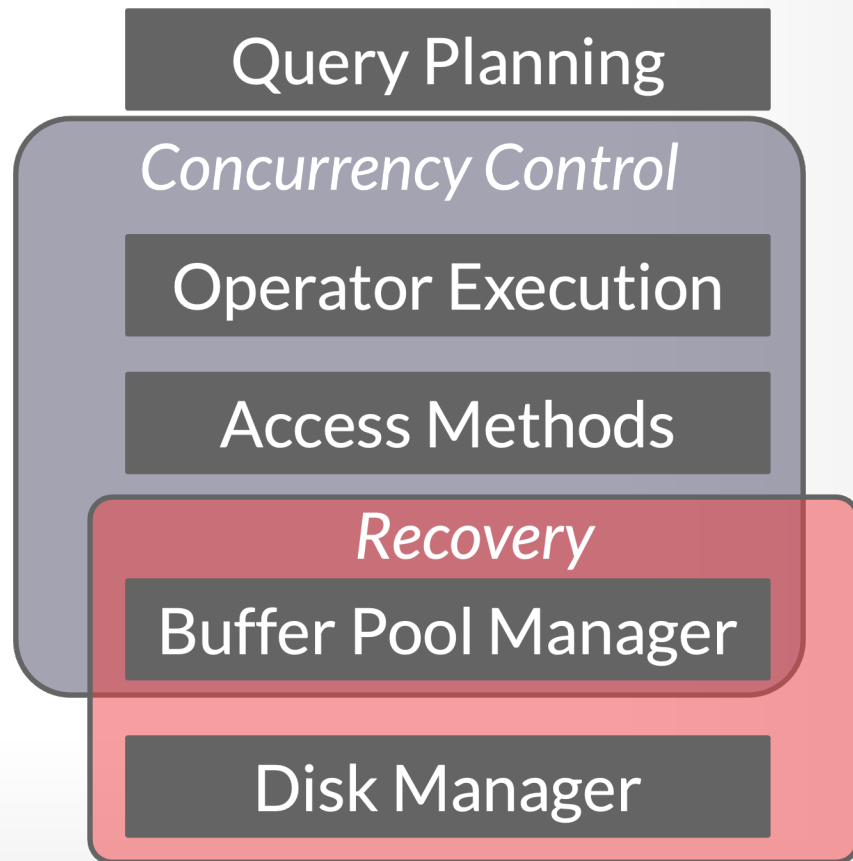
- Monday, April 21 @ 4:30pm
- EdgeQL with Gel
- Speaker: Michael Sullivan
- <https://cmu.zoom.us/j/93441451665>



COURSE STATUS

Databases are hard.

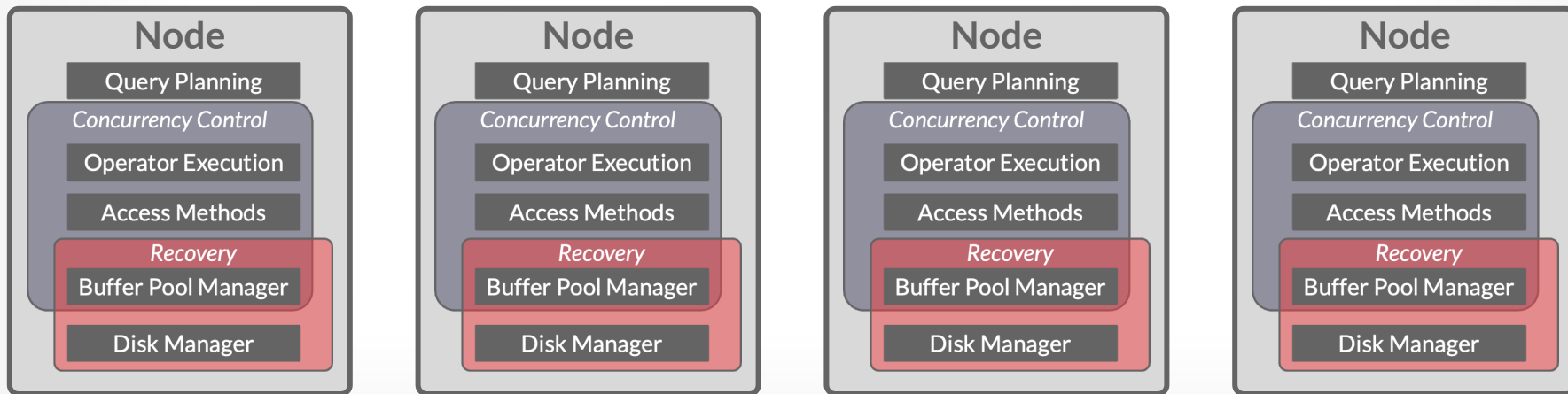
Distributed databases are harder.



COURSE STATUS

Databases are hard.

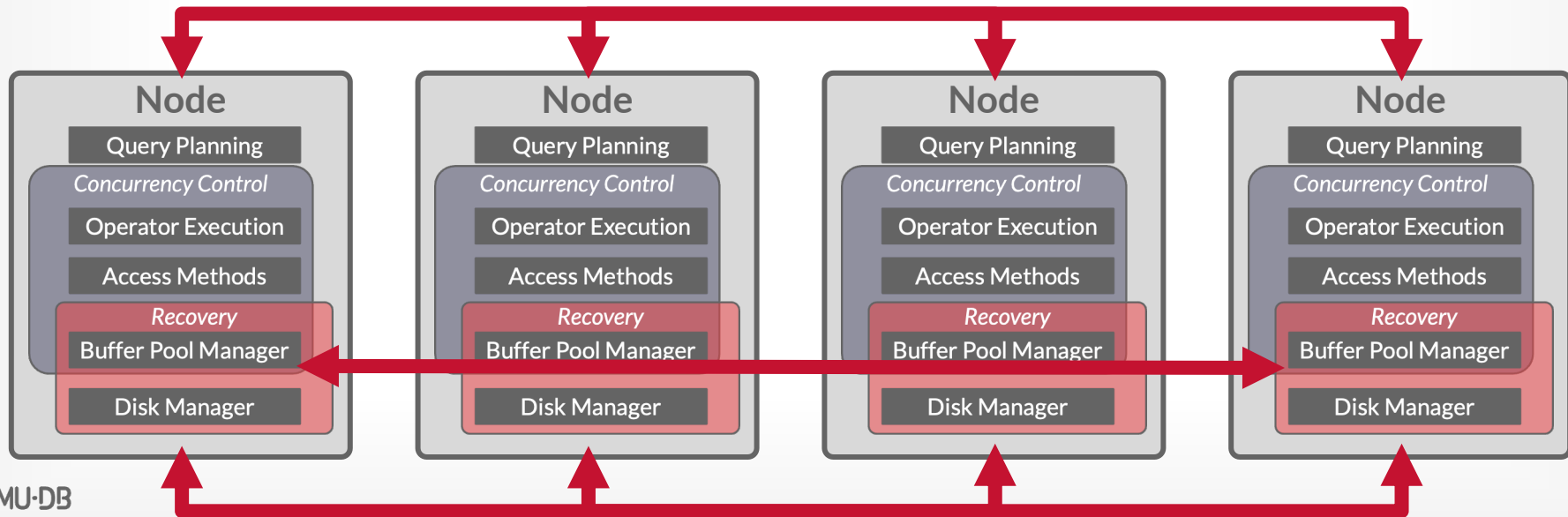
Distributed databases are harder.



COURSE STATUS

Databases are hard.

Distributed databases are harder.



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.

DISTRIBUTED DBMSs

Use the building blocks that we covered in single-node DBMSs to now support transaction processing and query execution in distributed environments.

- Optimization & Planning
- Concurrency Control
- Logging & Recovery

TODAY'S AGENDA

System Architectures

Design Issues

Partitioning Schemes

Distributed Concurrency Control

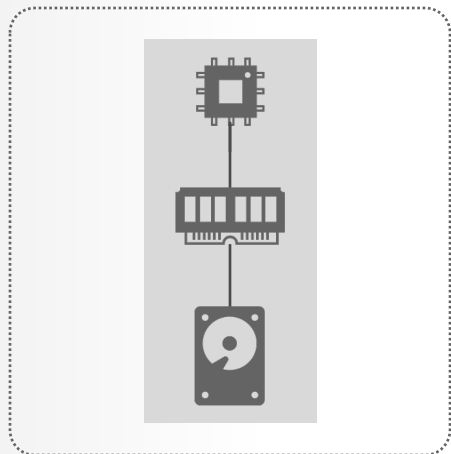
DB Flash Talk: [dbt Labs](#)

SYSTEM ARCHITECTURE

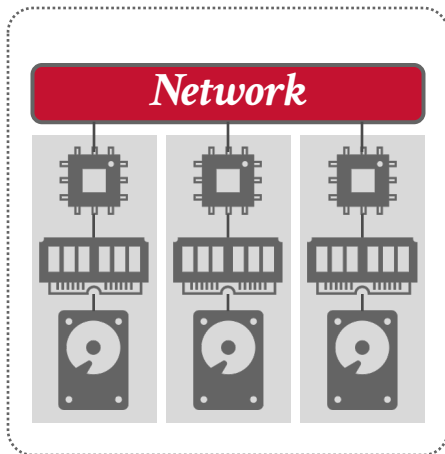
A distributed DBMS's system architecture specifies what shared resources are directly accessible to CPUs.

This affects how CPUs coordinate with each other and where they retrieve/store objects in the database.

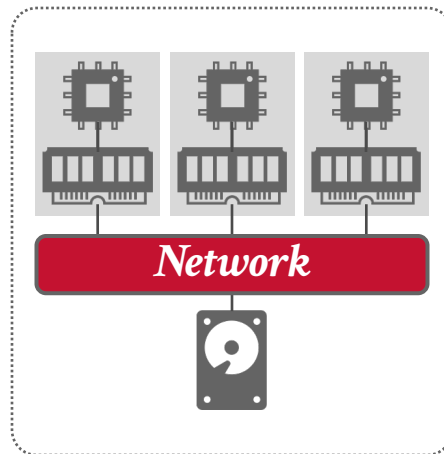
SYSTEM ARCHITECTURE



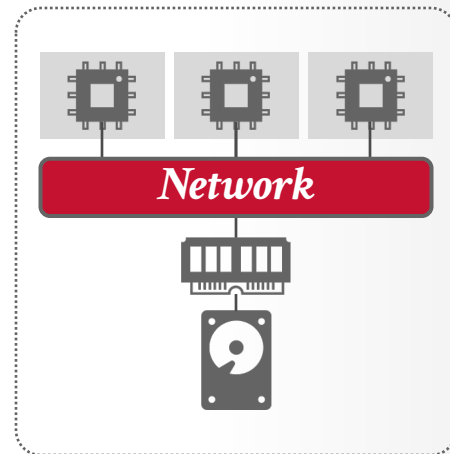
Shared
Everything



Shared
Nothing



Shared
Disk



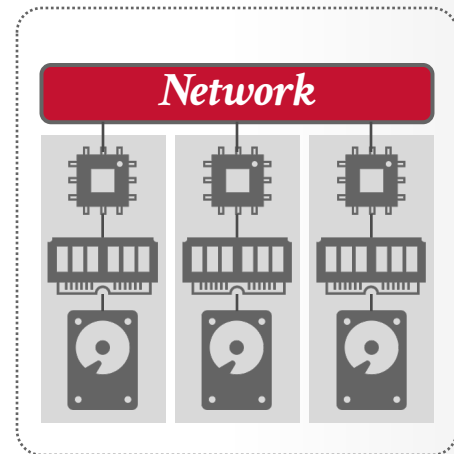
Shared
Memory

SHARED NOTHING

Each DBMS node has its own CPU, memory, and local disk.

Nodes only communicate with each other via network.

- Better performance & efficiency.
- Harder to scale capacity.
- Harder to ensure consistency.

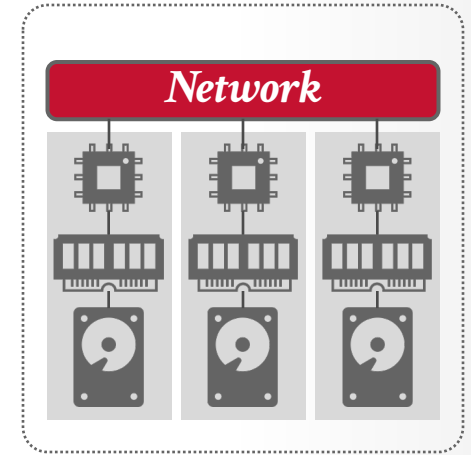


SHARED NOTHING

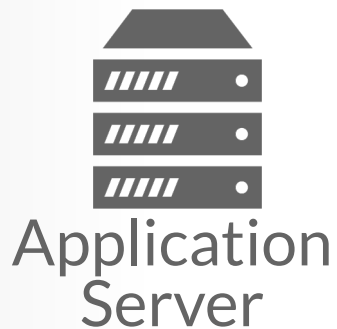
Each DBMS node has its own CPU, memory, and local disk.

Nodes only communicate with each other via network.

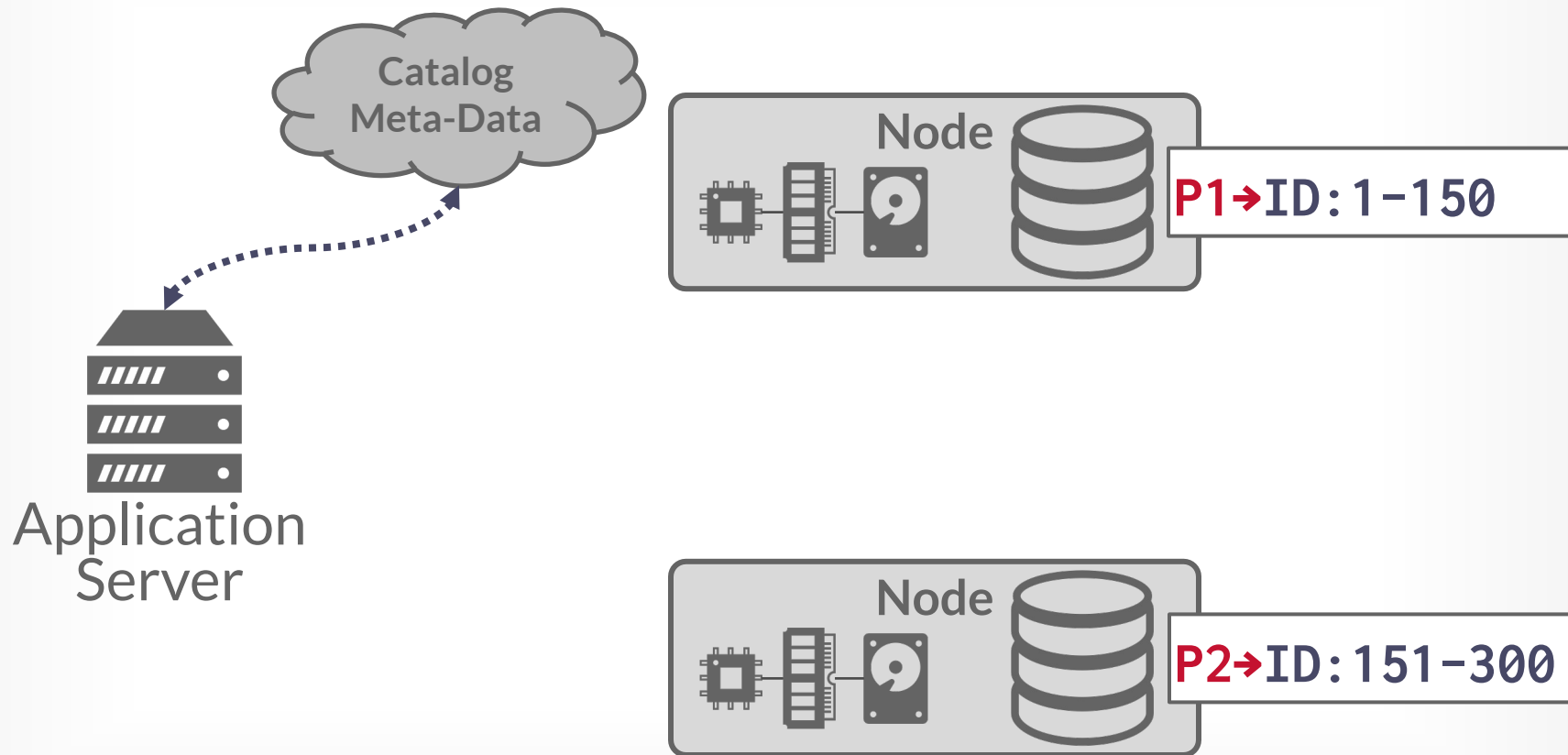
- Better performance & efficiency.
- Harder to scale capacity.
- Harder to ensure consistency.



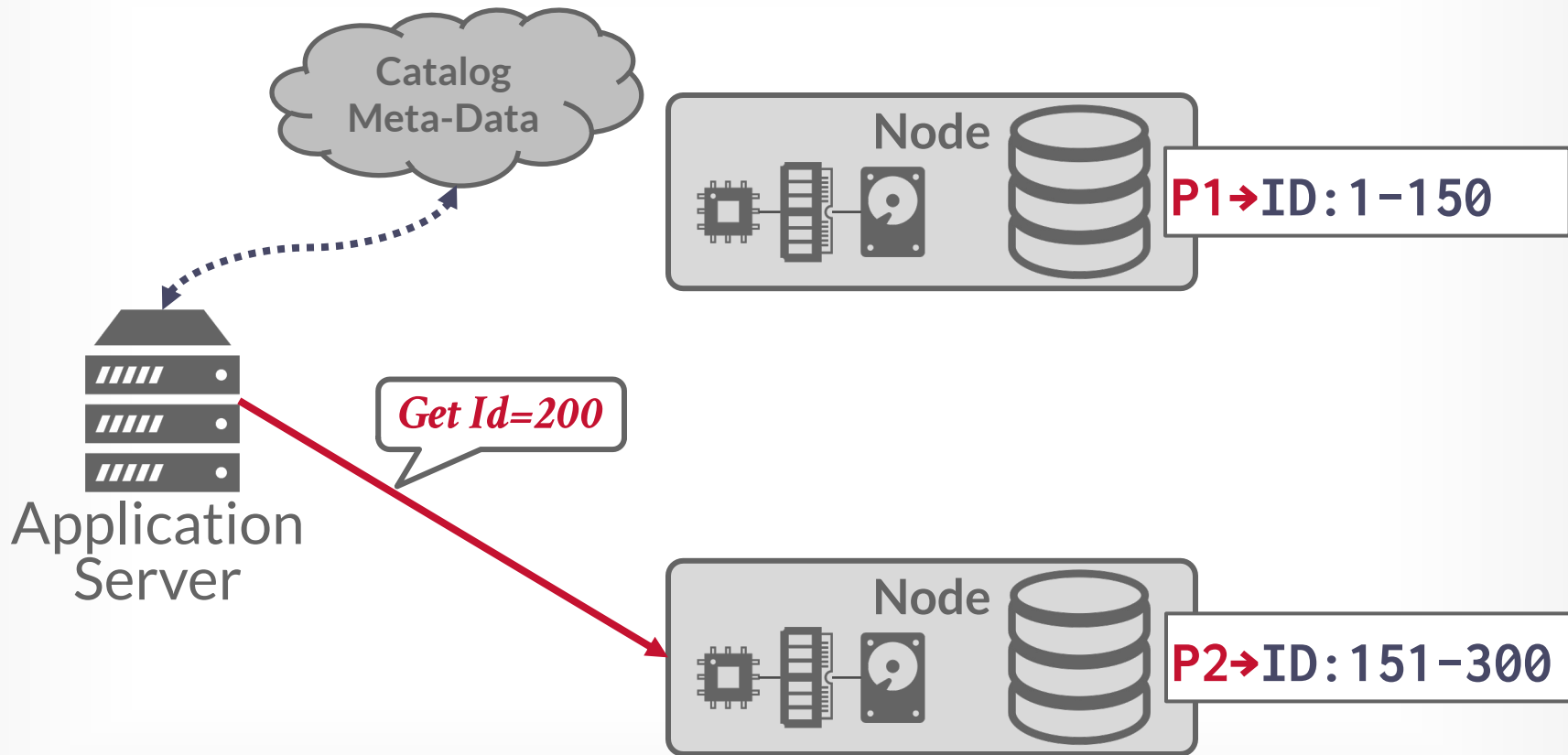
SHARED NOTHING EXAMPLE



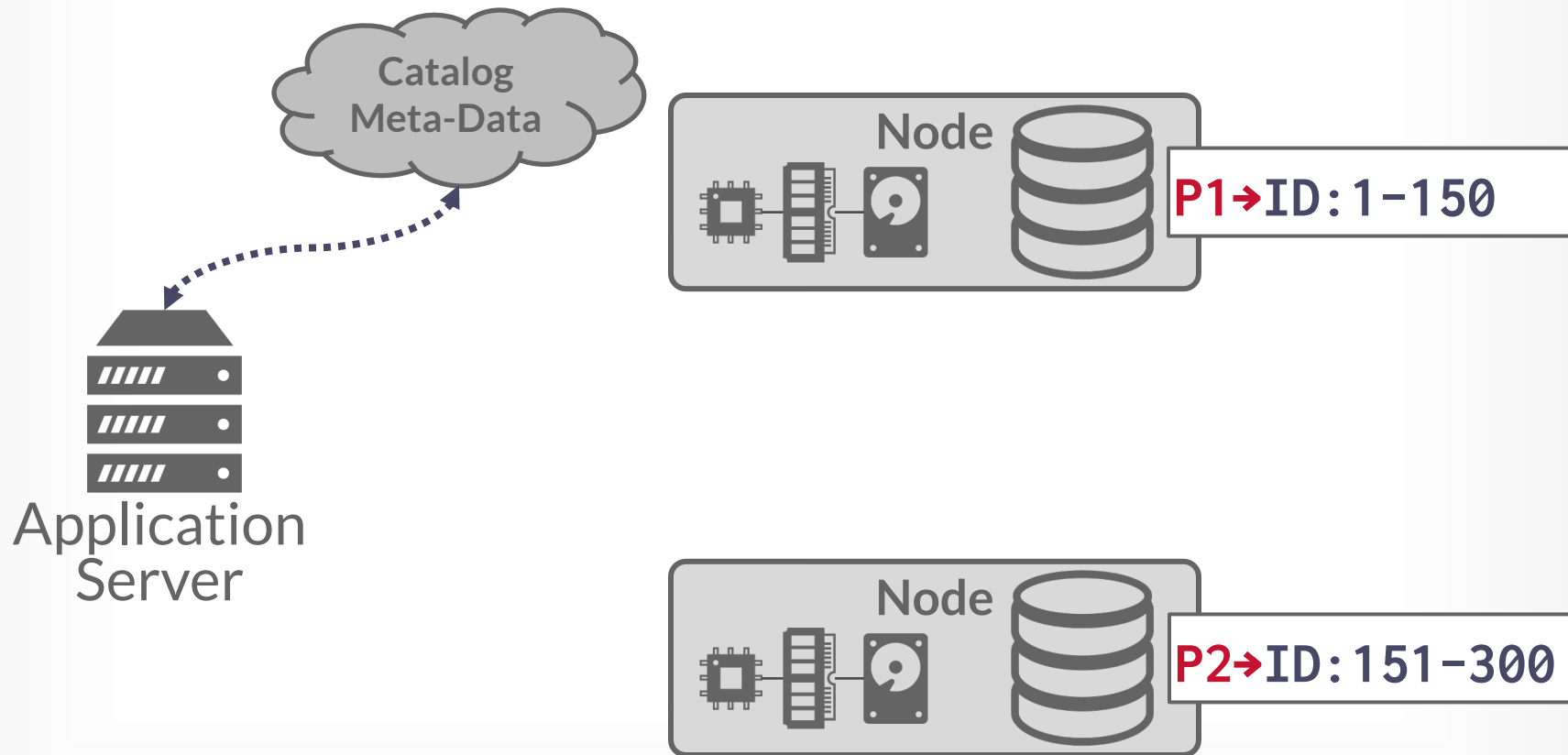
SHARED NOTHING EXAMPLE



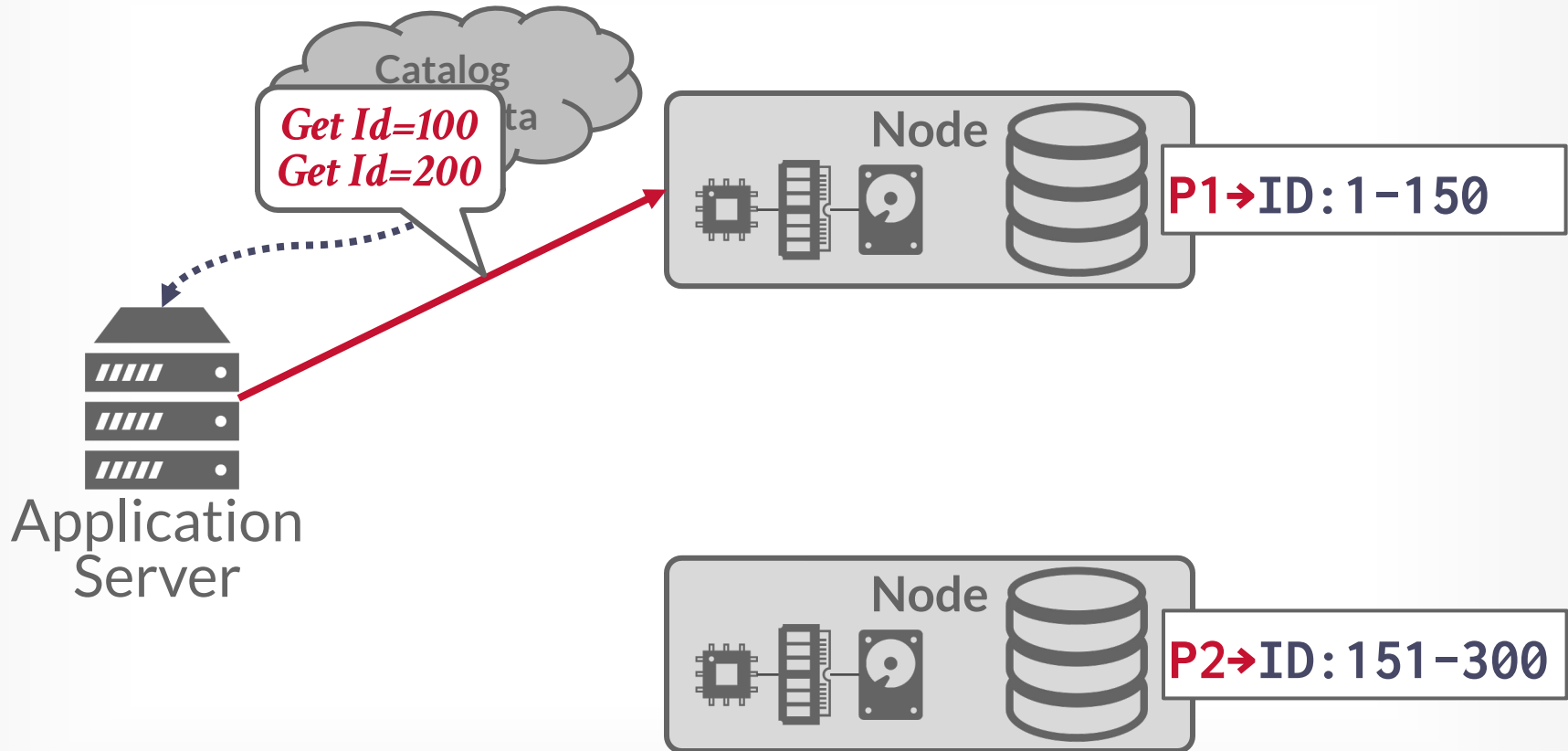
SHARED NOTHING EXAMPLE



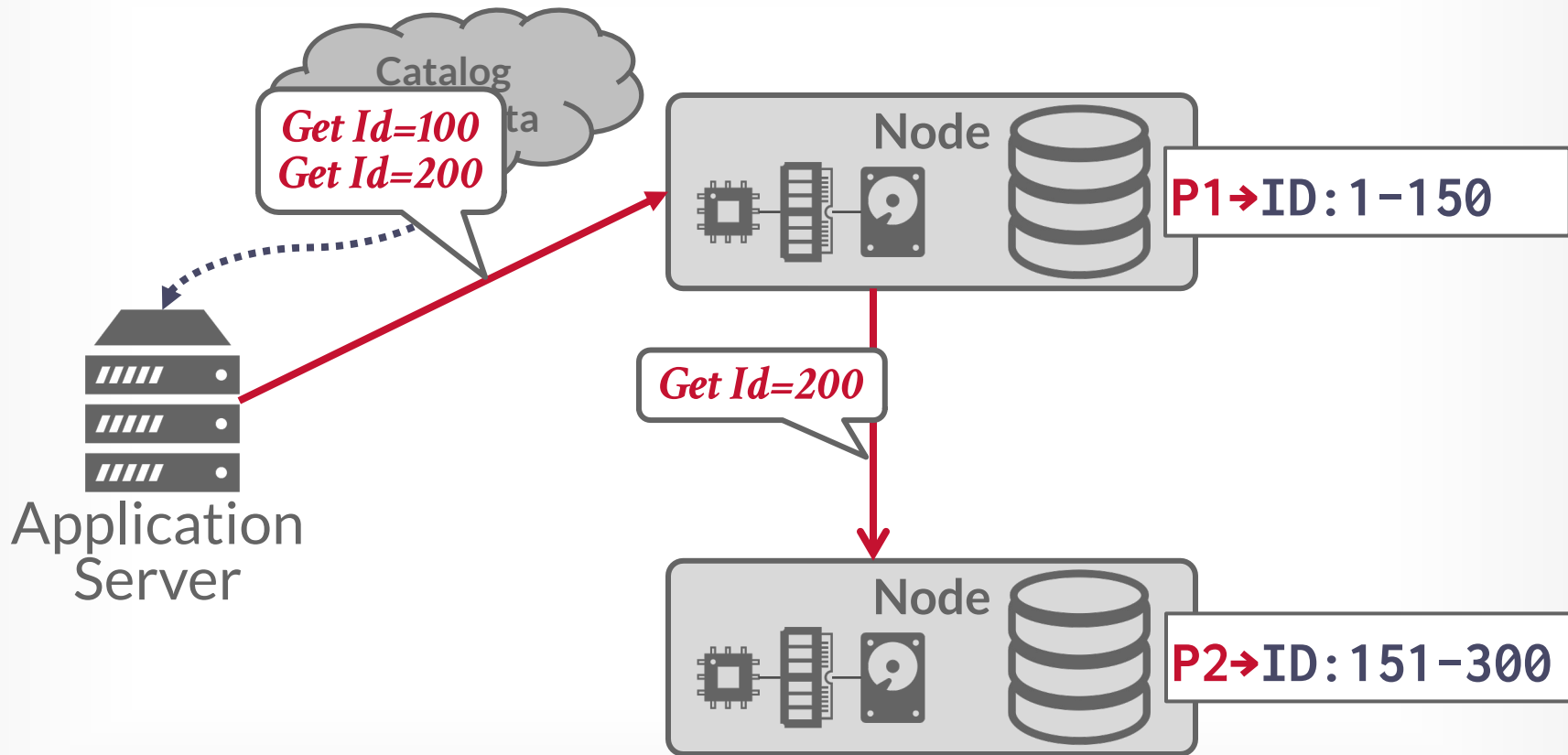
SHARED NOTHING EXAMPLE



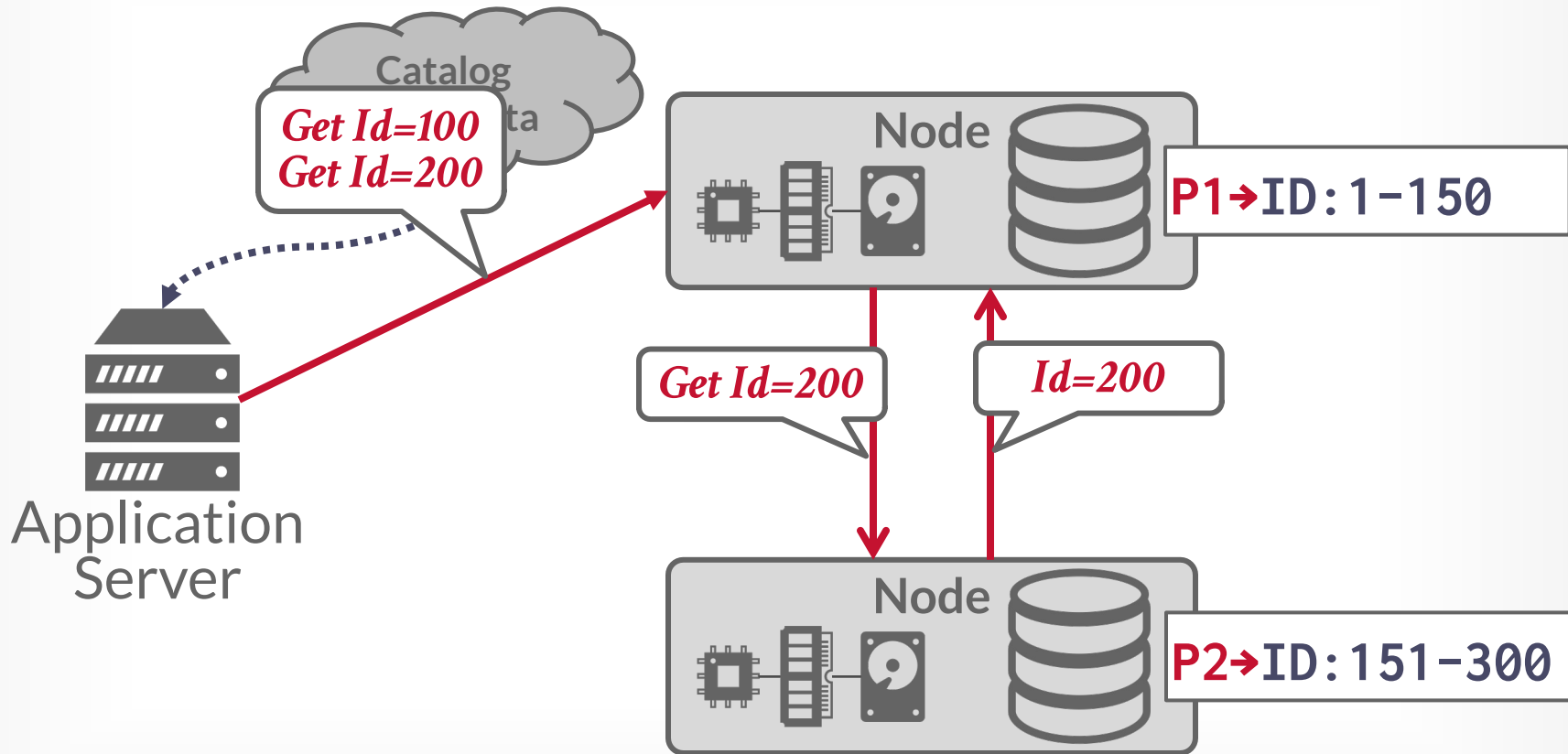
SHARED NOTHING EXAMPLE



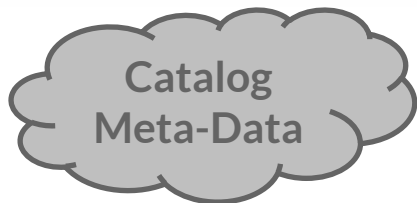
SHARED NOTHING EXAMPLE



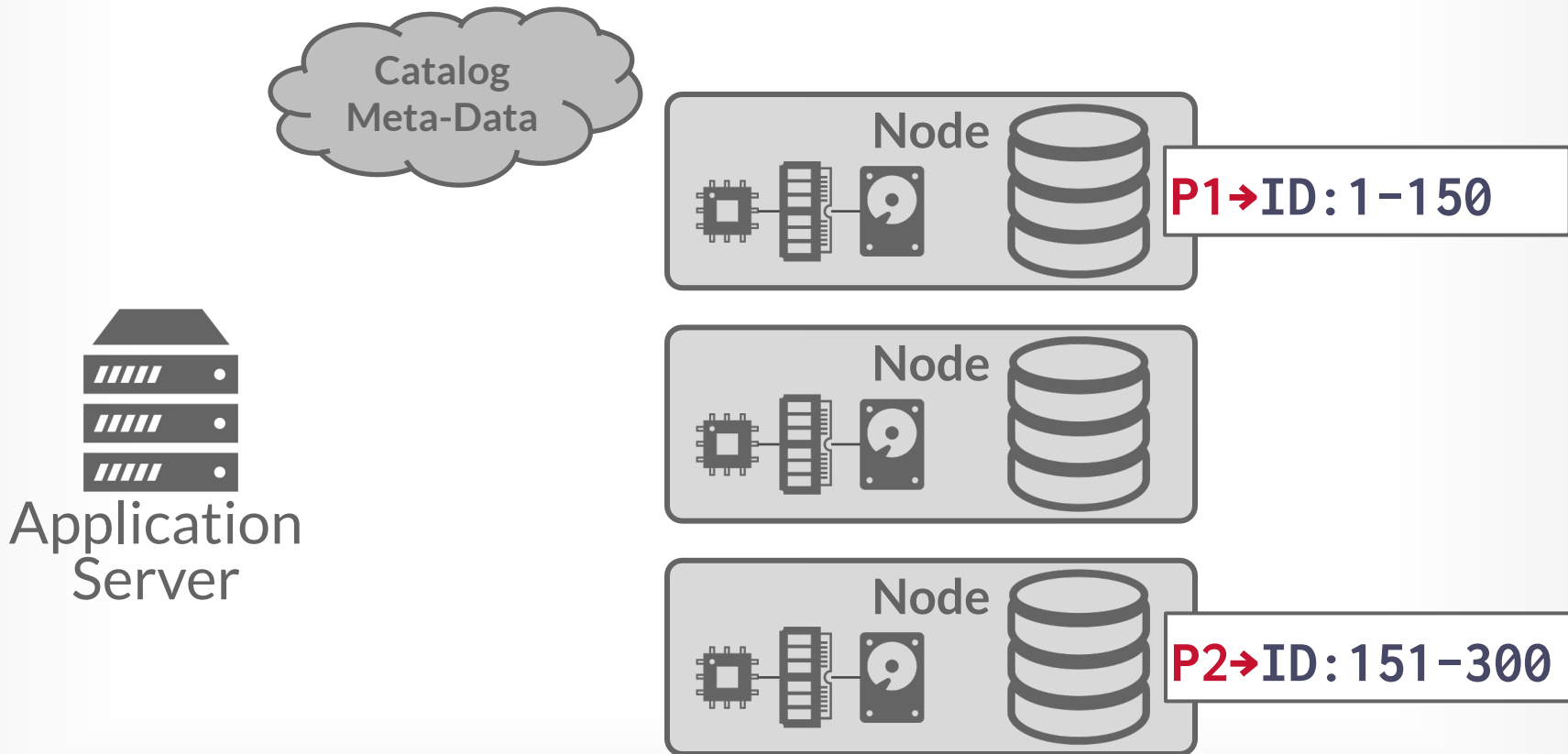
SHARED NOTHING EXAMPLE



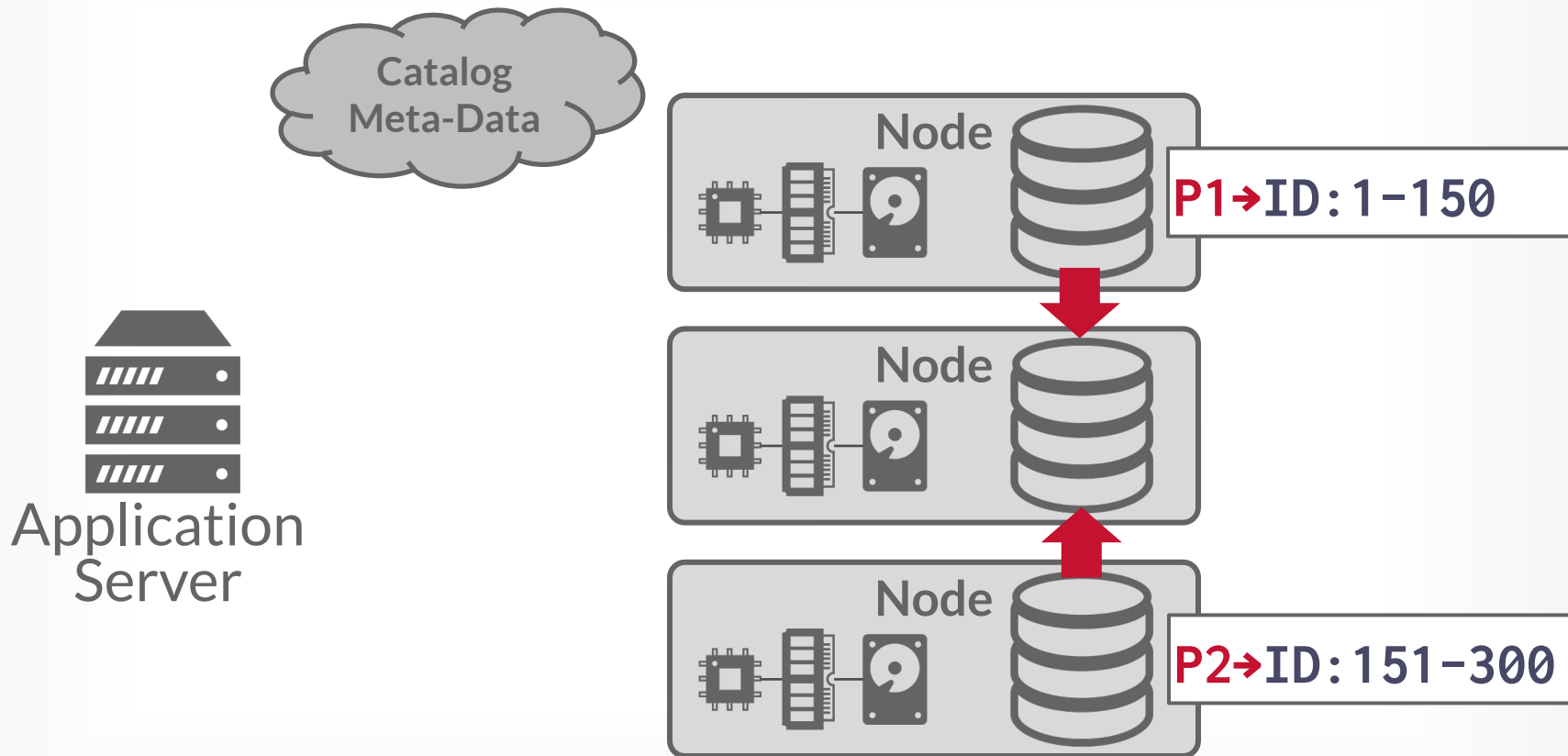
SHARED NOTHING EXAMPLE



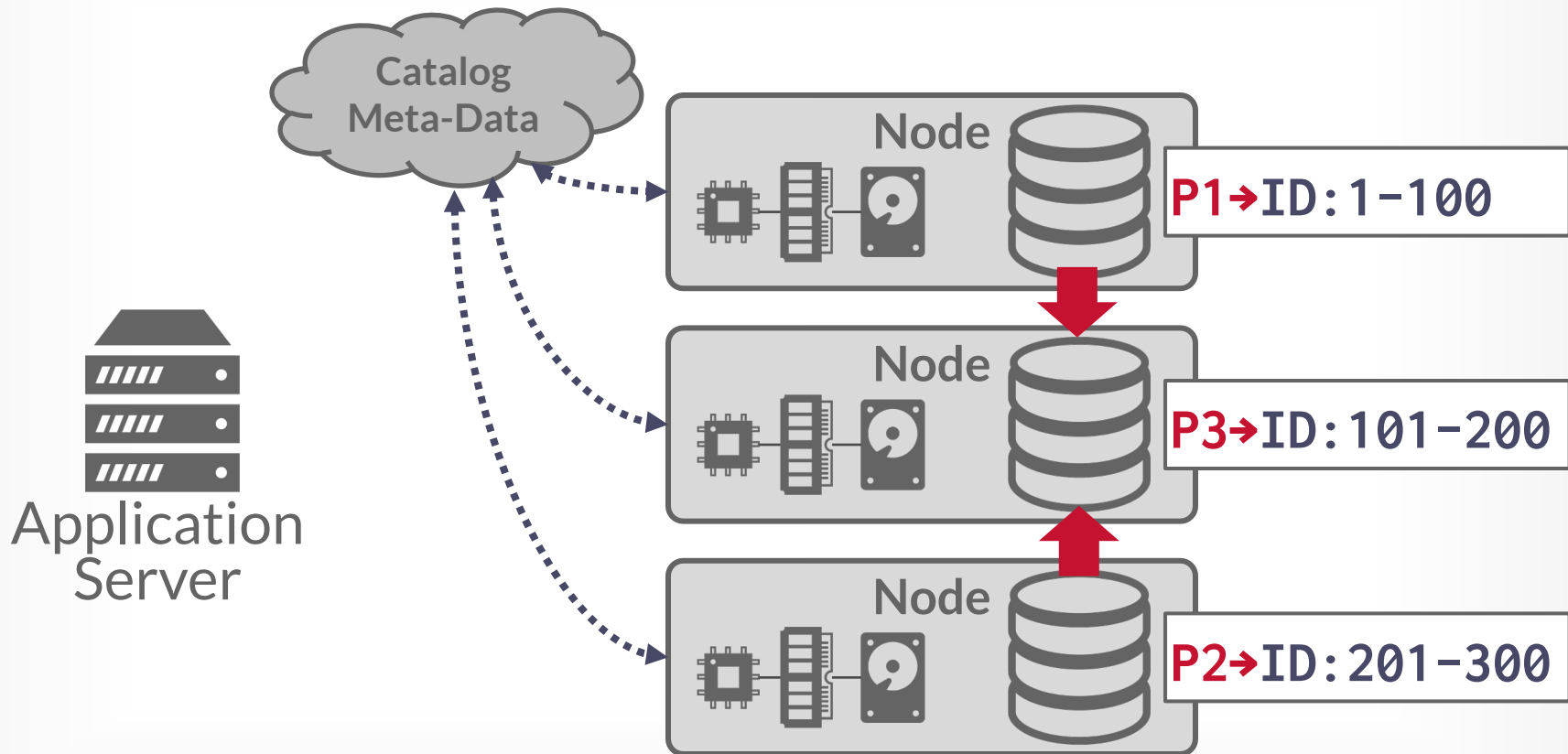
SHARED NOTHING EXAMPLE



SHARED NOTHING EXAMPLE



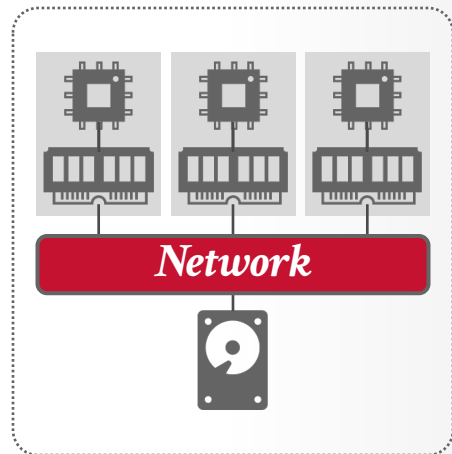
SHARED NOTHING EXAMPLE



SHARED DISK

Nodes access a single logical disk via an interconnect, but each have their own private memories.

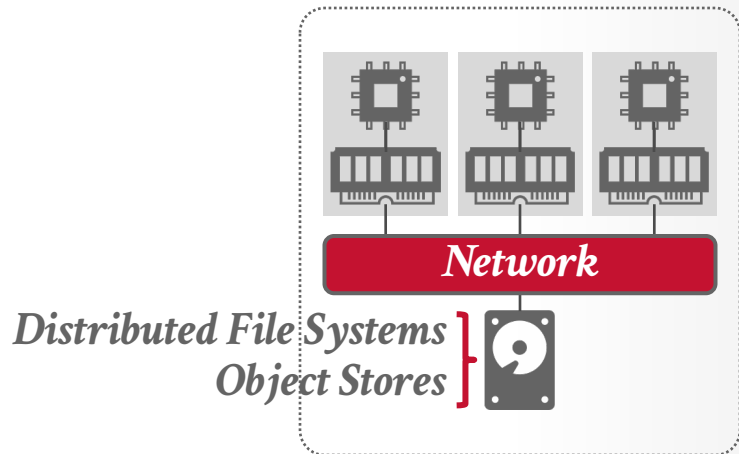
- Scale execution layer independently from the storage layer.
- Nodes can still use direct attached storage as a slower/larger cache.
- This architecture facilitates data lakes and serverless systems.



SHARED DISK

Nodes access a single logical disk via an interconnect, but each have their own private memories.

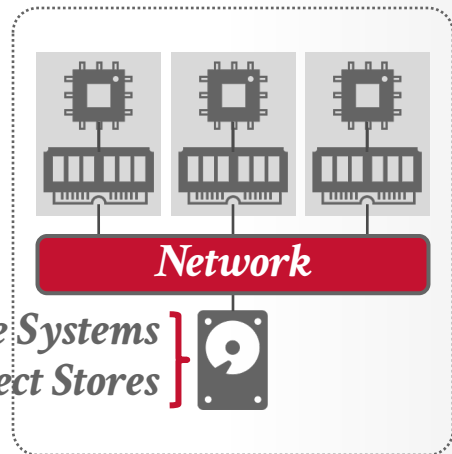
- Scale execution layer independently from the storage layer.
- Nodes can still use direct attached storage as a slower/larger cache.
- This architecture facilitates **data lakes** and **serverless** systems.



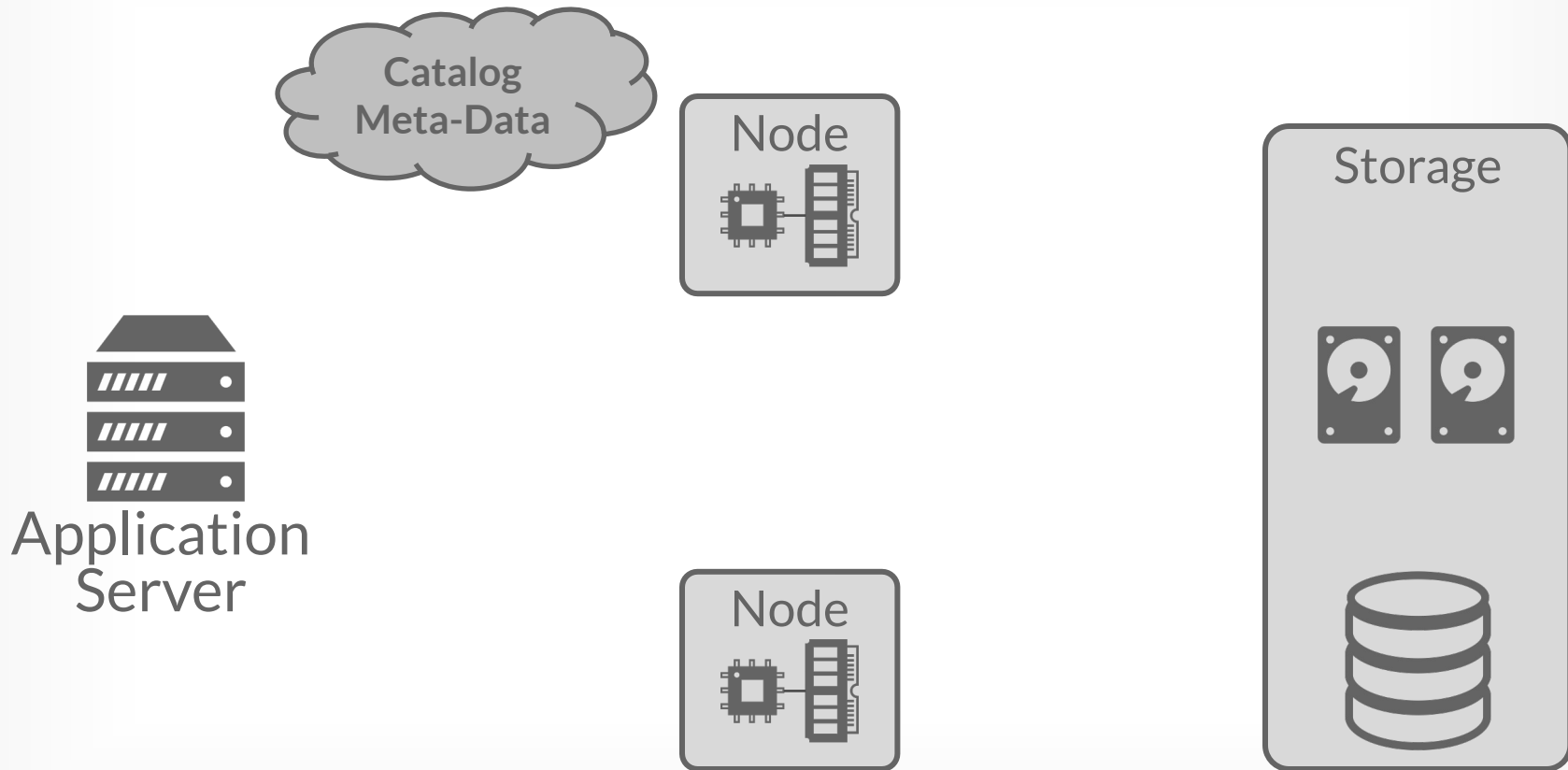
SHARED DISK

Nodes access a single logical disk via an interconnect, but each have their own private memories.

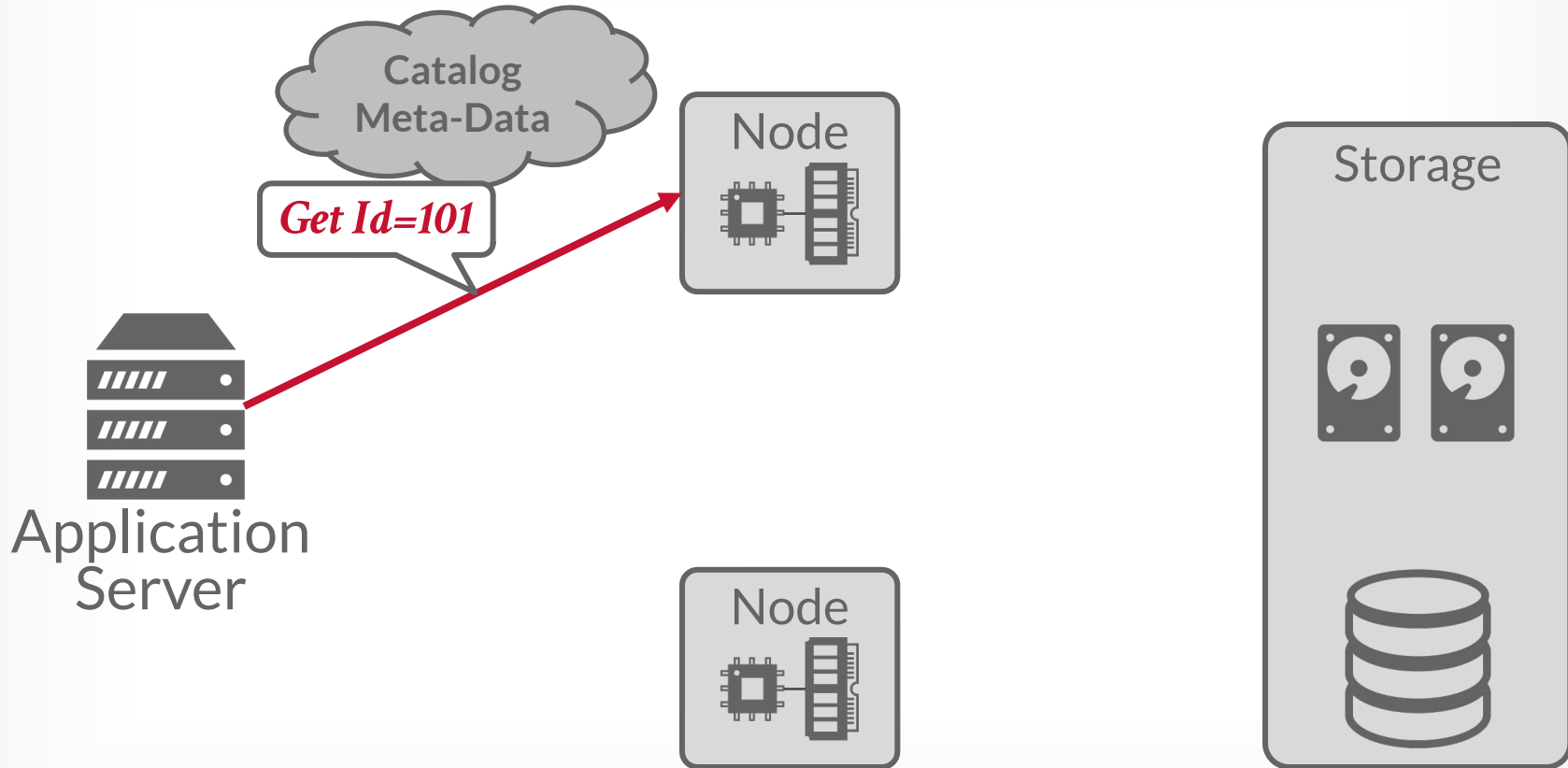
- Scale execution layer independently from the storage layer.
- Nodes can still use direct attached storage as a slower/larger cache.
- This architecture facilitates data lakes and serverless systems.



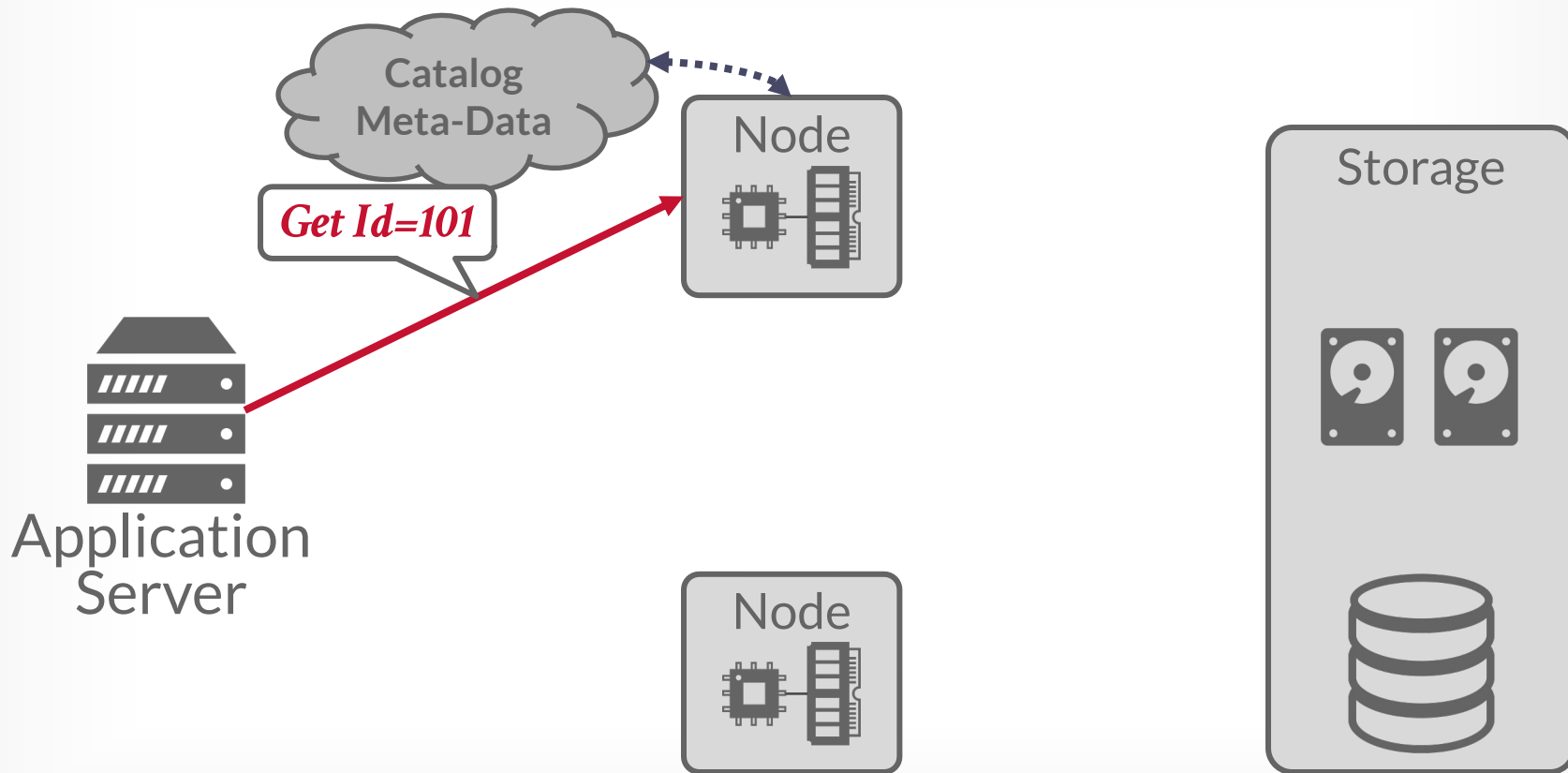
SHARED DISK EXAMPLE



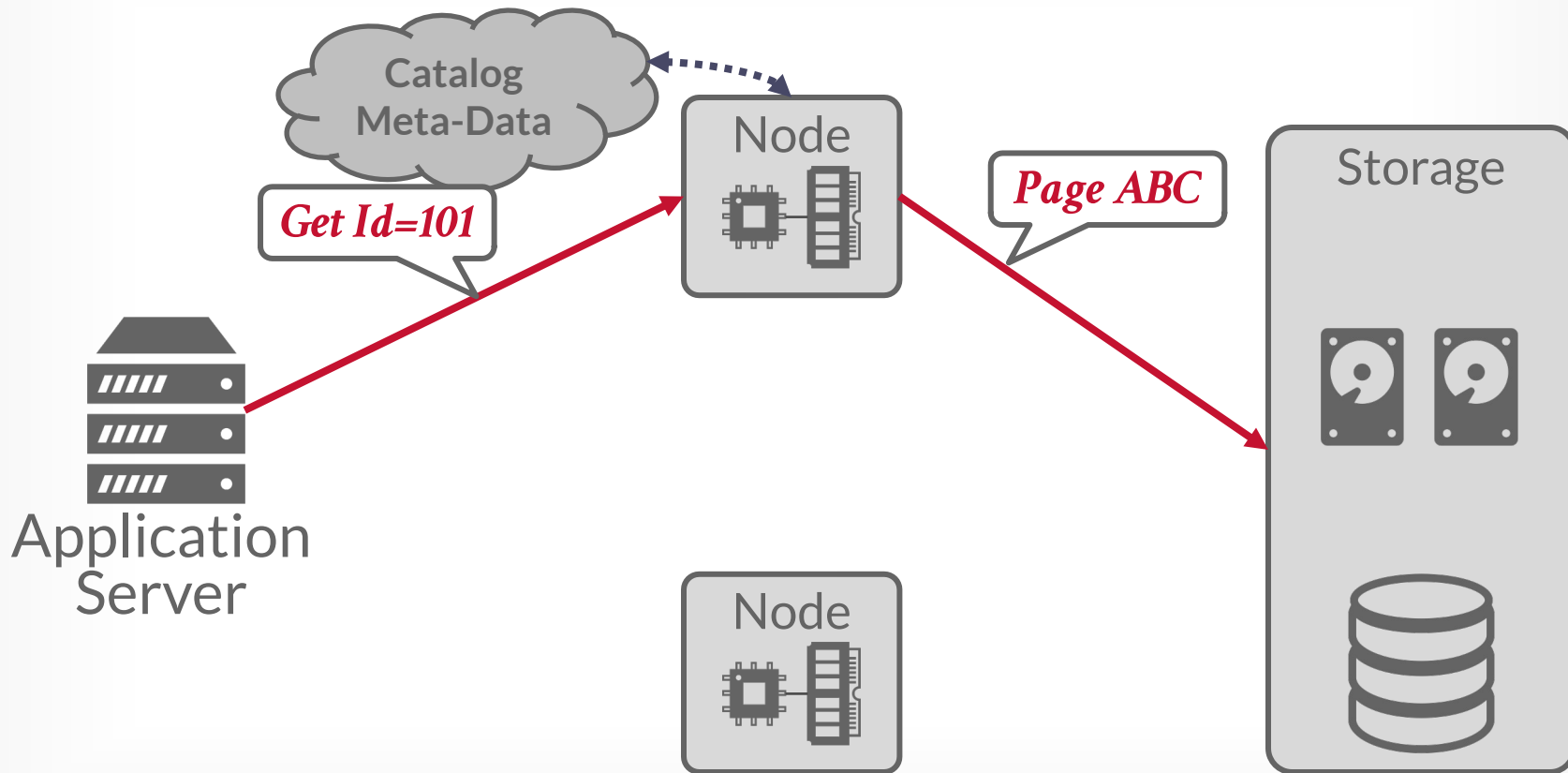
SHARED DISK EXAMPLE



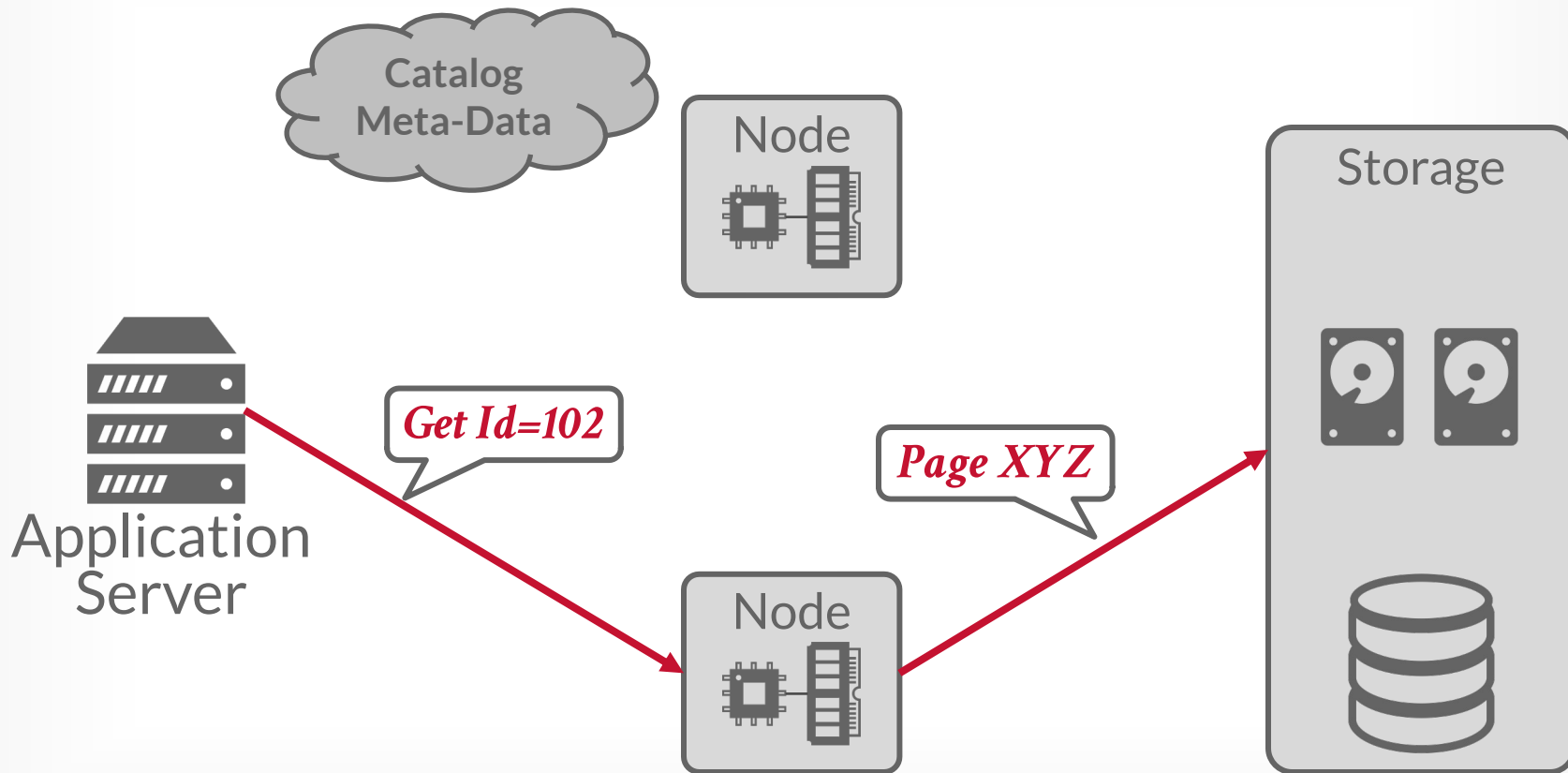
SHARED DISK EXAMPLE



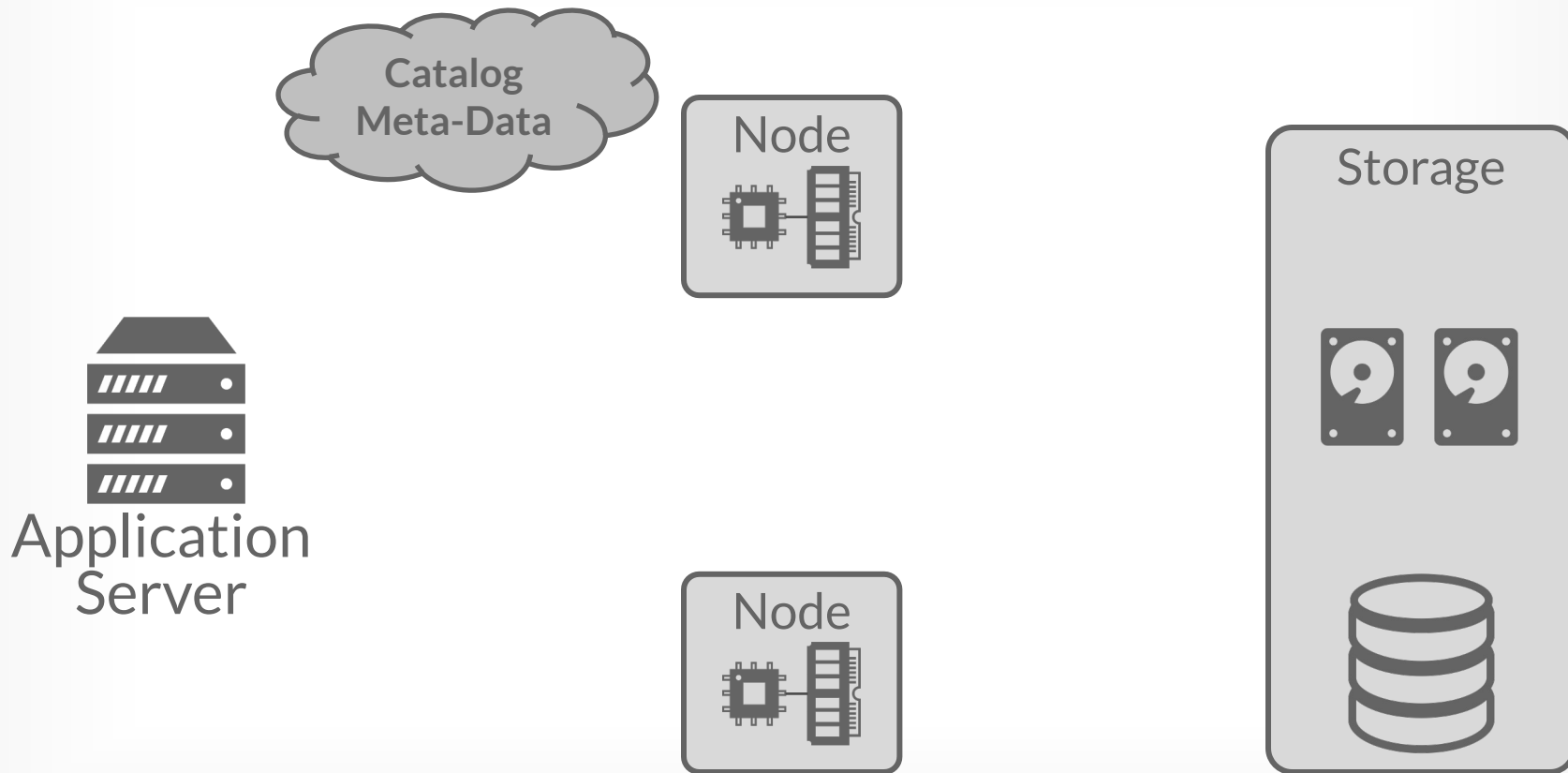
SHARED DISK EXAMPLE



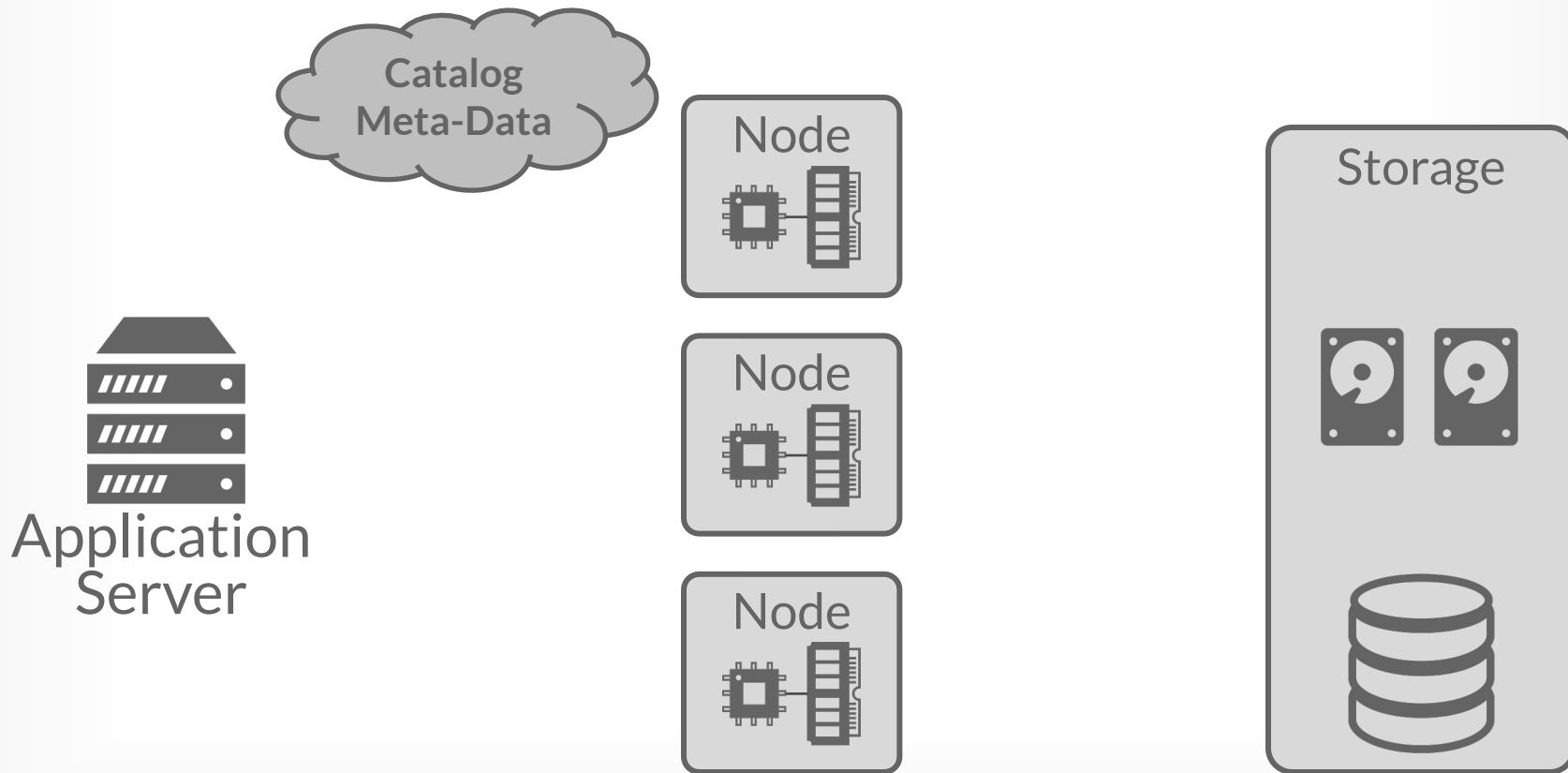
SHARED DISK EXAMPLE



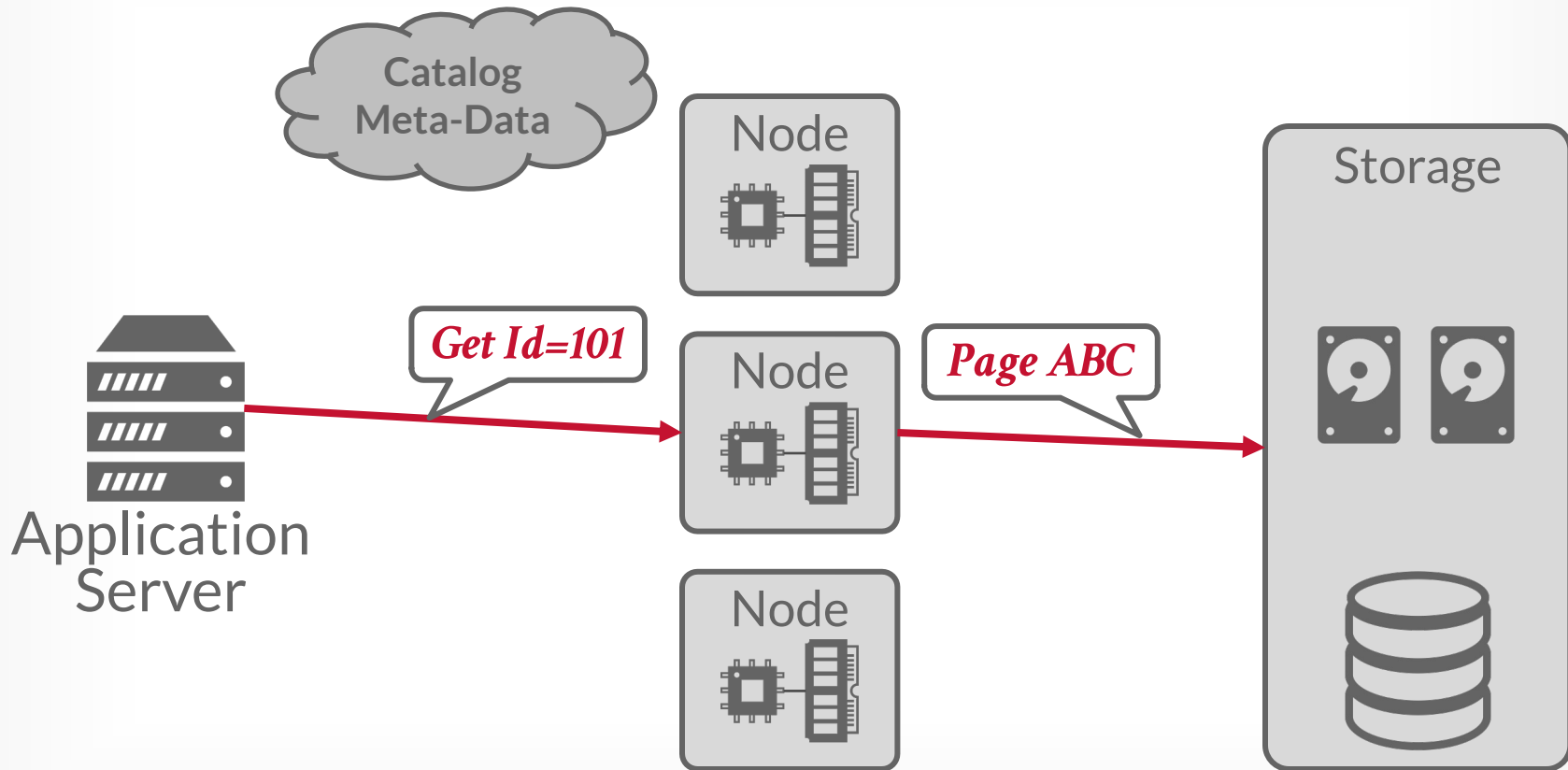
SHARED DISK EXAMPLE



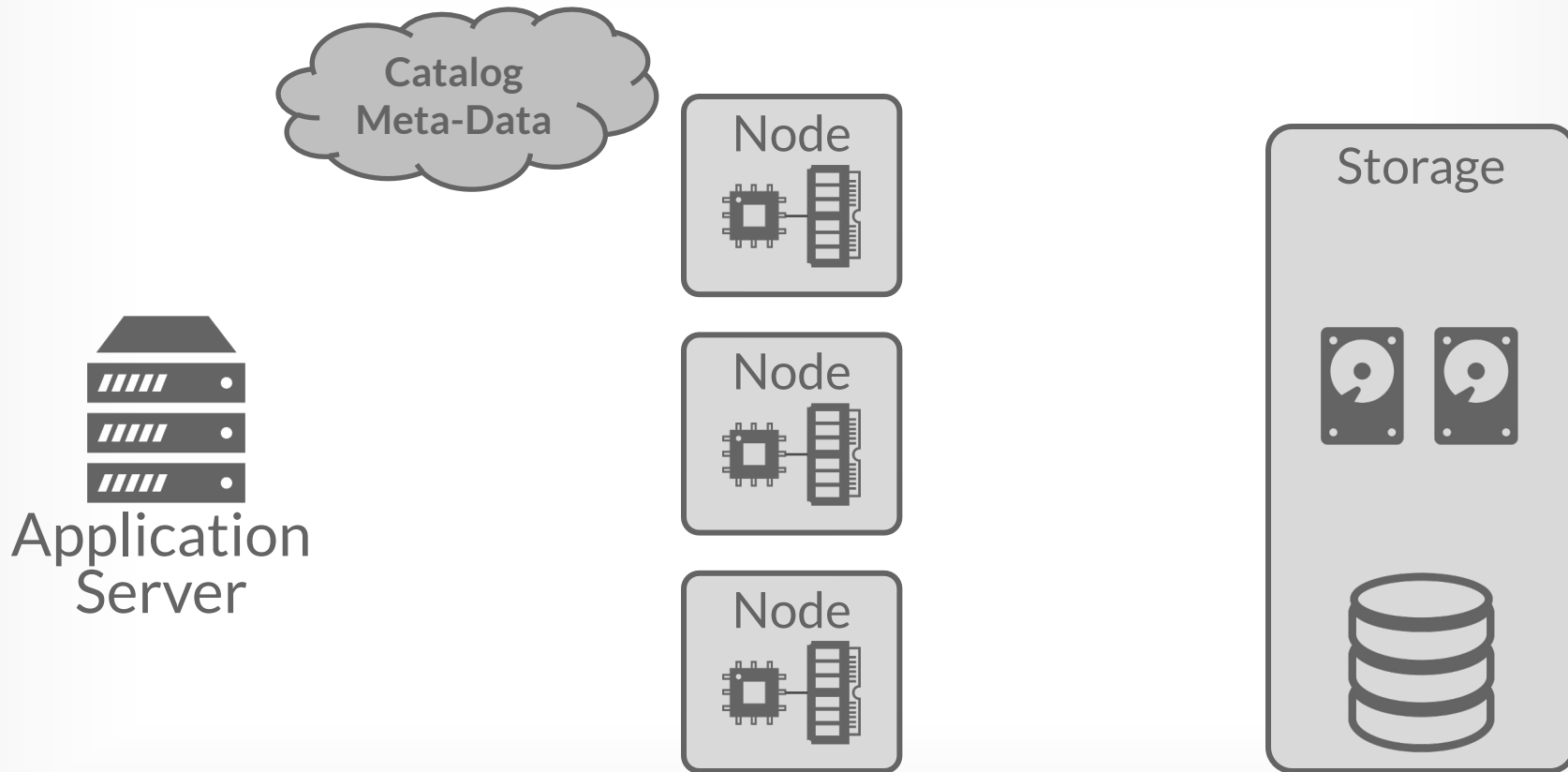
SHARED DISK EXAMPLE



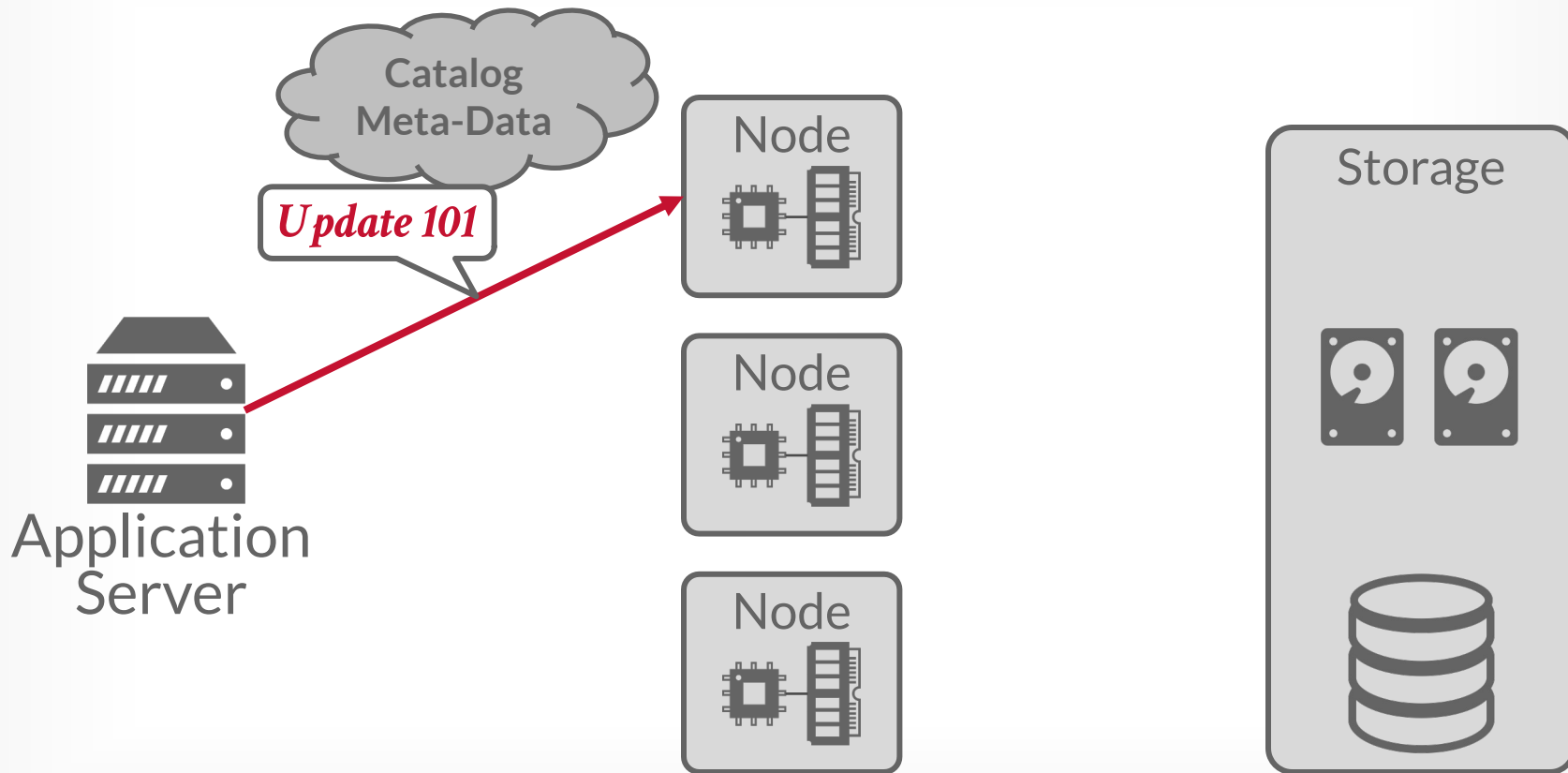
SHARED DISK EXAMPLE



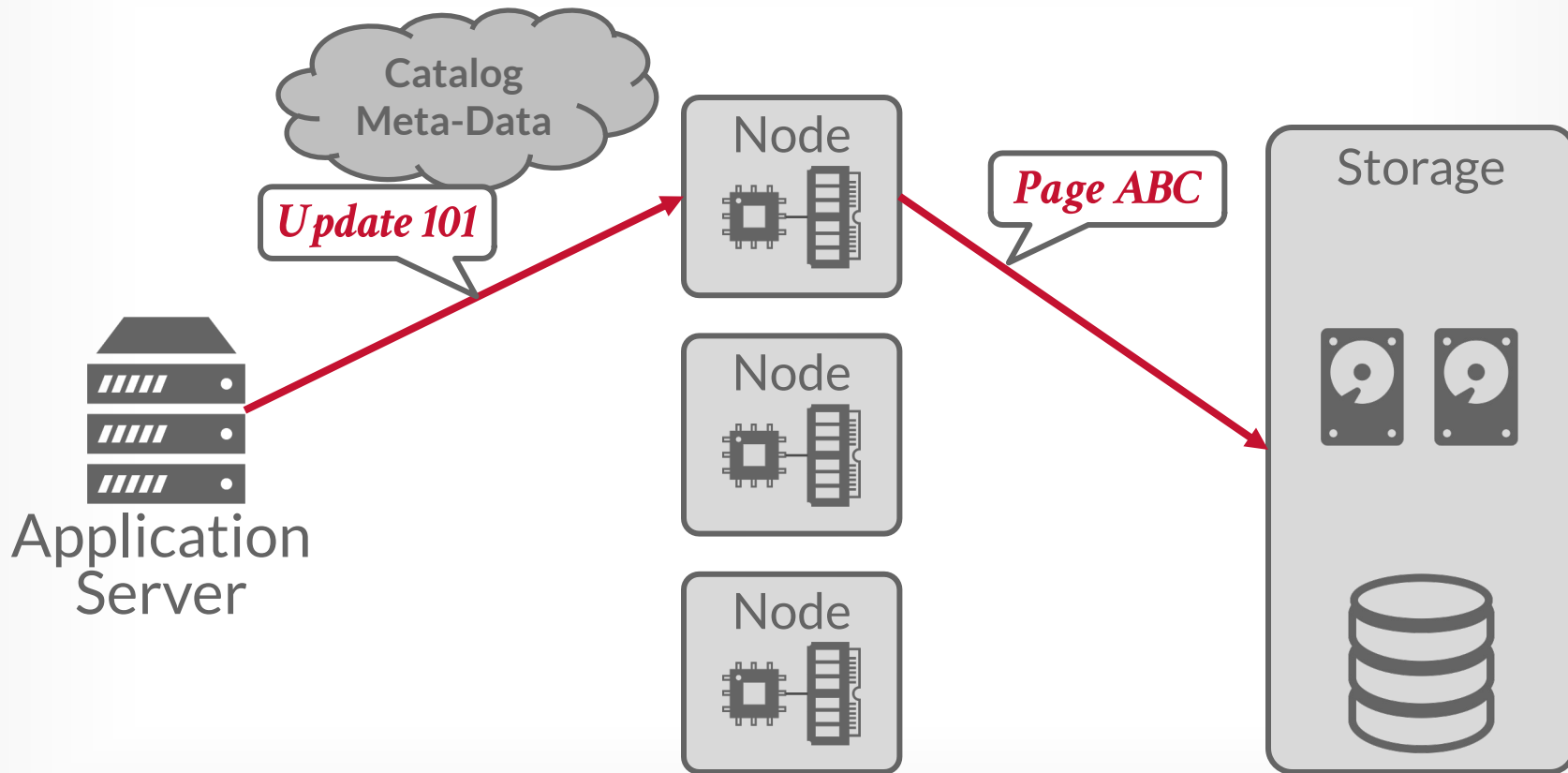
SHARED DISK EXAMPLE



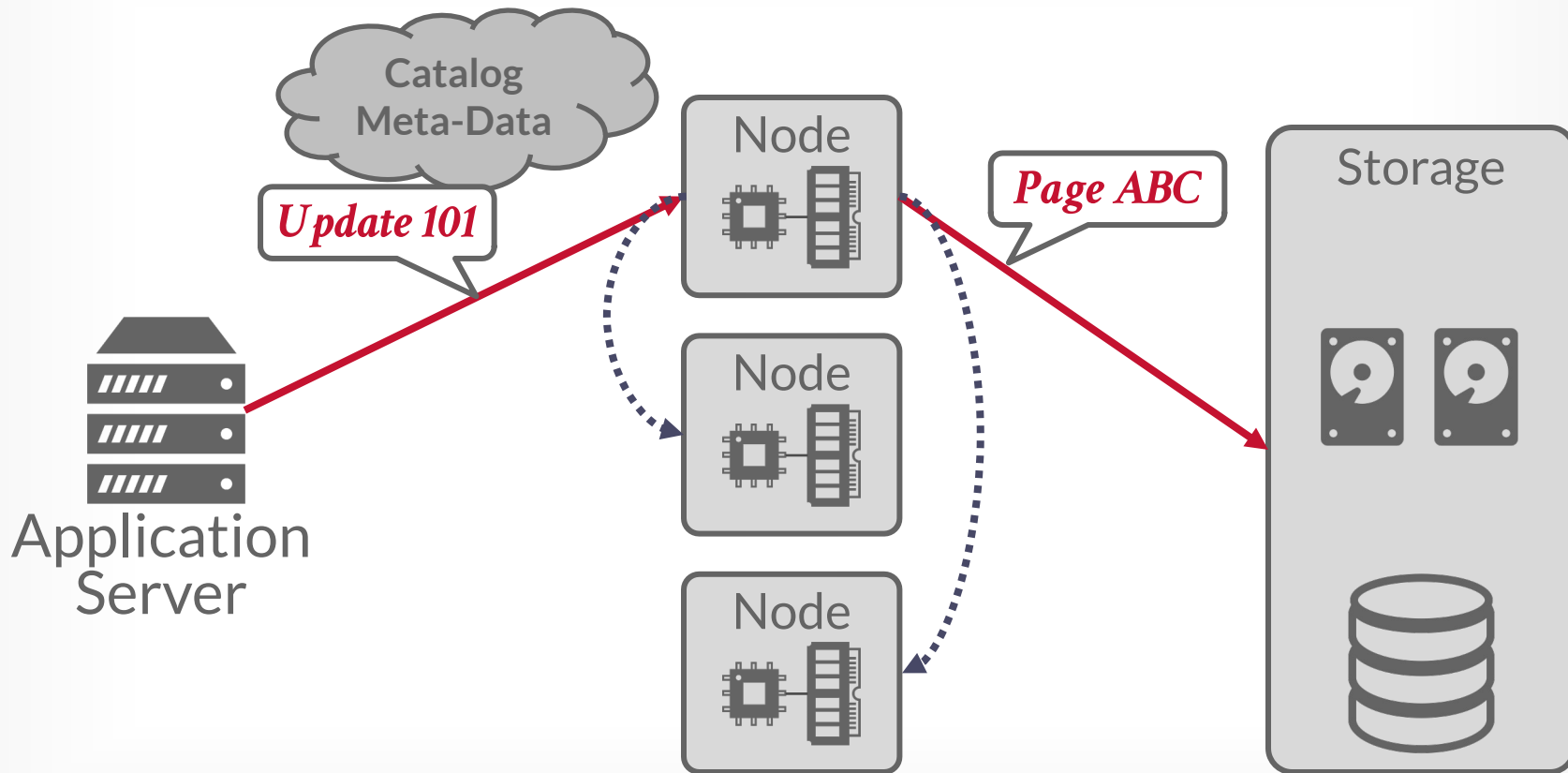
SHARED DISK EXAMPLE



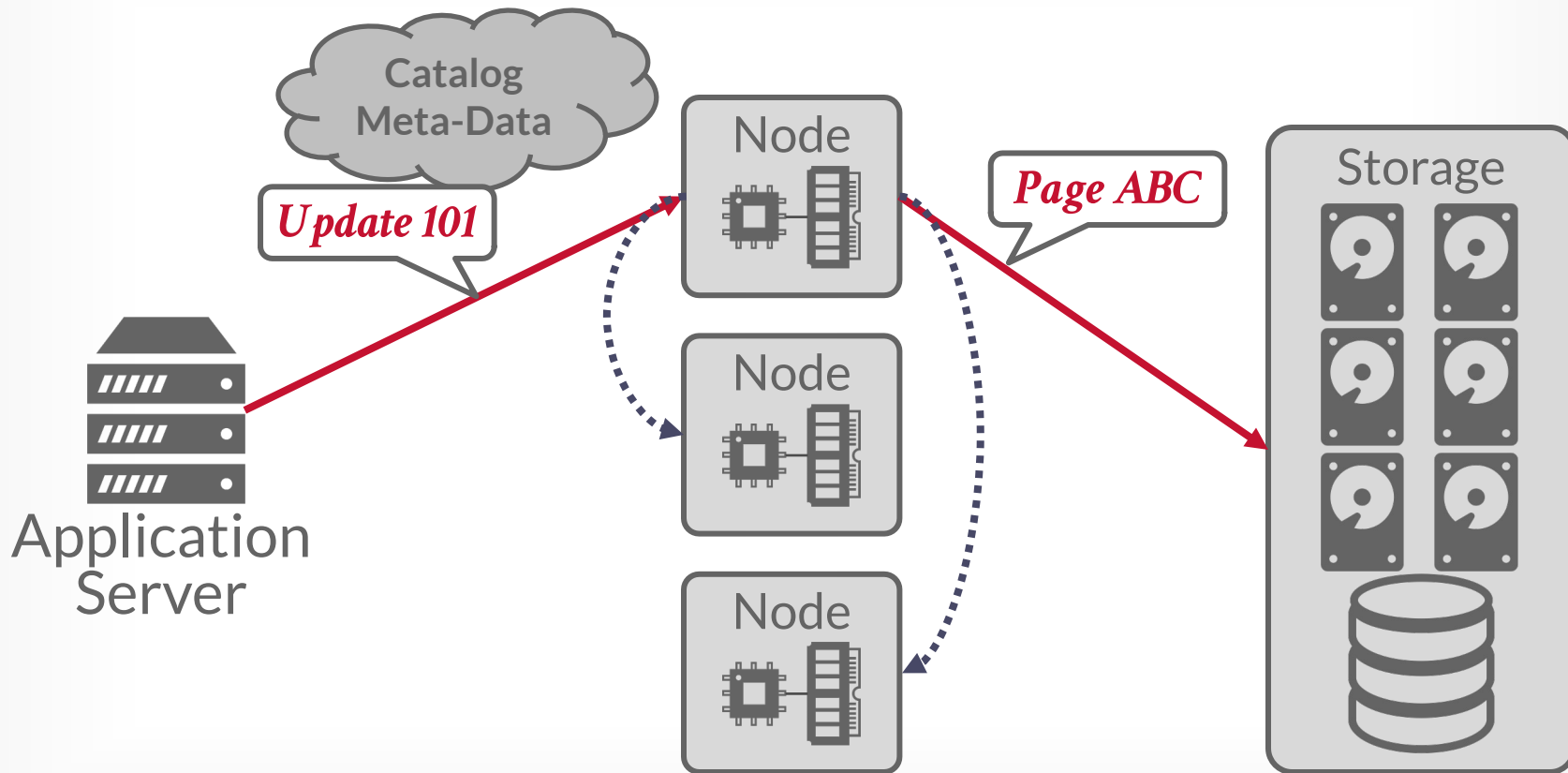
SHARED DISK EXAMPLE



SHARED DISK EXAMPLE



SHARED DISK EXAMPLE

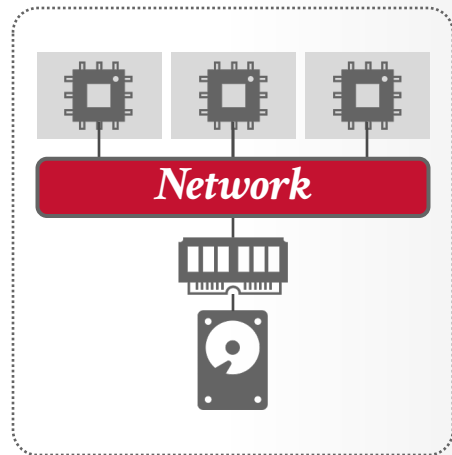


SHARED MEMORY

Nodes access a common memory address space via a fast interconnect.

- Each node has a global view of all the in-memory data structures.
- Can still use local memory / disk for intermediate results.

This looks a lot like shared-everything. Nobody does this.



DESIGN ISSUES

How does the application find data?

Where does the application send queries?

How to execute queries on distributed data?

→ Push query to data.

→ Pull data to query.

How do we divide the database across resources?

How does the DBMS ensure correctness?

DESIGN ISSUES

How does the application find data?

Where does the application send queries?

How to execute queries on distributed data?

→ Push query to data.

→ Pull data to query.

How do we divide the database across resources?

How does the DBMS ensure correctness? *Next Class*

DATA TRANSPARENCY

Applications should not be required to know where data is physically located in a distributed DBMS.

→ Any query that run on a single-node DBMS should produce the same result on a distributed DBMS.

In practice, developers need to be aware of the communication costs of queries to avoid excessively "expensive" data movement.

Chip -> Chip	Intra-rack	Planetary
~100 ns	~100 μ s	~100 ms

UPI ToR Switches Fiber optics

DATA TRANSPARENCY

Applications should not be required to know where data is physically located in a distributed DBMS.

→ Any query that run on a single-node DBMS should produce the same result on a distributed DBMS.

In practice, developers need to be aware of the communication costs of queries to avoid excessively "expensive" data movement.



DATABASE PARTITIONING

Split database across multiple resources:

→ Disks, nodes, processors.

→ Called "sharding" in NoSQL systems.

The DBMS executes query fragments on each partition and then combines the results to produce a single answer.

NAÏVE TABLE PARTITIONING

Assign an entire table to a single node.

Assumes that each node has enough storage space for an entire table.

Ideal if queries never join data across tables stored on different nodes and access patterns are uniform.

NAÏVE TABLE PARTITIONING

Table1

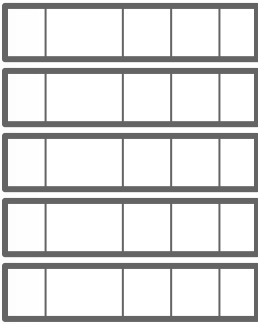
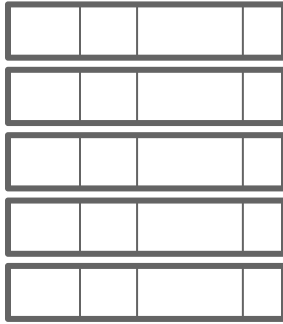


Table2



Partitions



Ideal Query:

```
SELECT * FROM table1
```


NAÏVE TABLE PARTITIONING

Table1

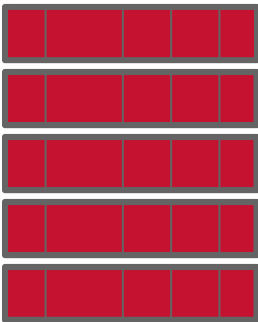
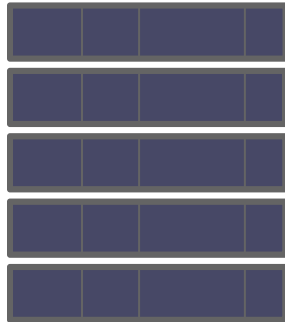
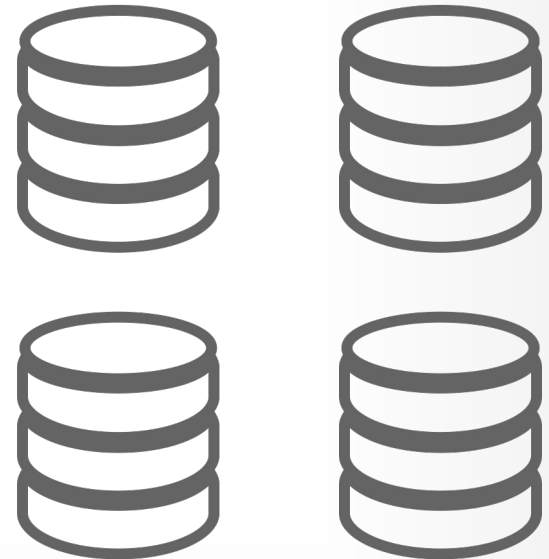


Table2



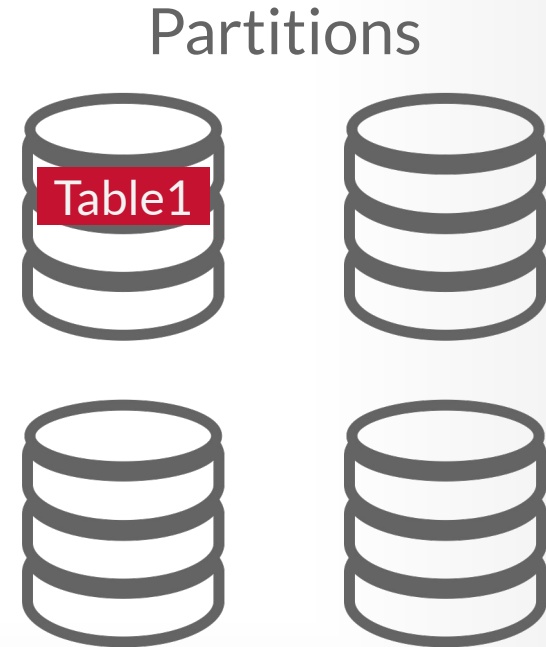
Partitions



Ideal Query:

```
SELECT * FROM table1
```

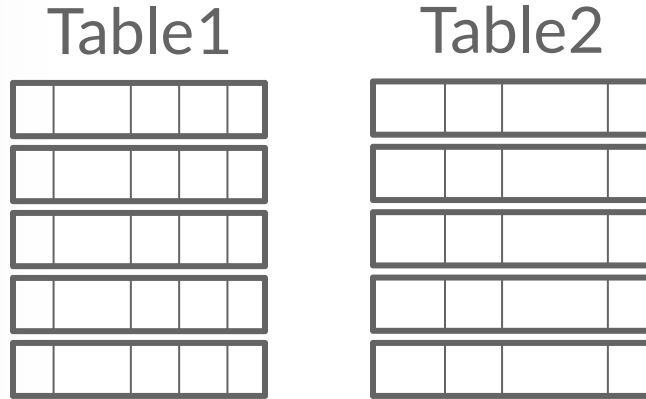
NAÏVE TABLE PARTITIONING



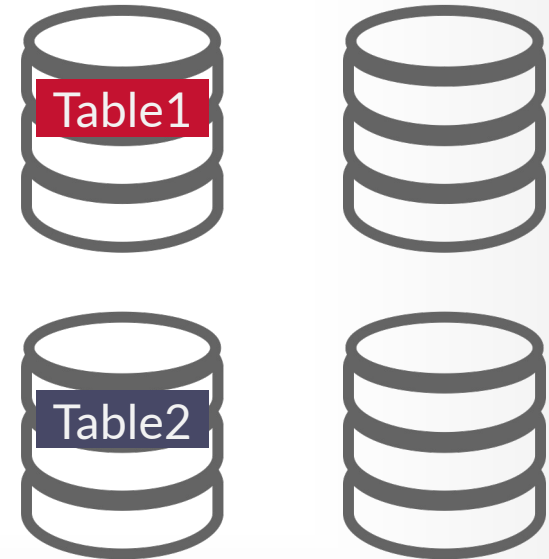
Ideal Query:

```
SELECT * FROM table1
```

NAÏVE TABLE PARTITIONING



Partitions



Ideal Query:

```
SELECT * FROM table1
```

VERTICAL PARTITIONING

Split a table's attributes into separate partitions.

Must store tuple information to reconstruct the original record.

```
CREATE TABLE foo (  
  attr1 INT,  
  attr2 INT,  
  attr3 INT,  
  attr4 TEXT  
);
```

Tuple#1	attr1	attr2	attr3	attr4
Tuple#2	attr1	attr2	attr3	attr4
Tuple#3	attr1	attr2	attr3	attr4
Tuple#4	attr1	attr2	attr3	attr4

VERTICAL PARTITIONING

Split a table's attributes into separate partitions.

Must store tuple information to reconstruct the original record.

```
CREATE TABLE foo (  
  attr1 INT,  
  attr2 INT,  
  attr3 INT,  
  attr4 TEXT  
);
```

Tuple#1	attr1	attr2	attr3	attr4
Tuple#2	attr1	attr2	attr3	attr4
Tuple#3	attr1	attr2	attr3	attr4
Tuple#4	attr1	attr2	attr3	attr4

VERTICAL PARTITIONING

Split a table's attributes into separate partitions.

Must store tuple information to reconstruct the original record.

```
CREATE TABLE foo (  
  attr1 INT,  
  attr2 INT,  
  attr3 INT,  
  attr4 TEXT  
);
```

Partition #1

Tuple#1	attr1	attr2	attr3
Tuple#2	attr1	attr2	attr3
Tuple#3	attr1	attr2	attr3
Tuple#4	attr1	attr2	attr3



Partition #2

Tuple#1	attr4
Tuple#2	attr4
Tuple#3	attr4
Tuple#4	attr4

HORIZONTAL PARTITIONING

Split a table's tuples into disjoint subsets based on some partitioning key and scheme.

→ Choose column(s) that divides the database equally in terms of size, load, or usage.

Partitioning Schemes:

→ Hashing

→ Ranges

→ Predicates

HORIZONTAL PARTITIONING

Table

101	a	XXX	2022-11-29
102	b	XXY	2022-11-28
103	c	XYZ	2022-11-29
104	d	XYX	2022-11-27
105	e	XYY	2022-11-29

Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

Partitions



HORIZONTAL PARTITIONING

Partitioning Key

Table

101	a	XXX	2022-11-29
102	b	XXY	2022-11-28
103	c	XYZ	2022-11-29
104	d	XYX	2022-11-27
105	e	XYY	2022-11-29

Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

Partitions



HORIZONTAL PARTITIONING

Partitioning Key

Table

101	a	XXX	2022-11-29	$hash(a)\%4 = P2$
102	b	XXY	2022-11-28	$hash(b)\%4 = P4$
103	c	XYZ	2022-11-29	$hash(c)\%4 = P3$
104	d	XYX	2022-11-27	$hash(d)\%4 = P2$
105	e	XYY	2022-11-29	$hash(e)\%4 = P1$

Partitions



Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

HORIZONTAL PARTITIONING

Partitioning Key

Table

101	a	XXX	2022-11-29	$\text{hash}(a)\%4 = P2$
102	b	XXV	2022-11-28	$\text{hash}(b)\%4 = P4$
103	c	XYZ	2022-11-29	$\text{hash}(c)\%4 = P3$
104	d	XYX	2022-11-27	$\text{hash}(d)\%4 = P2$
105	e	XYV	2022-11-29	$\text{hash}(e)\%4 = P1$

Partitions



Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

HORIZONTAL PARTITIONING

Partitioning Key

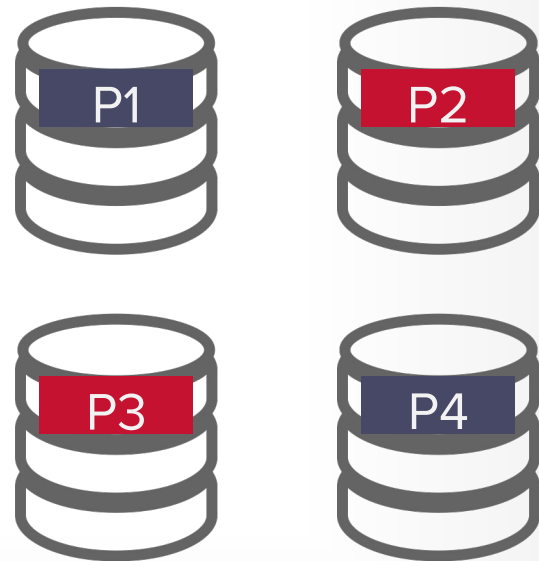
Table

101	a	XXX	2022-11-29	$\text{hash}(a)\%4 = P2$
102	b	XXY	2022-11-28	$\text{hash}(b)\%4 = P4$
103	c	XYZ	2022-11-29	$\text{hash}(c)\%4 = P3$
104	d	XYX	2022-11-27	$\text{hash}(d)\%4 = P2$
105	e	XYY	2022-11-29	$\text{hash}(e)\%4 = P1$

Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

Partitions



HORIZONTAL PARTITIONING

Partitioning Key

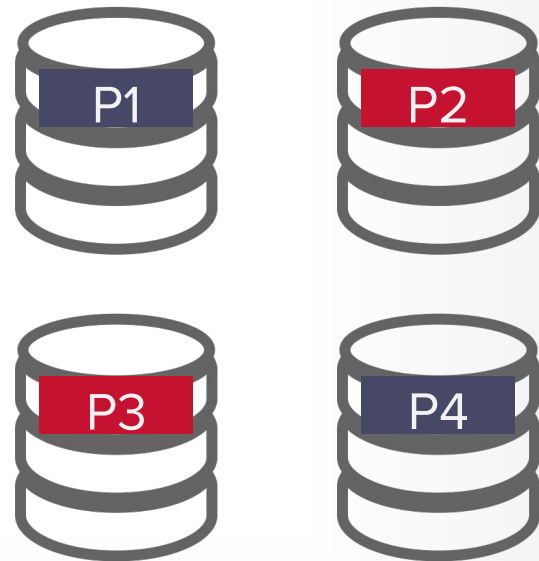
Table

101	a	XXX	2022-11-29	$hash(a)\%4 = P2$
102	b	XXY	2022-11-28	$hash(b)\%4 = P4$
103	c	XYZ	2022-11-29	$hash(c)\%4 = P3$
104	d	XYX	2022-11-27	$hash(d)\%4 = P2$
105	e	XYY	2022-11-29	$hash(e)\%4 = P1$

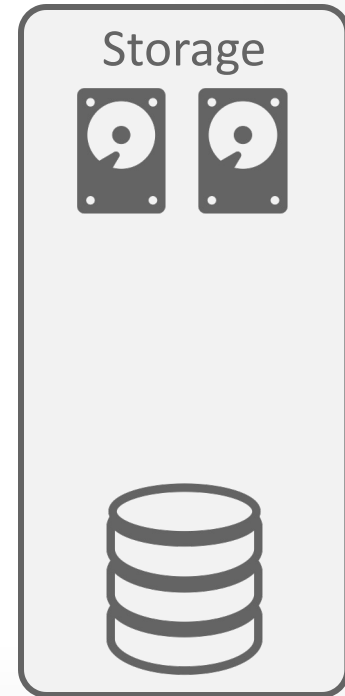
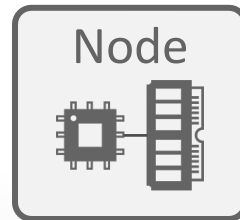
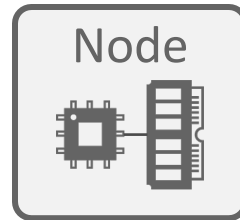
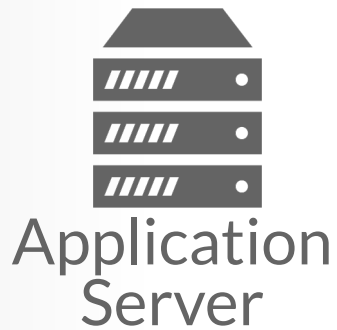
Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

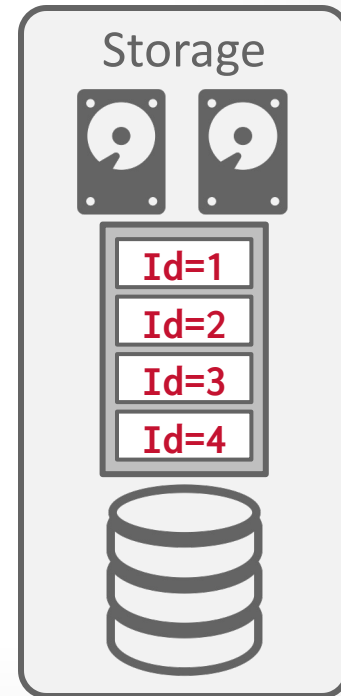
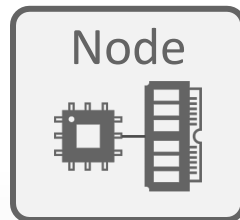
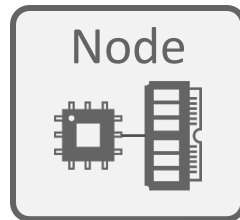
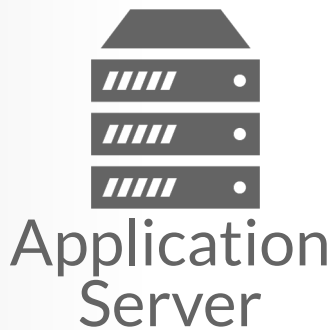
Partitions



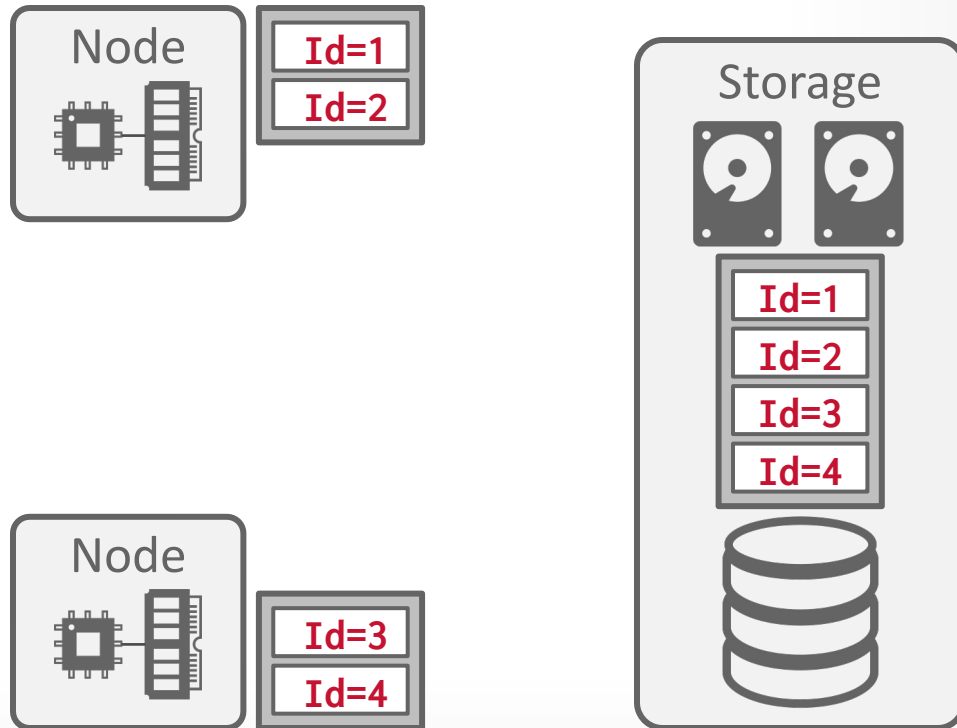
SHARED-DISK PARTITIONING



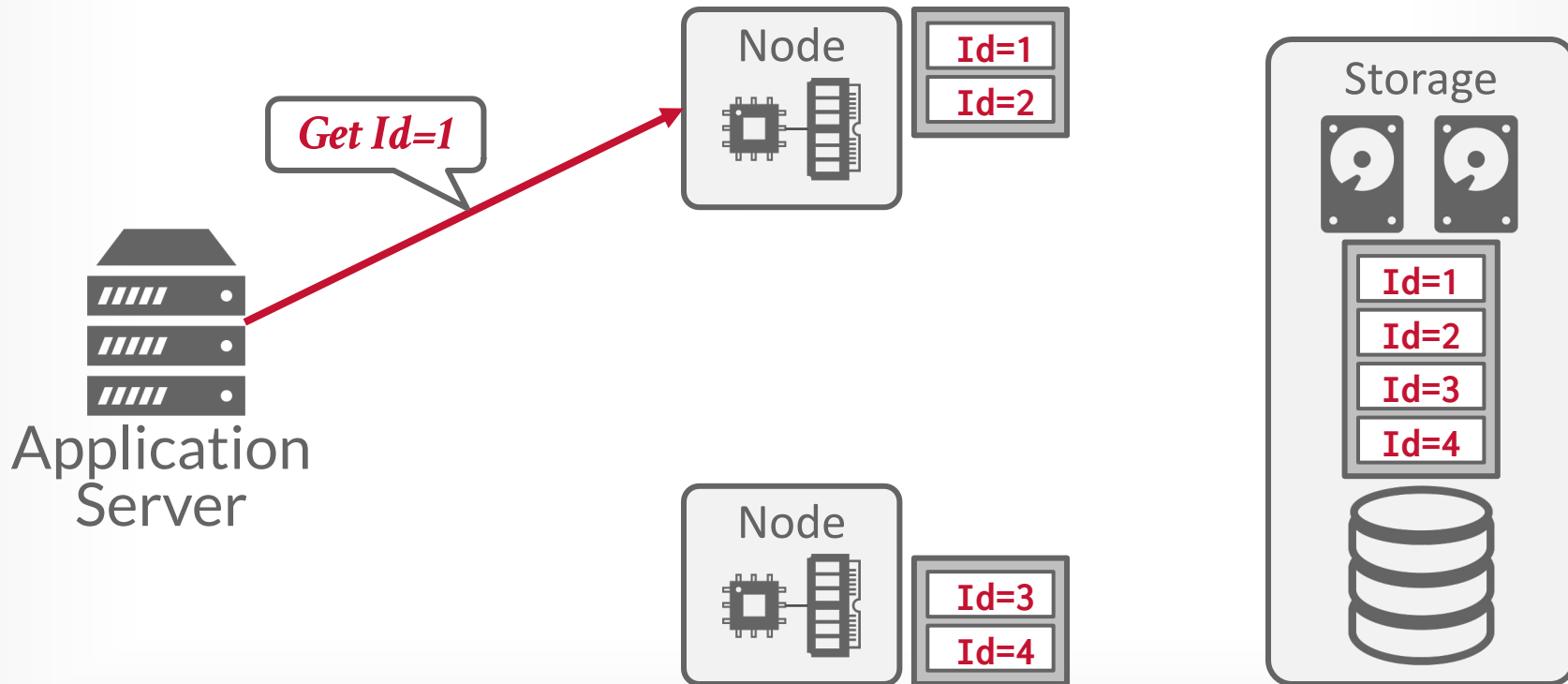
SHARED-DISK PARTITIONING



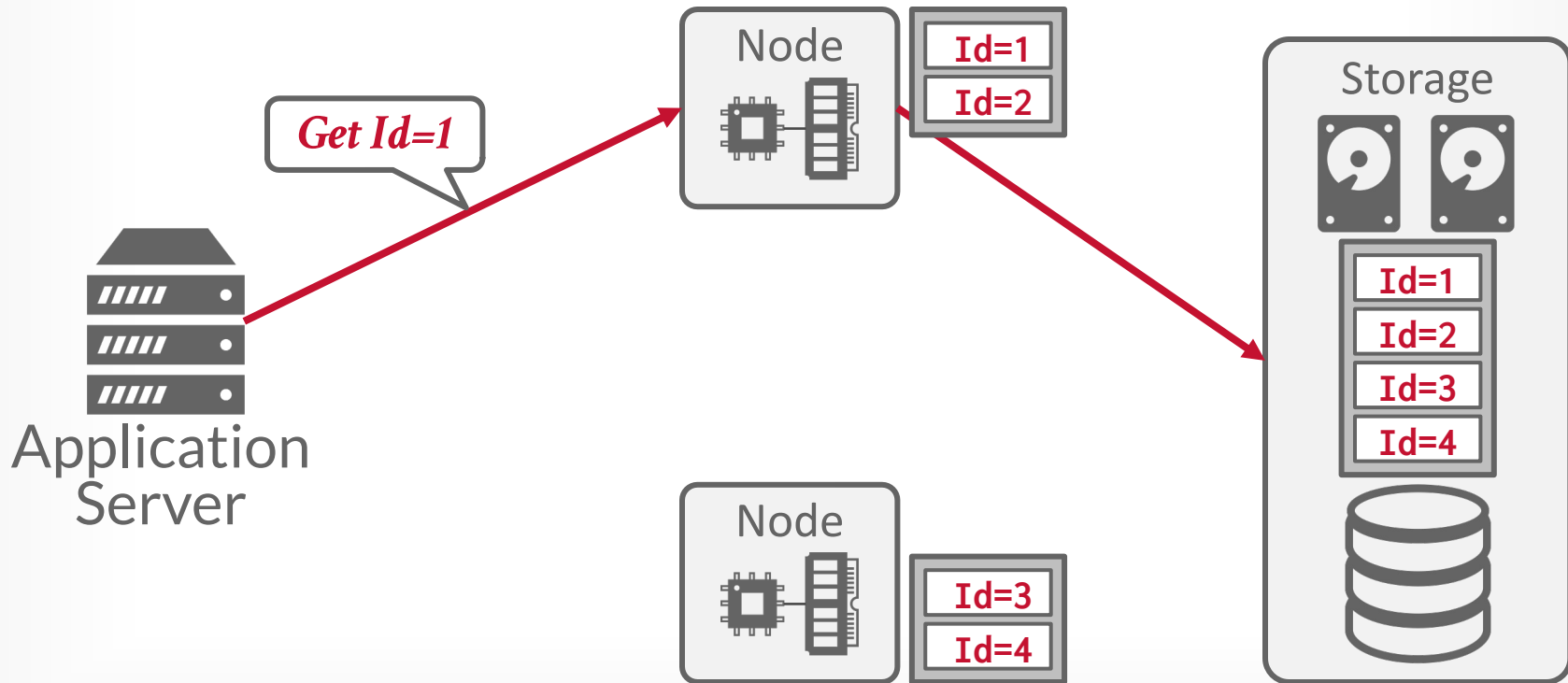
SHARED-DISK PARTITIONING



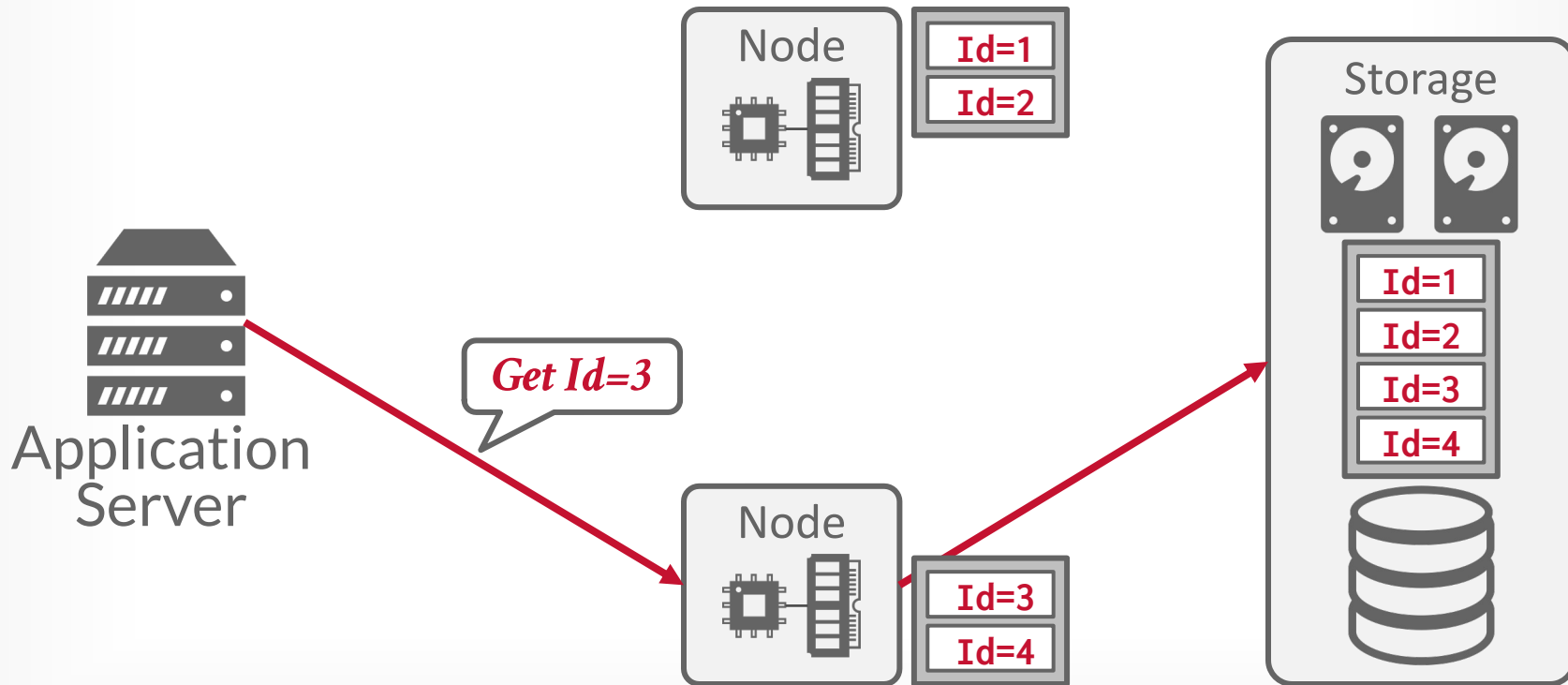
SHARED-DISK PARTITIONING



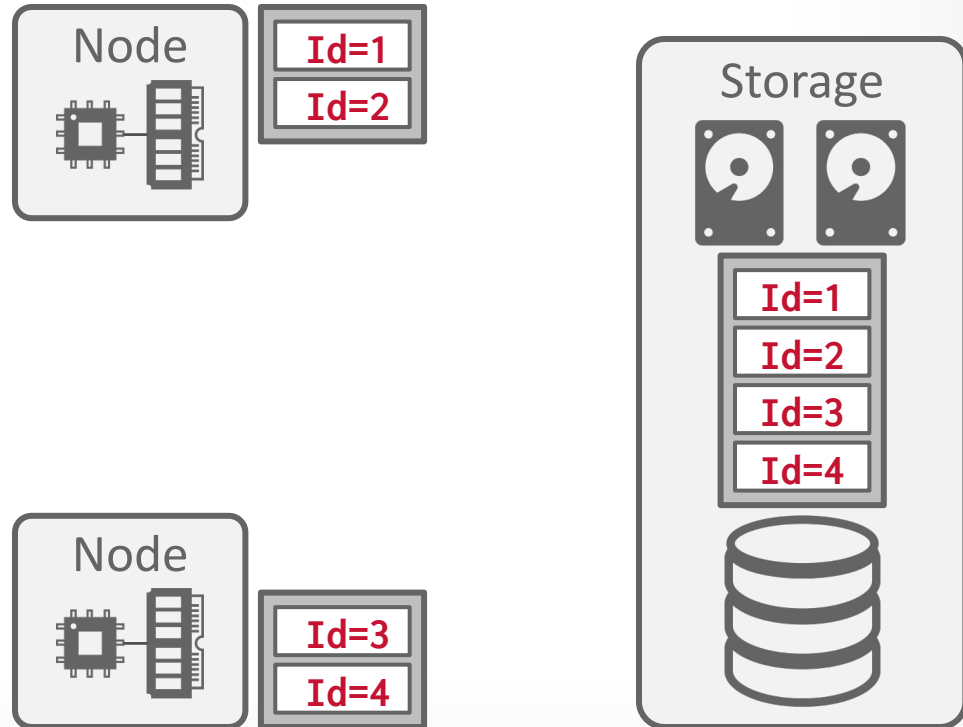
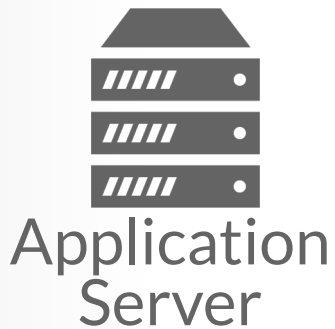
SHARED-DISK PARTITIONING



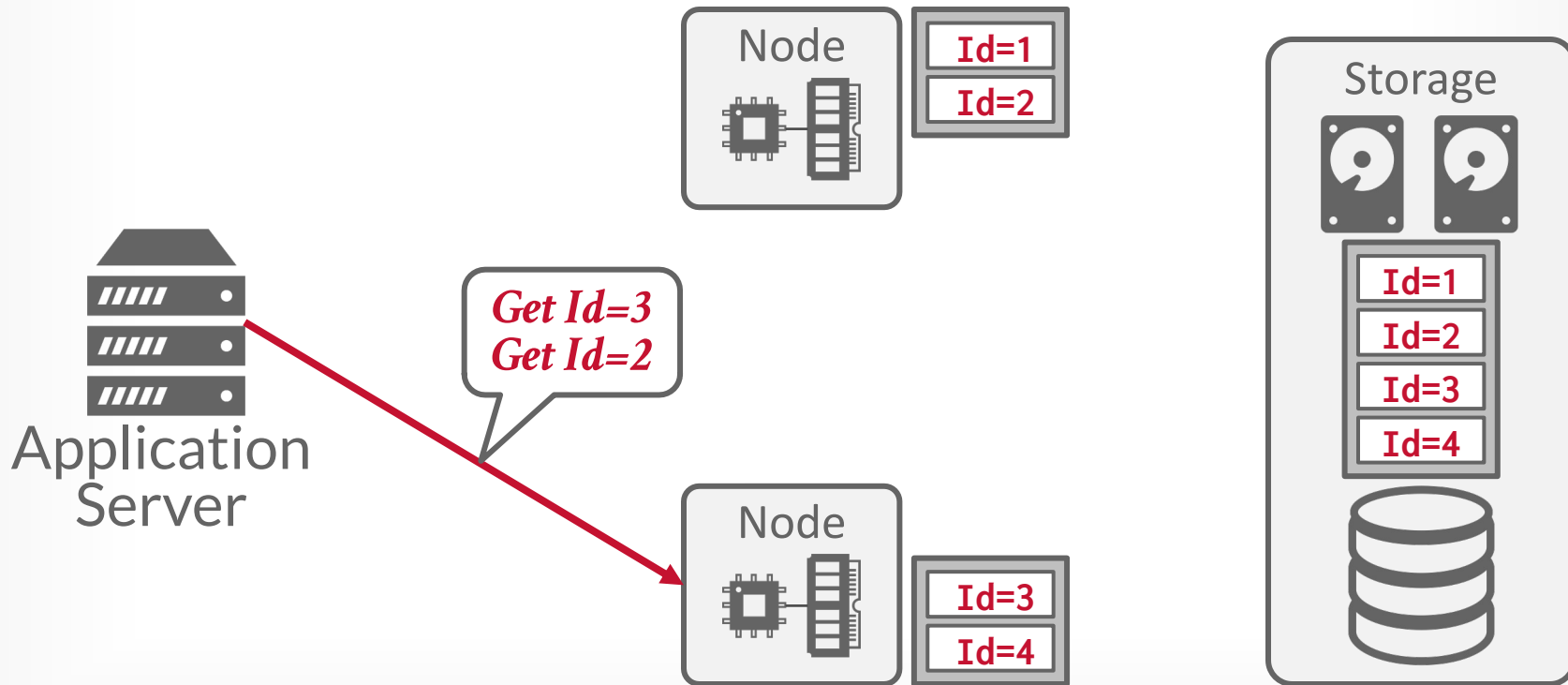
SHARED-DISK PARTITIONING



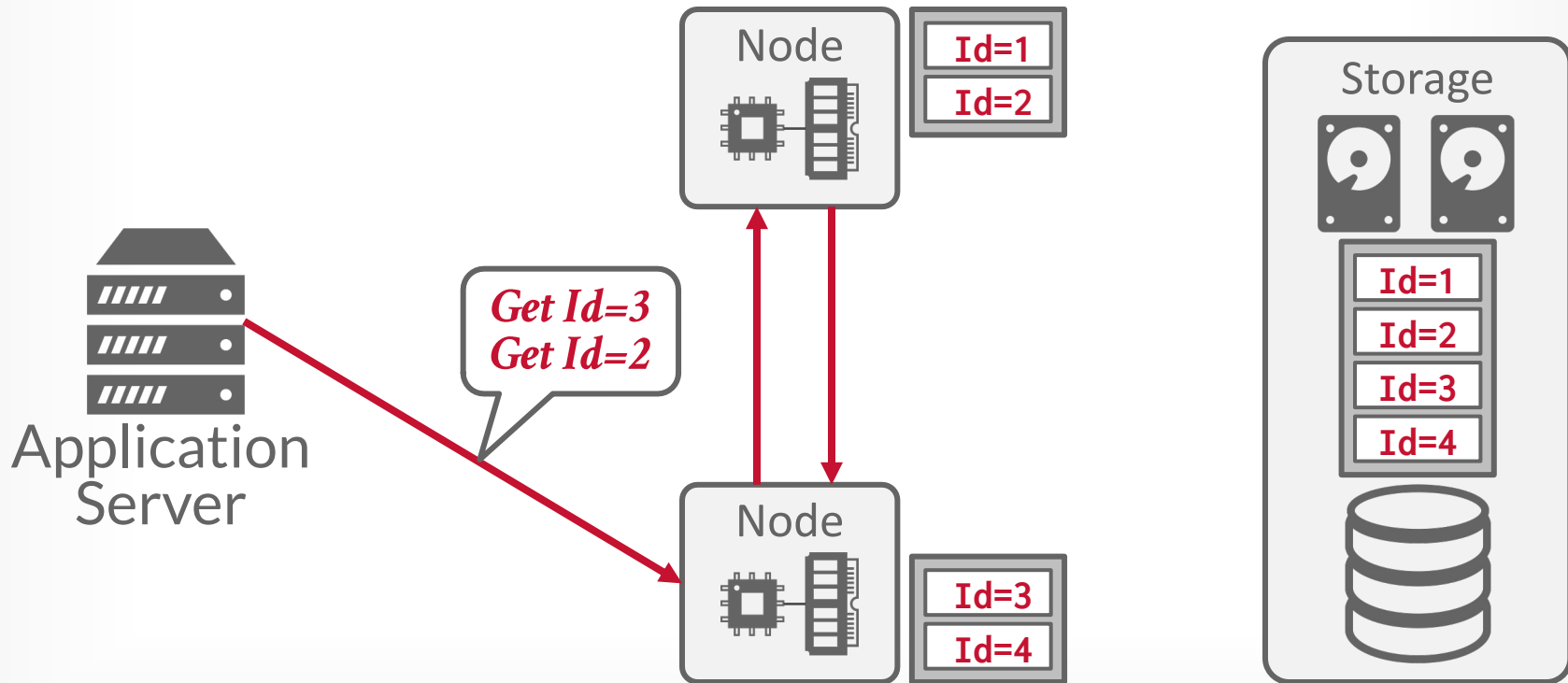
SHARED-DISK PARTITIONING



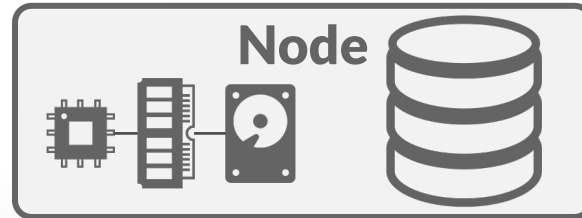
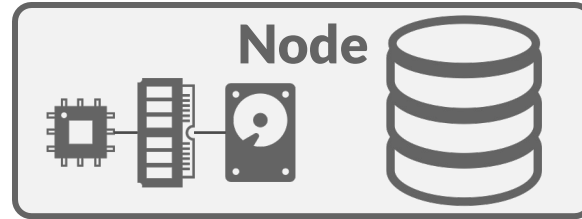
SHARED-DISK PARTITIONING



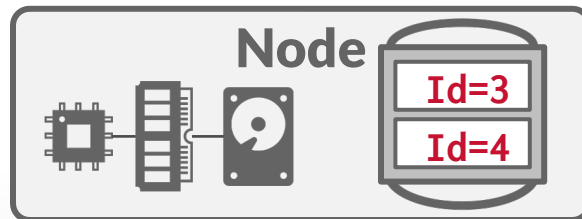
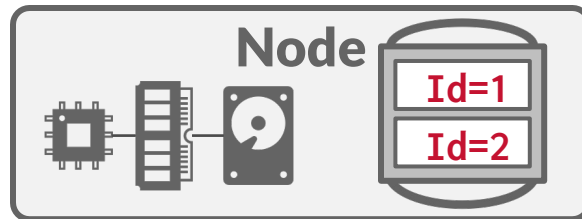
SHARED-DISK PARTITIONING



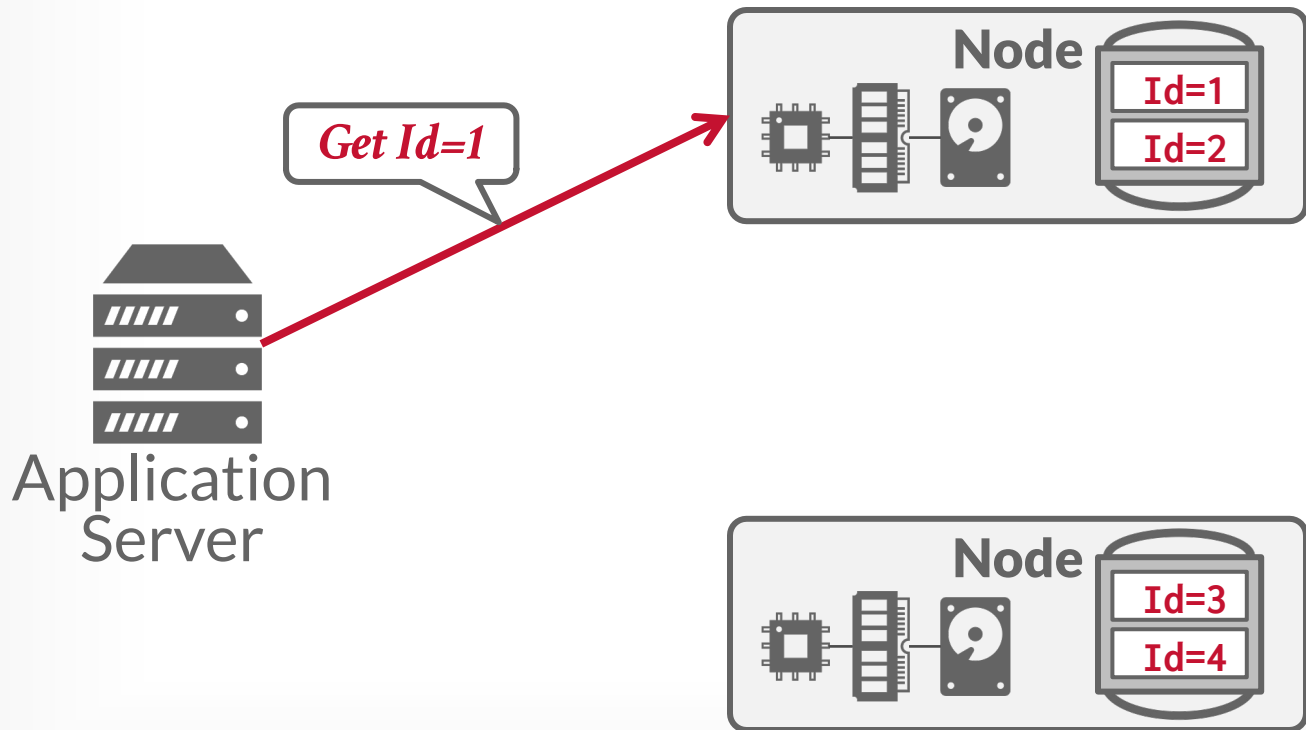
SHARED-NOTHING PARTITIONING



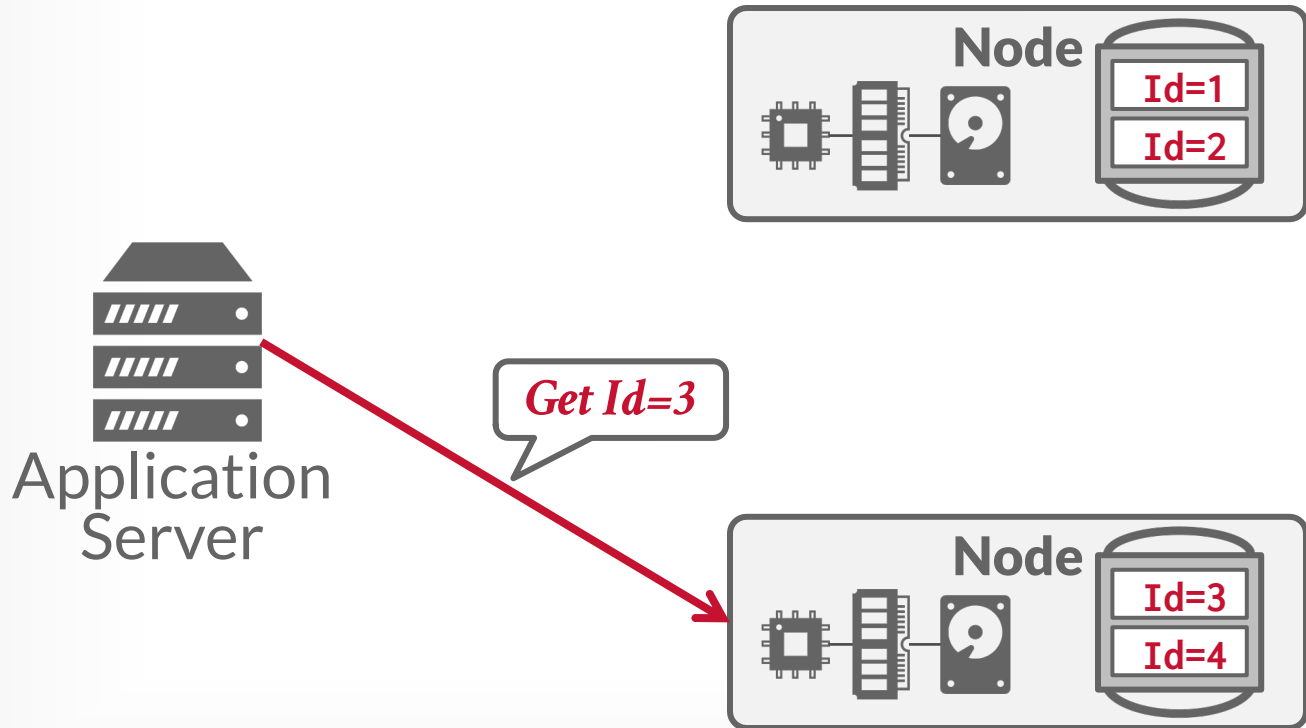
SHARED-NOTHING PARTITIONING



SHARED-NOTHING PARTITIONING



SHARED-NOTHING PARTITIONING



HORIZONTAL PARTITIONING

Partitioning Key

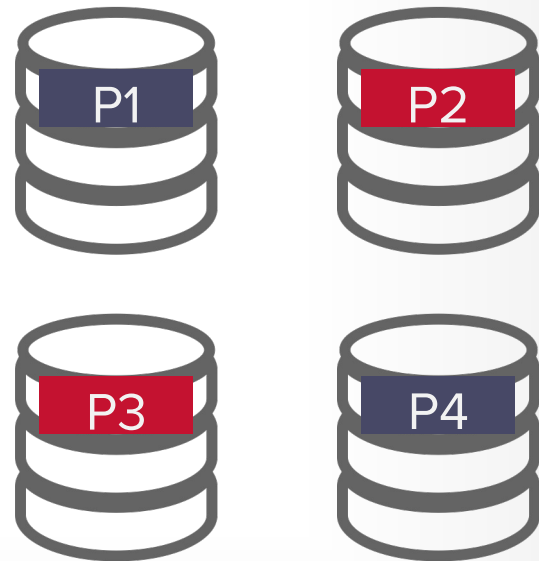
Table

101	a	XXX	2022-11-29	$\text{hash}(a)\%4 = P2$
102	b	XXY	2022-11-28	$\text{hash}(b)\%4 = P4$
103	c	XYZ	2022-11-29	$\text{hash}(c)\%4 = P3$
104	d	XYX	2022-11-27	$\text{hash}(d)\%4 = P2$
105	e	XYY	2022-11-29	$\text{hash}(e)\%4 = P1$

Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

Partitions



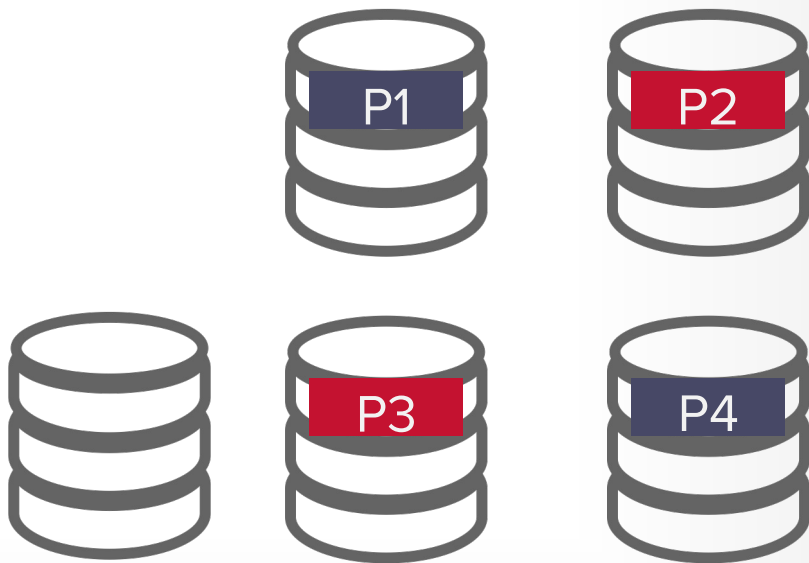
HORIZONTAL PARTITIONING

Partitioning Key

Table

101	a	XXX	2022-11-29	$hash(a)\%4 = P2$
102	b	XXY	2022-11-28	$hash(b)\%4 = P4$
103	c	XYZ	2022-11-29	$hash(c)\%4 = P3$
104	d	XYX	2022-11-27	$hash(d)\%4 = P2$
105	e	XYY	2022-11-29	$hash(e)\%4 = P1$

Partitions



Ideal Query:

```
SELECT * FROM table
WHERE partitionKey = ?
```

HORIZONTAL PARTITIONING

Partitioning Key

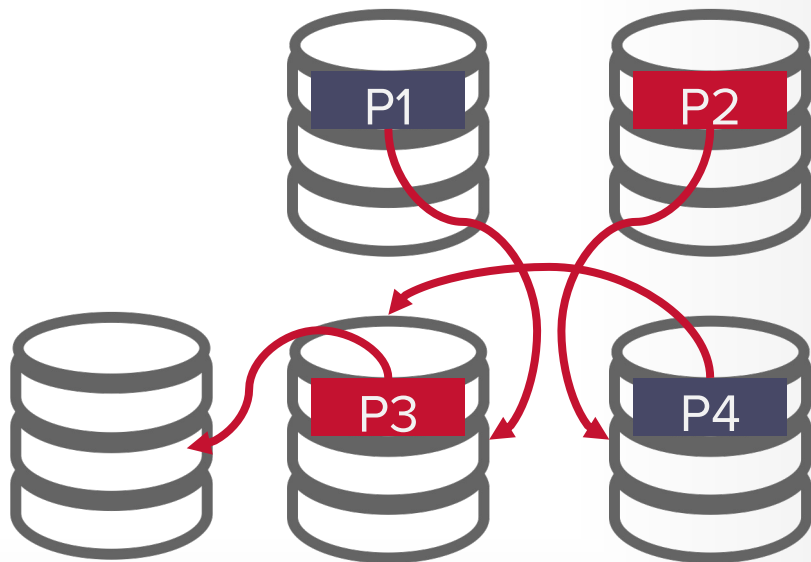
Table

101	a	XXX	2022-11-29	$\text{hash}(a)\%5 = P4$
102	b	XXY	2022-11-28	$\text{hash}(b)\%5 = P3$
103	c	XYZ	2022-11-29	$\text{hash}(c)\%5 = P5$
104	d	XYX	2022-11-27	$\text{hash}(d)\%5 = P1$
105	e	XYY	2022-11-29	$\text{hash}(e)\%5 = P3$

Ideal Query:

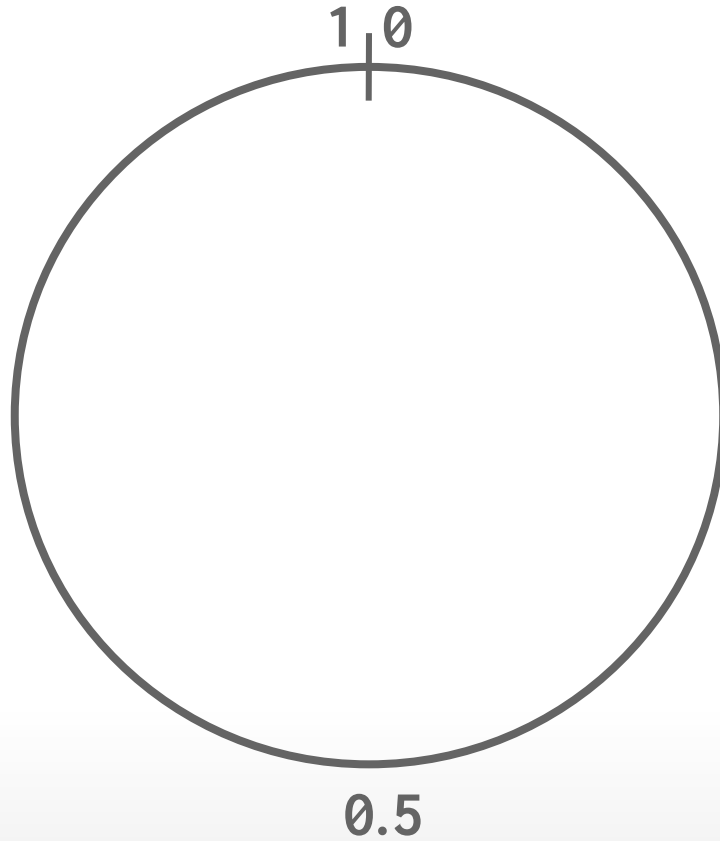
```
SELECT * FROM table
WHERE partitionKey = ?
```

Partitions

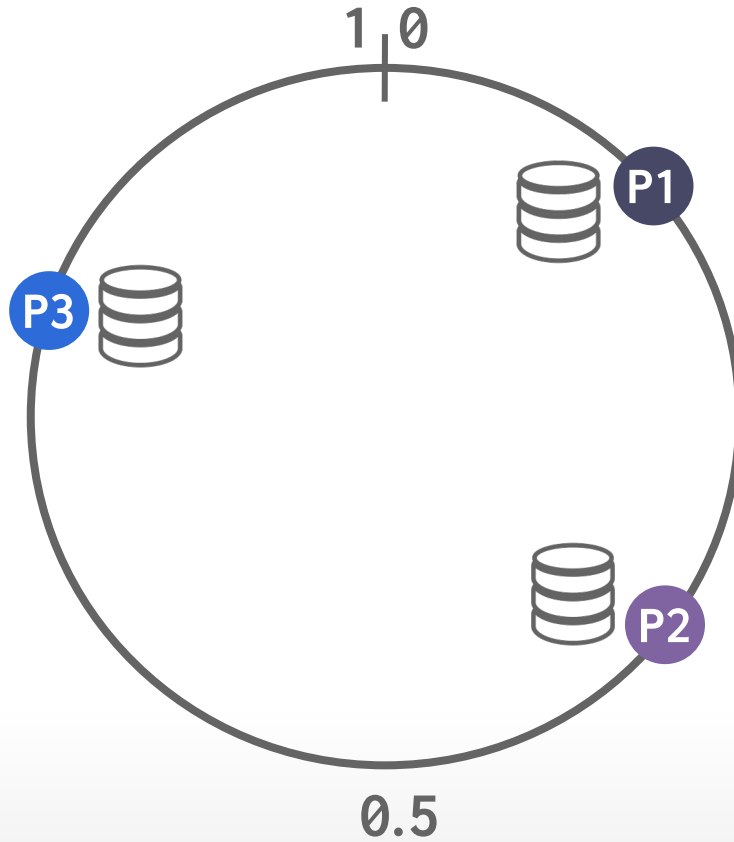


CONSISTENT HASHING

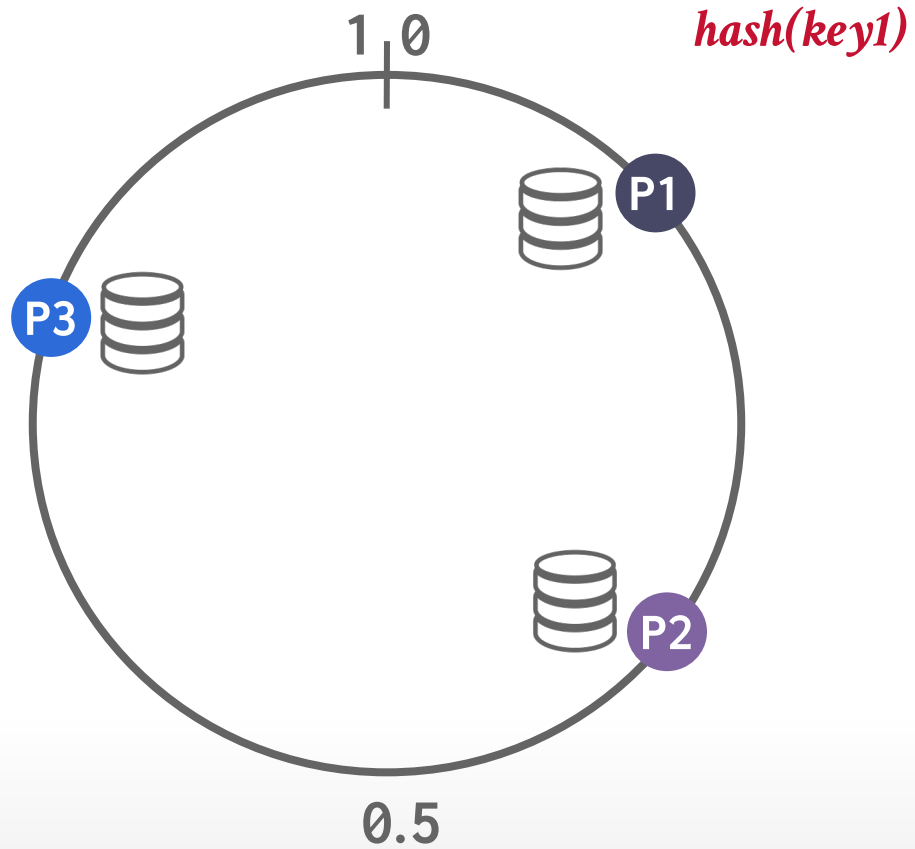
CONSISTENT HASHING



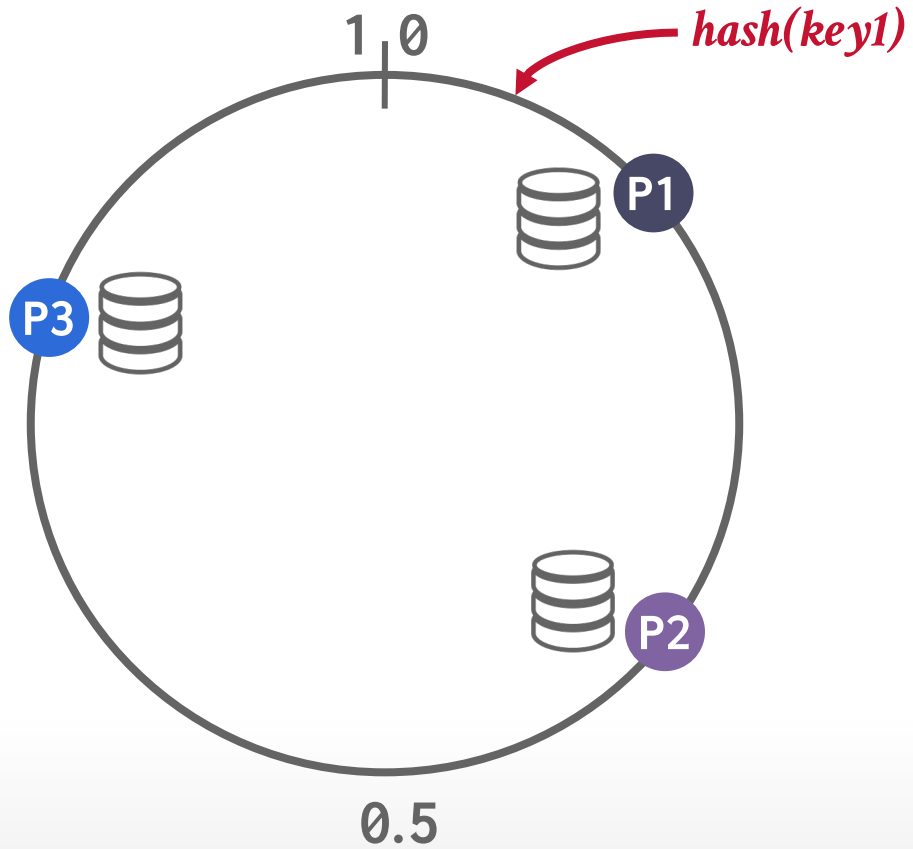
CONSISTENT HASHING



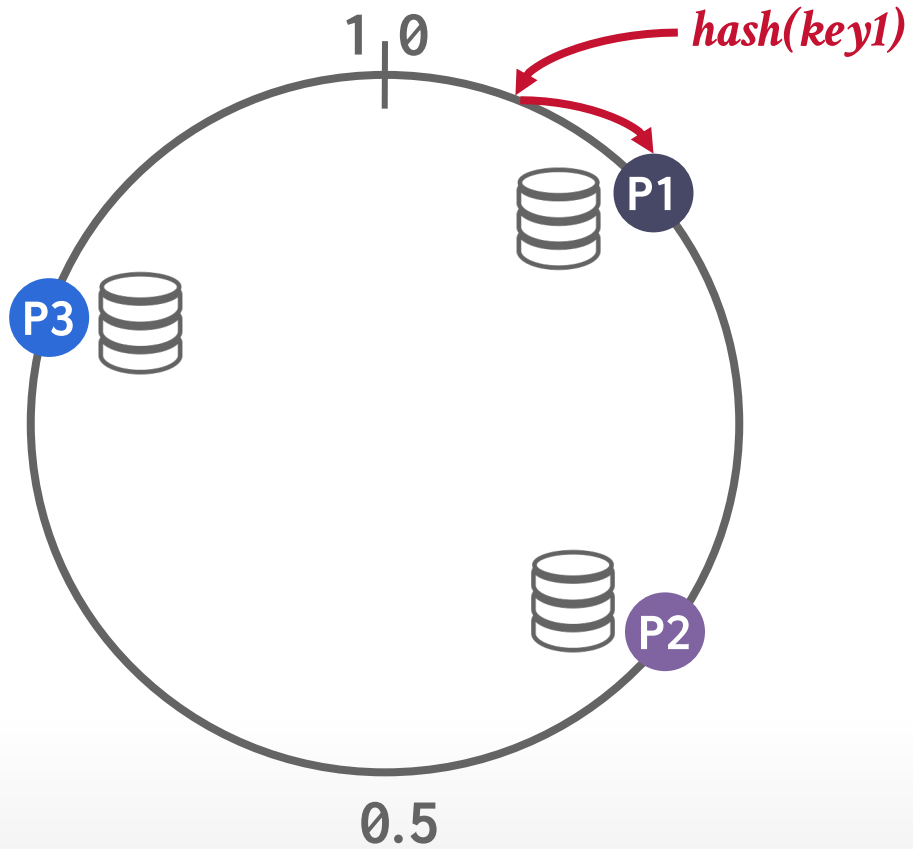
CONSISTENT HASHING



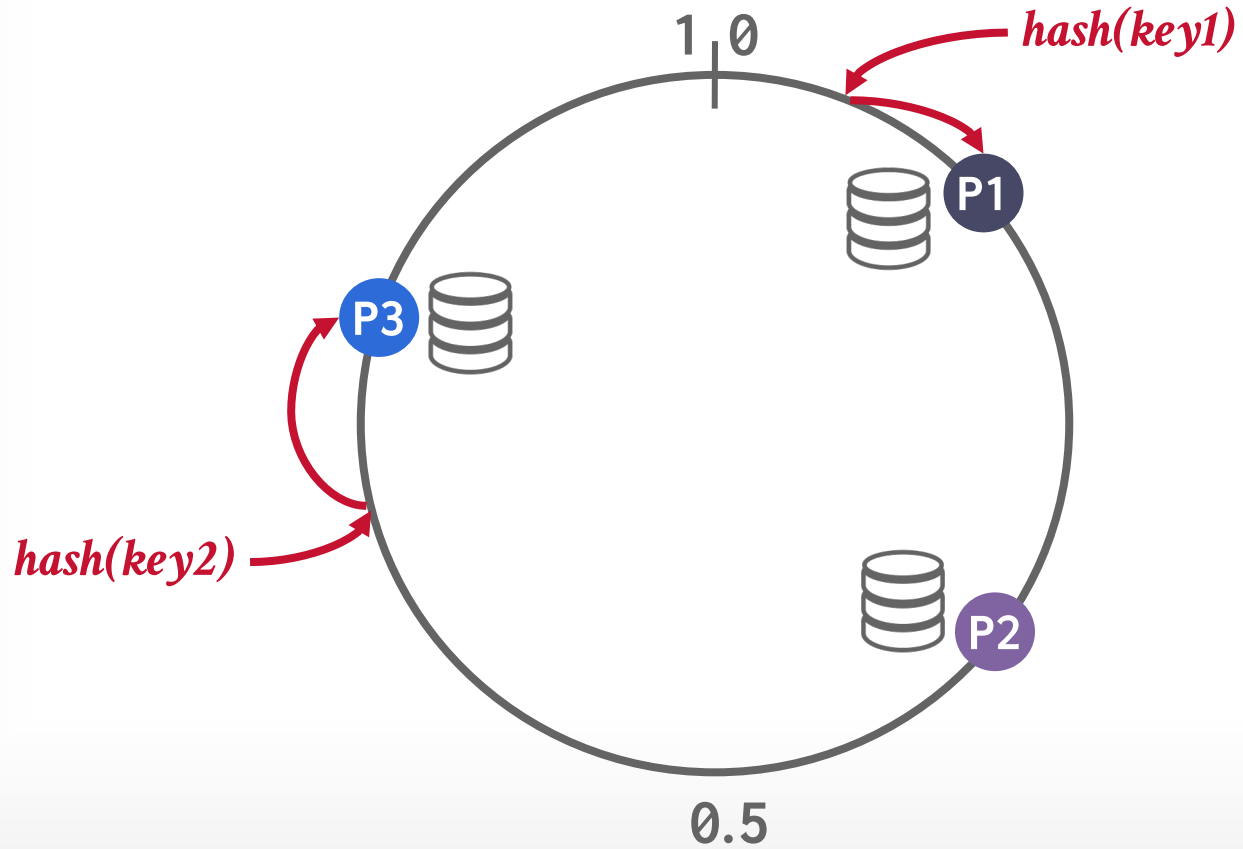
CONSISTENT HASHING



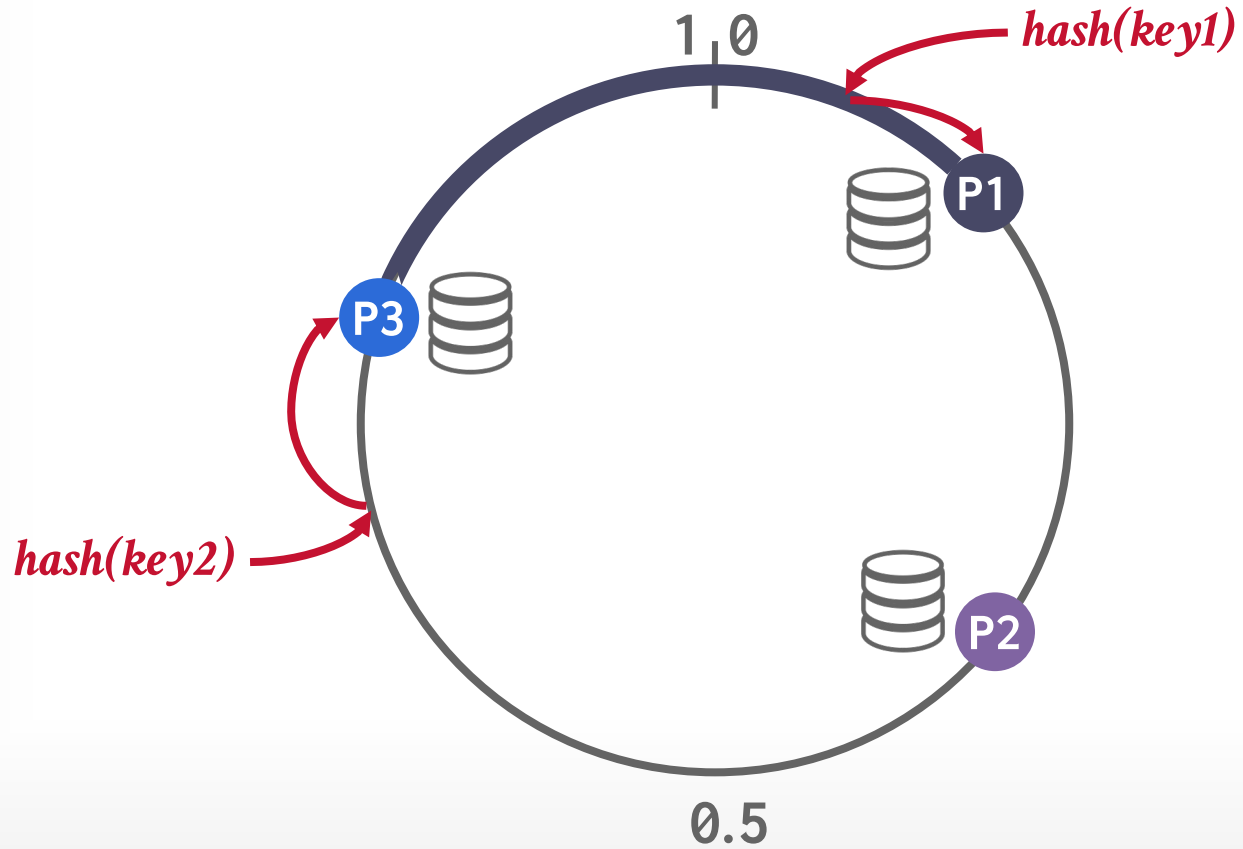
CONSISTENT HASHING



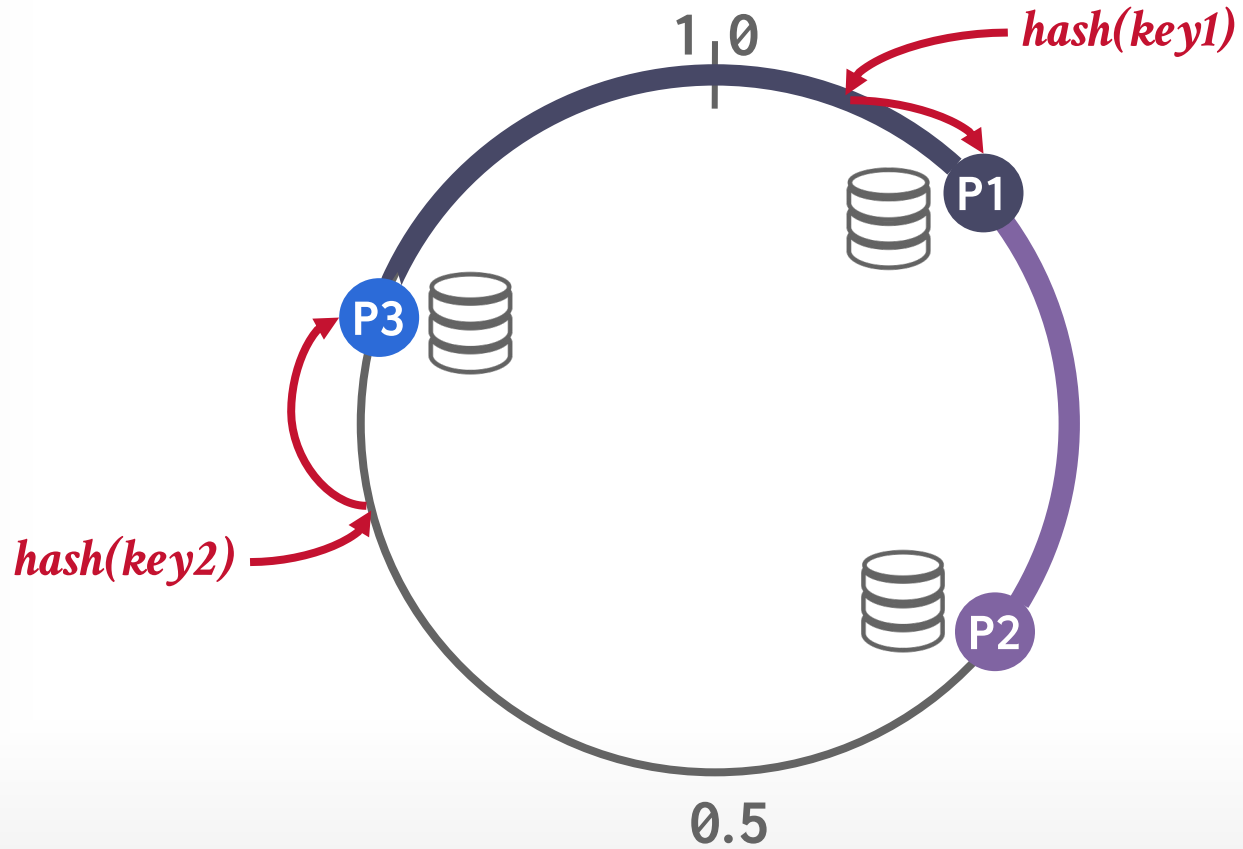
CONSISTENT HASHING



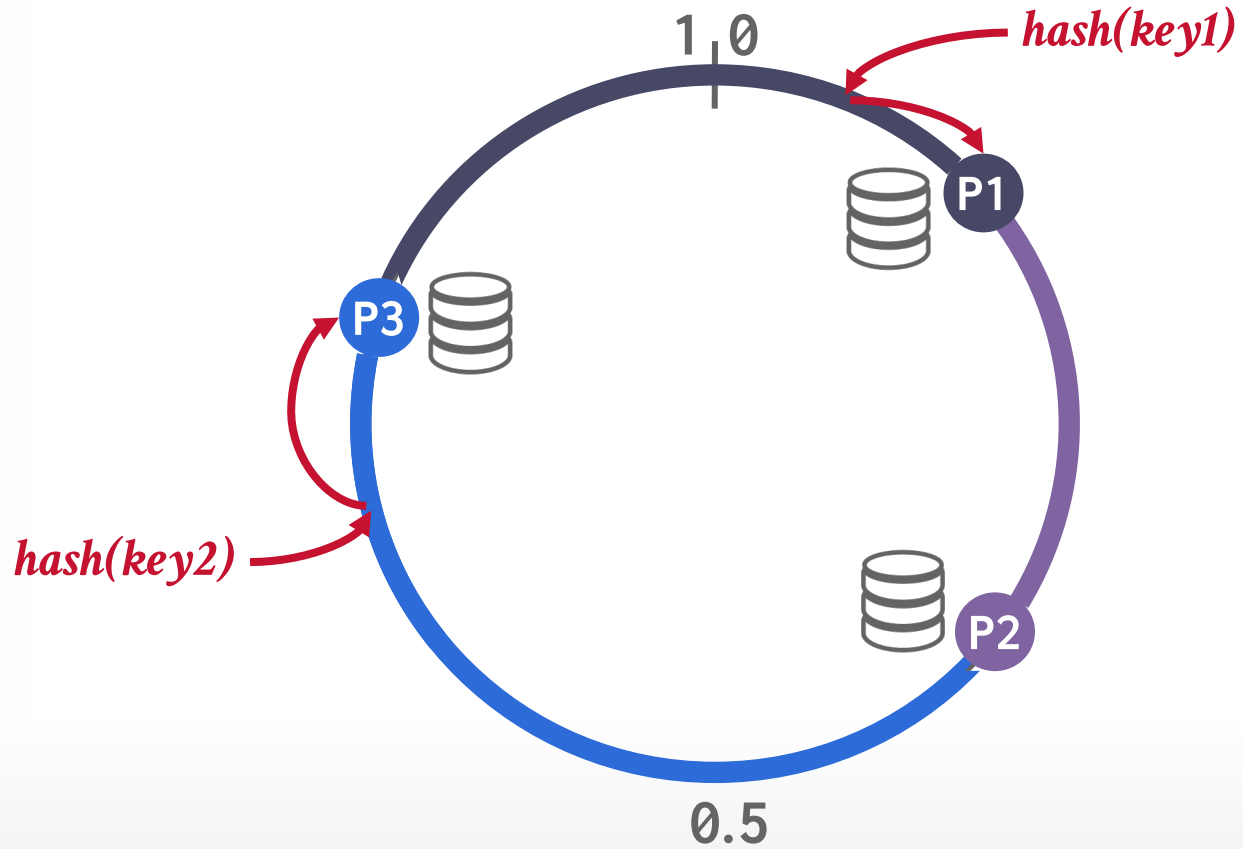
CONSISTENT HASHING



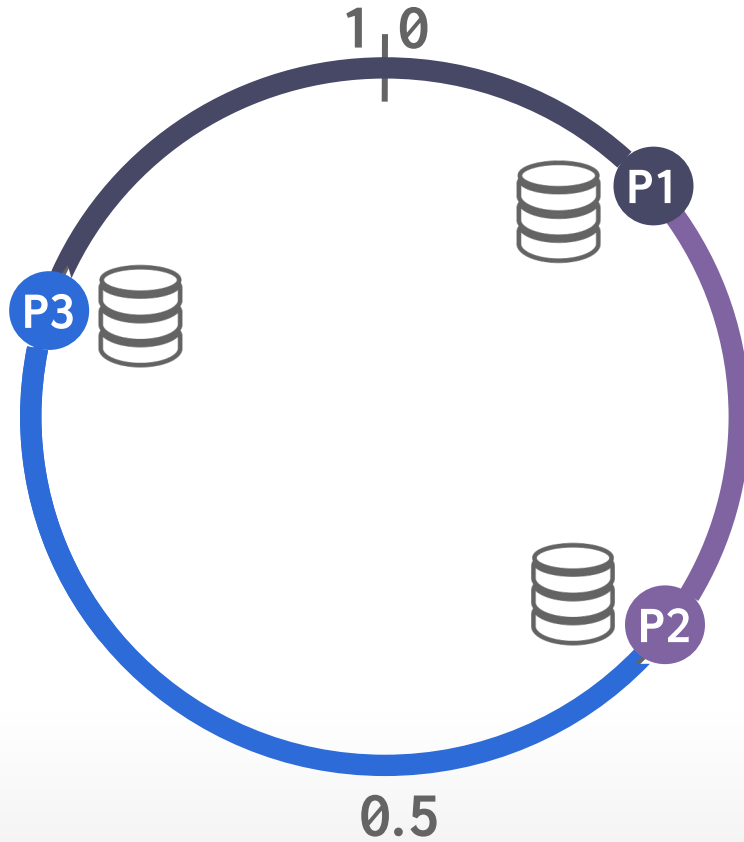
CONSISTENT HASHING



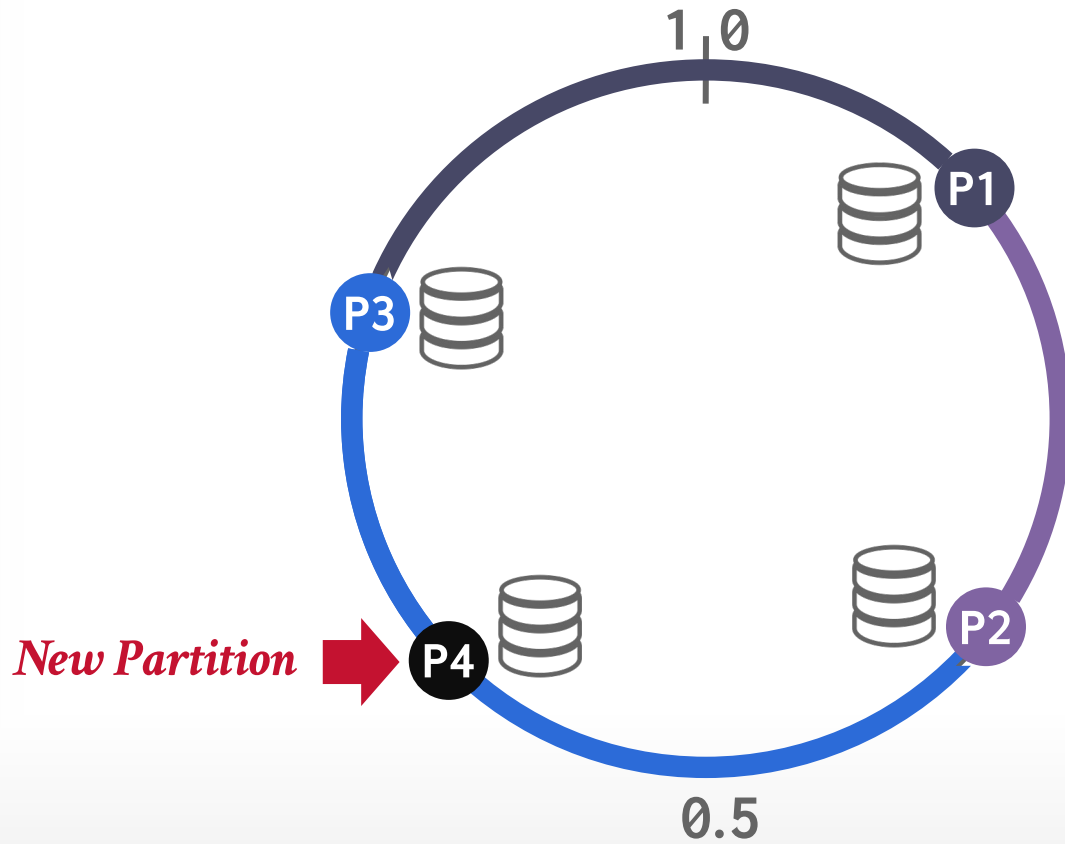
CONSISTENT HASHING



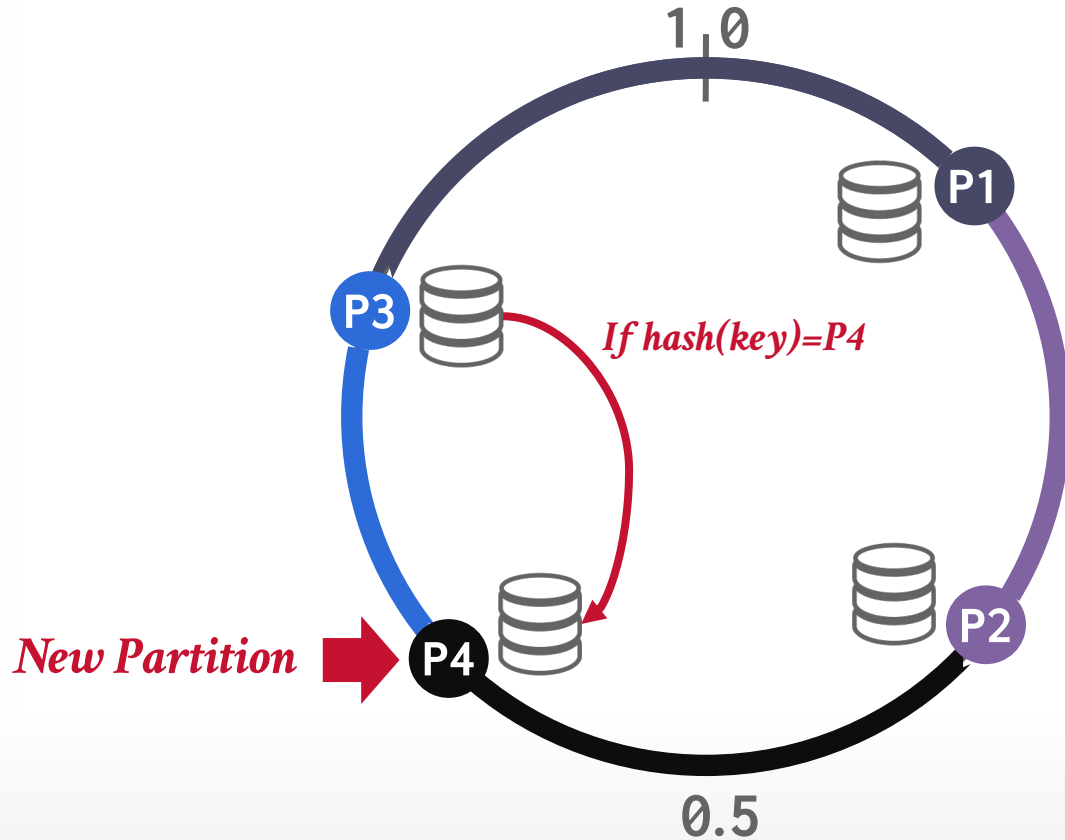
CONSISTENT HASHING



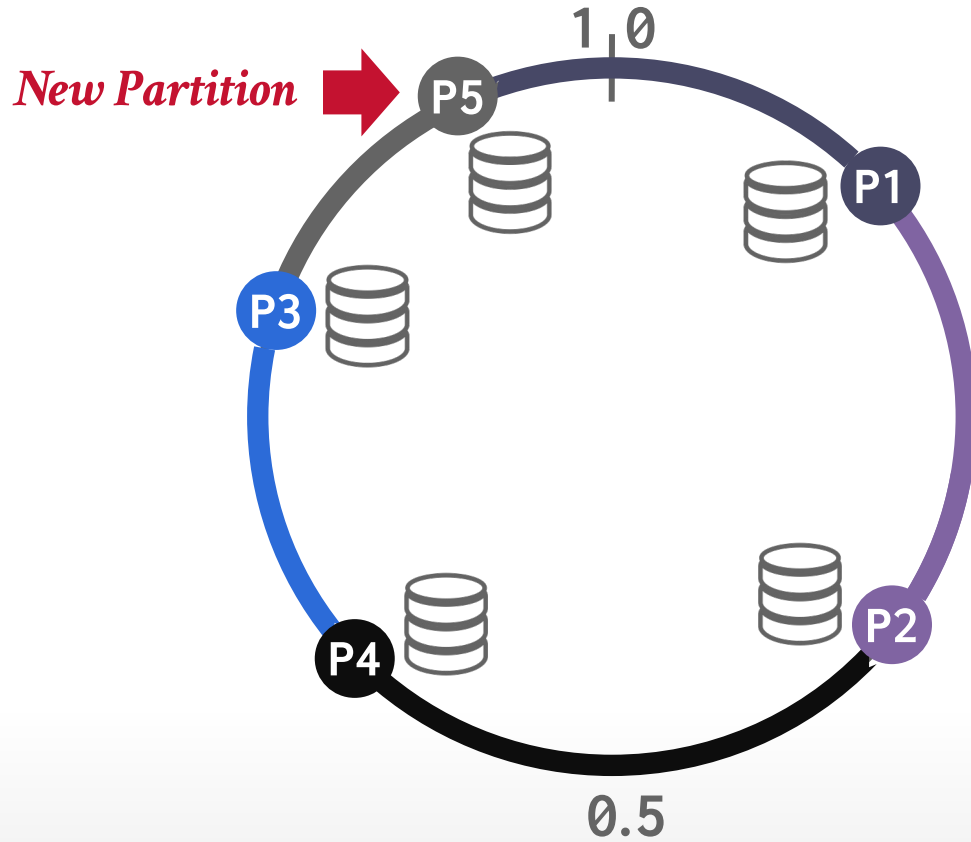
CONSISTENT HASHING



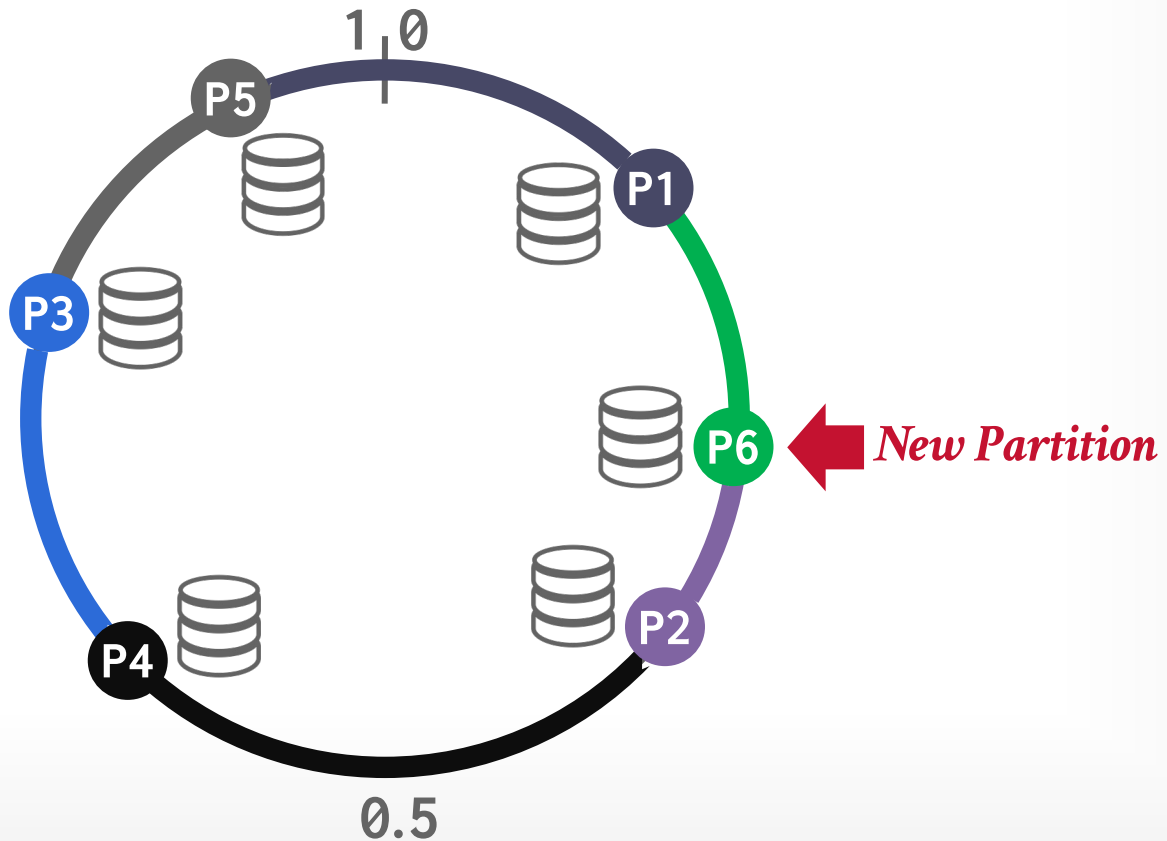
CONSISTENT HASHING



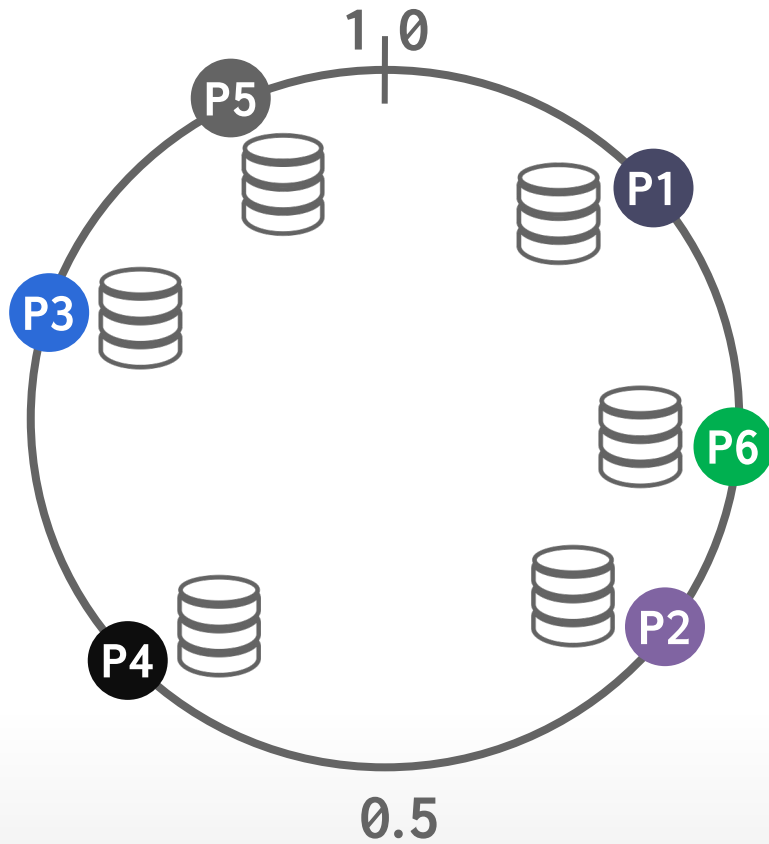
CONSISTENT HASHING



CONSISTENT HASHING

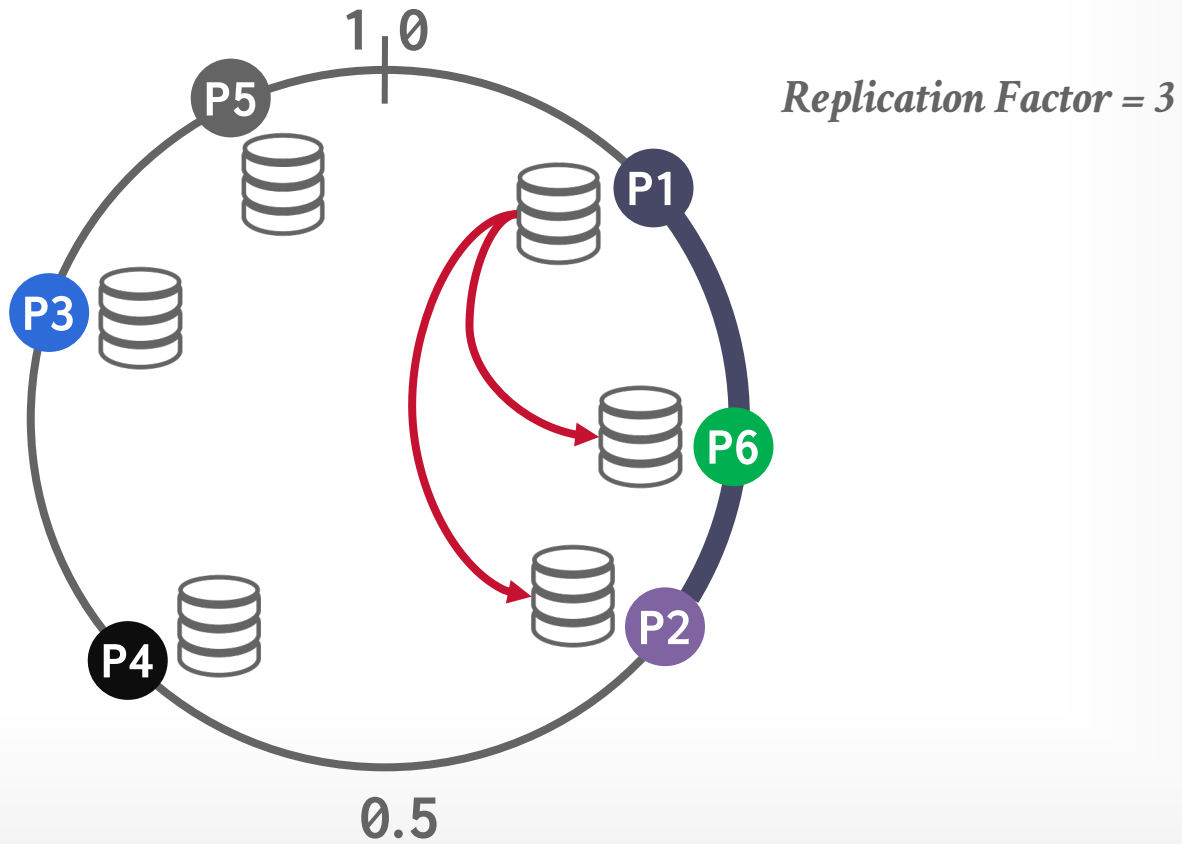


CONSISTENT HASHING

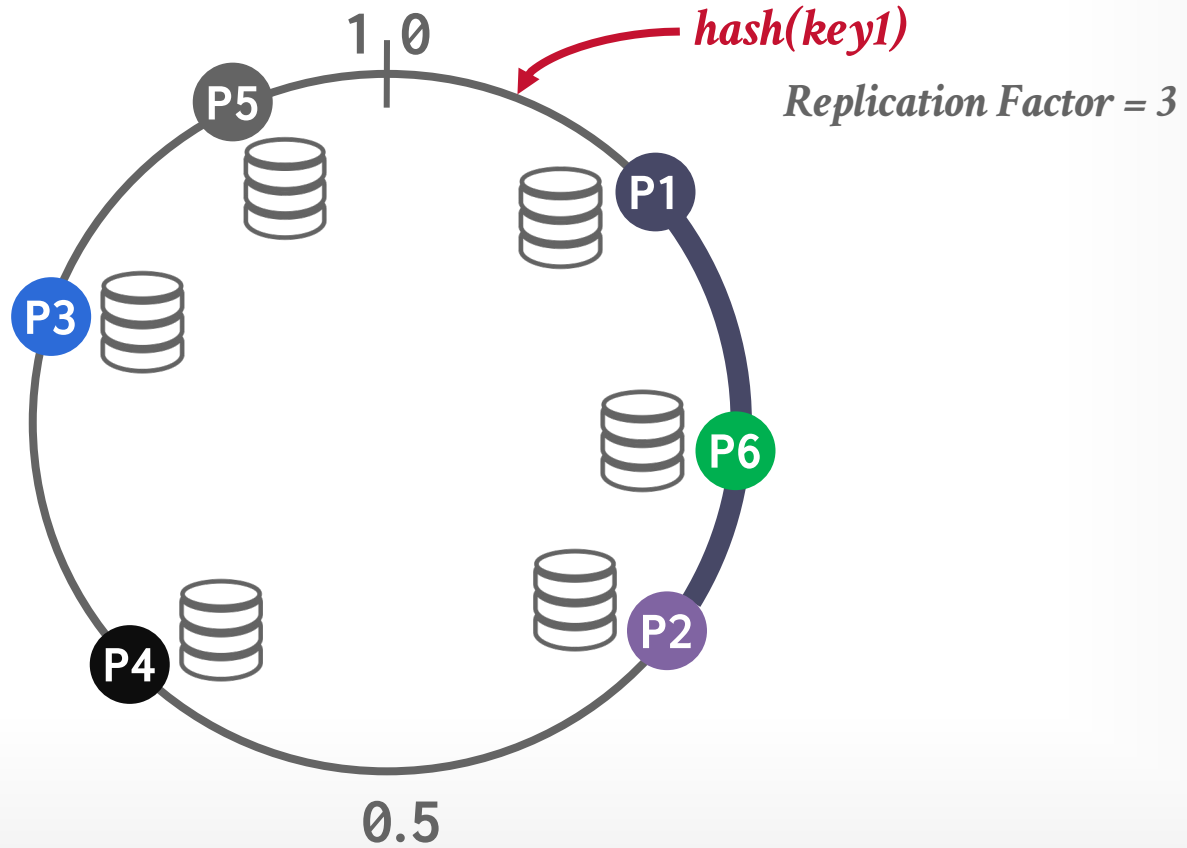


Replication Factor = 3

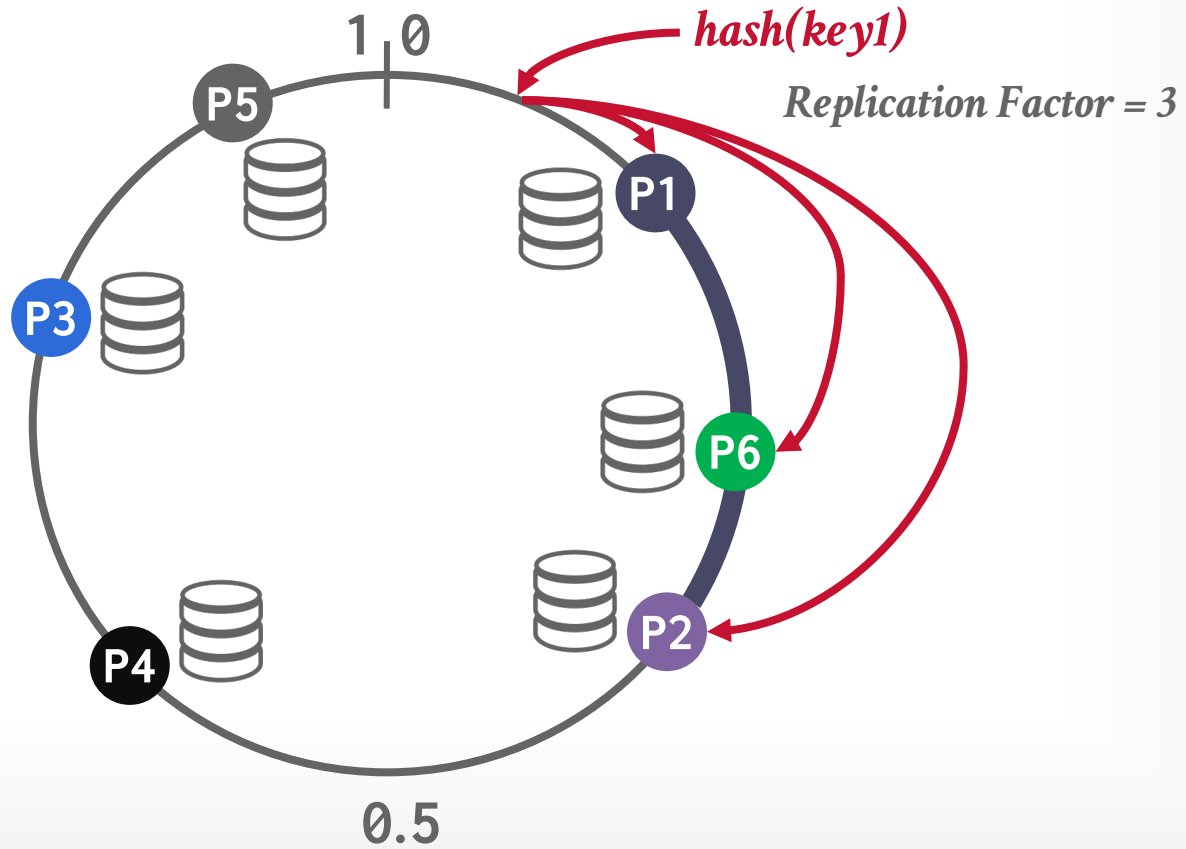
CONSISTENT HASHING



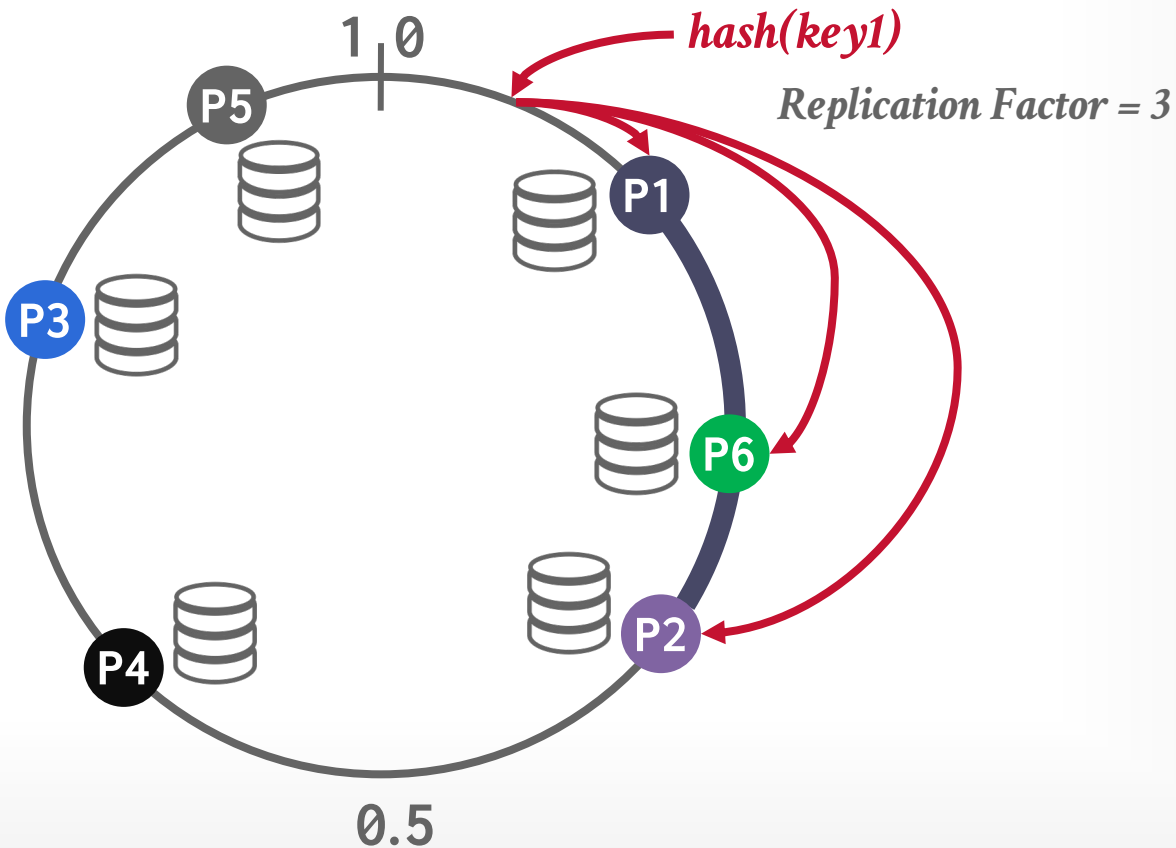
CONSISTENT HASHING



CONSISTENT HASHING



CONSISTENT HASHING



SINGLE-NODE VS. DISTRIBUTED

A single-node txn only accesses data that is contained on one partition.

→ The DBMS may not need check the behavior concurrent txns running on other nodes.

A distributed txn accesses data at one or more partitions.

→ Requires expensive coordination.

TRANSACTION COORDINATION

If our DBMS supports multi-operation and distributed txns, we need a way to coordinate their execution in the system.

Two different approaches:

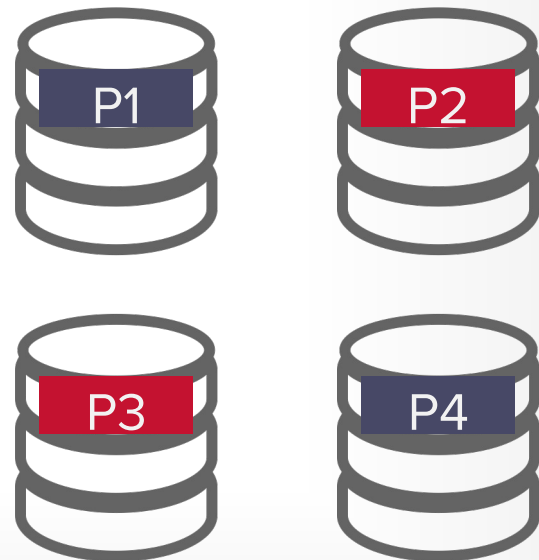
- **Centralized:** Global "traffic cop".
- **Decentralized:** Nodes organize themselves.

Most distributed DBMSs use a hybrid approach where they periodically elect some node to be a temporary coordinator.

CENTRALIZED COORDINATOR

Coordinator

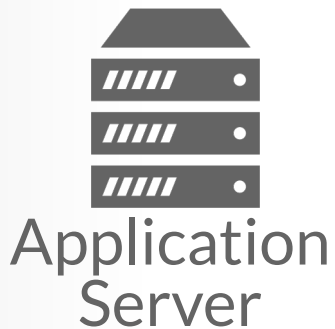
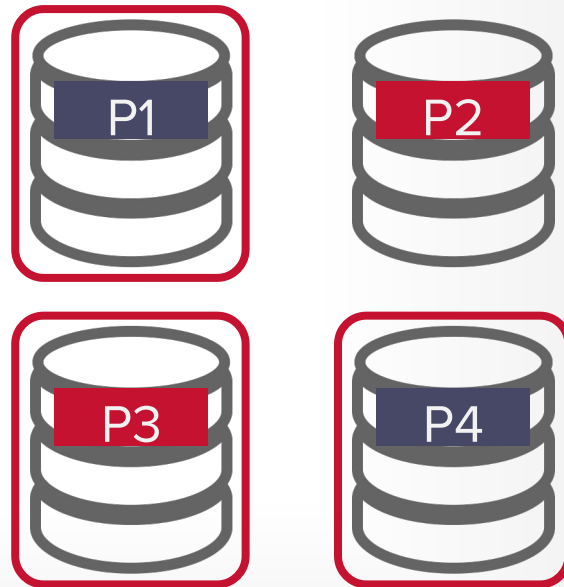
Partitions



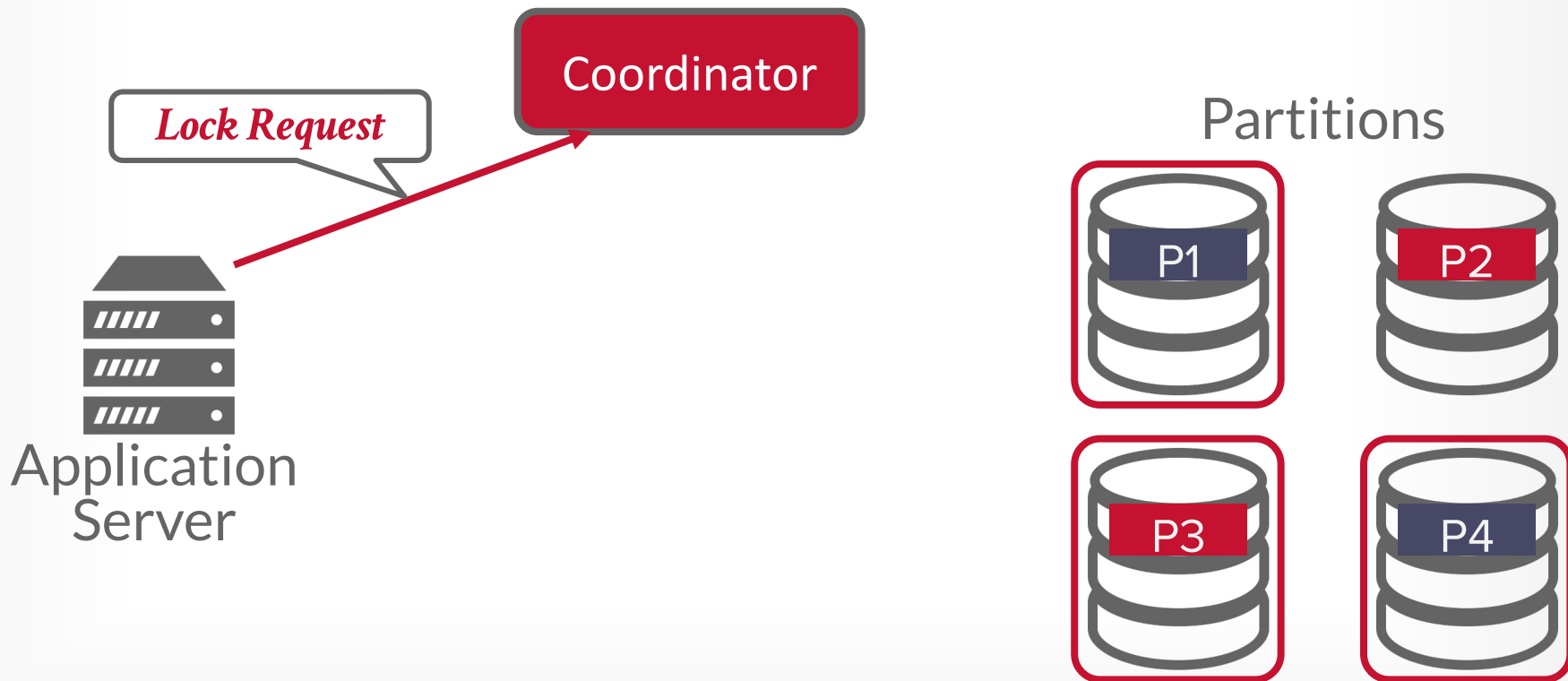
CENTRALIZED COORDINATOR

Coordinator

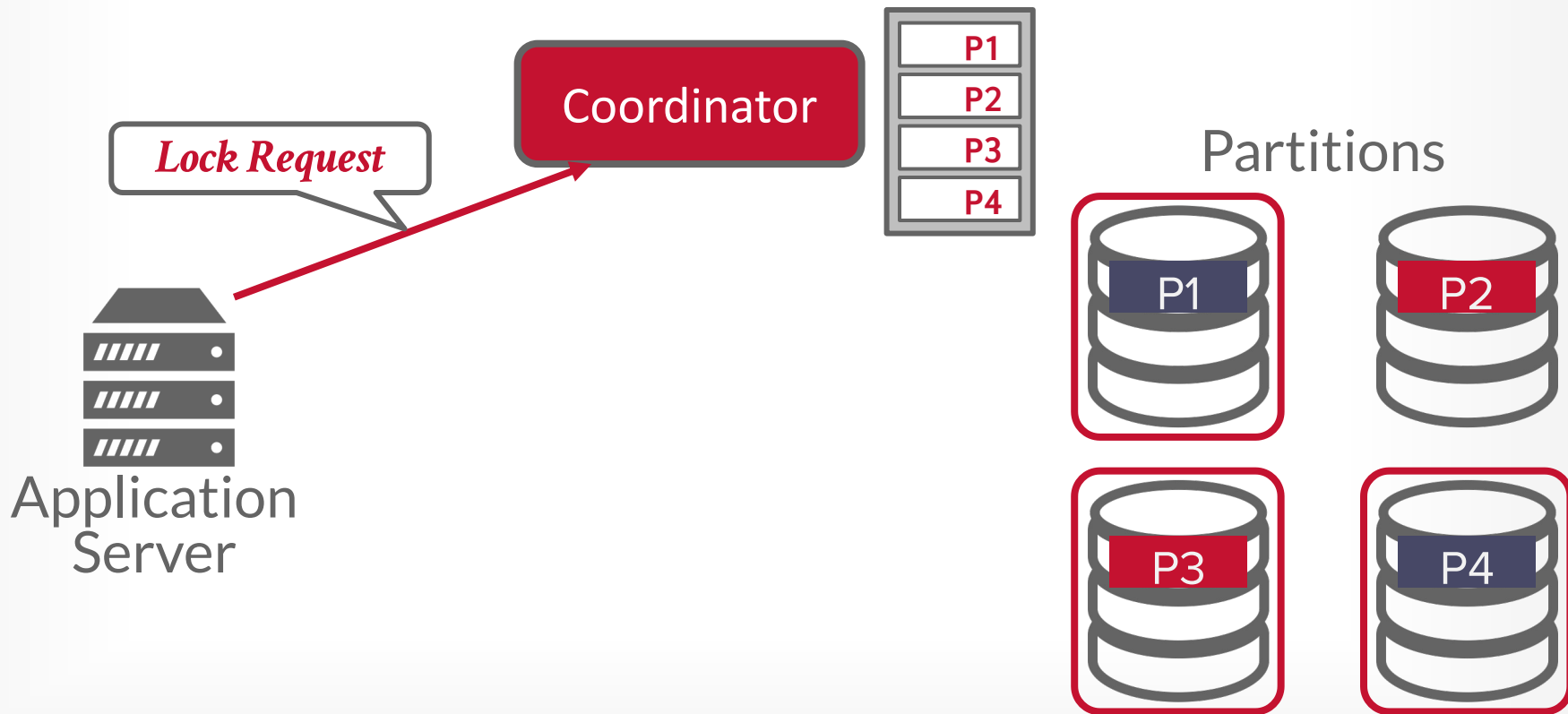
Partitions



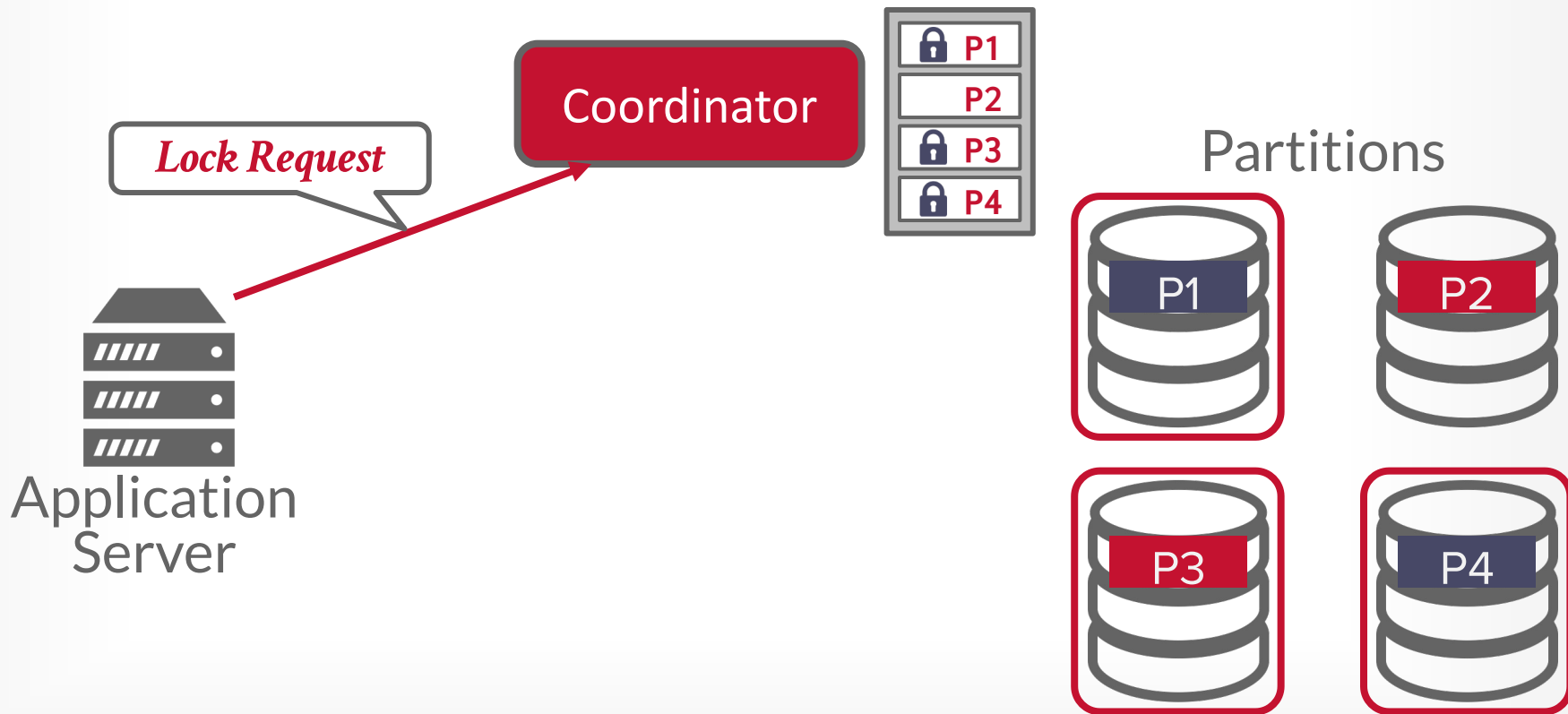
CENTRALIZED COORDINATOR



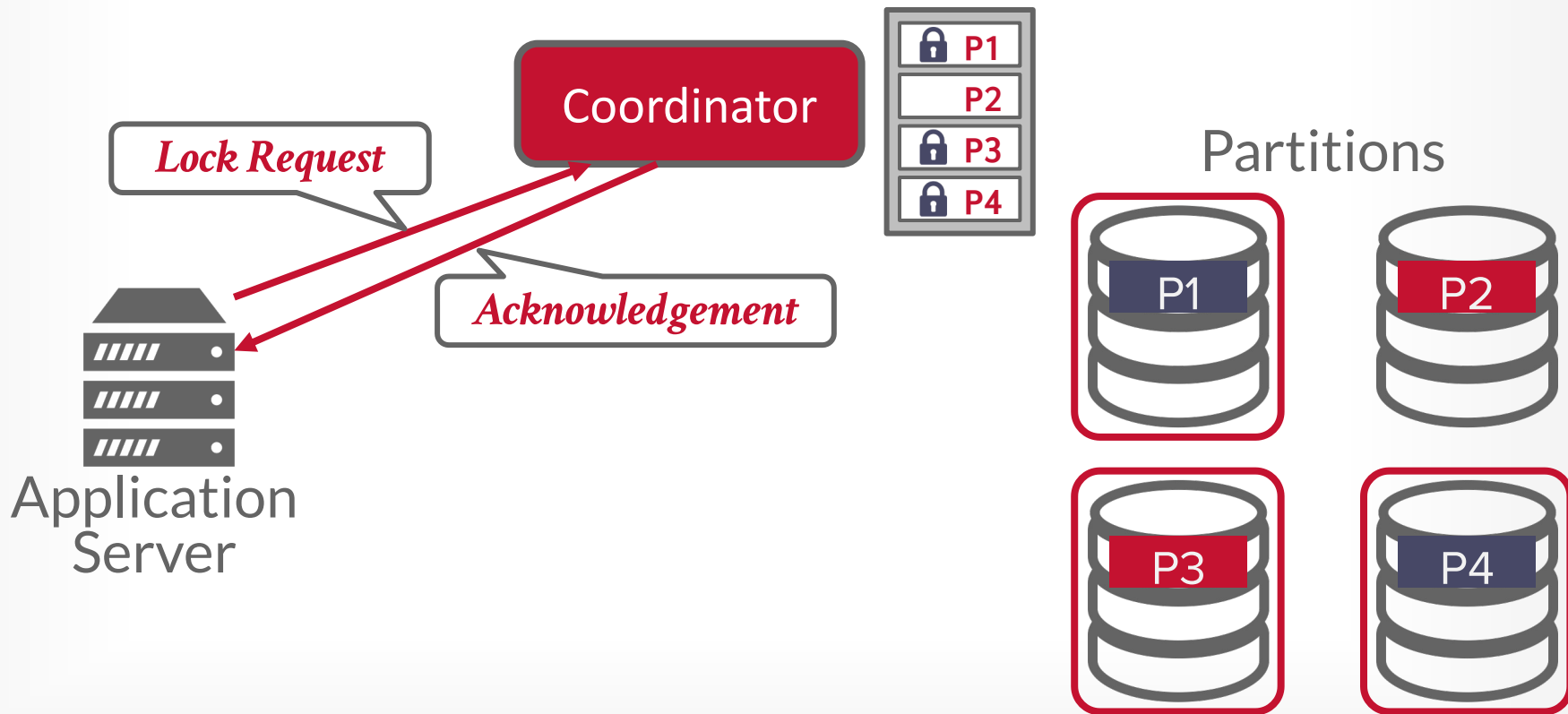
CENTRALIZED COORDINATOR



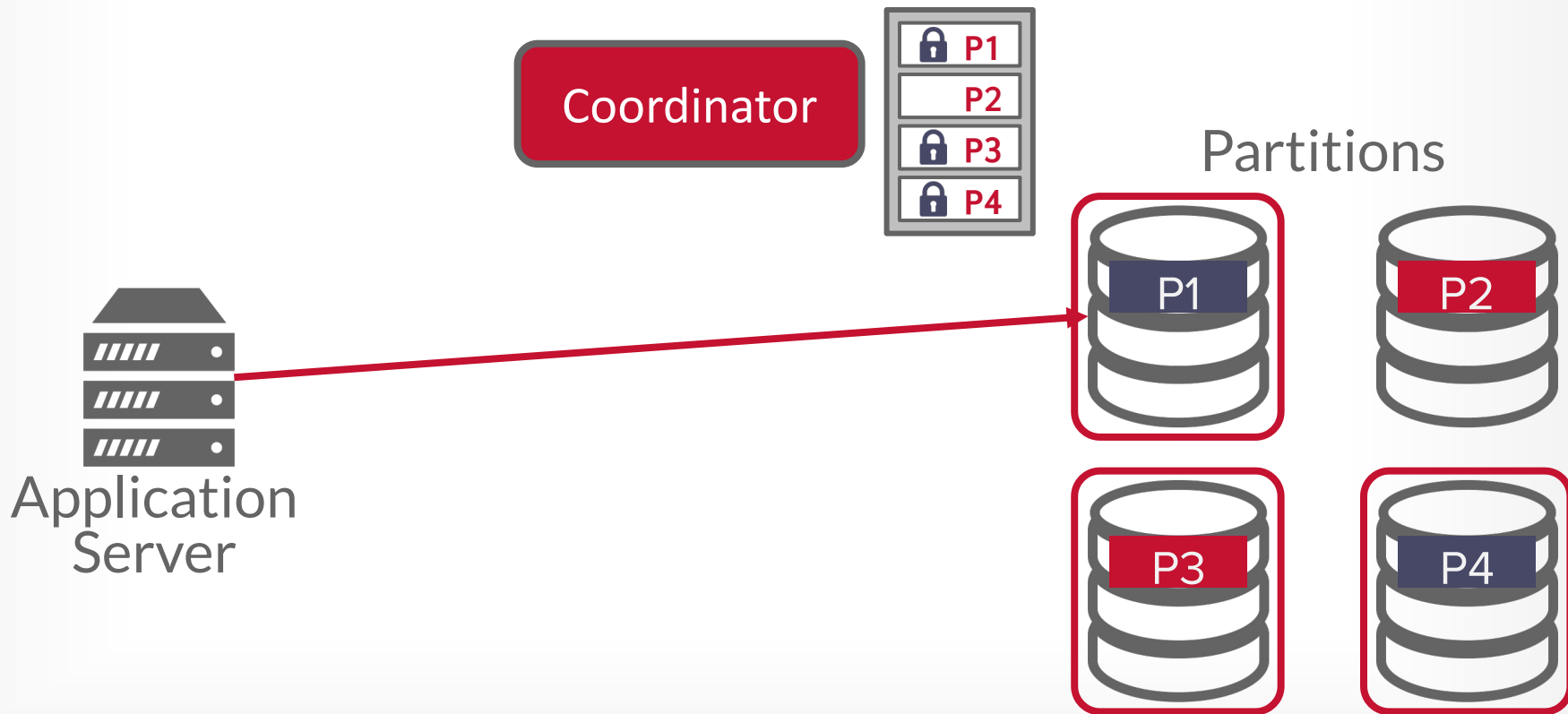
CENTRALIZED COORDINATOR



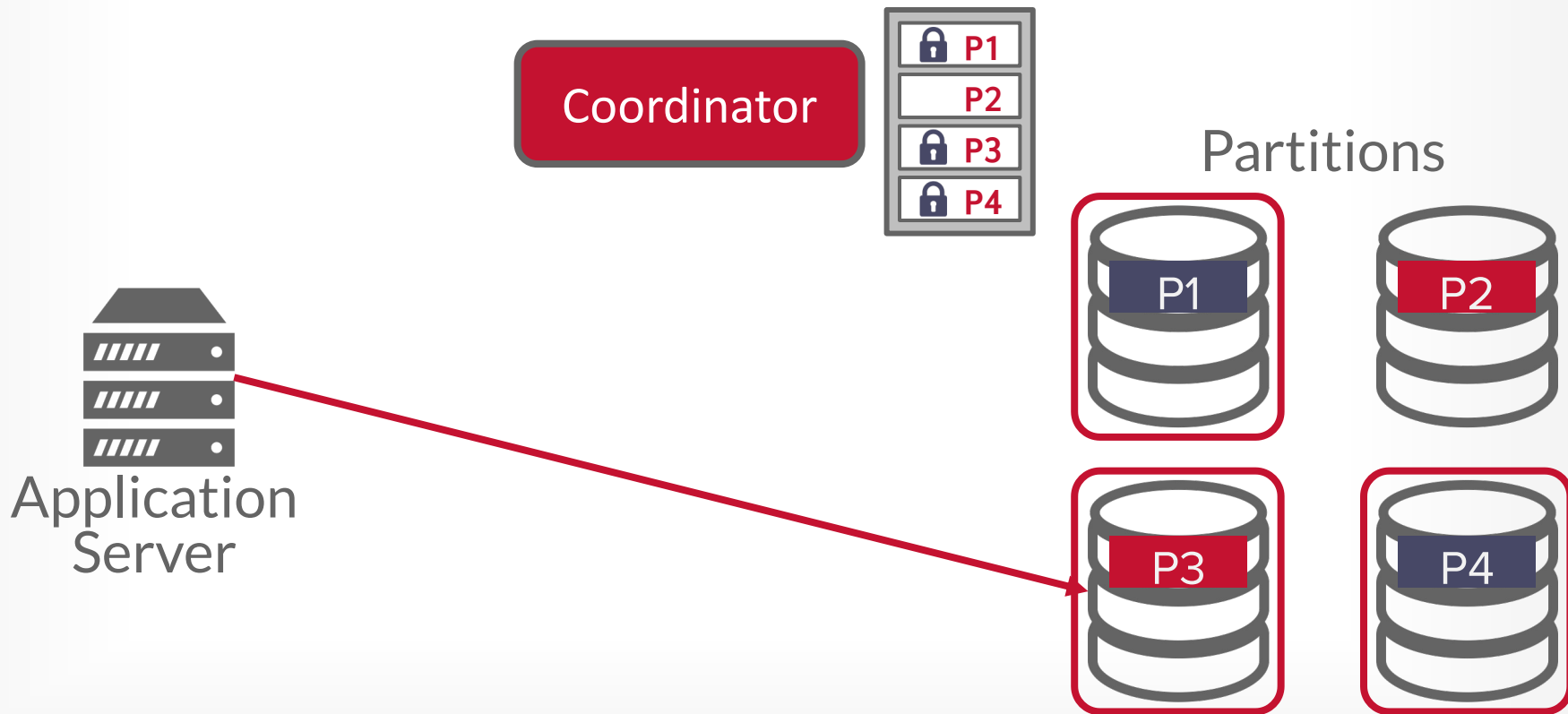
CENTRALIZED COORDINATOR



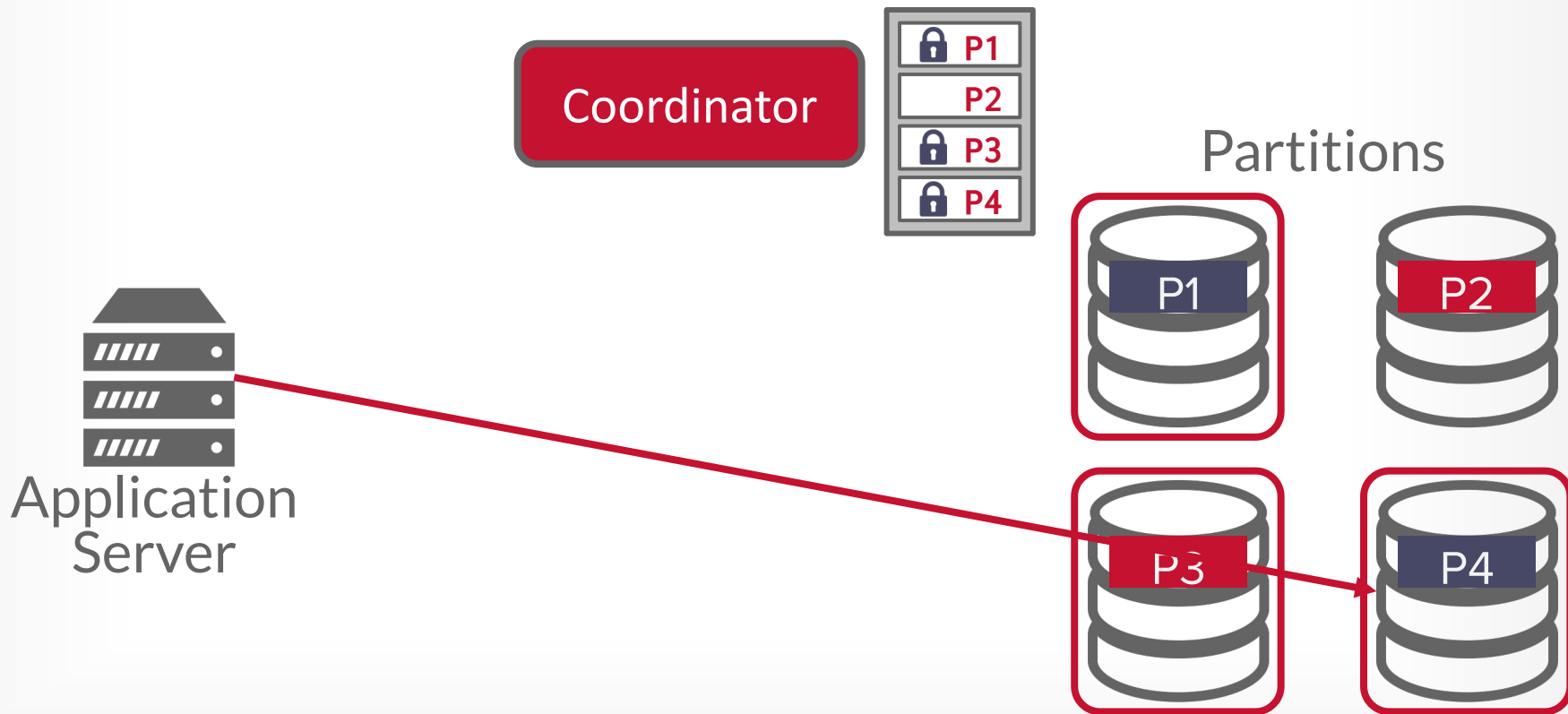
CENTRALIZED COORDINATOR



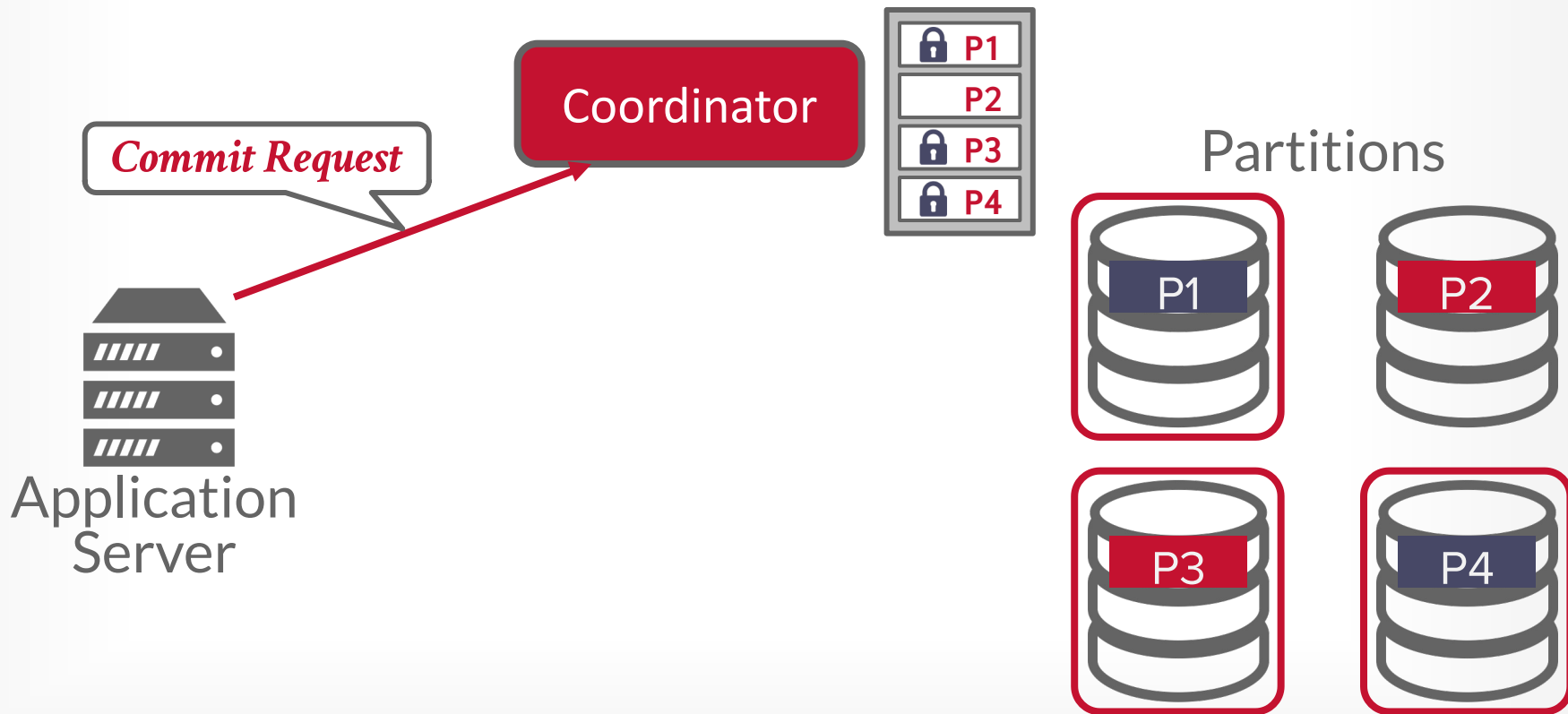
CENTRALIZED COORDINATOR



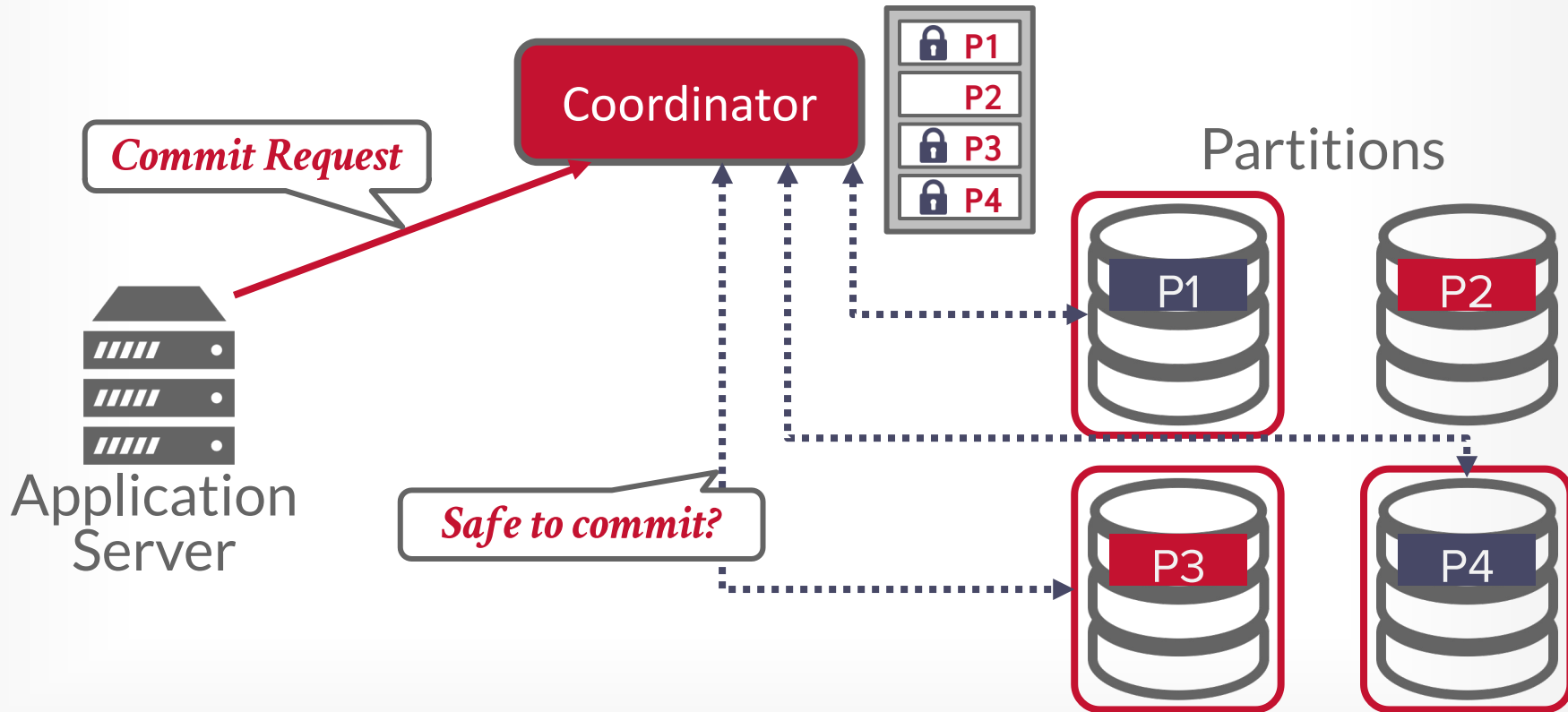
CENTRALIZED COORDINATOR



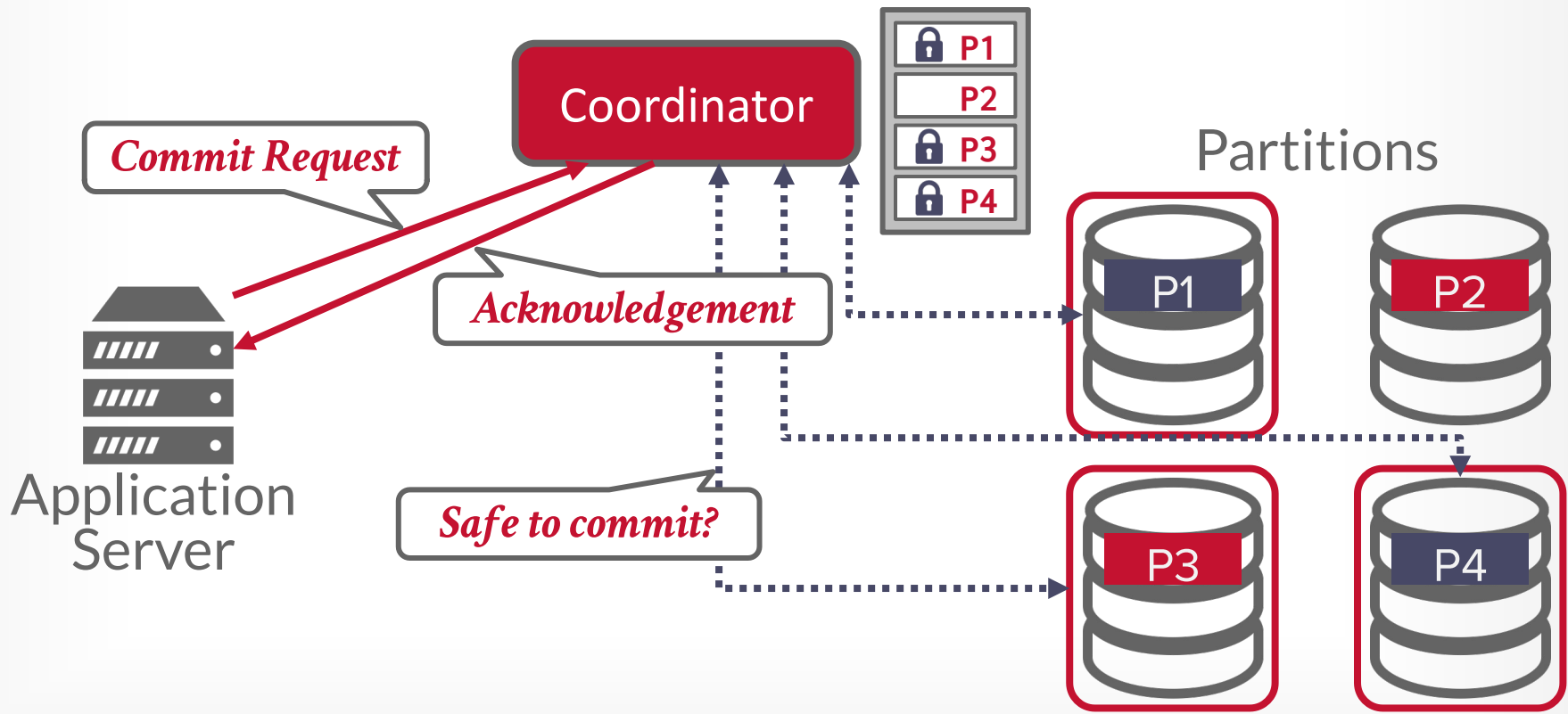
CENTRALIZED COORDINATOR



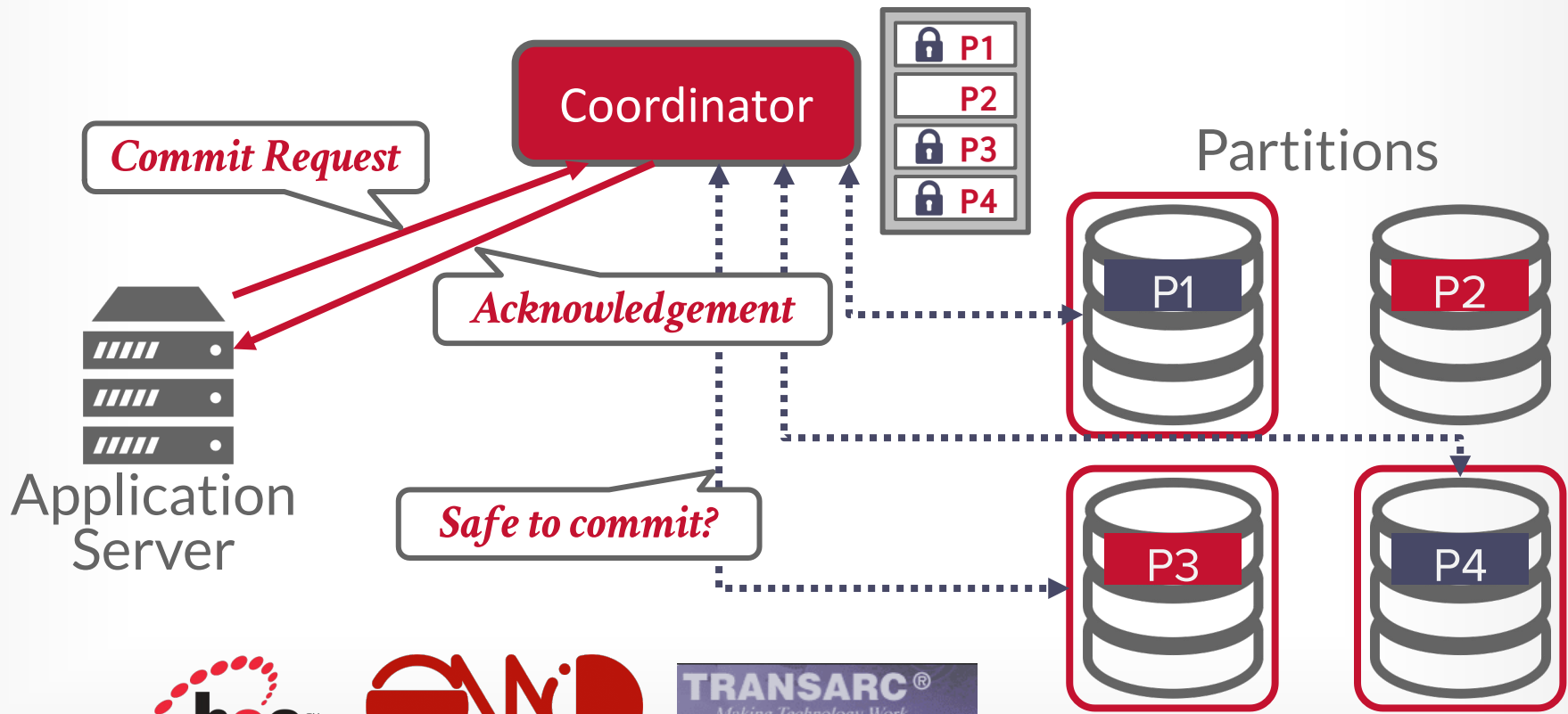
CENTRALIZED COORDINATOR



CENTRALIZED COORDINATOR



CENTRALIZED COORDINATOR



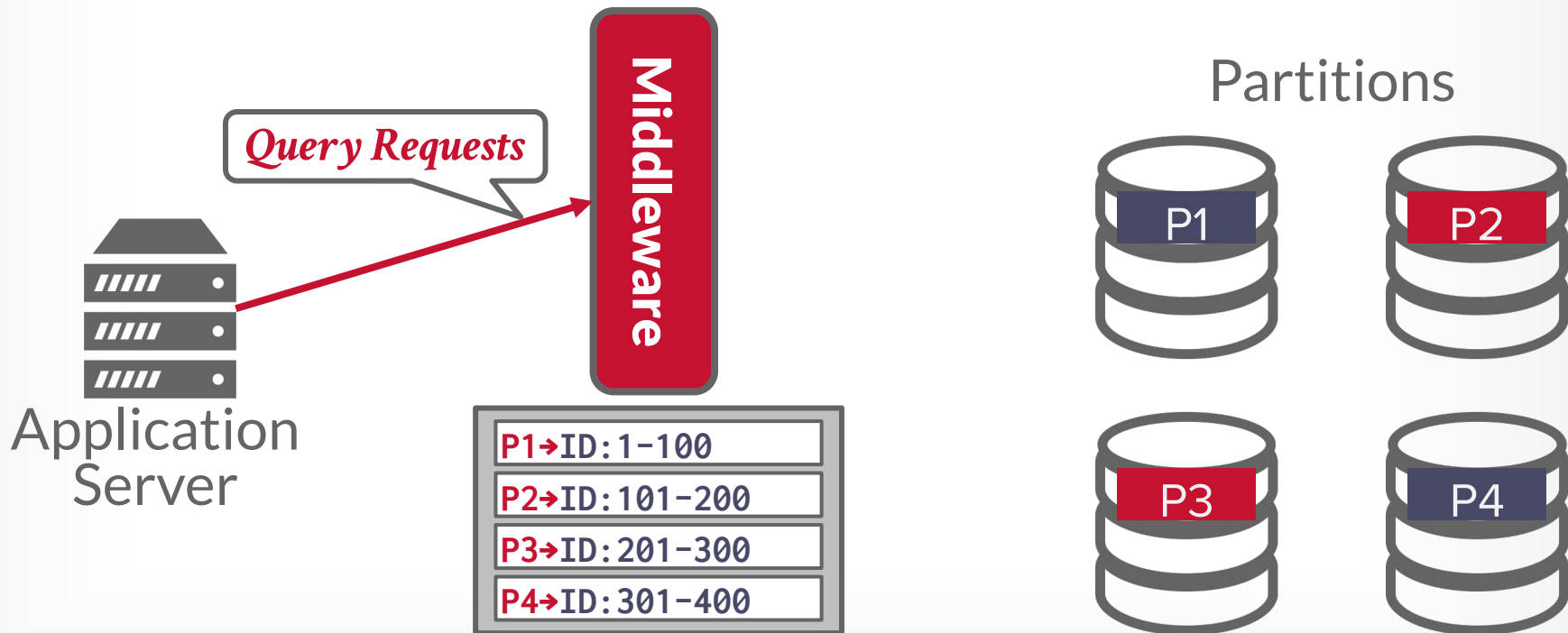
CENTRALIZED COORDINATOR



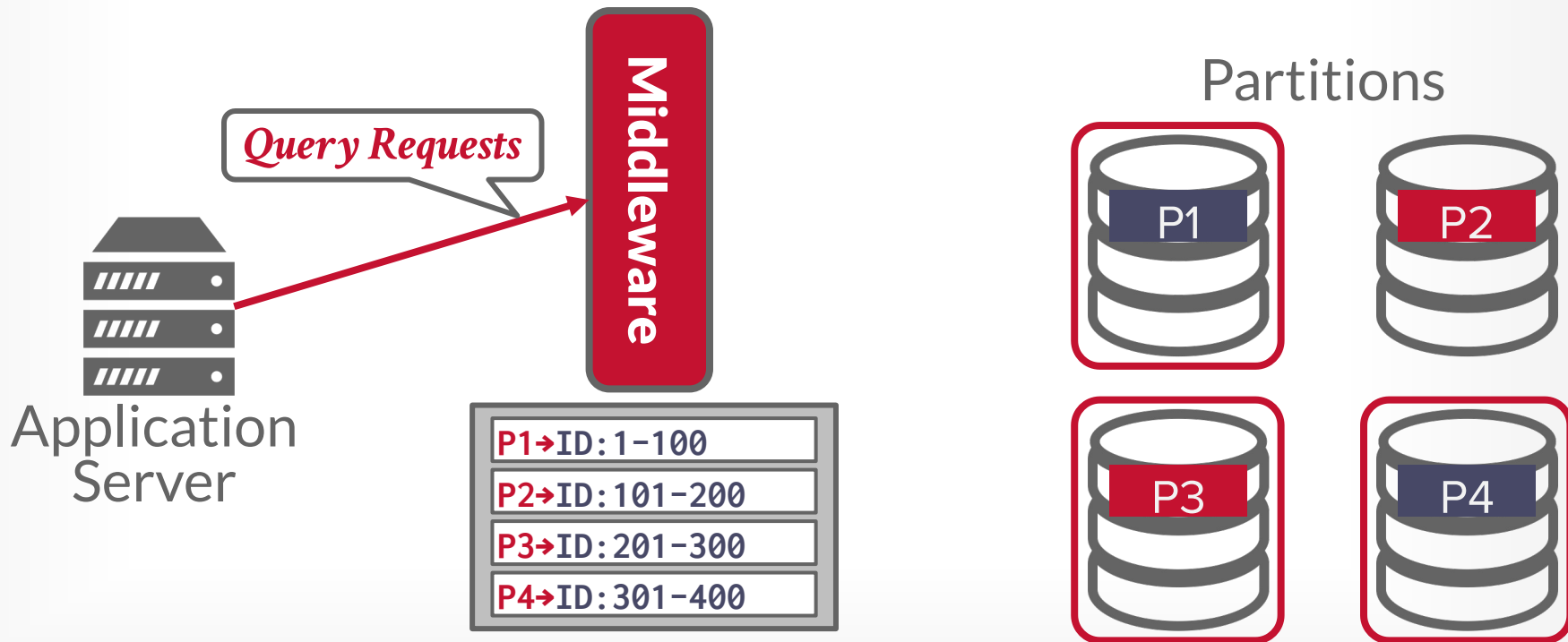
CENTRALIZED COORDINATOR



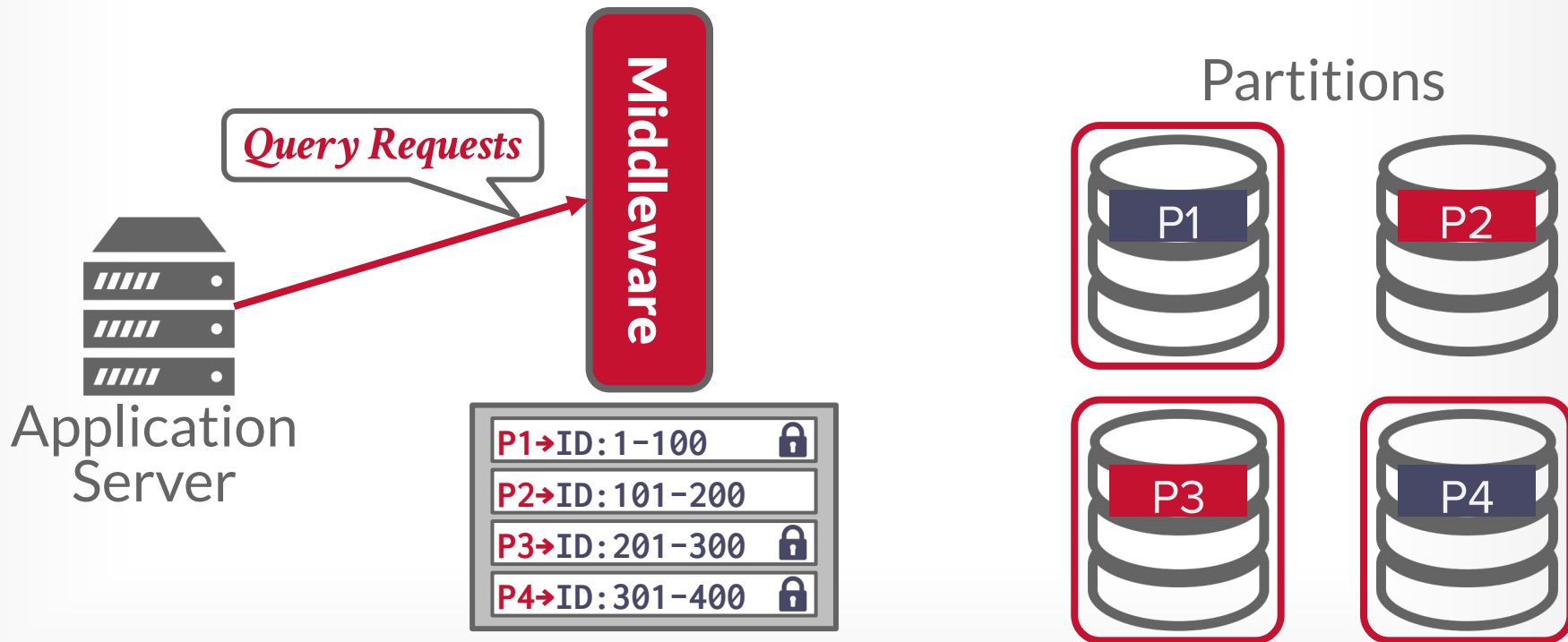
CENTRALIZED COORDINATOR



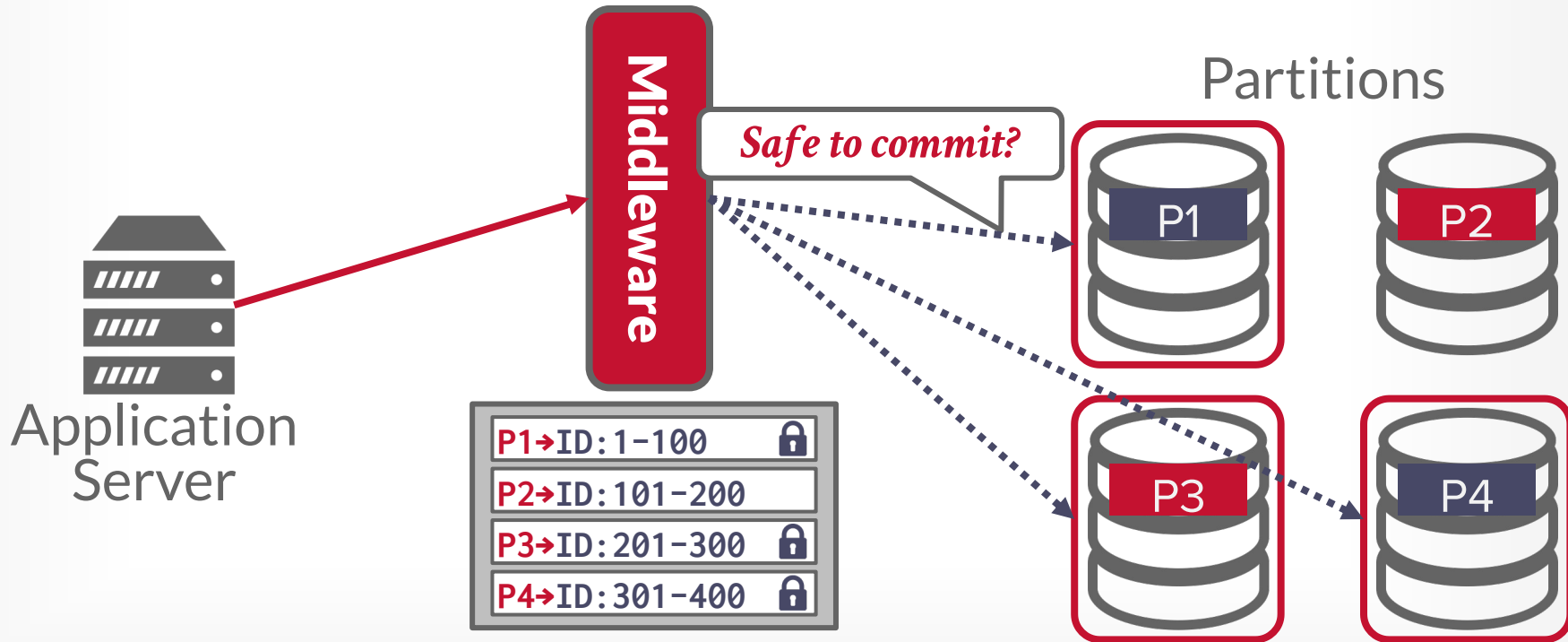
CENTRALIZED COORDINATOR



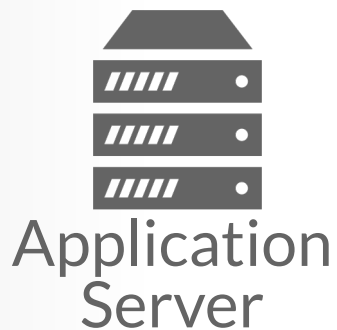
CENTRALIZED COORDINATOR



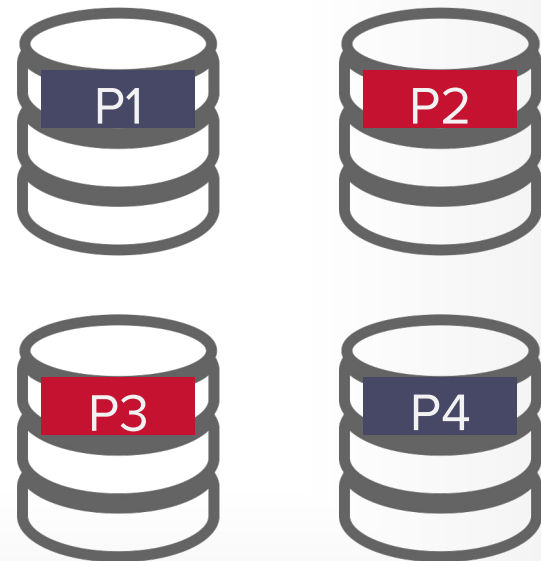
CENTRALIZED COORDINATOR



DECENTRALIZED COORDINATOR



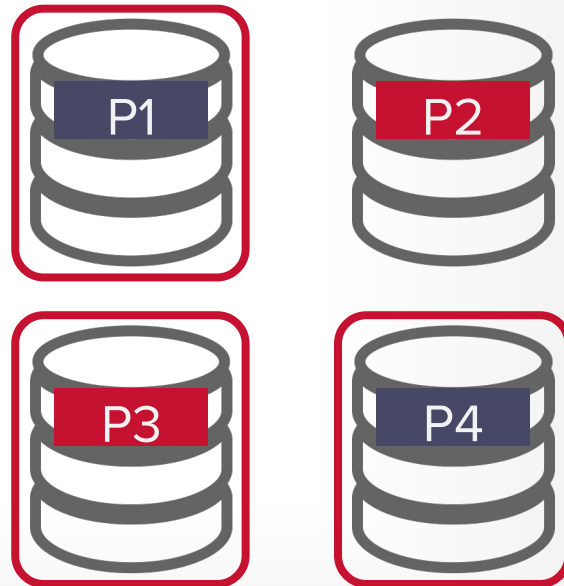
Partitions



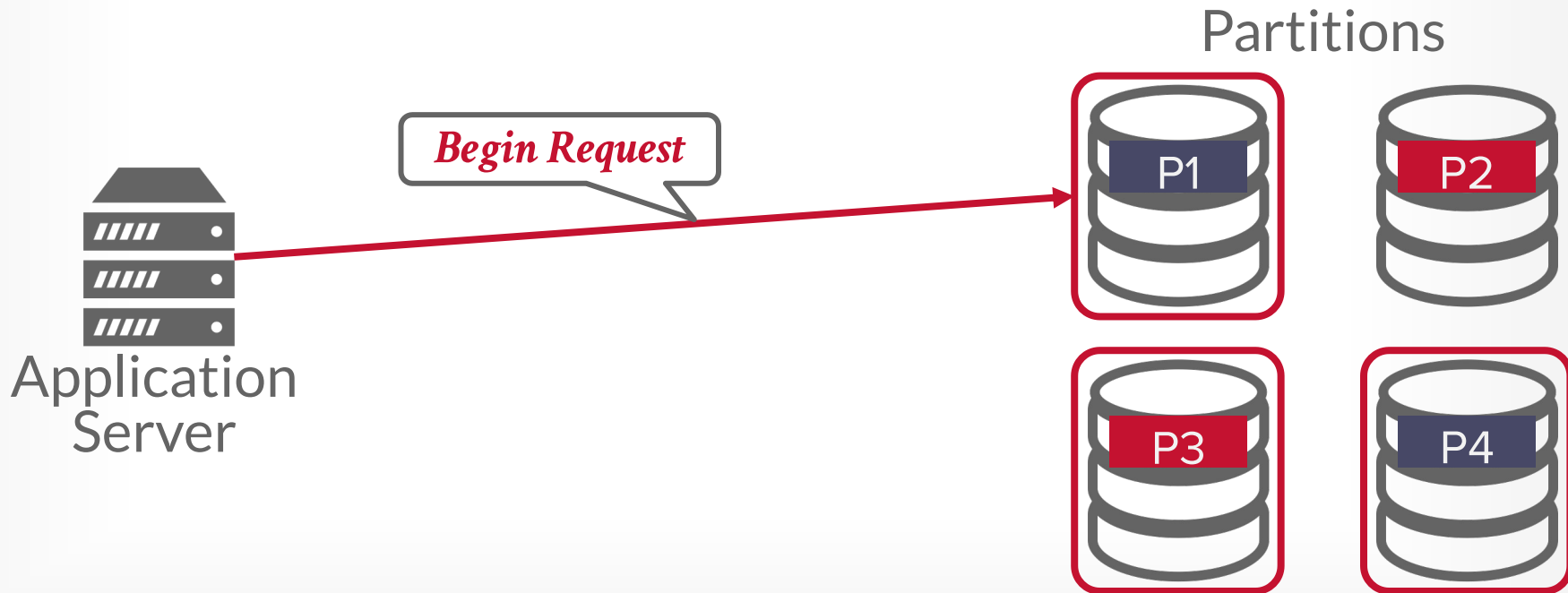
DECENTRALIZED COORDINATOR



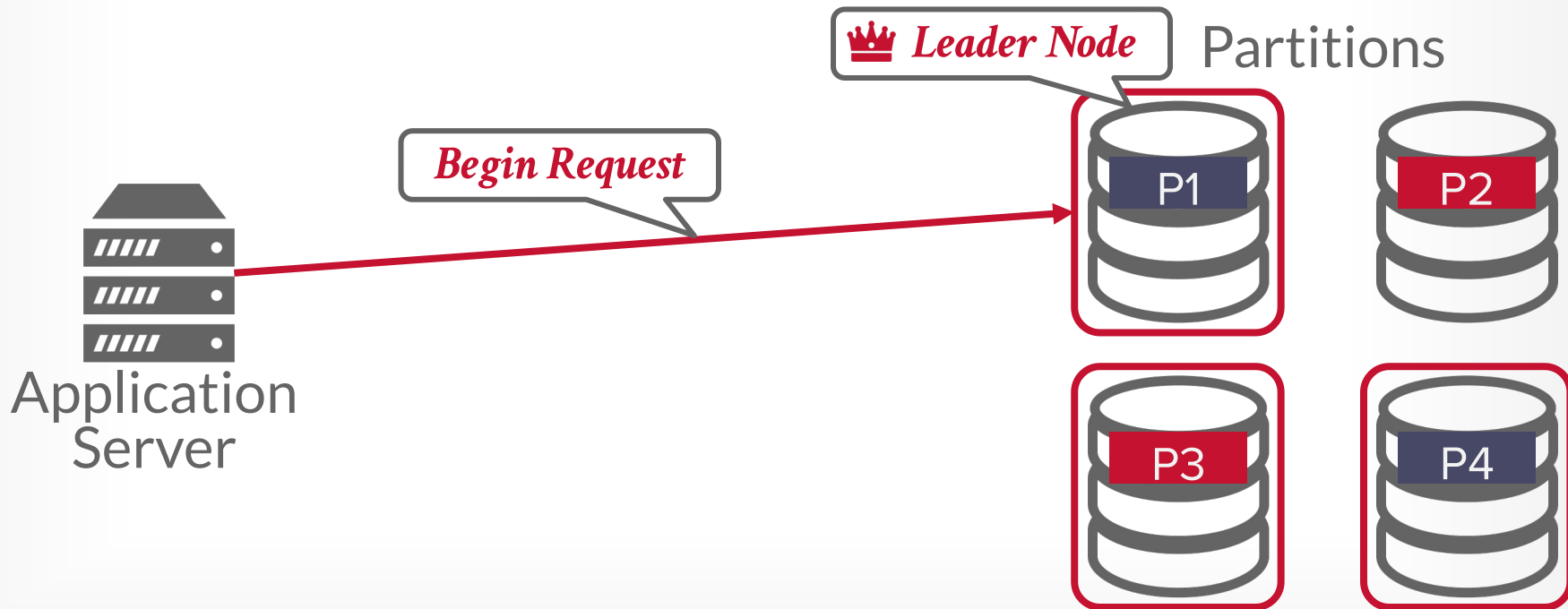
Partitions



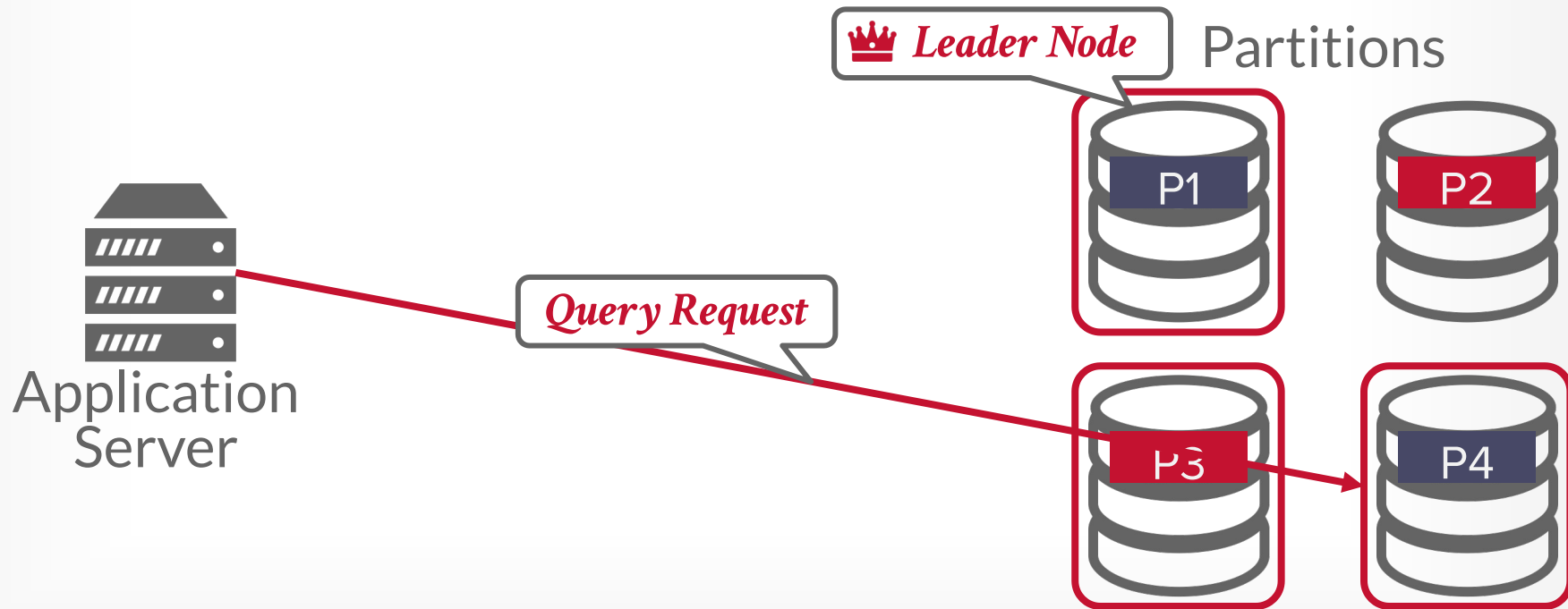
DECENTRALIZED COORDINATOR



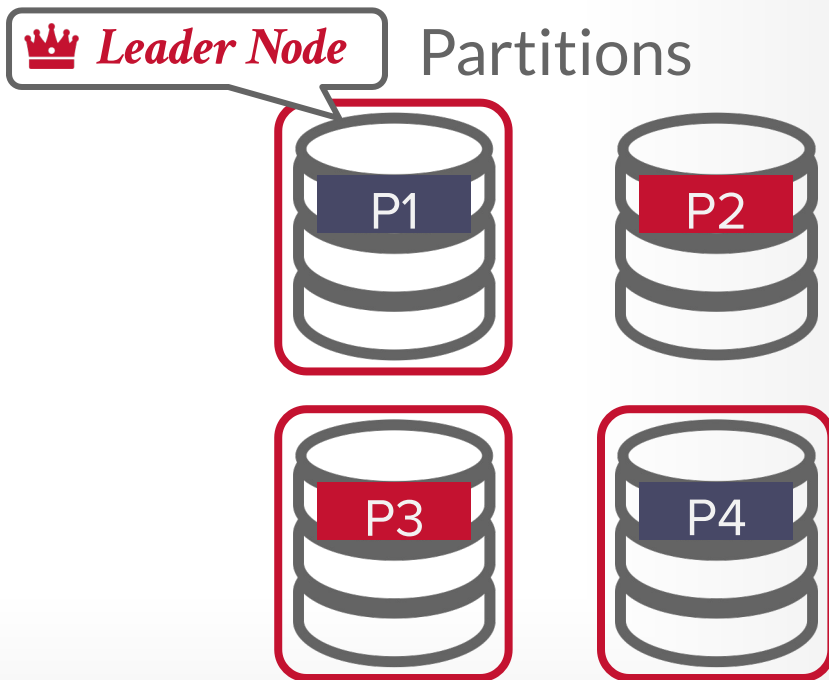
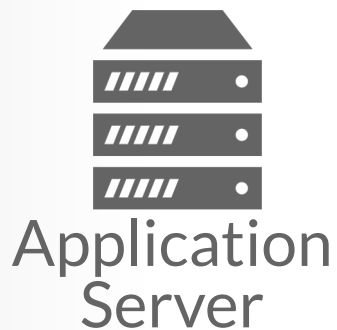
DECENTRALIZED COORDINATOR



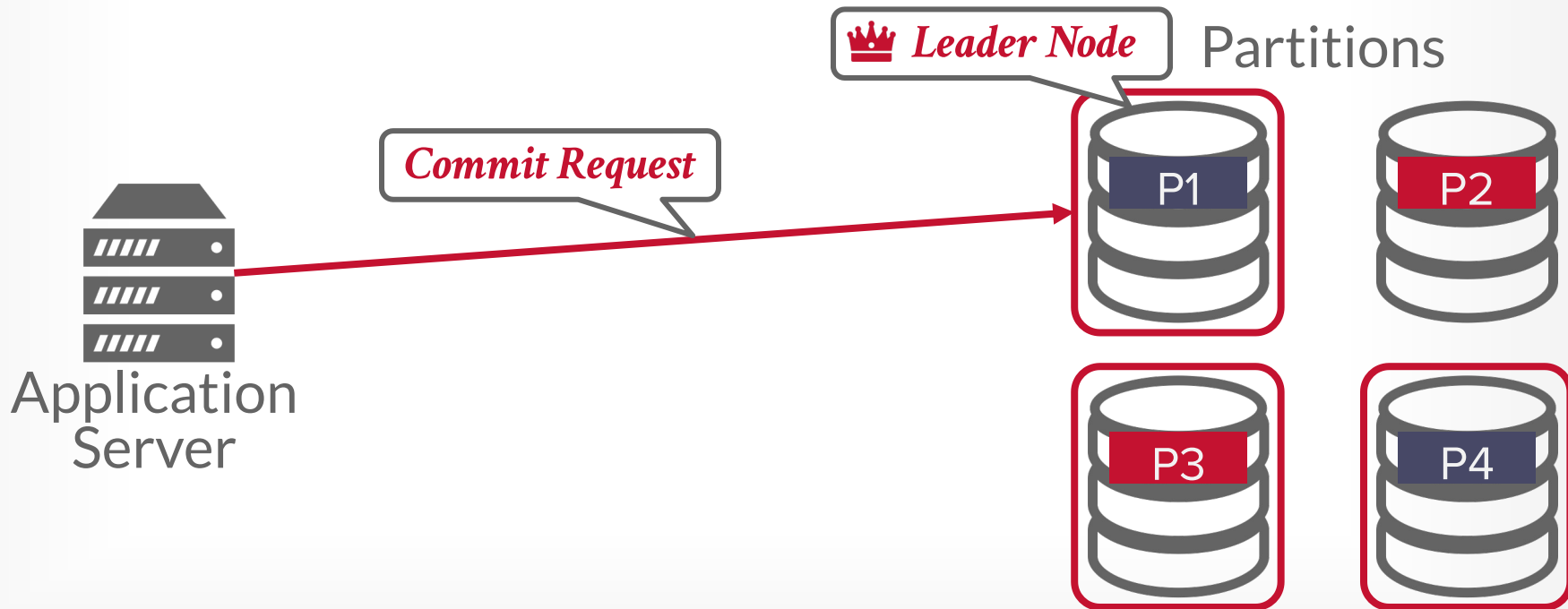
DECENTRALIZED COORDINATOR



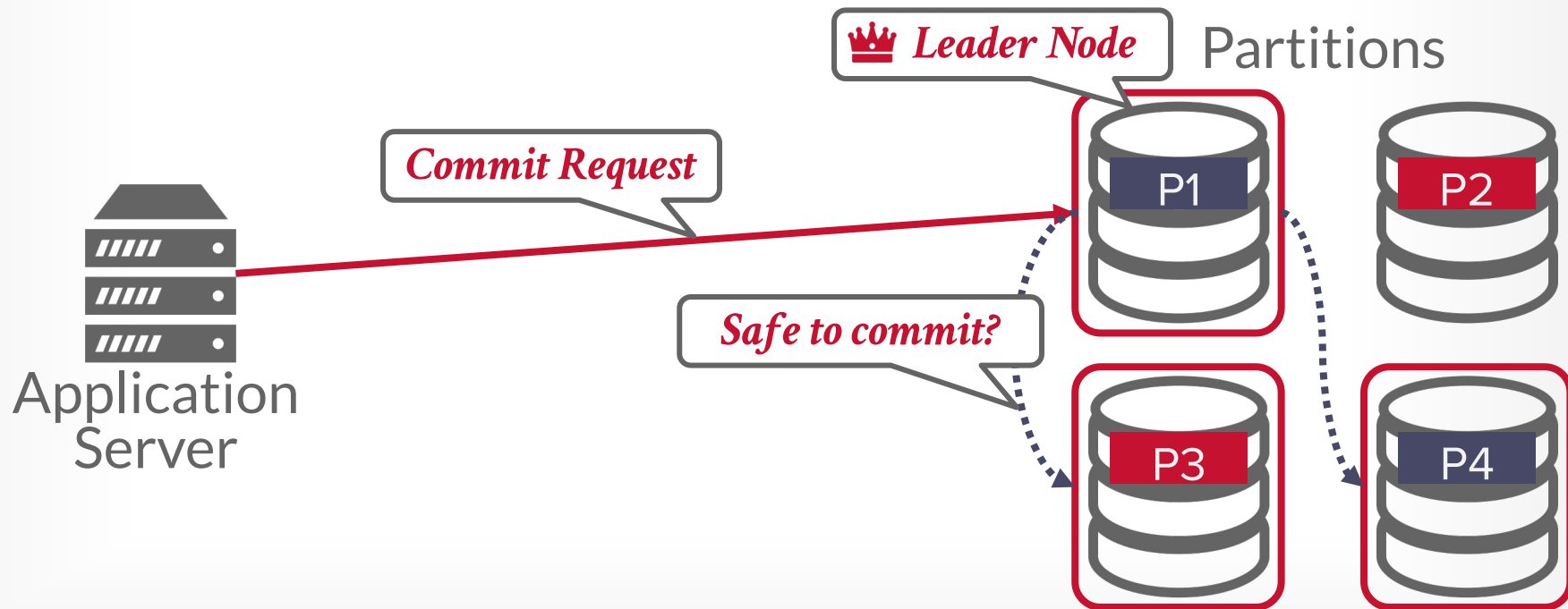
DECENTRALIZED COORDINATOR



DECENTRALIZED COORDINATOR



DECENTRALIZED COORDINATOR



OBSERVATION

We have assumed that the nodes in our distributed systems are running the same DBMS software.

But organizations often run many different DBMSs in their applications.

It would be nice if we could have a single interface for all our data.

FEDERATED DATABASES

Distributed architecture that connects disparate DBMSs into a single logical system.

→ Expose a single query interface that can access data at any location.

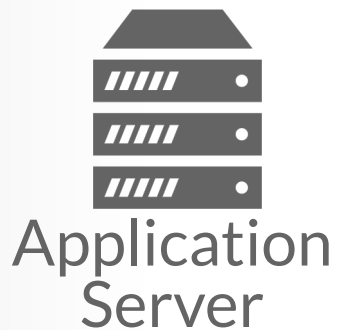
This is hard and nobody does it well

→ Different data models, query languages, limitations.

→ No easy way to optimize queries

→ Lots of data copying (bad).

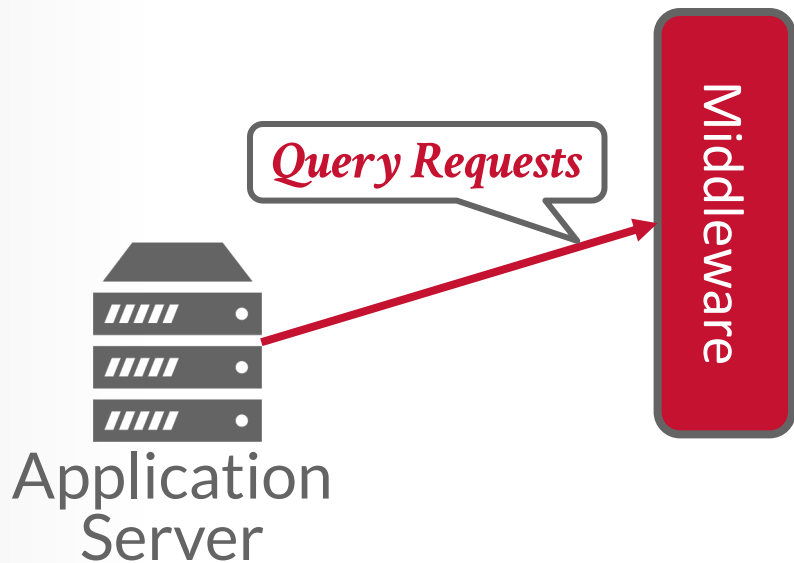
FEDERATED DATABASE EXAMPLE



Back-end DBMSs



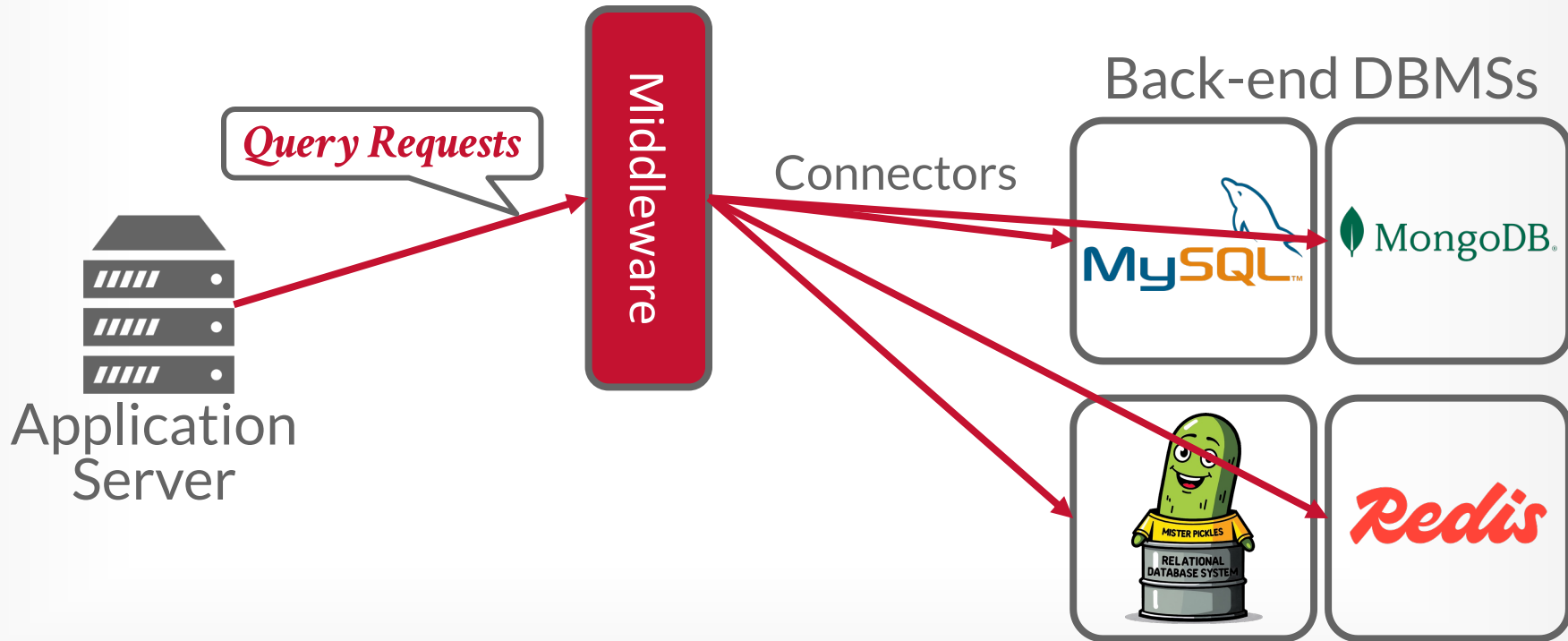
FEDERATED DATABASE EXAMPLE



Back-end DBMSs



FEDERATED DATABASE EXAMPLE



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.

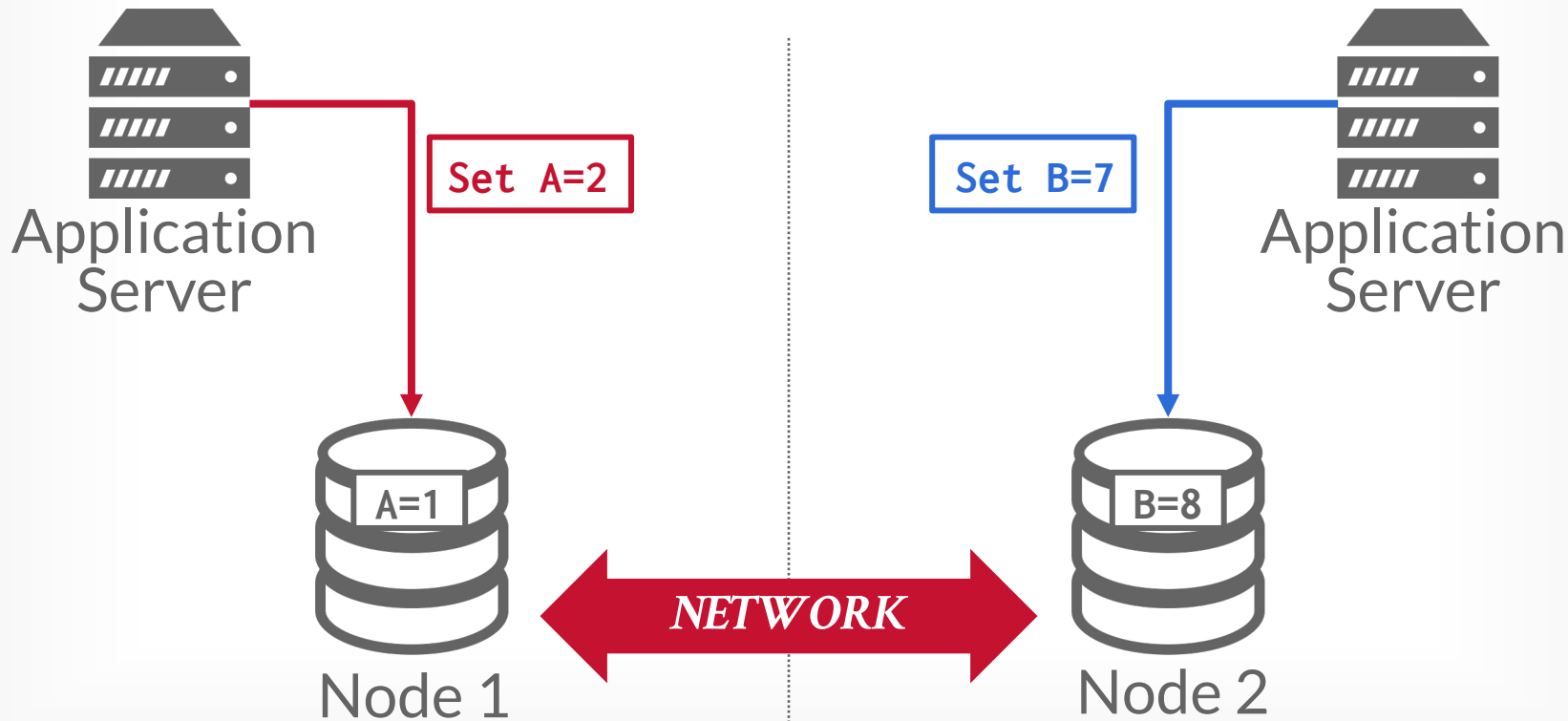
→ Node Failures (Permanent + Ephemeral).

→ Clock Skew.

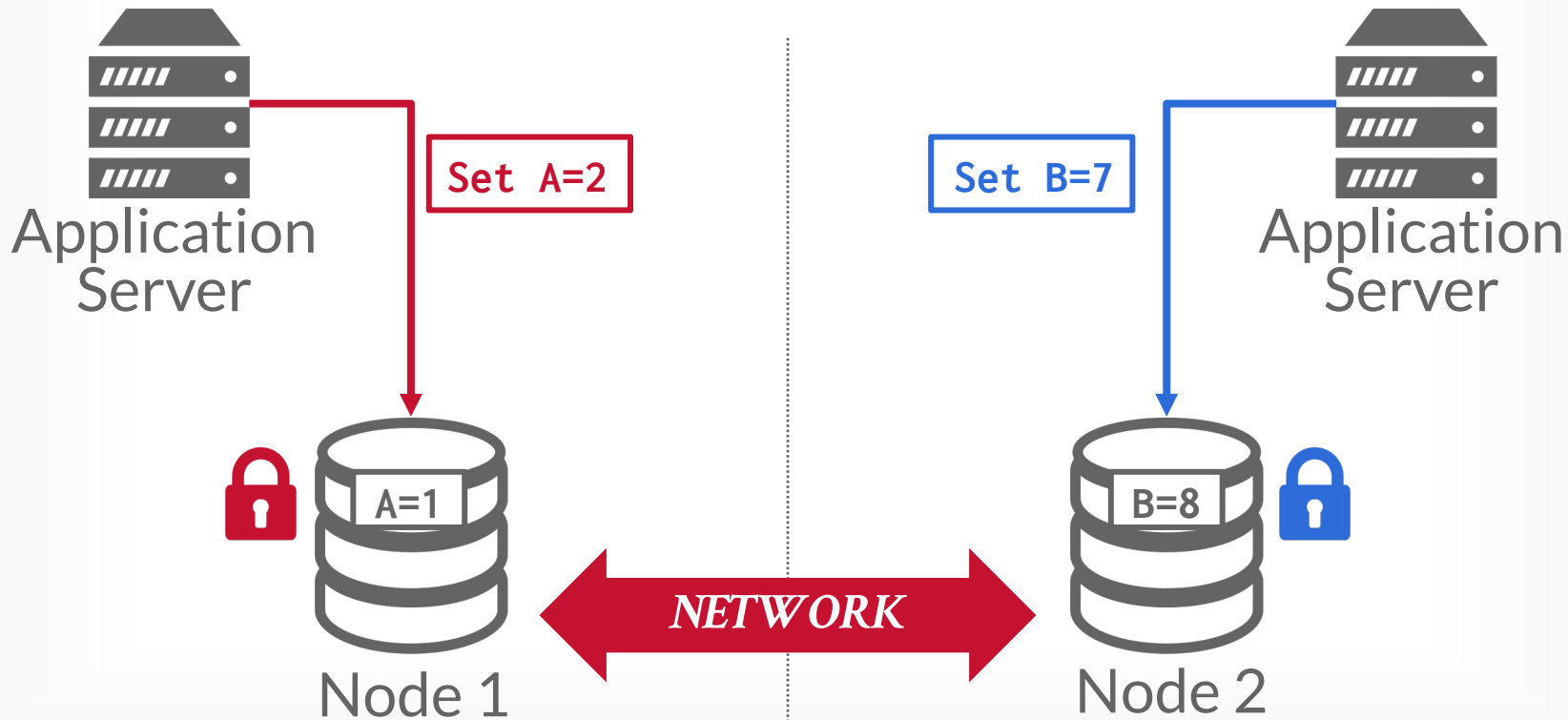
DISTRIBUTED 2PL



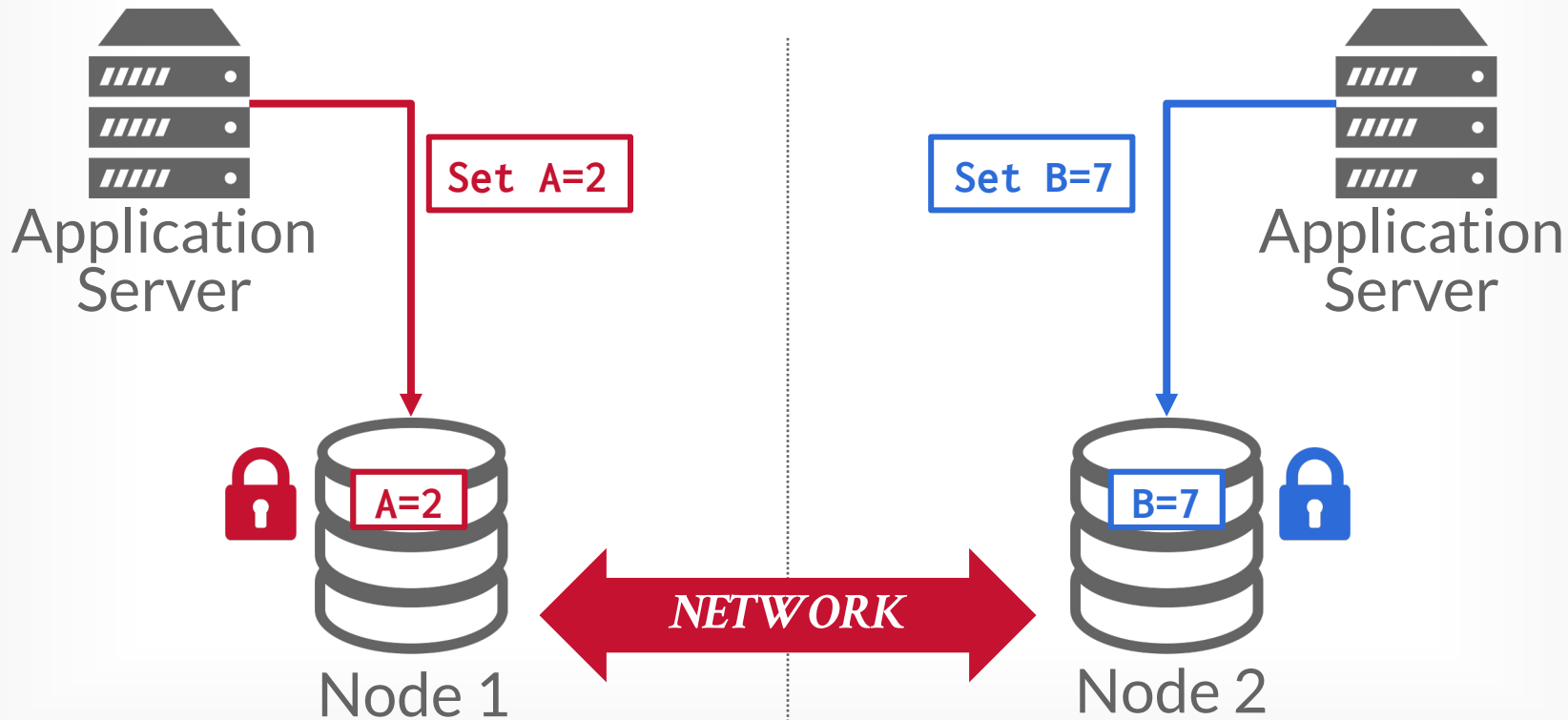
DISTRIBUTED 2PL



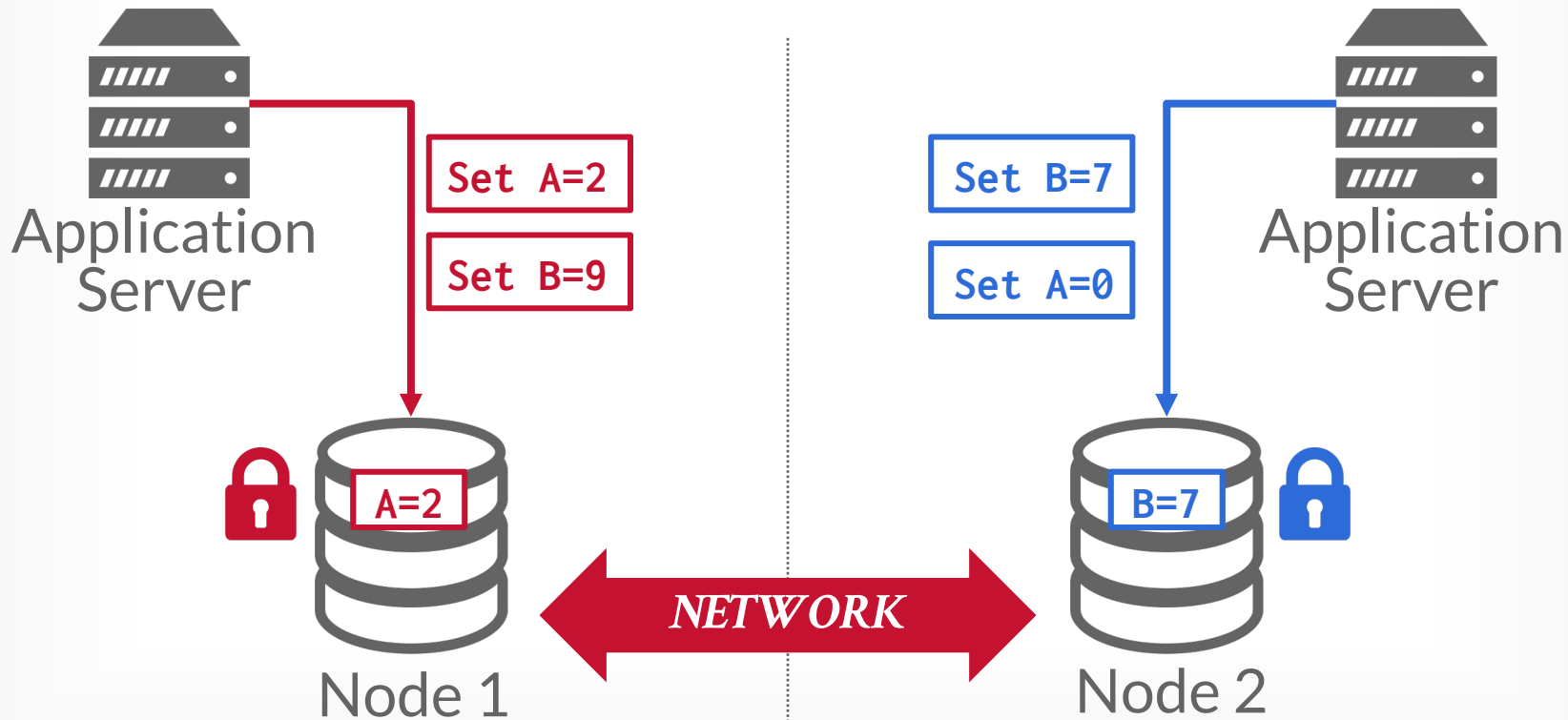
DISTRIBUTED 2PL



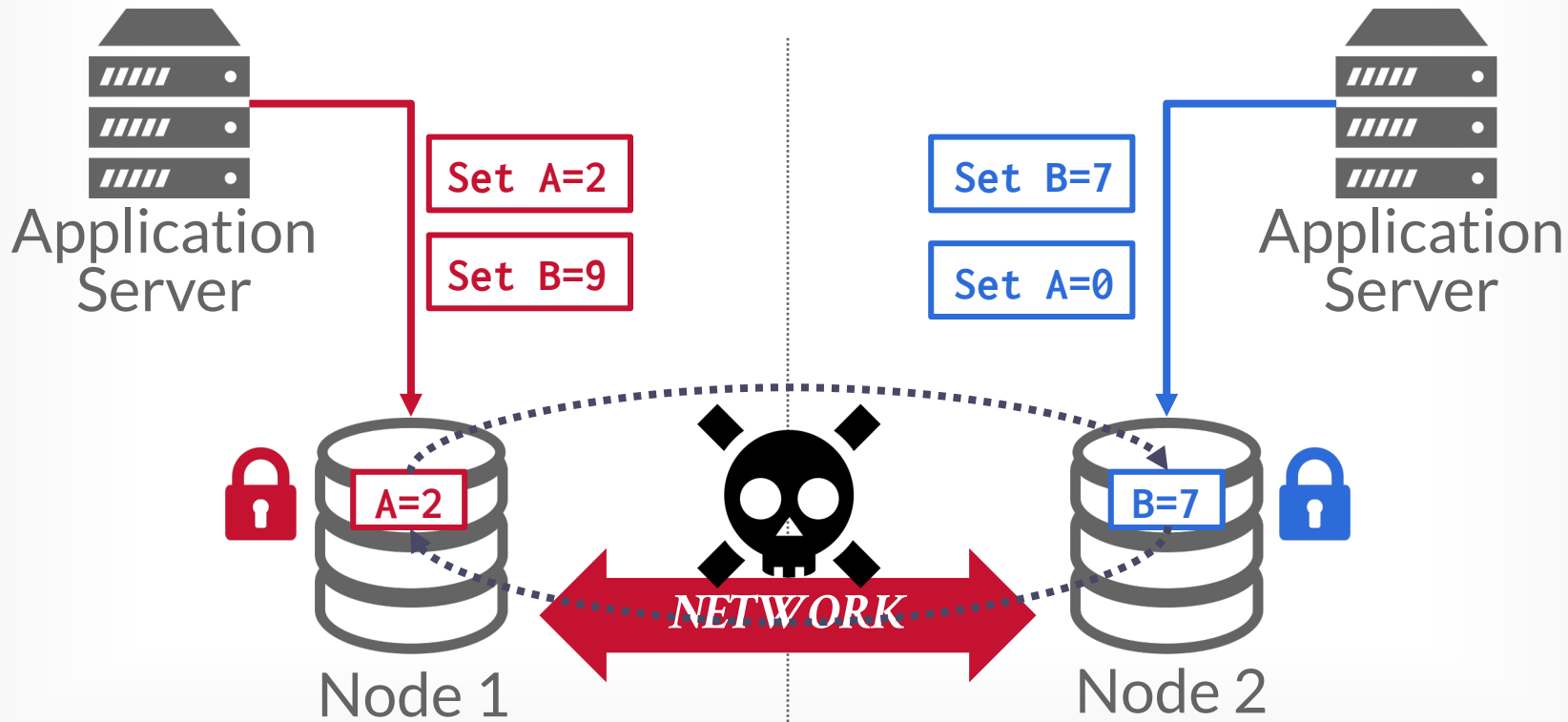
DISTRIBUTED 2PL



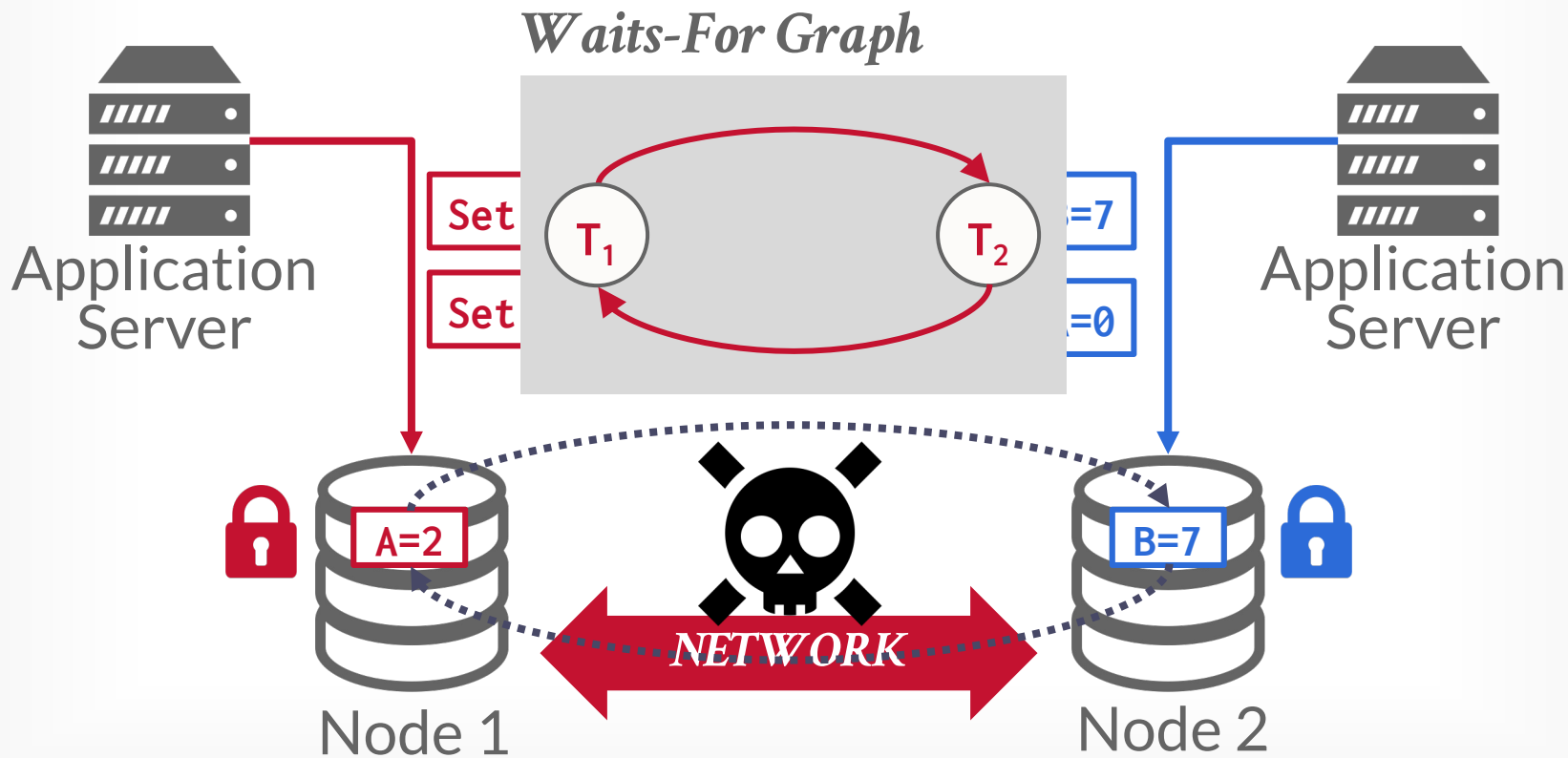
DISTRIBUTED 2PL



DISTRIBUTED 2PL



DISTRIBUTED 2PL



CONCLUSION

We have barely scratched the surface on distributed database systems...

It is hard to get this right.

NEXT CLASS

Distributed OLTP Systems

Replication

CAP Theorem

Real-World Examples