

# Fault-Tolerant Distributed Transactions on Blockchain

## *Introduction*



Suyash Gupta



Jelle Hellings



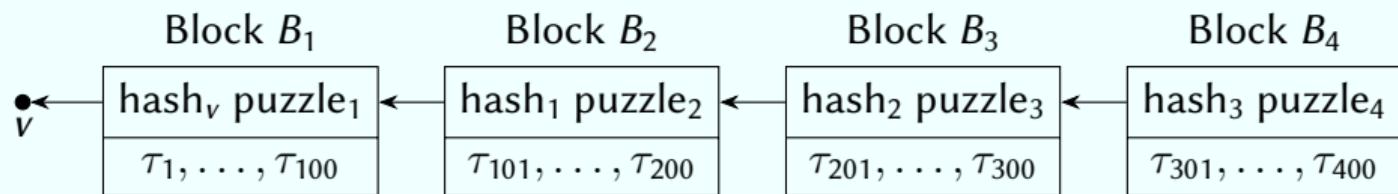
Mohammad Sadoghi



# A Quick Overview of Bitcoin

## What is Bitcoin?

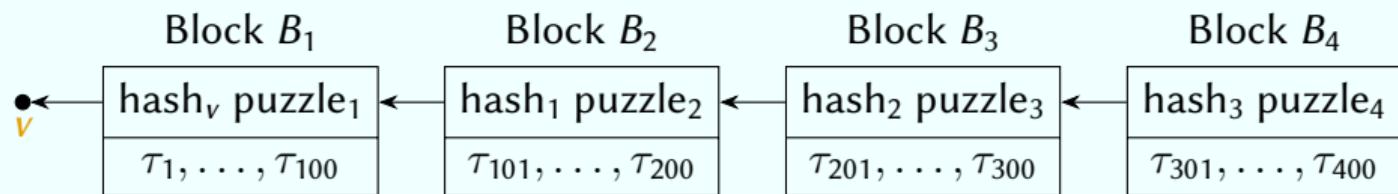
- ▶ A *monetary token*, the Bitcoin.
- ▶ An open and decentralized way to *transfer* these tokens.
- ▶ A representation of the history of these transfers: the *blockchain*:



# A Quick Overview of Bitcoin

## What is Bitcoin?

- ▶ A *monetary token*, the Bitcoin.
- ▶ An open and decentralized way to *transfer* these tokens.
- ▶ A representation of the history of these transfers: the *blockchain*:

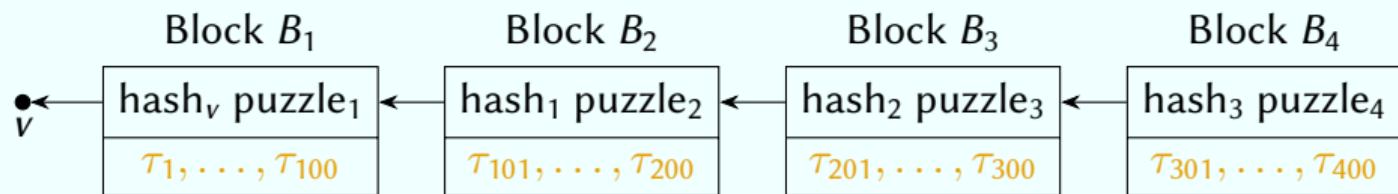


- ▶ A chain of blocks starting at a predefined initial *genesis block*  $v$ .

# A Quick Overview of Bitcoin

## What is Bitcoin?

- ▶ A *monetary token*, the Bitcoin.
- ▶ An open and decentralized way to *transfer* these tokens.
- ▶ A representation of the history of these transfers: the *blockchain*:

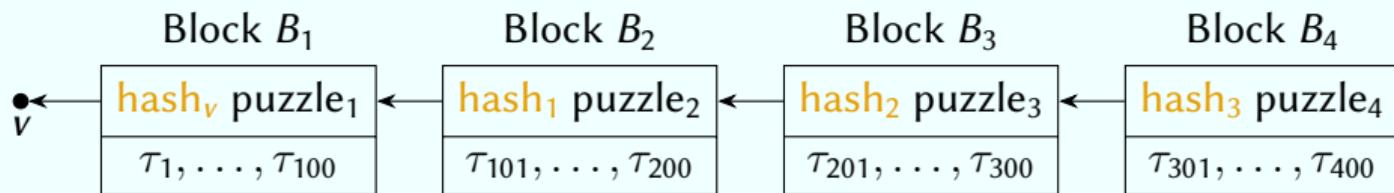


- ▶ Each block holds a *list of transactions* transferring Bitcoins.

# A Quick Overview of Bitcoin

## What is Bitcoin?

- ▶ A *monetary token*, the Bitcoin.
- ▶ An open and decentralized way to *transfer* these tokens.
- ▶ A representation of the history of these transfers: the *blockchain*:

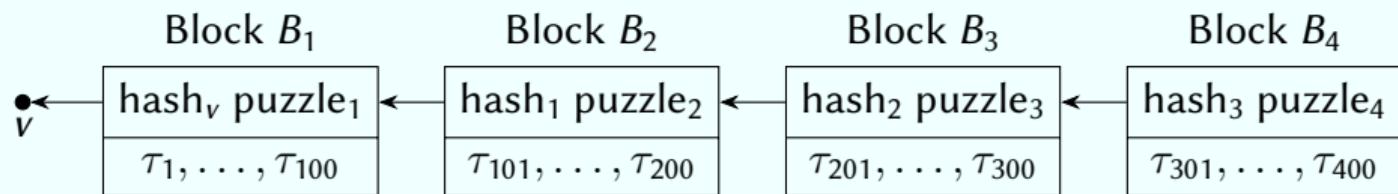


- ▶ Each block holds the *hash* of a previous block.

# A Quick Overview of Bitcoin

## What is Bitcoin?

- ▶ A *monetary token*, the Bitcoin.
- ▶ An open and decentralized way to *transfer* these tokens.
- ▶ A representation of the history of these transfers: the *blockchain*:

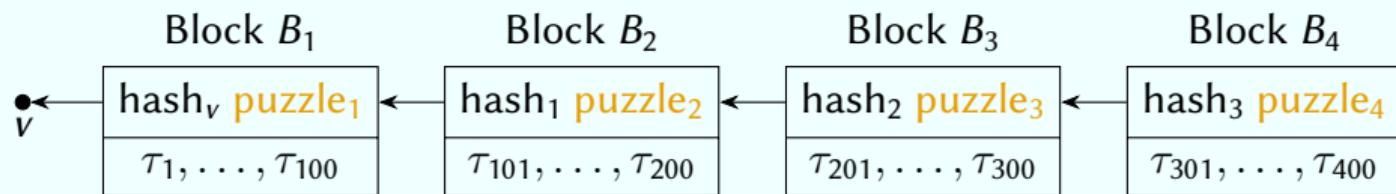


- ▶ **Incentive:** The creator of a block that is part of the chain gets a reward.

# A Quick Overview of Bitcoin

## What is Bitcoin?

- ▶ A *monetary token*, the Bitcoin.
- ▶ An open and decentralized way to *transfer* these tokens.
- ▶ A representation of the history of these transfers: the *blockchain*:

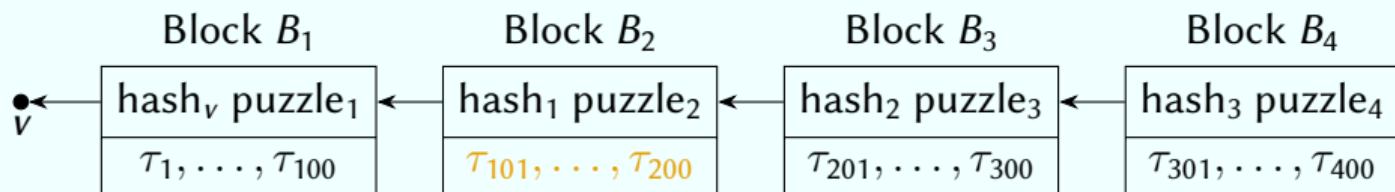


- ▶ Each block contains a solution to a computational *complex puzzle*, used to make the blockchain *tamper-proof*.

# The Tamper-Proof Design of Bitcoin

Incentive: The creator of a block that is part of the chain gets a reward.

Assumption: A malicious participant  $P$  wants to replace  $\tau_{123}$  by  $\tau'$

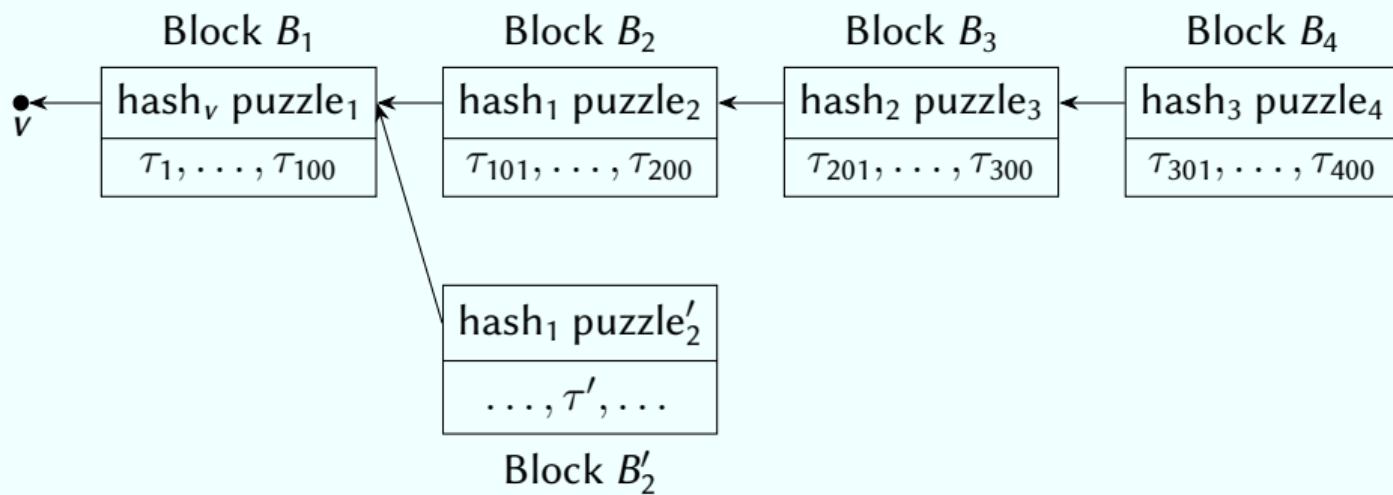


# The Tamper-Proof Design of Bitcoin

Incentive: The creator of a block that is part of the chain gets a reward.

Assumption: A malicious participant  $P$  wants to replace  $\tau_{123}$  by  $\tau'$

1.  $P$  creates a new block  $B'_2$  with  $\tau'$ .

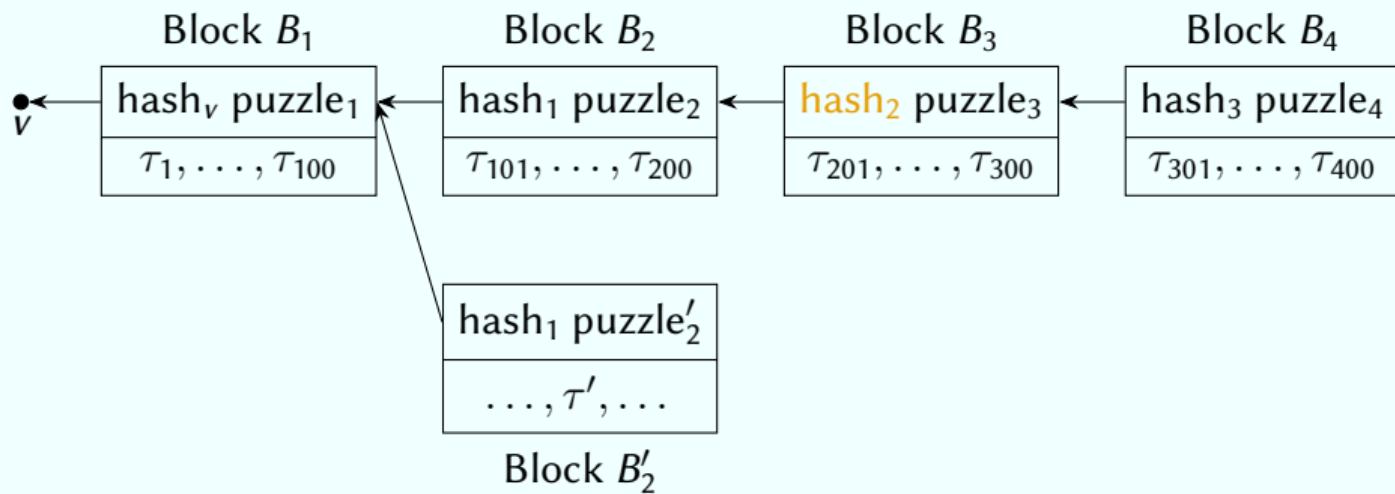


# The Tamper-Proof Design of Bitcoin

Incentive: The creator of a block that is part of the chain gets a reward.

Assumption: A malicious participant  $P$  wants to replace  $\tau_{123}$  by  $\tau'$

2.  $B_2 \neq B'_2$ . Hence,  $\text{hash}_2$  still points to  $B_2$ .

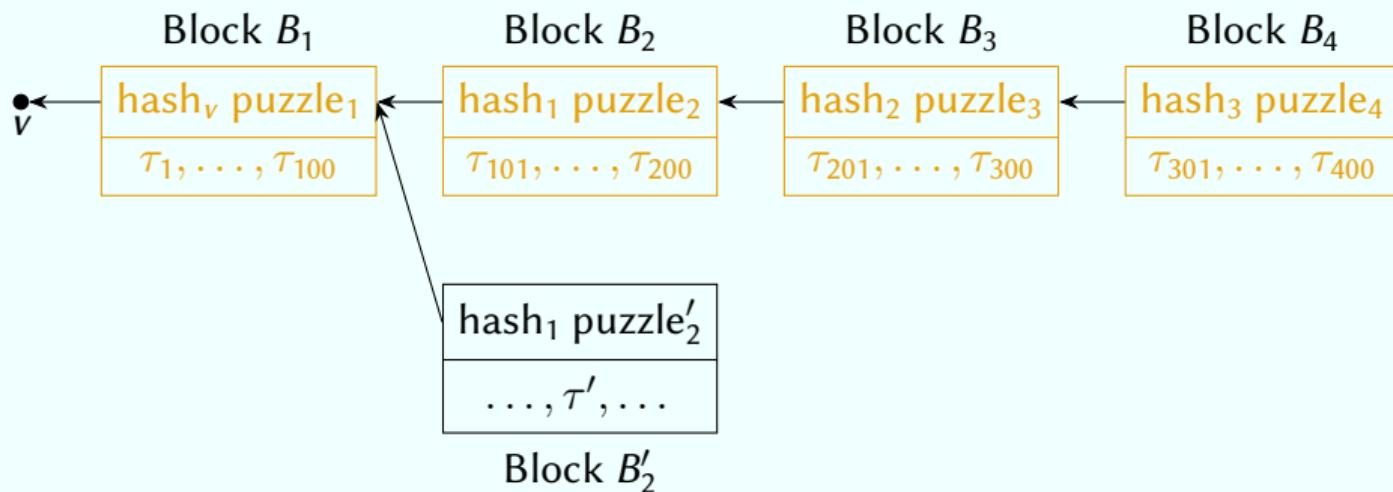


# The Tamper-Proof Design of Bitcoin

Incentive: The creator of a block that is part of the chain gets a reward.

Assumption: A malicious participant  $P$  wants to replace  $\tau_{123}$  by  $\tau'$

3. Good participants prefer to work on *long chains* over *short chains*.



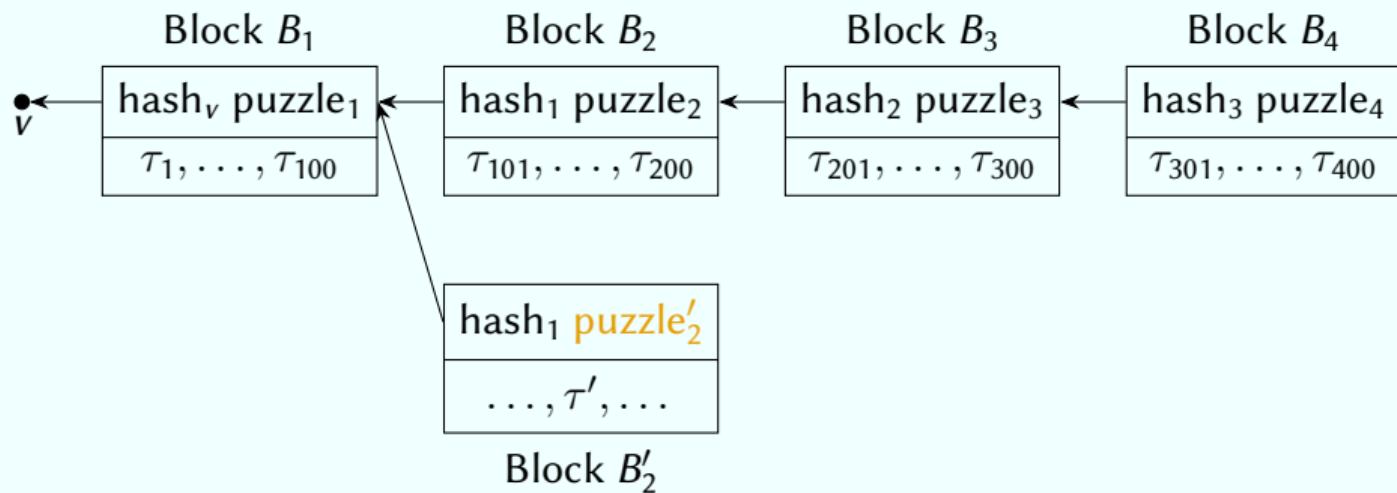
# The Tamper-Proof Design of Bitcoin

**Incentive:** The creator of a block that is part of the chain gets a reward.

**Assumption:** A malicious participant  $P$  wants to replace  $\tau_{123}$  by  $\tau'$

3. Good participants prefer to work on *long chains* over *short chains*.

Long chains give rewards to *more participants*: incentive to work on long chains.

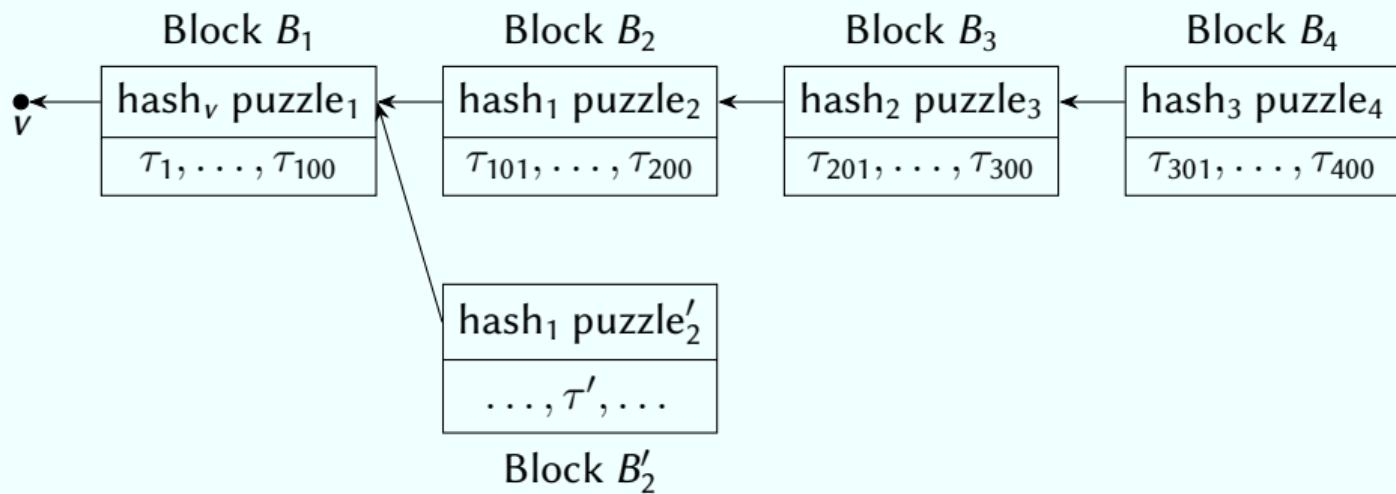


# The Tamper-Proof Design of Bitcoin

Incentive: The creator of a block that is part of the chain gets a reward.

Assumption: A malicious participant  $P$  wants to replace  $\tau_{123}$  by  $\tau'$

4. Complex puzzles prevent  $P$  from easily adding blocks to  $B'_2$ .



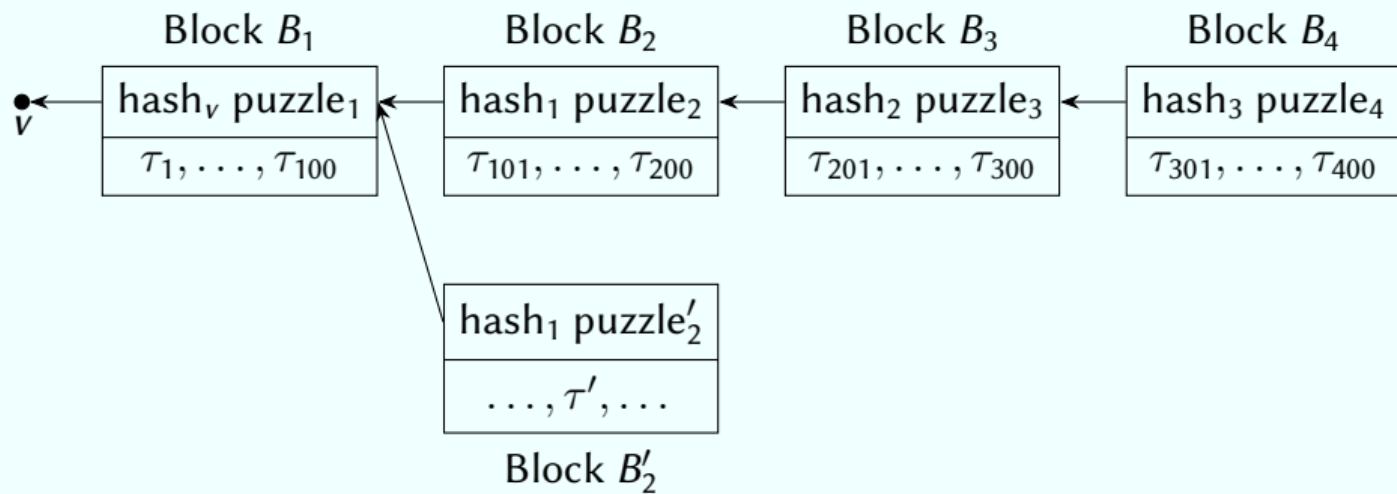
# The Tamper-Proof Design of Bitcoin

**Incentive:** The creator of a block that is part of the chain gets a reward.

**Assumption:** A malicious participant  $P$  wants to replace  $\tau_{123}$  by  $\tau'$

4. Complex puzzles prevent  $P$  from easily adding blocks to  $B'_2$ .

More incentives to continue from block  $B_4$  than block  $B'_2$ !



## Permissioned Blockchains

*Permissioned* blockchains have a well-defined set of *vetted* participants:  
the identities of participants is known to all!

## Permissioned Blockchains

*Permissioned* blockchains have a well-defined set of *vetted* participants:  
the identities of participants is known to all!

- ▶ ✗ No easy way to change the participants.
- ▶ ✗ Dozens of participants.
- ▶ ✓ Many thousands of transactions per second with low latencies.
- ▶ ✓ Strong consensus guarantees: no data inconsistencies.

## Permissioned Blockchains

*Permissioned* blockchains have a well-defined set of *vetted* participants:  
the identities of participants is known to all!

- ▶ ✗ No easy way to change the participants.
- ▶ ✗ Dozens of participants.
- ▶ ✓ Many thousands of transactions per second with low latencies.
- ▶ ✓ Strong consensus guarantees: no data inconsistencies.

Example: A e-health system managed by a consortium of health-care providers

- ▶ Each *vetted participant* manages their own systems (e.g., different software).
- ▶ Failure of individual participants should not break the system!
- ▶ Blockchains: *federated data management* and *resilience*.

## Permissioned Blockchains

*Permissioned* blockchains have a well-defined set of *vetted* participants:  
the identities of participants is known to all!

- ▶ ✗ No easy way to change the participants.
- ▶ ✗ Dozens of participants.
- ▶ ✓ Many thousands of transactions per second with low latencies.
- ▶ ✓ Strong consensus guarantees: no data inconsistencies.

Example: A e-health system managed by a consortium of health-care providers

- ▶ Each *vetted participant* manages their own systems (e.g., different software).
- ▶ Failure of individual participants should not break the system!
- ▶ Blockchains: *federated data management* and *resilience*.

Can be implemented using PBFT-style consensus.

## Permissionless Blockchains

*Permissionless* blockchains have *open membership*:  
anyone can participate, at any time, for any duration.

# Permissionless Blockchains

*Permissionless* blockchains have *open membership*:  
anyone can participate, at any time, for any duration.

- ▶ ✓ Participants are not known and can change.
- ▶ ✓ Thousands of participants.
- ▶ ✗ Few transactions per second with very high latencies.
- ▶ ✗ Weak consensus guarantees: data can be inconsistent (*forks*).

## Permissionless Blockchains

*Permissionless* blockchains have *open membership*:  
anyone can participate, at any time, for any duration.

- ▶ ✓ Participants are not known and can change.
- ▶ ✓ Thousands of participants.
- ▶ ✗ Few transactions per second with very high latencies.
- ▶ ✗ Weak consensus guarantees: data can be inconsistent (*forks*).

Example: Cryptocurrencies such as Bitcoin and Ethereum

# Permissionless Blockchains

*Permissionless* blockchains have *open membership*:  
anyone can participate, at any time, for any duration.

- ▶ ✓ Participants are not known and can change.
- ▶ ✓ Thousands of participants.
- ▶ ✗ Few transactions per second with very high latencies.
- ▶ ✗ Weak consensus guarantees: data can be inconsistent (*forks*).

Example: Cryptocurrencies such as Bitcoin and Ethereum

Challenge: *Sybil attacks*

A single entity can impersonate *many* participants to gain unfair control over the system.

# Permissionless Blockchains

*Permissionless* blockchains have *open membership*:  
anyone can participate, at any time, for any duration.

- ▶ ✓ Participants are not known and can change.
- ▶ ✓ Thousands of participants.
- ▶ ✗ Few transactions per second with very high latencies.
- ▶ ✗ Weak consensus guarantees: data can be inconsistent (*forks*).

Example: Cryptocurrencies such as Bitcoin and Ethereum

Challenge: *Sybil attacks*

A single entity can impersonate *many* participants to gain unfair control over the system.

Requires a consensus protocol that does not rely on *identities*.

# Distributed Systems

## Definition<sup>1</sup>

A distributed system is a collection of autonomous computing elements that appears to its users as a single coherent system.

---

<sup>1</sup>M. van Steen & A.S. Tanenbaum, *Distributed Systems*, 3rd ed., [distributed-systems.net](http://distributed-systems.net), 2017.

# Distributed Systems

## Definition<sup>1</sup>

A distributed system is a collection of autonomous computing elements that appears to its users as a single coherent system.

- ▶ Autonomous computing elements: nodes, replica, ...

---

<sup>1</sup>M. van Steen & A.S. Tanenbaum, *Distributed Systems*, 3rd ed., [distributed-systems.net](http://distributed-systems.net), 2017.

# Distributed Systems

## Definition<sup>1</sup>

A distributed system is a collection of autonomous computing elements that appears to its users as a single coherent system.

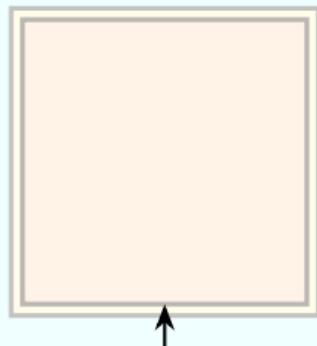
- ▶ Autonomous computing elements: nodes, replica, ....
- ▶ Single coherent system: “should behave as a single system”.
  - ▶ Users should see a single system.
  - ▶ Nodes must collaborate to provide that system.

---

<sup>1</sup>M. van Steen & A.S. Tanenbaum, *Distributed Systems*, 3rd ed., [distributed-systems.net](http://distributed-systems.net), 2017.

# Distributed Systems: Scalability

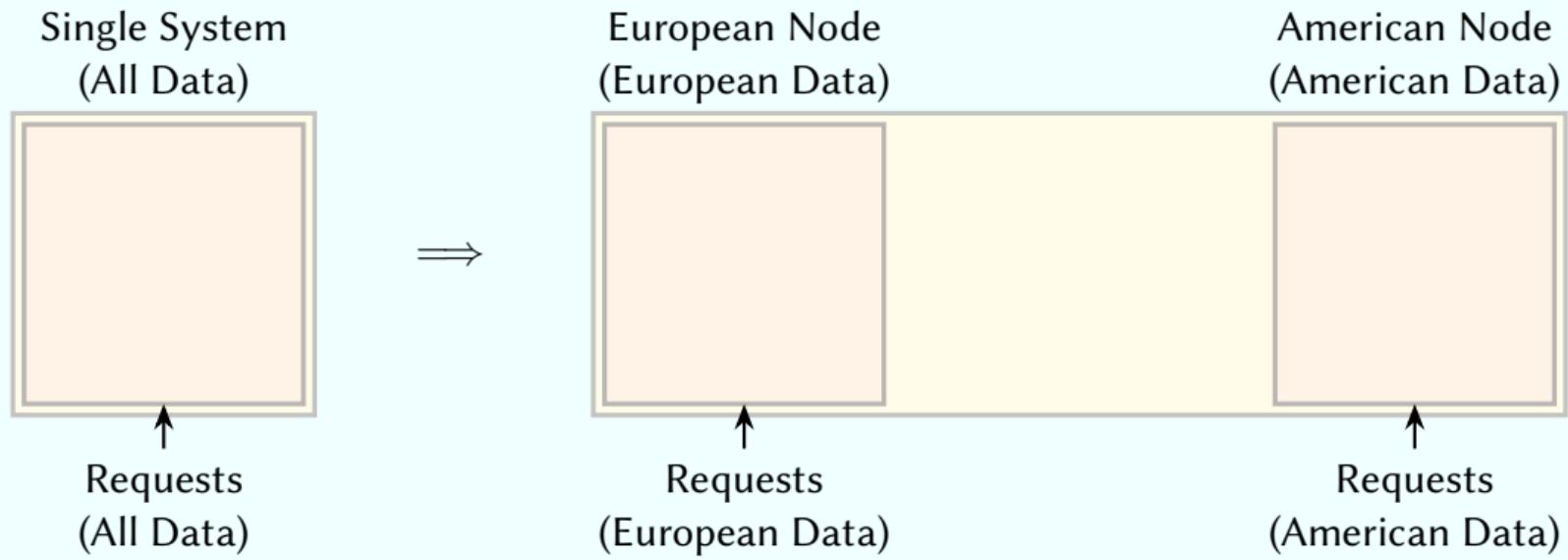
Single System  
(All Data)



Requests  
(All Data)

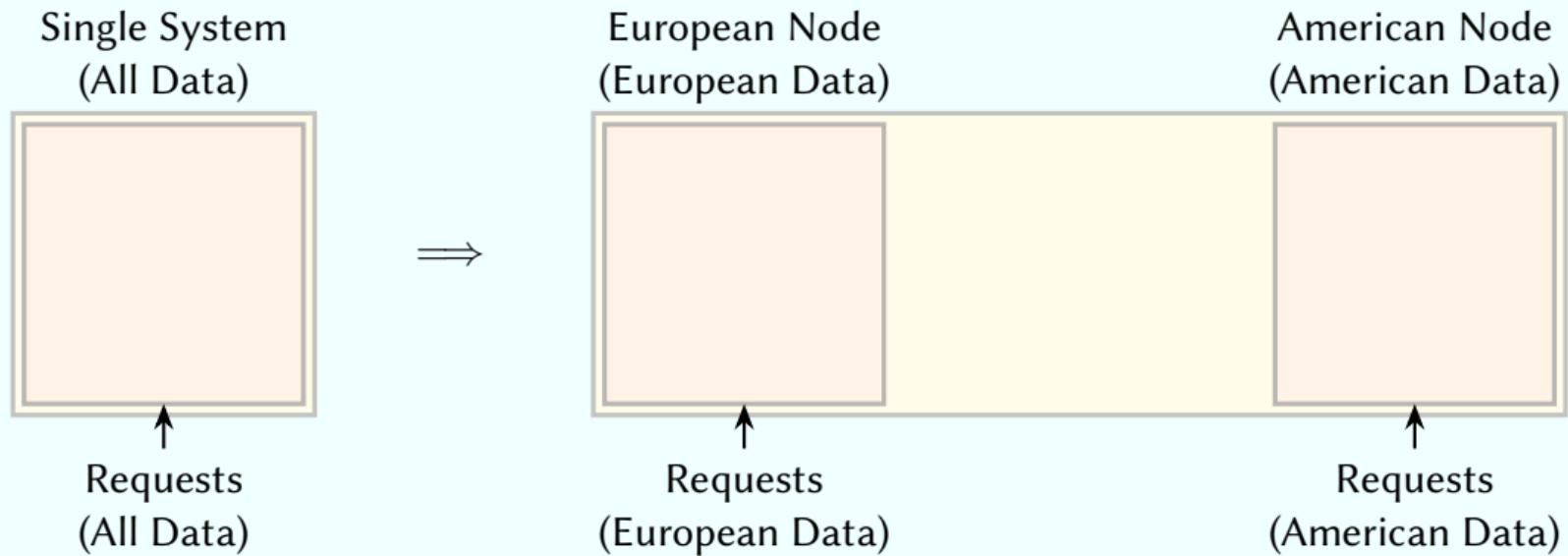
Single system: Storage and performance *bounded* by hardware.

# Distributed Systems: Scalability



Partition the system: More storage and *potentially* more performance.

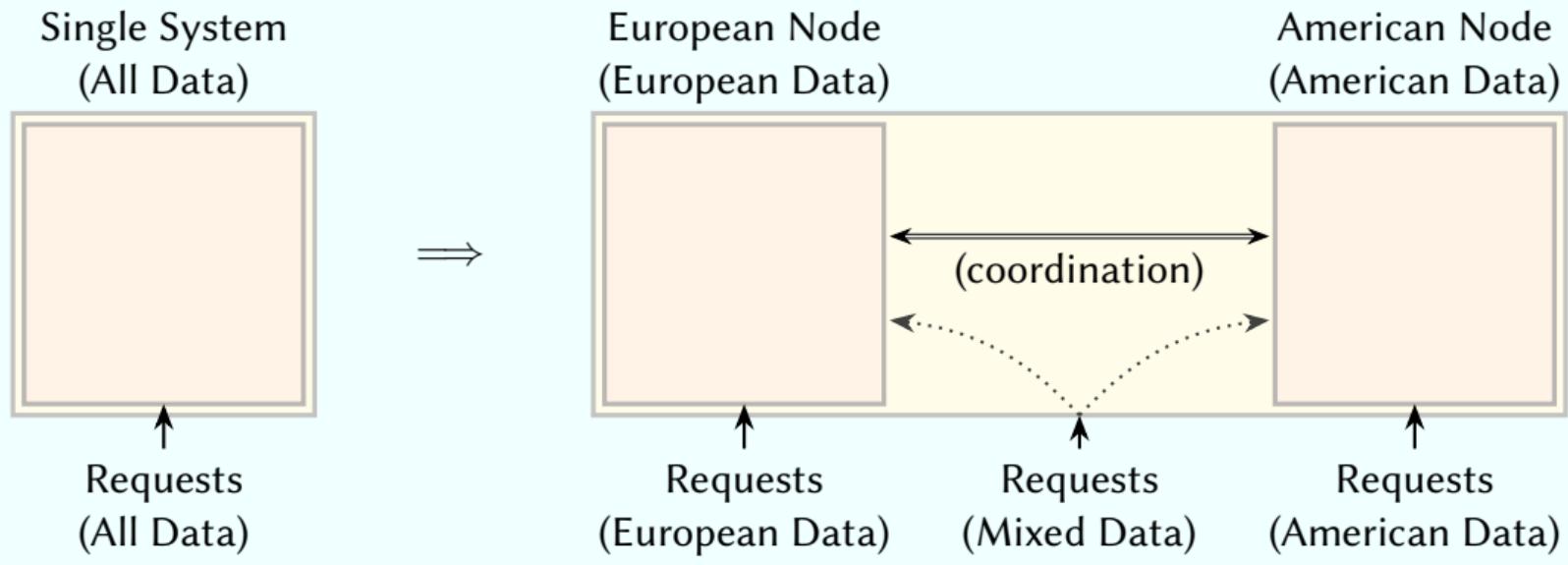
# Distributed Systems: Scalability



Partition the system: More storage and *potentially* more performance.

Potentially *lower latencies* if data ends up closer to users.

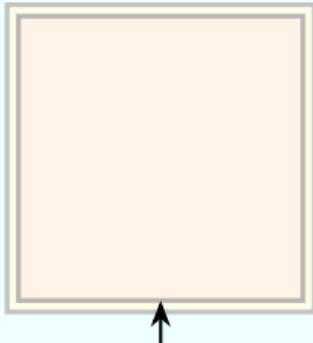
# Distributed Systems: Scalability



Complex requests: become more *costly* to answer!

# Distributed Systems: Specialization

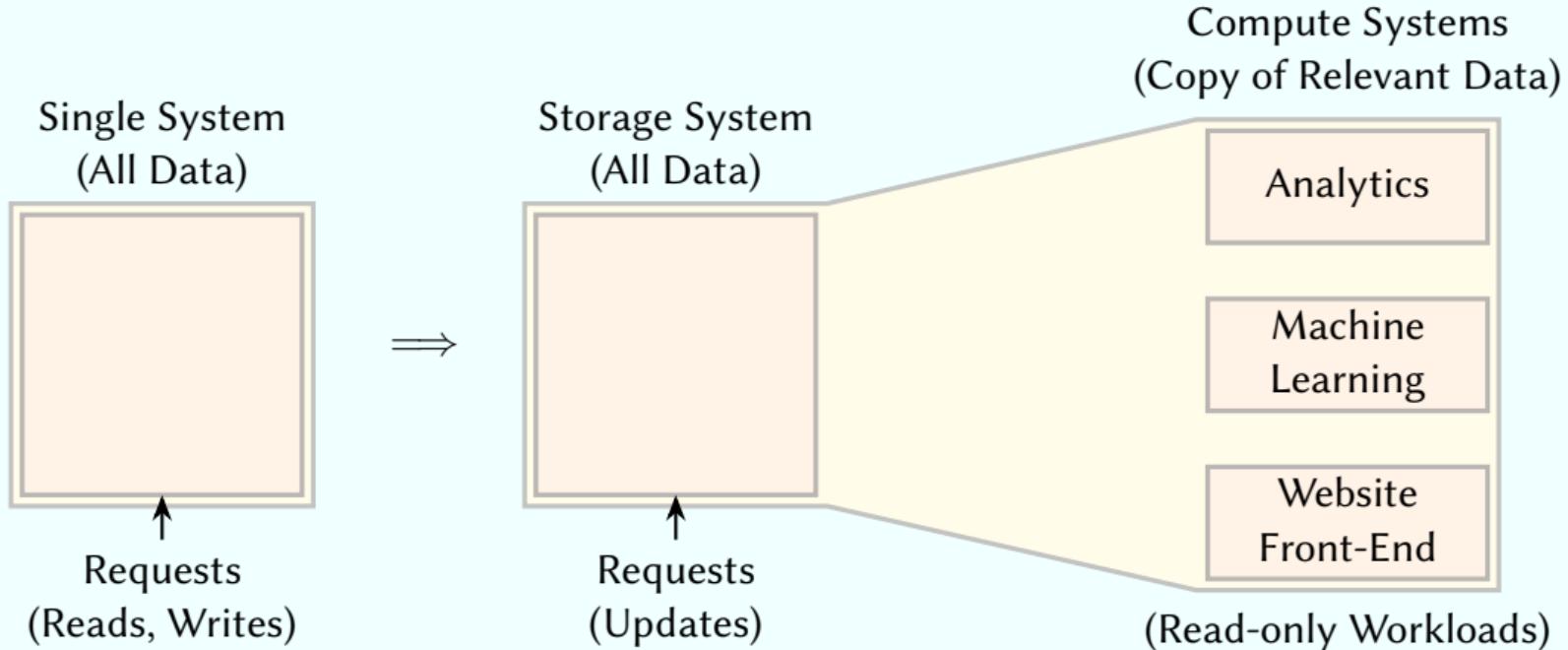
Single System  
(All Data)



Requests  
(Reads, Writes)

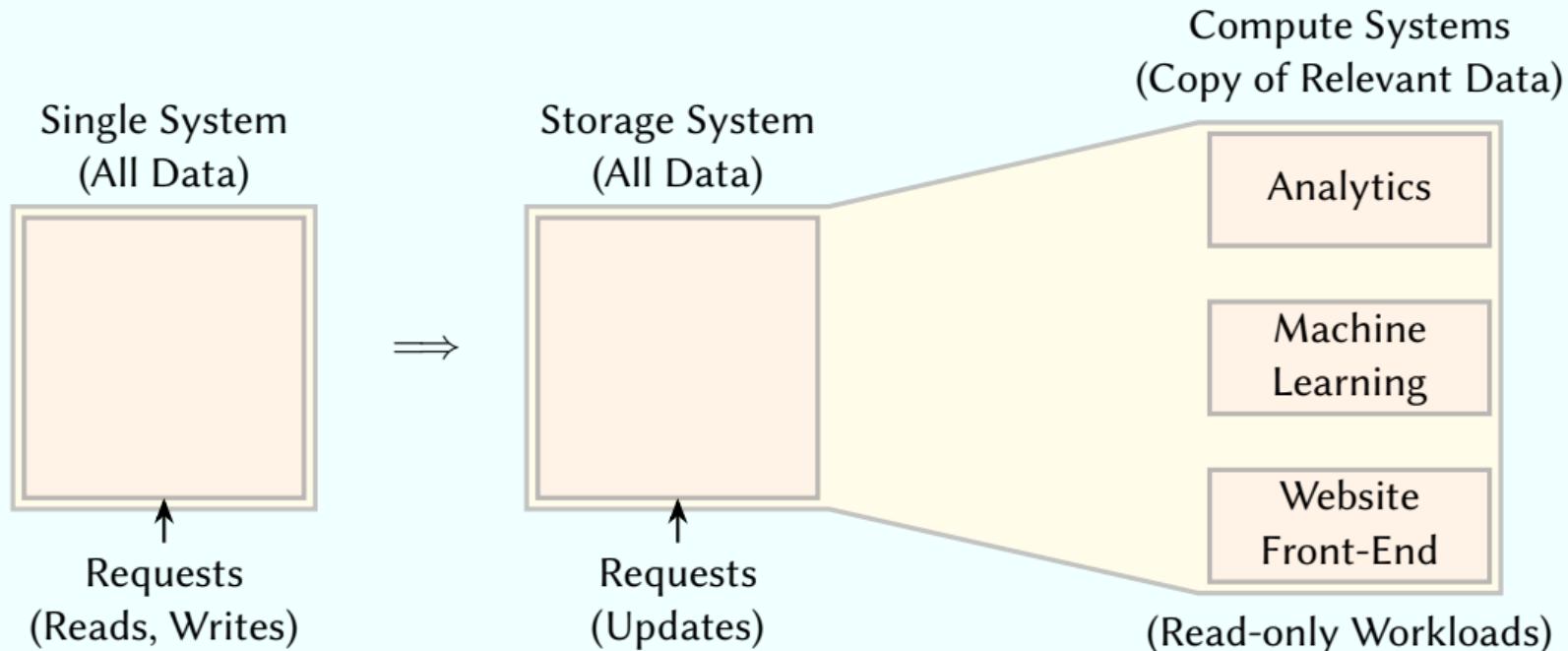
Single system: Compromise—cannot be optimized for *all* tasks.

# Distributed Systems: Specialization



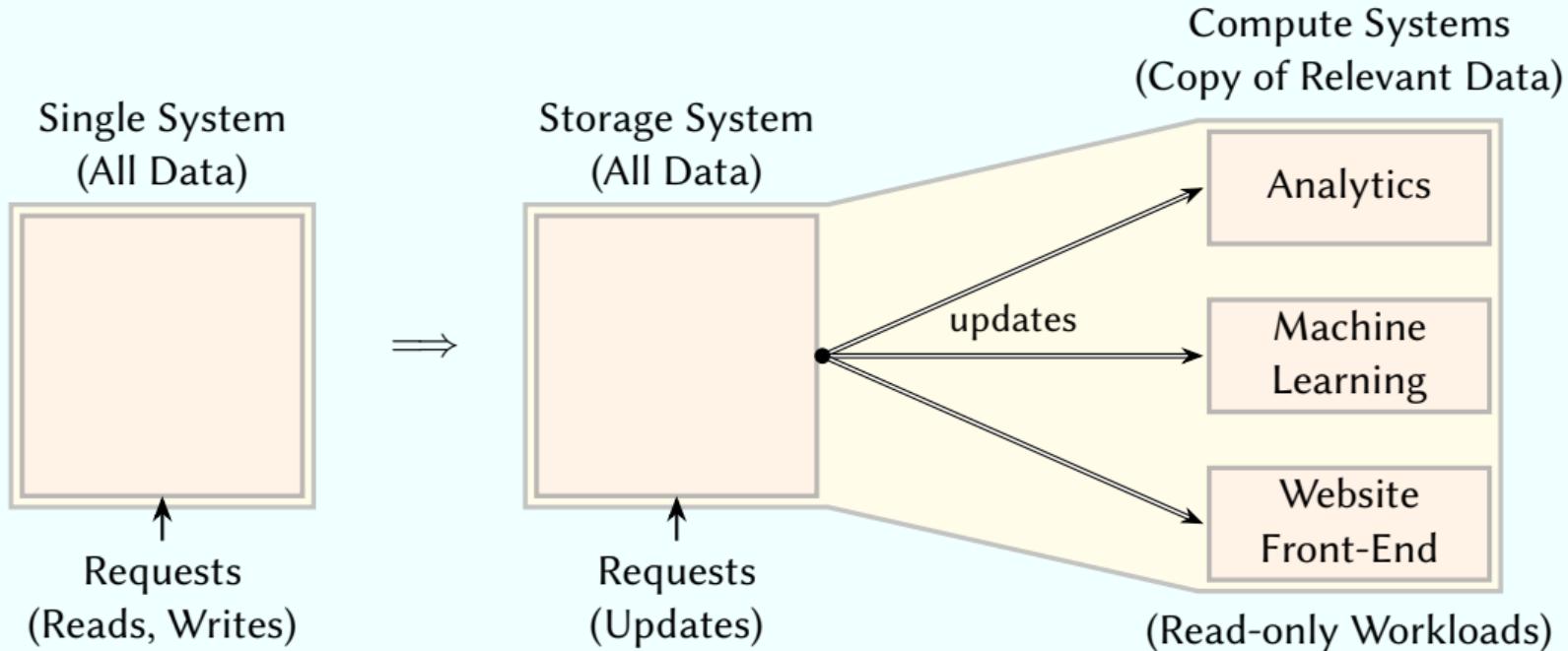
Specialize the system: Different nodes have distinct tasks.

# Distributed Systems: Specialization



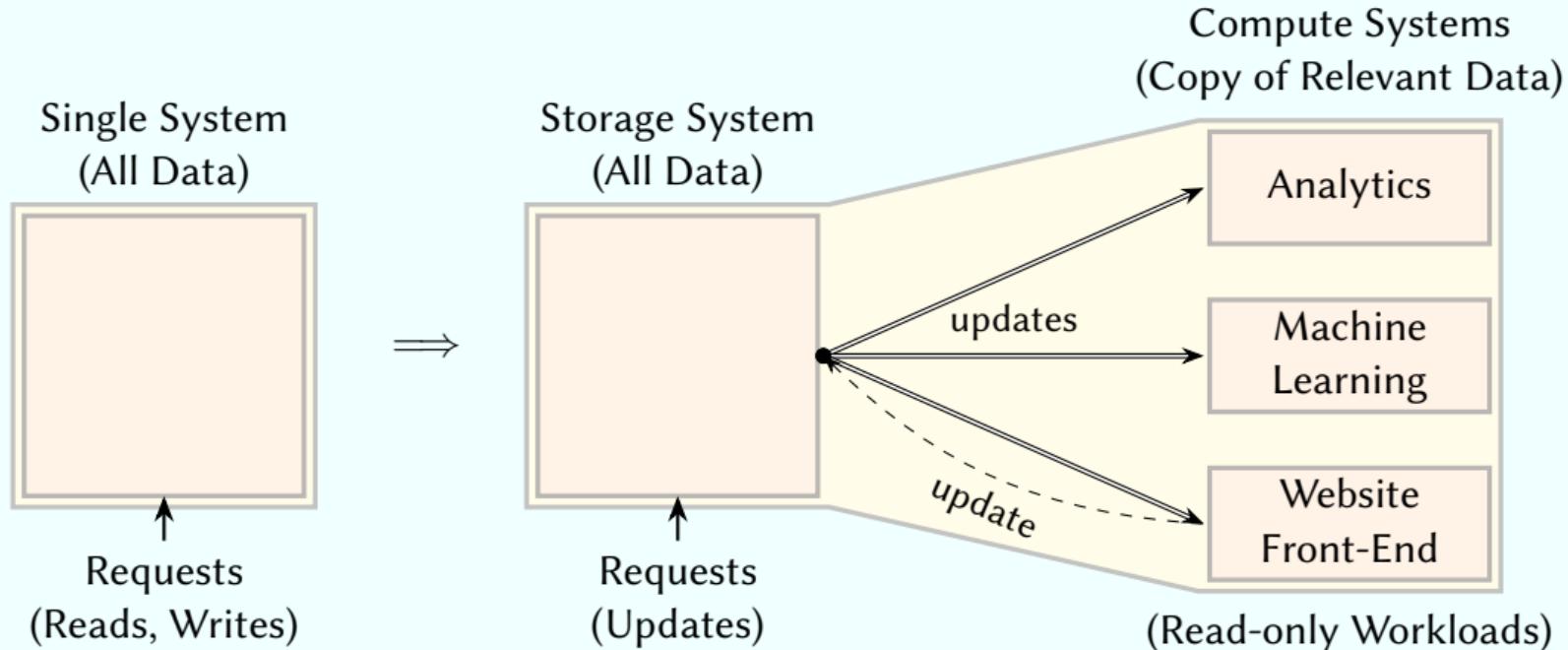
Specialize the system: Different nodes have distinct tasks.  
Specialized hardware and software *per* task.

# Distributed Systems: Specialization



Added cost: Keeping the compute systems *up-to-date*.

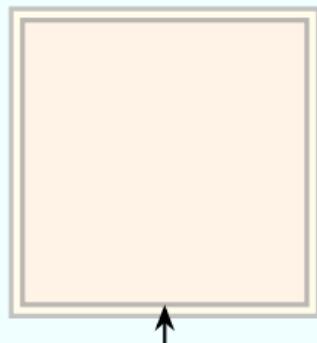
# Distributed Systems: Specialization



Design complexity: Updates from the compute systems?

# Distributed Systems: Reliability (Primary-Backup)

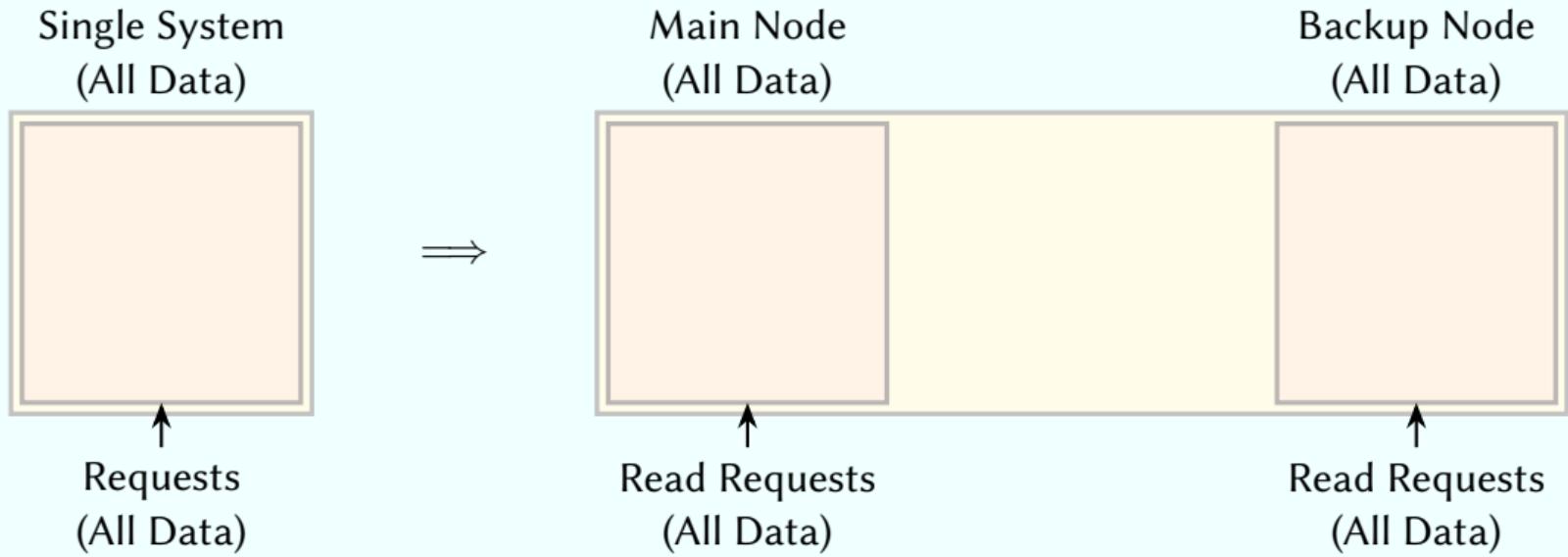
Single System  
(All Data)



Requests  
(All Data)

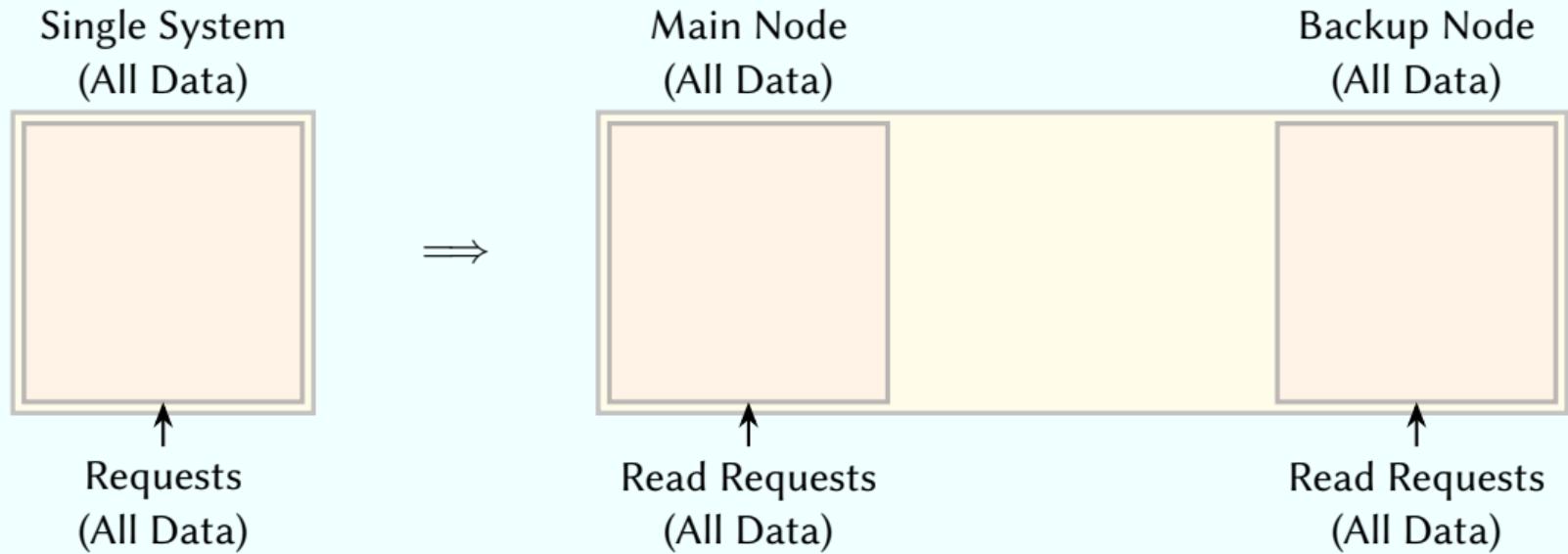
Single system: Single point of *failure*.

# Distributed Systems: Reliability (Primary-Backup)



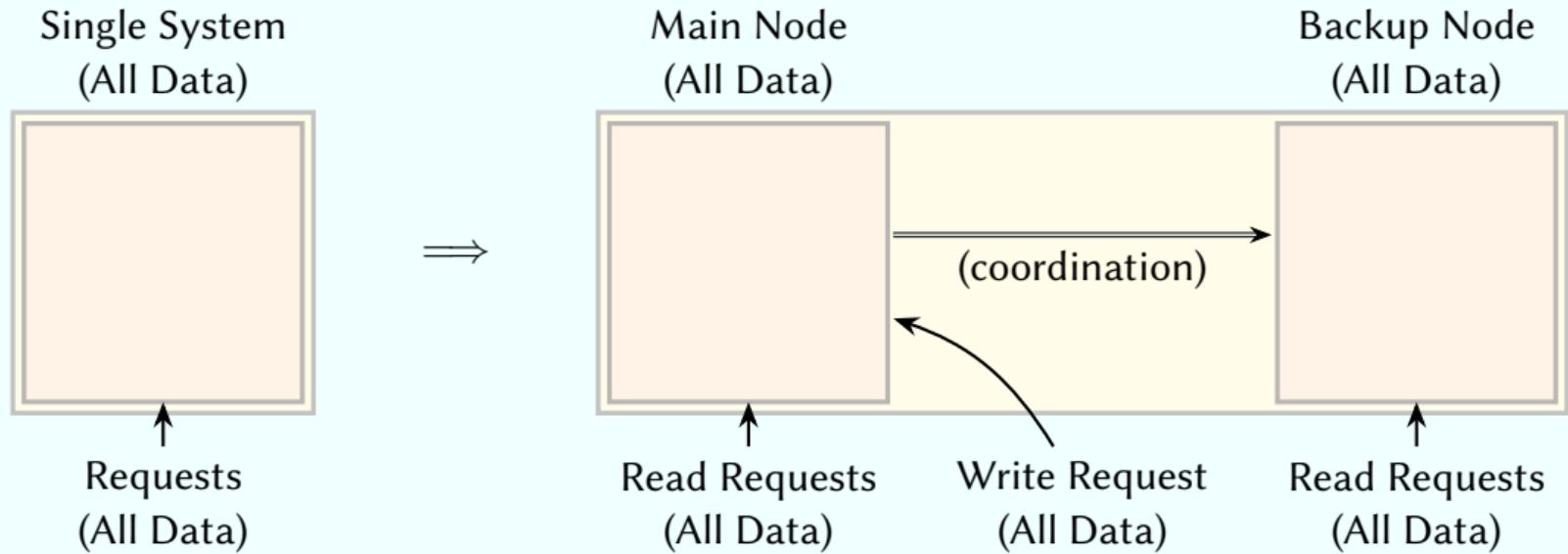
Multiple copies: Copies available after single failure.

# Distributed Systems: Reliability (Primary-Backup)



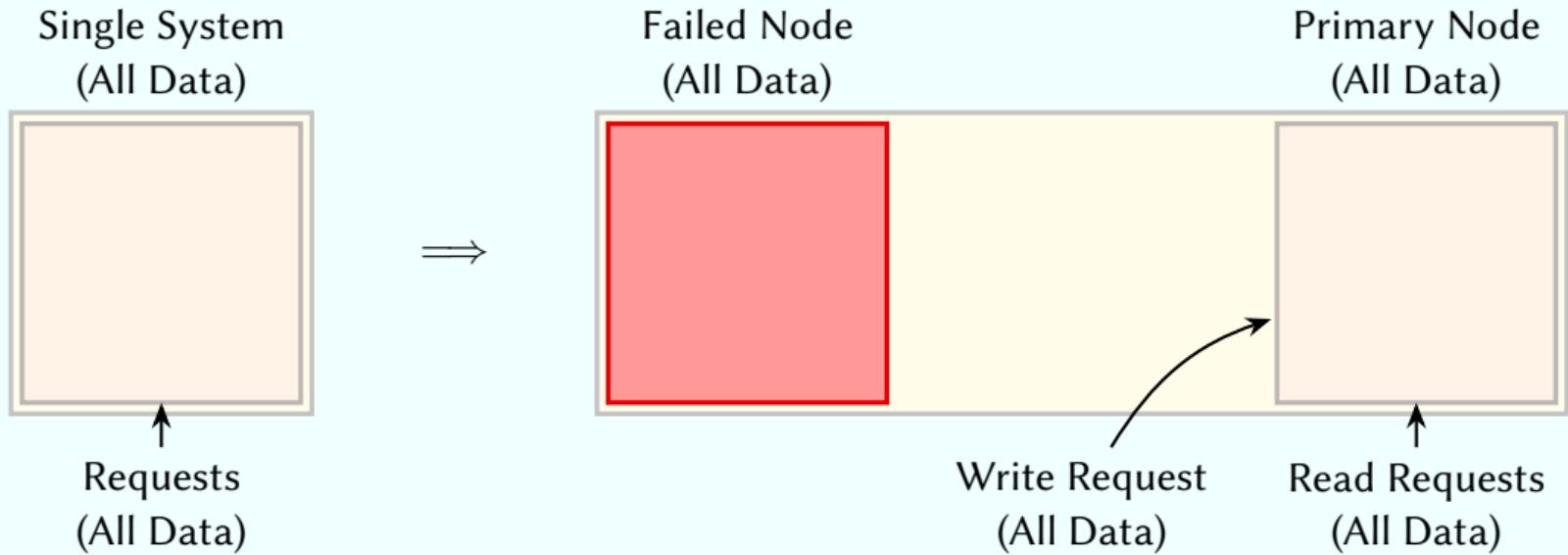
Multiple copies: Copies available after single failure.  
Potentially *lower latencies & more performance* when users read data.

# Distributed Systems: Reliability (Primary-Backup)



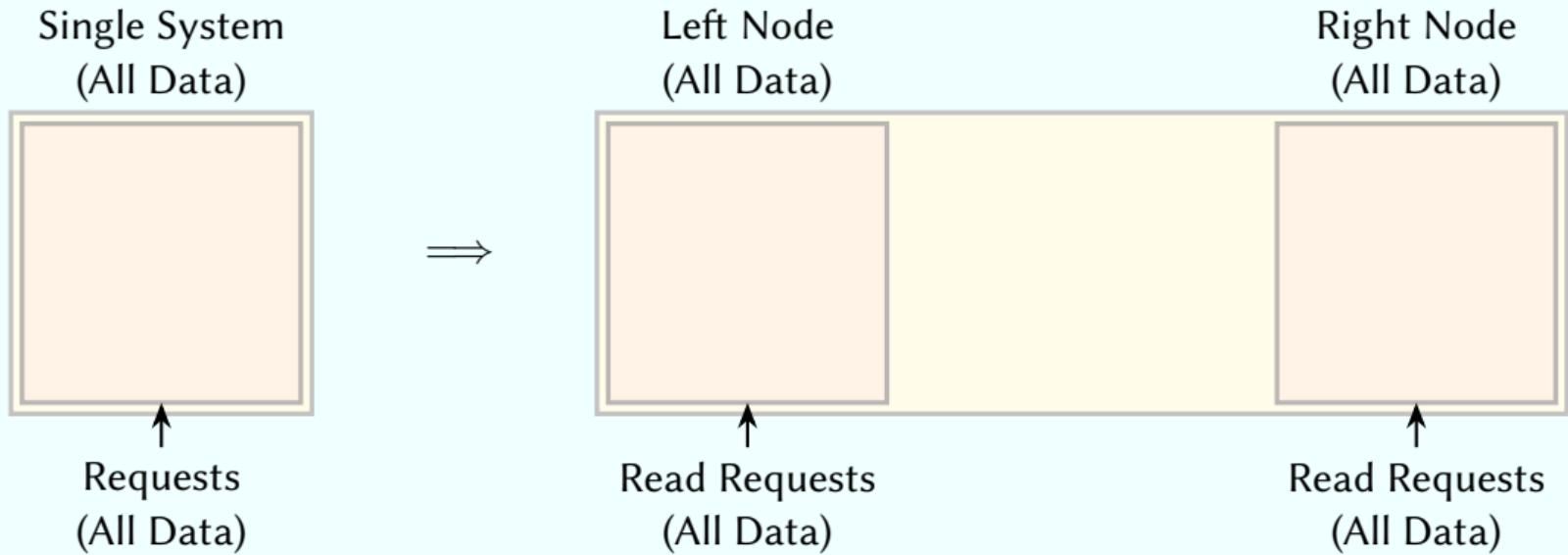
Changes to data requests: Primary *leads*, backups *follow*.

# Distributed Systems: Reliability (Primary-Backup)



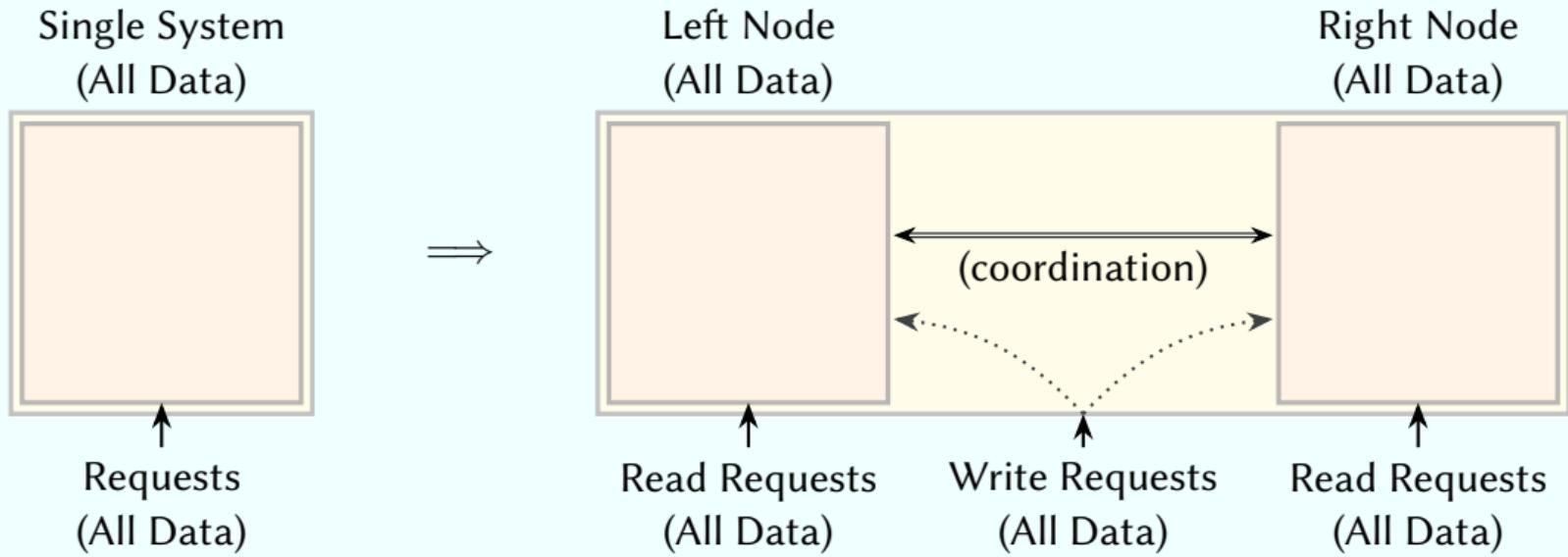
Failure: Recovery mechanisms—typically *complex*.

## Distributed Systems: Reliability (Decentralized)



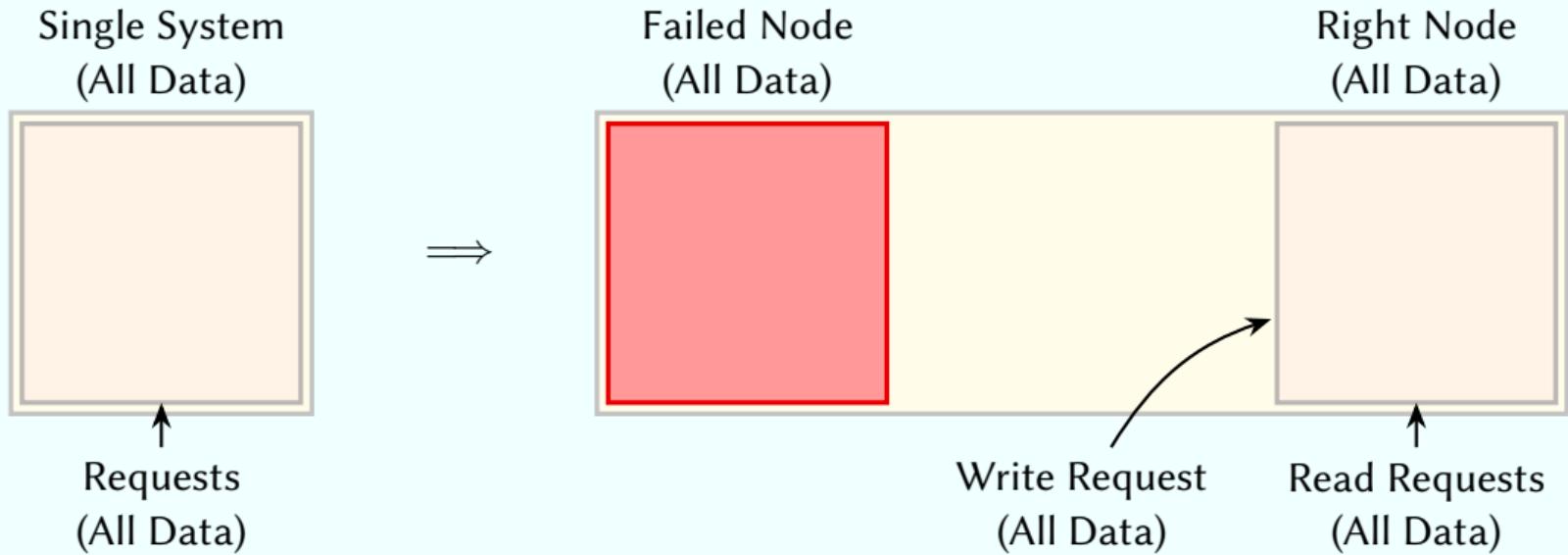
Multiple copies: Copies available after single failure.  
Potentially *lower latencies & more performance* when users read data.

# Distributed Systems: Reliability (Decentralized)



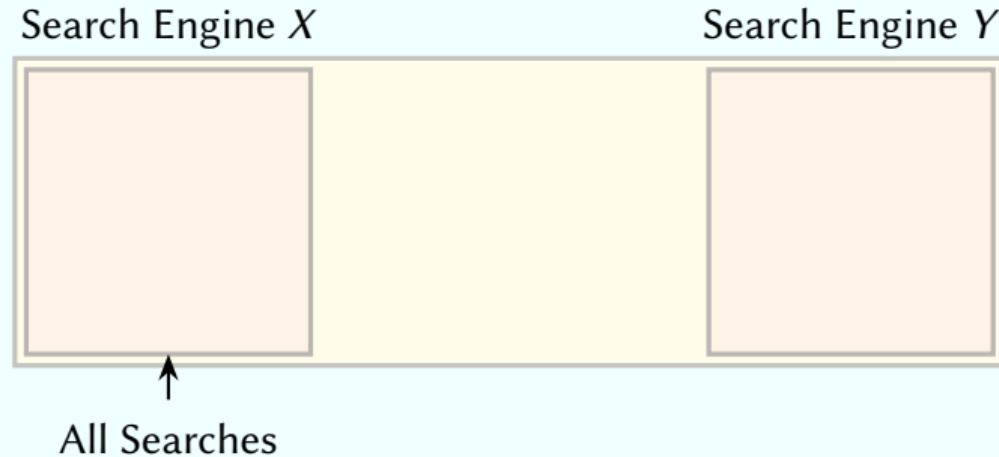
Changes to data requests: *Costly* coordinated decision among all active nodes.

## Distributed Systems: Reliability (Decentralized)



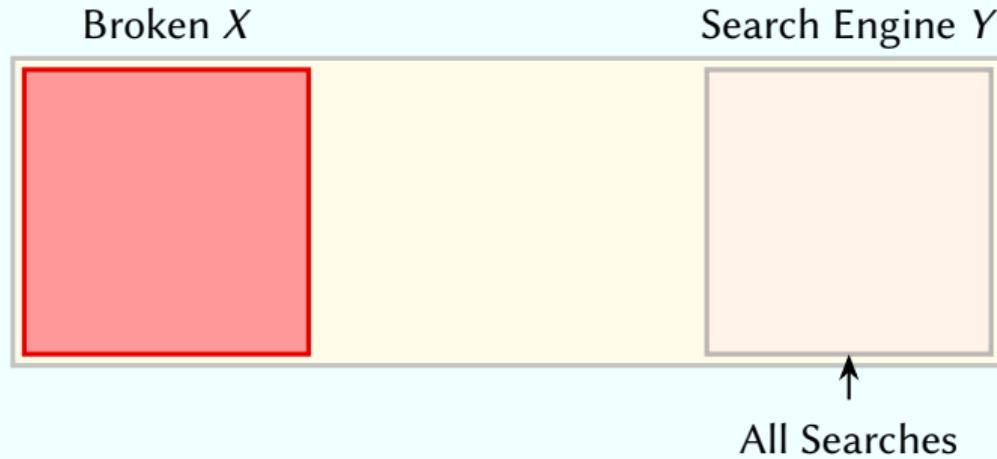
Failure: Recovery mechanisms—typically *easier* (or even *free*).

# Distributed System?



I perform all my searches on the web with Search Engine X.

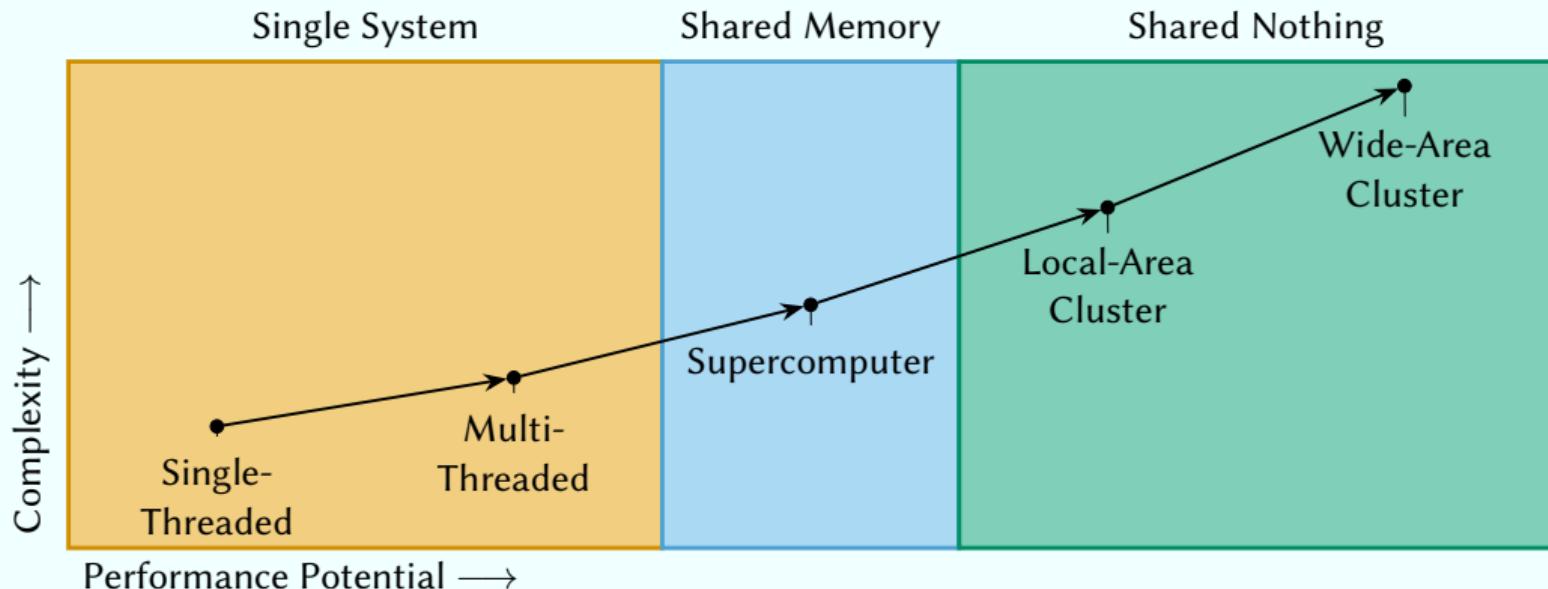
# Distributed System?



Today Search Engine  $X$  failed—I remembered the alternative Search Engine  $Y$ .

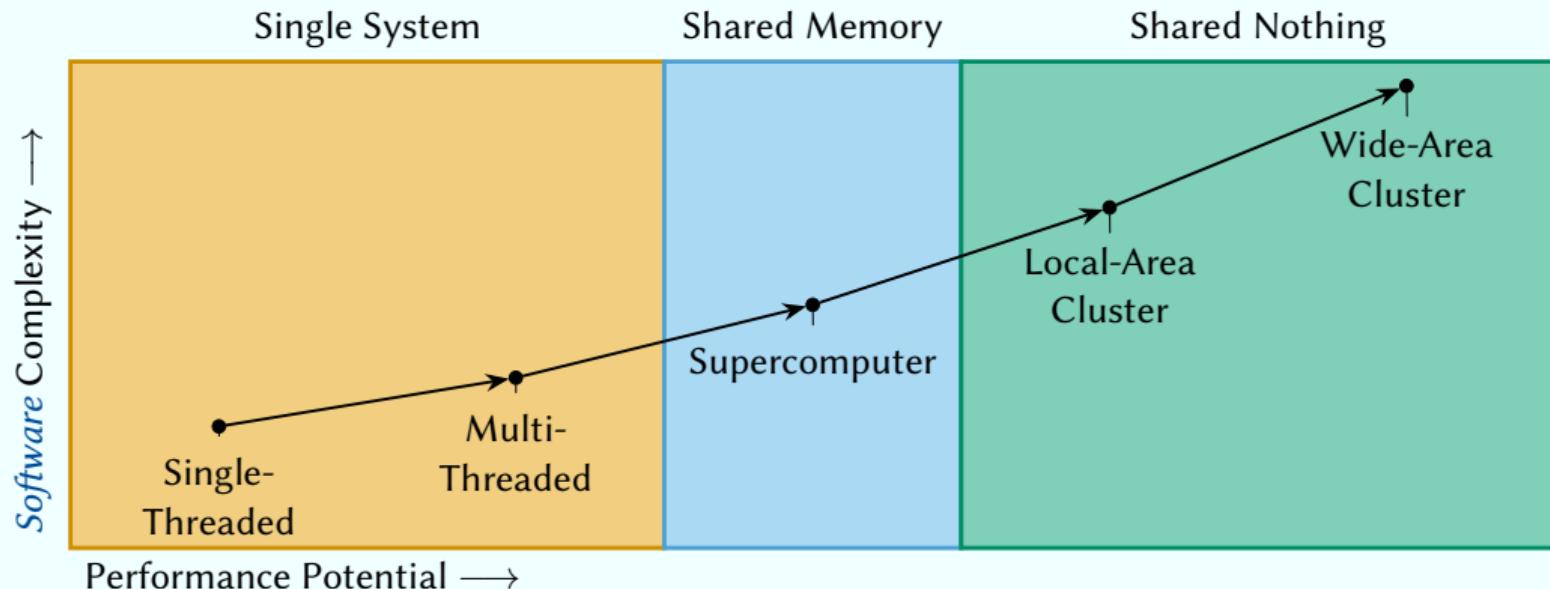
# Do we need Distributed Systems?

Distributed computing is *complex*



# Do we need Distributed Systems?

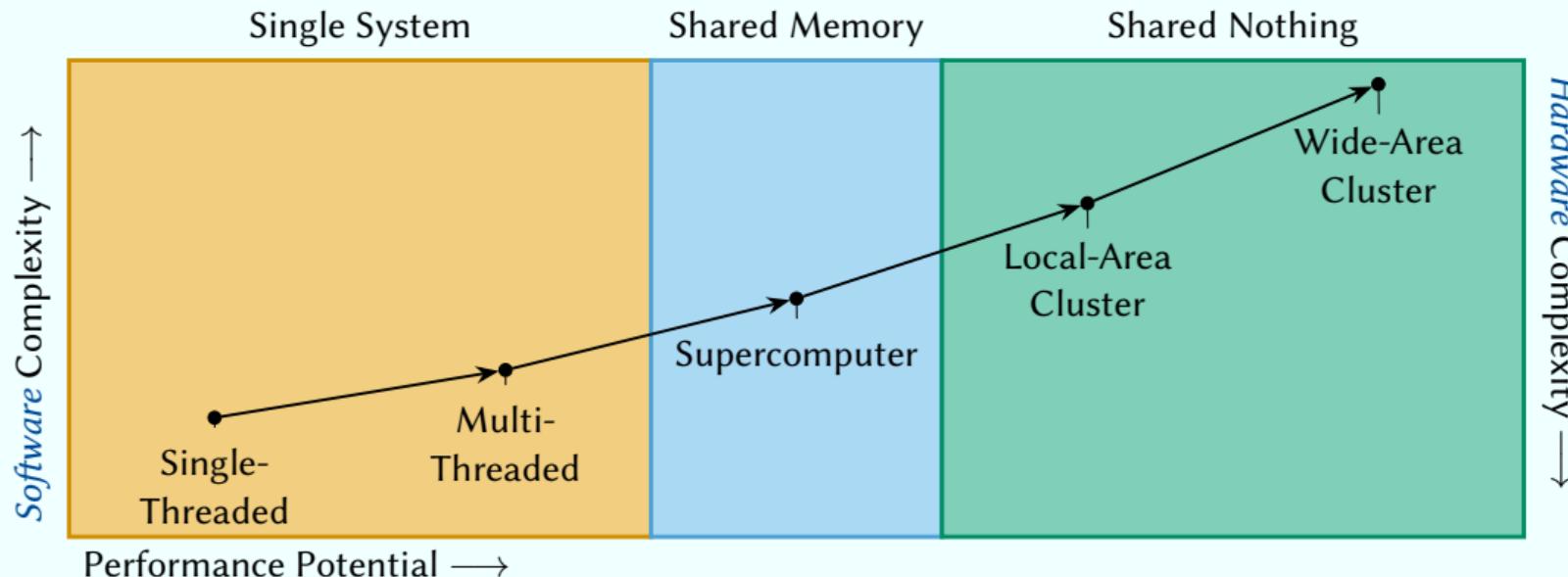
Distributed computing is *complex*



Complexity to write *efficient* software that solve a given problem.

# Do we need Distributed Systems?

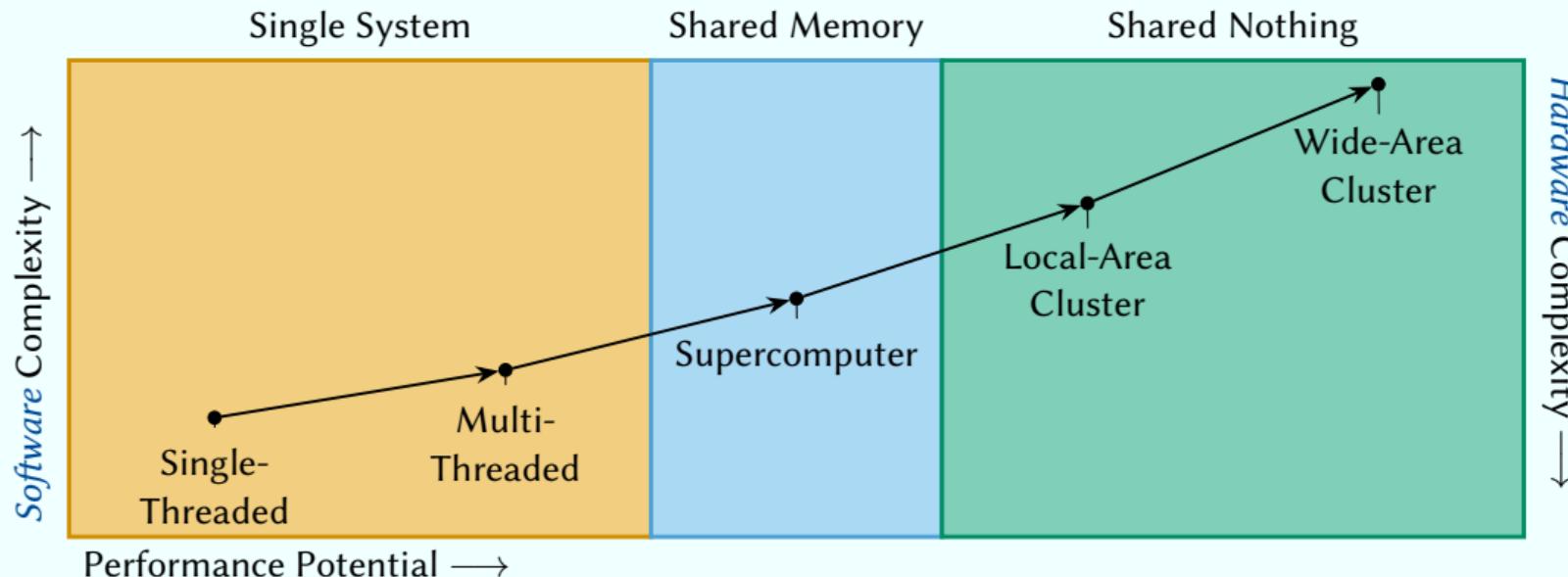
Distributed computing is *complex*



Complexity of *hardware* to solve a given problem in time  $t$ .

# Do we need Distributed Systems?

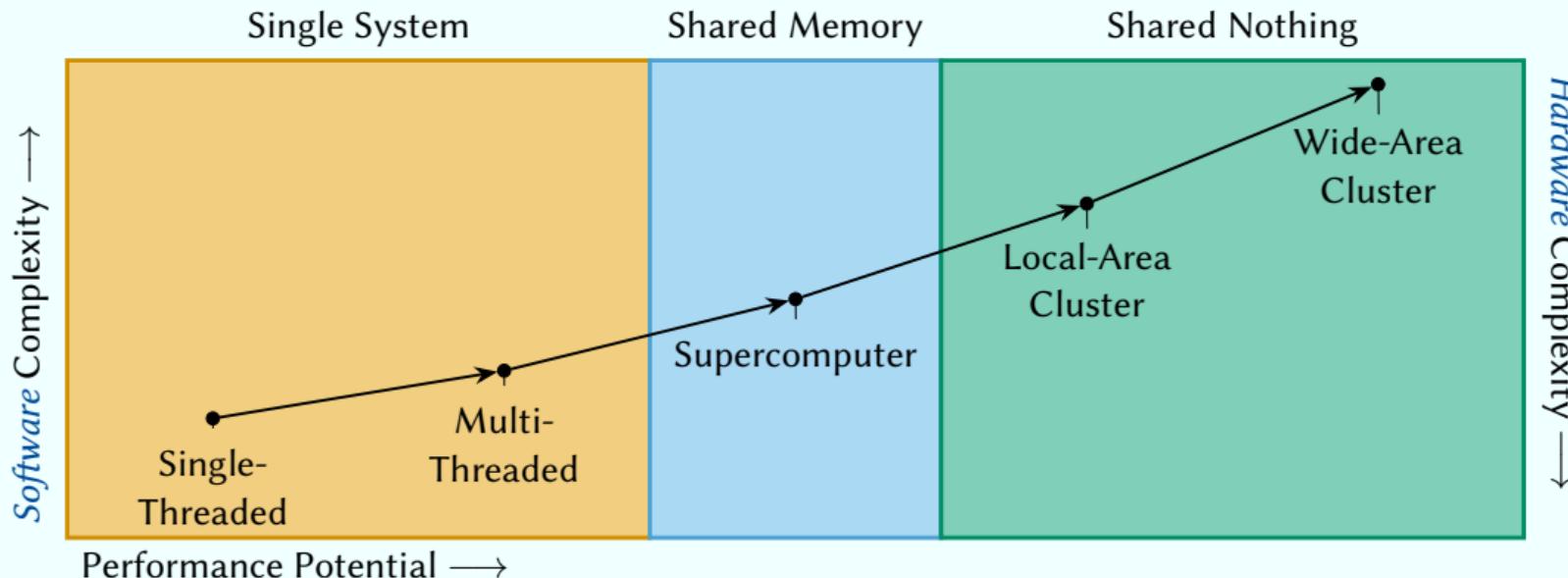
Distributed computing is *complex*



Increased software complexity is *fundamental*: Problems become theoretically harder.

# Do we need Distributed Systems?

Distributed computing is *complex*



Do we really need distributed systems?

## Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure
- ▶ Deployment failure
- ▶ Hardware failure
- ▶ Network failure
- ▶ Datacenter failure
- ▶ Malicious attacks

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure      Code practices? Regression tests? *Independent implementations?*
- ▶ Deployment failure
- ▶ Hardware failure
- ▶ Network failure
- ▶ Datacenter failure
- ▶ Malicious attacks

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure      Code practices? Regression tests? *Independent implementations?*
- ▶ Deployment failure    Deployment testing? *Staged deployment?*
- ▶ Hardware failure
- ▶ Network failure
- ▶ Datacenter failure
- ▶ Malicious attacks

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure      Code practices? Regression tests? *Independent implementations?*
- ▶ Deployment failure    Deployment testing? *Staged deployment?*
- ▶ Hardware failure      Use redundant, hot-swappable hardware? *Hot spares?*
- ▶ Network failure
- ▶ Datacenter failure
- ▶ Malicious attacks

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure      Code practices? Regression tests? *Independent implementations?*
- ▶ Deployment failure    Deployment testing? *Staged deployment?*
- ▶ Hardware failure     Use redundant, hot-swappable hardware? *Hot spares?*
- ▶ Network failure       Network redundancy? *Remote hot spares?*
- ▶ Datacenter failure
- ▶ Malicious attacks

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure      Code practices? Regression tests? *Independent implementations?*
- ▶ Deployment failure    Deployment testing? *Staged deployment?*
- ▶ Hardware failure     Use redundant, hot-swappable hardware? *Hot spares?*
- ▶ Network failure       Network redundancy? *Remote hot spares?*
- ▶ Datacenter failure    *Remote hot spares?*
- ▶ Malicious attacks

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure      Code practices? Regression tests? *Independent implementations?*
- ▶ Deployment failure    Deployment testing? *Staged deployment?*
- ▶ Hardware failure     Use redundant, hot-swappable hardware? *Hot spares?*
- ▶ Network failure       Network redundancy? *Remote hot spares?*
- ▶ Datacenter failure    *Remote hot spares?*
- ▶ Malicious attacks    *Decentralized & independent implementations?*

# Distributed System Need: Resilience

Consider: a *service* with a very high SLA (may almost never stop).

## Causes of failure

- ▶ Software failure      Code practices? Regression tests? *Independent implementations?*
- ▶ Deployment failure    Deployment testing? *Staged deployment?*
- ▶ Hardware failure     Use redundant, hot-swappable hardware? *Hot spares?*
- ▶ Network failure       Network redundancy? *Remote hot spares?*
- ▶ Datacenter failure    *Remote hot spares?*
- ▶ Malicious attacks    *Decentralized & independent implementations?*

*Distributed* designs can provide resilience.

# Desirable Properties of a Distributed System

Consider a simple & minimalistic distributed system

Fully-replicated: each node (replica) holds the *same data*.

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# Desirable Properties of a Distributed System

Consider a simple & minimalistic distributed system

Fully-replicated: each node (replica) holds the *same data*.

The CAP Properties<sup>2</sup>

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# Desirable Properties of a Distributed System

Consider a simple & minimalistic distributed system

Fully-replicated: each node (replica) holds the *same data*.

## The CAP Properties<sup>2</sup>

**Consistency** All replicas have a single up-to-date copy of the data.  
(All replicas always have the same date).

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# Desirable Properties of a Distributed System

Consider a simple & minimalistic distributed system

Fully-replicated: each node (replica) holds the *same data*.

## The CAP Properties<sup>2</sup>

**Consistency** All replicas have a single up-to-date copy of the data.  
(All replicas always have the same date).

**Availability** All data is always accessible & updateable.  
(The system can always provide services to clients).

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# Desirable Properties of a Distributed System

Consider a simple & minimalistic distributed system

Fully-replicated: each node (replica) holds the *same data*.

## The CAP Properties<sup>2</sup>

**Consistency** All replicas have a single up-to-date copy of the data.  
(All replicas always have the same date).

**Availability** All data is always accessible & updateable.  
(The system can always provide services to clients).

**Partitioning** The system can tolerate network partitions.  
(The system can cope with network failures).

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# Desirable Properties of a Distributed System

Consider a simple & minimalistic distributed system

Fully-replicated: each node (replica) holds the *same data*.

## The CAP Properties<sup>2</sup>

**Consistency** All replicas have a single up-to-date copy of the data.  
(All replicas always have the same date).

**Availability** All data is always accessible & updateable.  
(The system can always provide services to clients).

**Partitioning** The system can tolerate network partitions.  
(The system can cope with network failures).

## Theorem

A *fully-replicated* system can only provide two-out-of-three CAP properties.

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# The CAP Theorem: Revisited<sup>2</sup>

The CAP Theorem lacks nuance.

**Consistency** All replicas have a single up-to-date copy of the data.

**Availability** All data is always accessible & updateable.

**Partitioning** The system can tolerate network partitions.

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# The CAP Theorem: Revisited<sup>2</sup>

The CAP Theorem lacks nuance.

**Consistency** All replicas have a single up-to-date copy of the data.

**Nuance:** Perfect consistency is a very high bar & will limit performance.

**Alternatives:** eventual consistency & state convergence.

**Availability** All data is always accessible & updateable.

**Partitioning** The system can tolerate network partitions.

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# The CAP Theorem: Revisited<sup>2</sup>

The CAP Theorem lacks nuance.

**Consistency** All replicas have a single up-to-date copy of the data.

**Nuance:** Perfect consistency is a very high bar & will limit performance.

**Alternatives:** eventual consistency & state convergence.

**Availability** All data is always accessible & updateable.

**Nuance:** What about latency?

**Alternatives:** Read-only mode (e.g., offline mode), flexible response times.

**Partitioning** The system can tolerate network partitions.

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# The CAP Theorem: Revisited<sup>2</sup>

The CAP Theorem lacks nuance.

**Consistency** All replicas have a single up-to-date copy of the data.

**Nuance:** Perfect consistency is a very high bar & will limit performance.

**Alternatives:** eventual consistency & state convergence.

**Availability** All data is always accessible & updateable.

**Nuance:** What about latency?

**Alternatives:** Read-only mode (e.g., offline mode), flexible response times.

**Partitioning** The system can tolerate network partitions.

**Nuance:** For non-mobile systems, partition are rare.

**Alternatives:** aim at CA & compensation after partition recovery.

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# The CAP Theorem: Revisited<sup>2</sup>

The CAP Theorem lacks nuance.

**Consistency** All replicas have a single up-to-date copy of the data.

**Nuance:** Perfect consistency is a very high bar & will limit performance.

**Alternatives:** eventual consistency & state convergence.

**Availability** All data is always accessible & updateable.

**Nuance:** What about latency?

**Alternatives:** Read-only mode (e.g., offline mode), flexible response times.

**Partitioning** The system can tolerate network partitions.

**Nuance:** For non-mobile systems, partition are rare.

**Alternatives:** aim at CA & compensation after partition recovery.

There is a huge design space—CAP only covers a tiny part.

---

<sup>2</sup>E. Brewer, *CAP Twelve Years Later: How the “Rules” Have Changed*, 2012.

# Distributed Databases

## Definition<sup>3</sup>

**Distributed database** A collection of multiple, logically interrelated databases located at the nodes of a distributed system.

**Distributed DBMS** the software system that manages the distributed database & makes the distribution transparent to the users

---

<sup>3</sup>M.T. Özsu & P. Valduriez, *Principles of Distributed Database Systems*, 4th ed., 2020.

# Resilient Systems

## Definition

A resilient system is a system that can tolerate failures while continuously providing service.

# Resilient Systems

## Definition

A resilient system is a system that can **tolerate failures** while continuously providing service.

- ▶ *Failure*: hardware, software, network, malicious attacks, ....  
Crash fault tolerance versus Byzantine fault tolerance.

# Resilient Systems

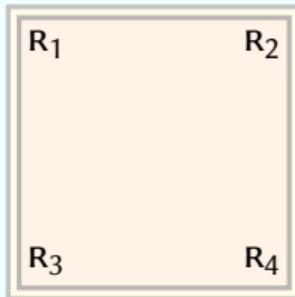
## Definition

A resilient system is a system that can **tolerate failures** while **continuously providing service**.

- ▶ *Failure*: hardware, software, network, malicious attacks, ....  
Crash fault tolerance versus Byzantine fault tolerance.
- ▶ *Continuous services*: No downtime, manual intervention, restarts, ....

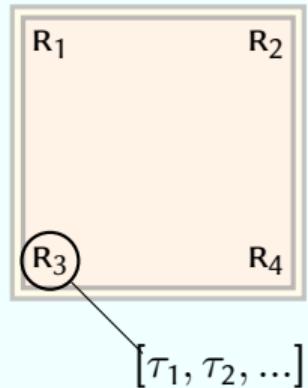
# Outline of a Blockchain-Based Resilient System

## 1. Replicas.



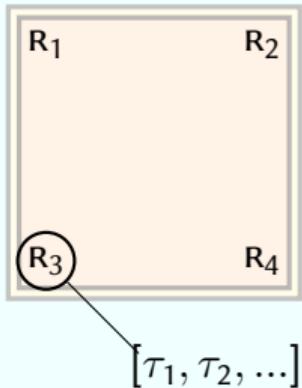
# Outline of a Blockchain-Based Resilient System

1. Replicas.
2. Holding a *ledger* of transactions.



# Outline of a Blockchain-Based Resilient System

1. Replicas.
2. Holding a *ledger* of transactions.



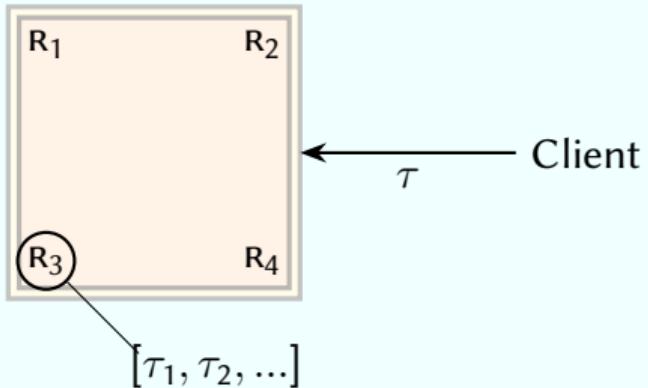
Assumption: Deterministic execution

The ledger provides *consistency*:

Good replicas execute the same operations in the same order and get the same results.

# Outline of a Blockchain-Based Resilient System

1. Replicas.
2. Holding a *ledger* of transactions.
3. Clients request new transactions.



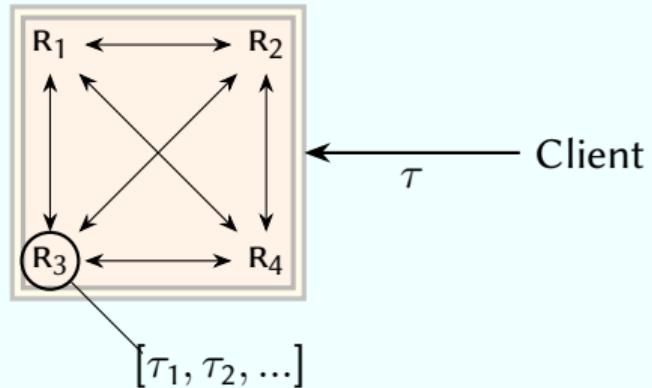
Assumption: Deterministic execution

The ledger provides *consistency*:

Good replicas execute the same operations in the same order and get the same results.

# Outline of a Blockchain-Based Resilient System

1. Replicas.
2. Holding a *ledger* of transactions.
3. Clients request new transactions.
4. Transaction agreement via *consensus*.



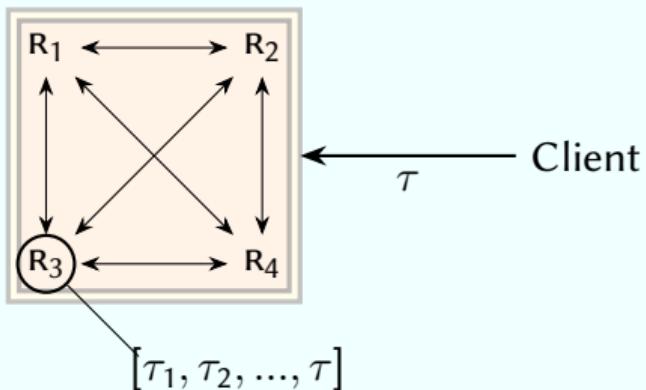
Assumption: Deterministic execution

The ledger provides *consistency*:

Good replicas execute the same operations in the same order and get the same results.

# Outline of a Blockchain-Based Resilient System

1. Replicas.
2. Holding a *ledger* of transactions.
3. Clients request new transactions.
4. Transaction agreement via *consensus*.
5. Append-only updates to ledger.



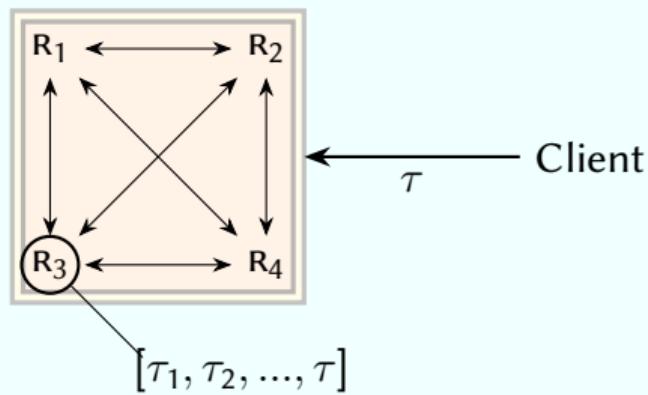
Assumption: Deterministic execution

The ledger provides *consistency*:

Good replicas execute the same operations in the same order and get the same results.

# Outline of a Blockchain-Based Resilient System

1. Replicas.
2. Holding a *ledger* of transactions.
3. Clients request new transactions.
4. Transaction agreement via *consensus*.
5. Append-only updates to ledger.
6. Cryptography.



Assumption: Deterministic execution

The ledger provides *consistency*:

Good replicas execute the same operations in the same order and get the same results.

# Classifying Resilient Systems

The complexity of operating a resilient system depends on many factors:

**Failures.** In what ways can replicas fail?

**Communication.** What assumptions are made on communication?

**Authentication.** How are messages and their senders identified and verified?

## Failure Models

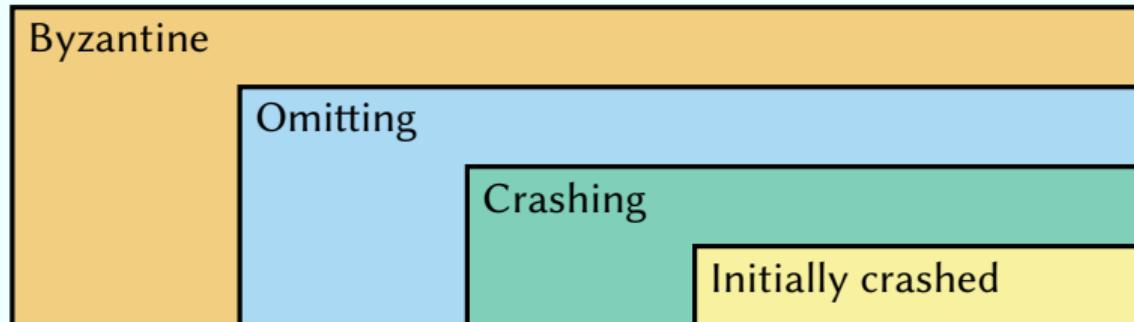
We say that a replica is

- initially crashed** if it will never do anything;
- crashed** if it stops doing anything at some point;
- omitting** if it can omit coordination steps;
- Byzantine** if it can behave *arbitrary*  
(e.g., omitting steps, performing the wrong steps.).

# Failure Models

We say that a replica is  
**initially crashed** if it will never do anything;  
**crashed** if it stops doing anything at some point;

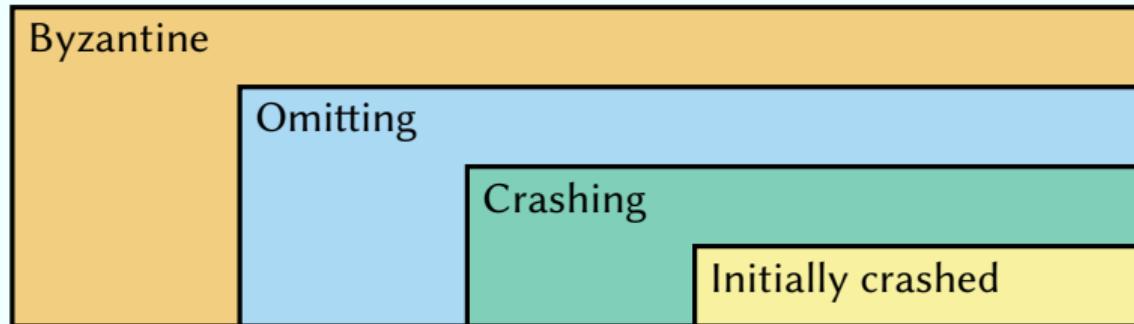
**omitting** if it can omit coordination steps;  
**Byzantine** if it can behave *arbitrary*  
(e.g., omitting steps, performing the wrong steps.).



# Failure Models

We say that a replica is  
**initially crashed** if it will never do anything;  
**crashed** if it stops doing anything at some point;

**omitting** if it can omit coordination steps;  
**Byzantine** if it can behave *arbitrary*  
(e.g., omitting steps, performing the wrong steps.).



Both *omitting* and *Byzantine* replicas can be **malicious**:  
these replicas can **coordinate** among themselves in attempts to disrupt the system.

# Communication Models

We say that communication is

**synchronous** if every message sent will arrive only at its destination,  
will do so exactly once within some *known delay*.  
Communication can be modeled in *rounds*.

**asynchronous** if messages can be arbitrarily delayed, duplicated, or *dropped*.  
Arbitrary delay: message can arrive *out of order*.

# Communication Models

We say that communication is

**synchronous** if every message sent will arrive only at its destination,  
will do so exactly once within some *known delay*.  
Communication can be modeled in *rounds*.

**asynchronous** if messages can be arbitrarily delayed, duplicated, or *dropped*.  
Arbitrary delay: message can arrive *out of order*.

- ▶ Synchronous: *assume* reliable communication, no partitions.

# Communication Models

We say that communication is

**synchronous** if every message sent will arrive only at its destination,  
will do so exactly once within some *known delay*.  
Communication can be modeled in *rounds*.

**asynchronous** if messages can be arbitrarily delayed, duplicated, or *dropped*.  
Arbitrary delay: message can arrive *out of order*.

- ▶ Synchronous: *assume* reliable communication, no partitions.    *unrealistic*

# Communication Models

We say that communication is

**synchronous** if every message sent will arrive only at its destination,  
will do so exactly once within some *known delay*.  
Communication can be modeled in *rounds*.

**asynchronous** if messages can be arbitrarily delayed, duplicated, or *dropped*.  
Arbitrary delay: message can arrive *out of order*.

- ▶ Synchronous: *assume* reliable communication, no partitions.    *unrealistic*
- ▶ Asynchronous: *worst-case* assumptions.

# Communication Models

We say that communication is

**synchronous** if every message sent will arrive only at its destination,  
will do so exactly once within some *known delay*.  
Communication can be modeled in *rounds*.

**asynchronous** if messages can be arbitrarily delayed, duplicated, or *dropped*.  
Arbitrary delay: message can arrive *out of order*.

- ▶ Synchronous: *assume* reliable communication, no partitions.      *unrealistic*
- ▶ Asynchronous: *worst-case* assumptions.      *complex*

## Byzantine Fault Tolerance: Authentication

Problem: Impersonation

Byzantine replicas can act arbitrarily and malicious:  
they can try to impersonate many good replicas!

# Byzantine Fault Tolerance: Authentication

## Problem: Impersonation

Byzantine replicas can act arbitrarily and malicious:  
they can try to impersonate many good replicas!

## Solution: Authenticated communication

If replica  $Q$  determines that  $m$  was sent by replica  $R$  then:

- ▶  $m$  must have been sent by  $R$  if  $R$  is good; and
- ▶  $m$  must be sent by some faulty replica if  $R$  is faulty.  
Faulty replicas can only impersonate each other.

# Byzantine Fault Tolerance: Authentication

## Problem: Impersonation

Byzantine replicas can act arbitrarily and malicious:  
they can try to impersonate many good replicas!

## Solution: Authenticated communication

If replica Q determines that  $m$  was sent by replica R then:

- ▶  $m$  must have been sent by R if R is good; and
- ▶  $m$  must be sent by some faulty replica if R is faulty.  
Faulty replicas can only impersonate each other.

Implementation: *message authentication codes* (MACs).

(*cheap* symmetric-key cryptography).

# Byzantine Fault Tolerance: Authentication

## Problem: Impersonation

Byzantine replicas can act arbitrarily and malicious:  
they can try to impersonate many good replicas!

## Solution: Authenticated communication

If replica Q determines that  $m$  was sent by replica R then:

- ▶  $m$  must have been sent by R if R is good; and
- ▶  $m$  must be sent by some faulty replica if R is faulty.  
Faulty replicas can only impersonate each other.

Implementation: *message authentication codes* (MACs).  
(*cheap* symmetric-key cryptography).

## Second Problem: Message corruption

Faulty replicas can corrupt any *forwarded messages*.

# Byzantine Fault Tolerance: Authentication

## Problem: Impersonation

Byzantine replicas can act arbitrarily and malicious:  
they can try to impersonate many good replicas!

## Second Problem: Message corruption

Faulty replicas can corrupt any *forwarded messages*.

## Solution: Digital signatures

- ▶ Each replica  $R$  can *sign* any message  $m$ , yielding  $\text{cert}(m, R)$ .
- ▶ A signed message  $\text{cert}(m, R)$  is non-forgeable without the help of  $R$ .  
Faulty replicas can only corrupt messages from each other when forwarding.

# Byzantine Fault Tolerance: Authentication

## Problem: Impersonation

Byzantine replicas can act arbitrarily and malicious:  
they can try to impersonate many good replicas!

## Second Problem: Message corruption

Faulty replicas can corrupt any *forwarded messages*.

## Solution: Digital signatures

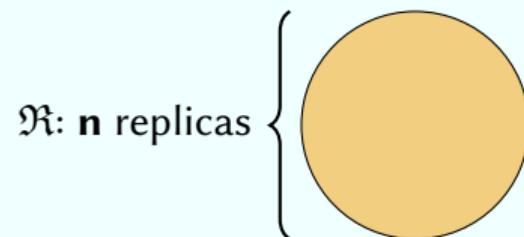
- ▶ Each replica  $R$  can *sign* any message  $m$ , yielding  $\text{cert}(m, R)$ .
- ▶ A signed message  $\text{cert}(m, R)$  is non-forgeable without the help of  $R$ .  
Faulty replicas can only corrupt messages from each other when forwarding.

Implementation: *costly* public-key cryptography.

# Notations for Resilient Systems

A resilient system is a *system* with

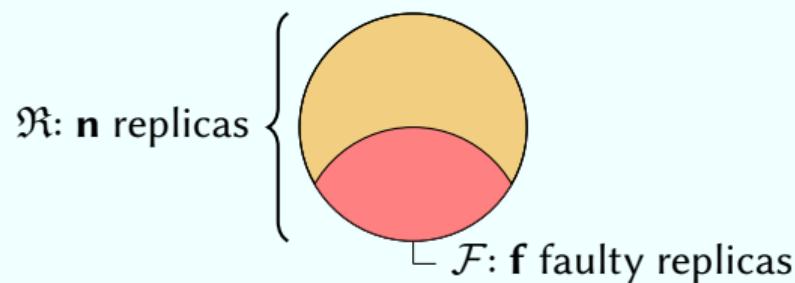
- ▶ A set of  $n$  *replicas*  $\mathfrak{R}$ .



# Notations for Resilient Systems

A resilient system is a *system* with

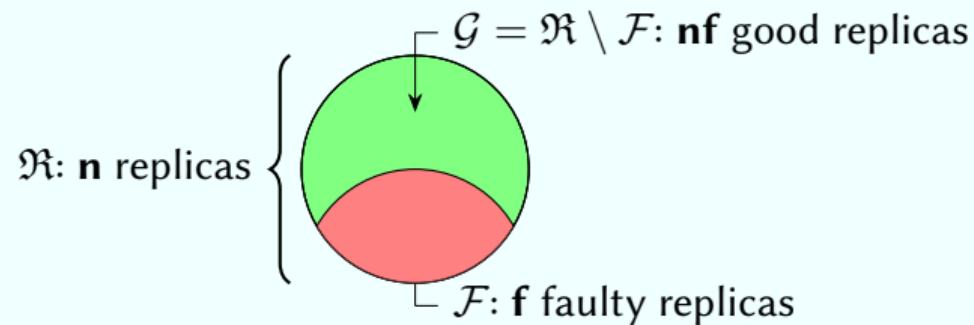
- ▶ A set of  $n$  *replicas*  $\mathfrak{R}$ .
- ▶ A set of  $f$  such replicas  $\mathcal{F} \subseteq \mathfrak{R}$  are *faulty*.



# Notations for Resilient Systems

A resilient system is a *system* with

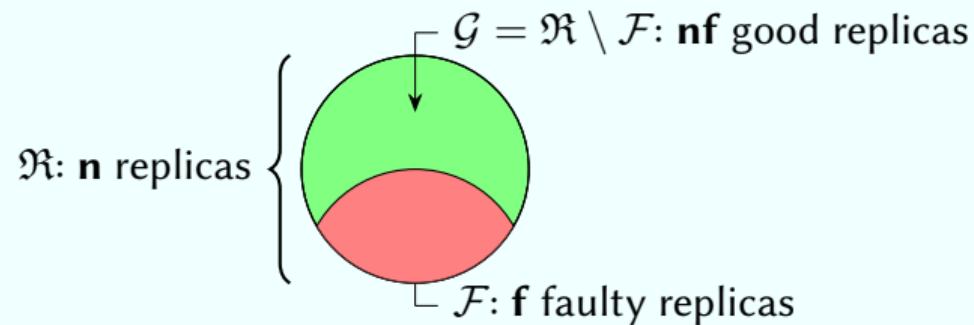
- ▶ A set of  $n$  *replicas*  $\mathfrak{R}$ .
- ▶ A set of  $f$  such replicas  $\mathcal{F} \subseteq \mathfrak{R}$  are *faulty*.
- ▶ We have  $n-f$  *good (non-faulty) replicas*  $\mathcal{G} = \mathfrak{R} \setminus \mathcal{F}$ .



# Notations for Resilient Systems

A resilient system is a *system* with

- ▶ A set of  $n$  *replicas*  $\mathfrak{R}$ .
- ▶ A set of  $f$  such replicas  $\mathcal{F} \subseteq \mathfrak{R}$  are *faulty*.
- ▶ We have  $nf$  *good (non-faulty) replicas*  $\mathcal{G} = \mathfrak{R} \setminus \mathcal{F}$ .
- ▶ Each replica  $r \in \mathfrak{R}$  has a *unique identifier*  $0 \leq id(r) < n$ .



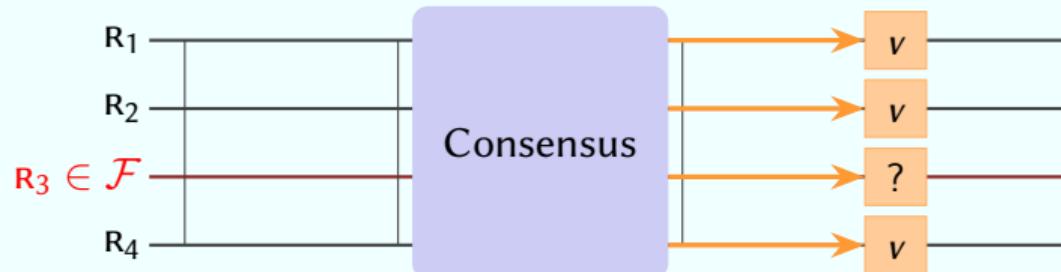
# Coordination in Resilient Systems: Consensus

A protocol provides *consensus* if upon completion of the protocol:

**Termination** Eventually, each good replica  $r \in \mathcal{G}$  decides on a value  $v(r)$ .

**Non-divergence** All good replicas decide on the same value.

Hence, if  $r_1, r_2 \in \mathcal{G}$  decide  $v(r_1)$  and  $v(r_2)$ , then  $v(r_1) = v(r_2)$ .



# Coordination in Resilient Systems: Consensus

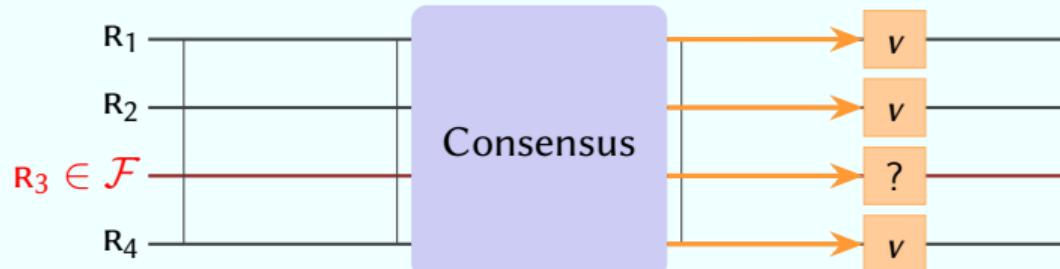
A protocol provides *consensus* if upon completion of the protocol:

**Termination** Eventually, each good replica  $r \in \mathcal{G}$  decides on a value  $v(r)$ .

**Non-divergence** All good replicas decide on the same value.

Hence, if  $r_1, r_2 \in \mathcal{G}$  decide  $v(r_1)$  and  $v(r_2)$ , then  $v(r_1) = v(r_2)$ .

**Non-triviality** In different runs of the protocol, replicas can decide on different values.



Excludes trivial solutions: e.g., good replicas always deciding a pre-defined value.

# Coordination in Resilient Systems: Consensus

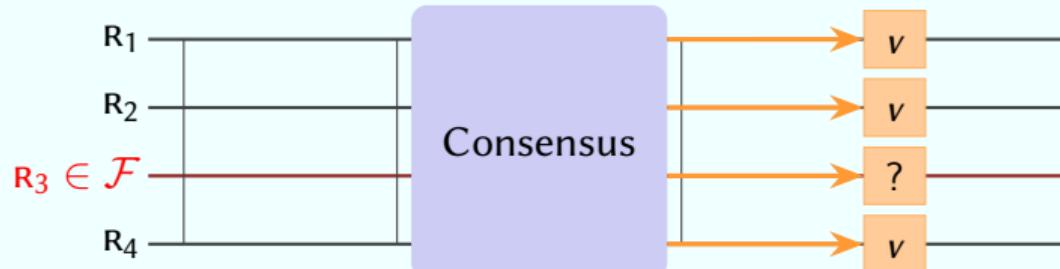
A protocol provides *consensus* if upon completion of the protocol:

**Termination** Eventually, each good replica  $r \in \mathcal{G}$  decides on a value  $v(r)$ . *Liveness*

**Non-divergence** All good replicas decide on the same value. *Safety*

Hence, if  $r_1, r_2 \in \mathcal{G}$  decide  $v(r_1)$  and  $v(r_2)$ , then  $v(r_1) = v(r_2)$ .

**Non-triviality** In different runs of the protocol, replicas can decide on different values.



Excludes trivial solutions: e.g., good replicas always deciding a pre-defined value.

## Consensus: From Formalization to Practice

Formal consensus is *abstract*: Where do the decided values come from?

# Consensus: From Formalization to Practice

Formal consensus is *abstract*: Where do the decided values come from?

Towards consensus in practice

Expand *non-triviality* by putting application-specific requirements on decided values.

## Consensus: From Formalization to Practice

Formal consensus is *abstract*: Where do the decided values come from?

Towards consensus in practice

Expand *non-triviality* by putting application-specific requirements on decided values.

Example: System processing client transactions

Decided-upon values that are not-yet executed client-requested transactions.

# Coordination in Resilient Systems: Interactive Consistency

Assumption: Each replica  $R$  holds an initial value  $v(R)$

A protocol provides *interactive consistency* if upon completion of the protocol:

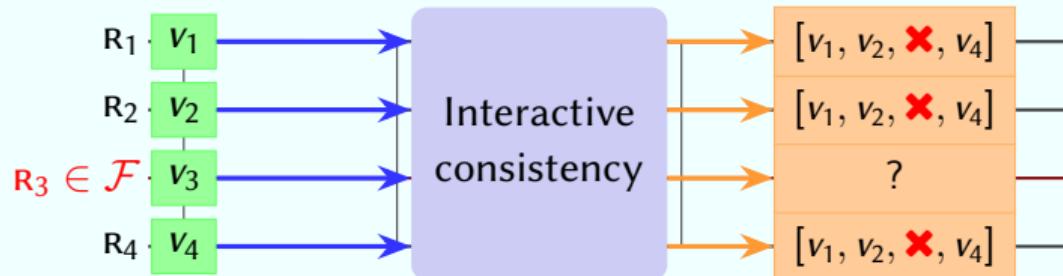
**Termination** Eventually, each good replica  $R \in \mathcal{G}$  decides on a list  $L(R)$  of  $n$  values.

**Non-divergence** All good replicas decide on the same list.

Hence, if  $R_1, R_2 \in \mathcal{G}$  decide  $L(R_1)$  and  $L(R_2)$ , then  $L(R_1) = L(R_2)$ .

**Dependence** Let  $R \in \mathcal{G}$ . Good replicas will have  $v(R)$  as the  $\text{id}(R)$ -th list value.

Hence, if replica  $Q \in \mathcal{G}$  decided on  $L(Q)$ , then  $L(Q)[\text{id}(R)] = v(R)$ .



## From Interactive Consistency to Consensus

Example: System processing client transactions

1. Each  $r \in \mathcal{G}$  chooses as  $v(r)$  a new client transaction.

## From Interactive Consistency to Consensus

Example: System processing client transactions

1. Each  $r \in \mathcal{G}$  chooses as  $v(r)$  a new client transaction.
2. All replicas participate in *interactive consistency*.

## From Interactive Consistency to Consensus

Example: System processing client transactions

1. Each  $r \in \mathcal{G}$  chooses as  $v(r)$  a new client transaction.
2. All replicas participate in *interactive consistency*.  
Each good replica obtains the same list  $L$  of  $n$  values.

## From Interactive Consistency to Consensus

Example: System processing client transactions

1. Each  $r \in \mathcal{G}$  chooses as  $v(r)$  a new client transaction.
2. All replicas participate in *interactive consistency*.  
Each good replica obtains the same list  $L$  of  $n$  values.
3. At least  $nf = n - f$  values in  $L$  are *valid new client transactions*.

# From Interactive Consistency to Consensus

## Example: System processing client transactions

1. Each  $r \in \mathcal{G}$  chooses as  $v(r)$  a new client transaction.
2. All replicas participate in *interactive consistency*.  
Each good replica obtains the same list  $L$  of  $n$  values.
3. At least  $nf = n - f$  values in  $L$  are *valid new client transactions*.
4. Good replicas use a *deterministic method* to choose a valid transaction from  $L$ . E.g.,
  - ▶ the first valid transaction in  $L$ ; or
  - ▶ all valid transactions in  $L$  (optimization).

# Coordination in Resilient Systems: Byzantine Broadcast

Assumption: Some replica  $g$  holds an initial value  $w$

A protocol provides *byzantine broadcast* if upon completion of the protocol:

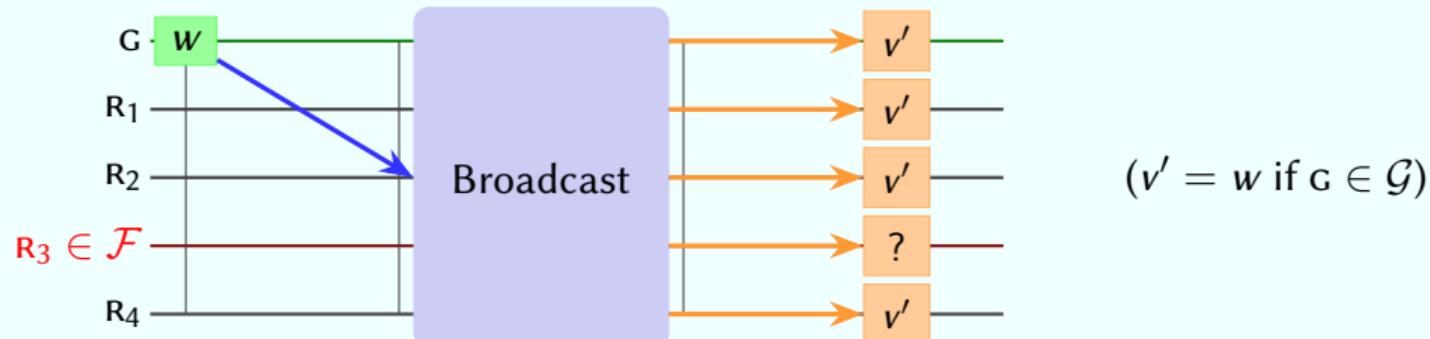
**Termination** Eventually, each good replica  $r \in \mathcal{G}$  decides on a value  $v(r)$ .

**Non-divergence** All good replicas decide on the same value.

Hence, if  $r_1, r_2 \in \mathcal{G}$  decide  $v(r_1)$  and  $v(r_2)$ , then  $v(r_1) = v(r_2)$ .

**Dependence** If  $g$  is good, then good replicas decide  $w$ .

Hence, if  $g \in \mathcal{G}$  and  $q \in \mathcal{G}$ , then  $v(q) = w$ .



## From Byzantine Broadcasts to Interactive Consistency

Assumption: Each replica  $r$  holds an initial value  $v(r)$

1. Each replica  $r$  performs Byzantine broadcast of  $v(r)$ .

## From Byzantine Broadcasts to Interactive Consistency

Assumption: Each replica  $r$  holds an initial value  $v(r)$

1. Each replica  $r$  performs Byzantine broadcast of  $v(r)$ .
2. What about faulty replicas? They might not broadcast!

# From Byzantine Broadcasts to Interactive Consistency

Assumption: Each replica  $r$  holds an initial value  $v(r)$

1. Each replica  $r$  performs Byzantine broadcast of  $v(r)$ .
2. What about faulty replicas? They might not broadcast!
3. Each good replica  $q \in \mathfrak{R}$  constructs  $L(q) = [v_0, \dots, v_{n-1}]$   
in which  $v_i$ ,  $0 \leq i < n$ , is the value broadcasted by the  $i$ -th replica.

# From Byzantine Broadcasts to Interactive Consistency

Assumption: Each replica  $r$  holds an initial value  $v(r)$

1. Each replica  $r$  performs Byzantine broadcast of  $v(r)$ .
2. What about faulty replicas? They might not broadcast!
3. Each good replica  $q \in \mathfrak{R}$  constructs  $L(q) = [v_0, \dots, v_{n-1}]$  in which  $v_i$ ,  $0 \leq i < n$ , is the value broadcasted by the  $i$ -th replica.

Interactive consistency and Byzantine broadcasts solve the same problem.

Both can be used to provide practical consensus.

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

Asynchronous	Orange	Light Blue	Dark Green	Yellow	Dark Blue	Orange
Synchronous	Light Orange	Light Blue	Light Green	Light Yellow	Light Blue	Light Orange
	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

Asynchronous	Orange	Light Blue	Dark Green	Yellow	Dark Blue	Orange
Synchronous	$n > f$					
	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$					
Synchronous	$n > f$					

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$					
Synchronous	$n > f$	$n > f$				

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$	Impossible <sup>4</sup>				
Synchronous	$n > f$	$n > f$				

<sup>4</sup>The FLP Impossibility Theorem.

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$	Impossible <sup>4</sup>	Impossible	Impossible	Impossible	Impossible
Synchronous	$n > f$	$n > f$				

<sup>4</sup>The FLP Impossibility Theorem.

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	$n > 2f$	Impossible <sup>4</sup>	Impossible	Impossible	Impossible	Impossible
Asynchronous						
Synchronous	$n > f$	$n > f$	$n > f$			
	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)

<sup>4</sup>The FLP Impossibility Theorem.

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	$n > 2f$	Impossible <sup>4</sup>	Impossible	Impossible	Impossible	Impossible
Asynchronous						
Synchronous	$n > f$	$n > f$	$n > f$	Impossible		
	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)

<sup>4</sup>The FLP Impossibility Theorem.

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	$n > 2f$	Impossible <sup>4</sup>	Impossible	Impossible	Impossible	Impossible
Asynchronous						
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	
	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)

<sup>4</sup>The FLP Impossibility Theorem.

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	$n > 2f$	Impossible <sup>4</sup>	Impossible	Impossible	Impossible	Impossible
Asynchronous						
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	$n > f$
	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)

<sup>4</sup>The FLP Impossibility Theorem.

# On the Complexity of Consensus

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$	Impossible <sup>4</sup>	Impossible	Impossible	Impossible	Impossible
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	$n \geq f$ $n > 2f$

<sup>4</sup>The FLP Impossibility Theorem.

## Consensus in Practice: Impossible?

### Theorem (FLP Impossibility Theorem)

*The consensus problem cannot be solved for systems that*

- ▶ *uses asynchronous communication; and*
- ▶ *allows faulty replicas to crash.*

Real-world networks: Asynchronous, but mostly well-behaved

## Consensus in Practice: Impossible?

### Theorem (FLP Impossibility Theorem)

*The consensus problem cannot be solved for systems that*

- ▶ *uses asynchronous communication; and*
- ▶ *allows faulty replicas to crash.*

Real-world networks: Asynchronous, but mostly well-behaved

Do not solve consensus: weaken the *termination* guarantee:

Each good replica  $r \in \mathcal{G}$  decides on a value  $v(r)$ .

## Consensus in Practice: Impossible?

### Theorem (FLP Impossibility Theorem)

*The consensus problem cannot be solved for systems that*

- ▶ *uses asynchronous communication; and*
- ▶ *allows faulty replicas to crash.*

**Real-world networks: Asynchronous, but mostly well-behaved**

Do not solve consensus: weaken the *termination* guarantee:

Each good replica  $r \in \mathcal{G}$  decides on a value  $v(r)$ .

Weak termination Good replicas decide *when communication is well-behaved*.

Probabilistic termination Good replicas decide *with high probability*.

# On the Complexity of Consensus in Practice

Consider a system with  $n$  replicas of which  $f$  are faulty.

The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$			Impossible		
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	$n \geq f$ $n > 2f$

# On the Complexity of Consensus in Practice

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$	$n > 2f^5$		Impossible		
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	$n \geq f$ $n > 2f$

<sup>5</sup>Weak or probabilistic termination.

# On the Complexity of Consensus in Practice

Consider a system with  $n$  replicas of which  $f$  are faulty.

The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$	$n > 2f^5$	$n > 2f^5$	Impossible		
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	$n \geq f$ $n > 2f$

<sup>5</sup>Weak or probabilistic termination.

# On the Complexity of Consensus in Practice

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$	$n > 2f^5$	$n > 2f^5$	Impossible	$n > 3f^5$	
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	$n \geq f$ $n > 2f$

<sup>5</sup>Weak or probabilistic termination.

# On the Complexity of Consensus in Practice

Consider a system with  $n$  replicas of which  $f$  are faulty.

## The consensus problem (informal)

Can the good replicas reach *agreement* on a value?

	Initially Crashed	Crashed	Omitting	Byzantine	Byzantine (with MACs)	Byzantine (with DSs)
Asynchronous	$n > 2f$	$n > 2f^5$	$n > 2f^5$	Impossible	$n > 3f^5$	$n > 3f^5$
Synchronous	$n > f$	$n > f$	$n > f$	Impossible	$n > 3f$	$n \geq f$ $n > 2f$

<sup>5</sup>Weak or probabilistic termination.

# The Complexity and Cost of Consensus

Many other limitations on consensus are known

Communication phases a worst-case approach in at-least  $f + 1$  phases;  
at-least  $t + 2$  phases if  $t \leq f$  replicas behave faulty (optimistic);

Communication cost an exchange of at-least  $nf$  signatures;  
an exchange of at least  $n + f^2$  messages;

Network structure at-least  $2f + 1$  disjoint communication paths  
between every pair of replicas.