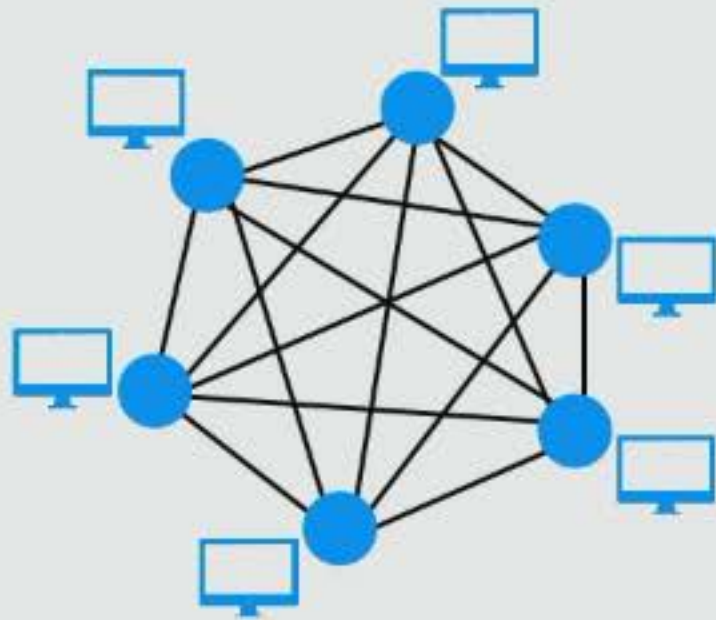# RDMA:
# Provably More Powerful Communication
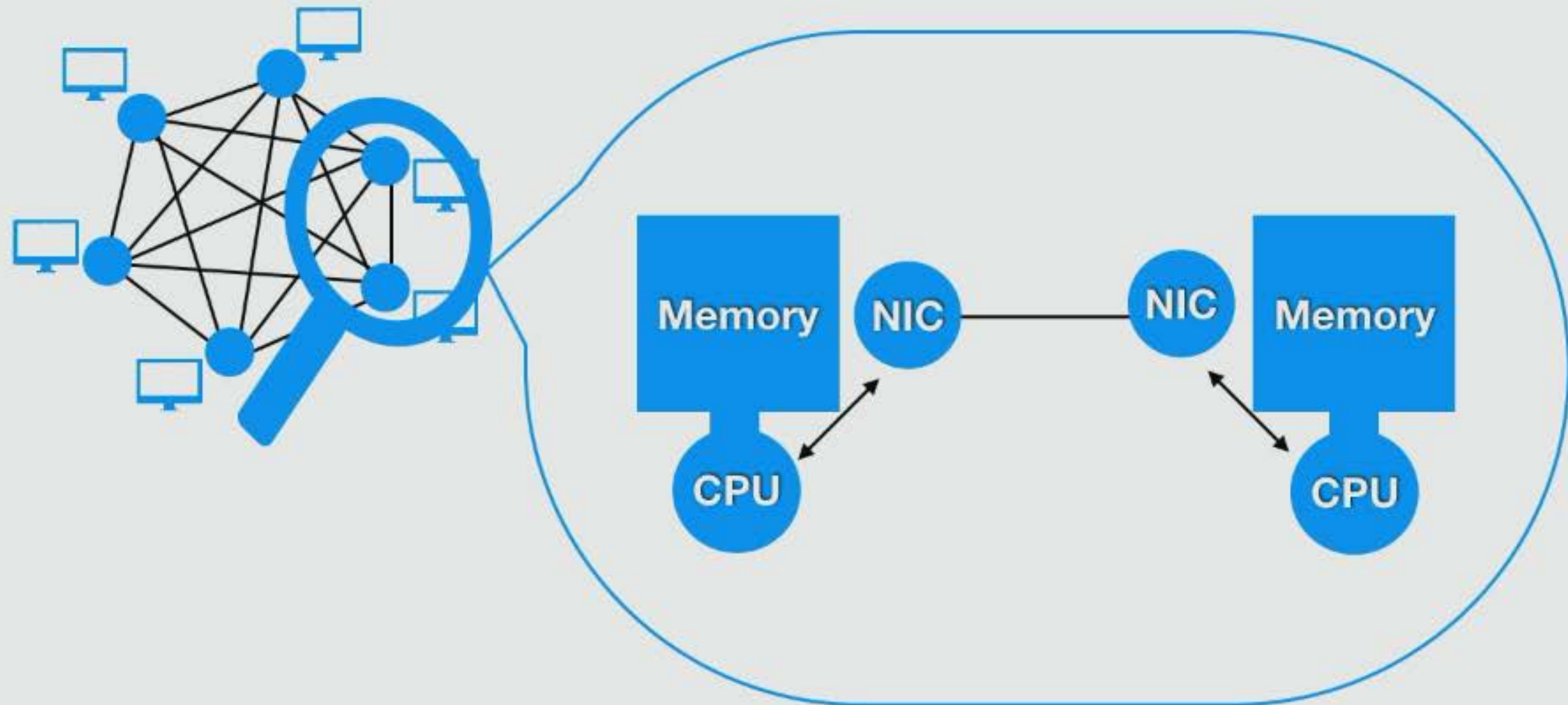
Naama Ben-David (CMU)

PODC'18, PODC'19

Marcos Aguilera, Irina Calciu, Rachid Guerraoui, Virendra Marathe,
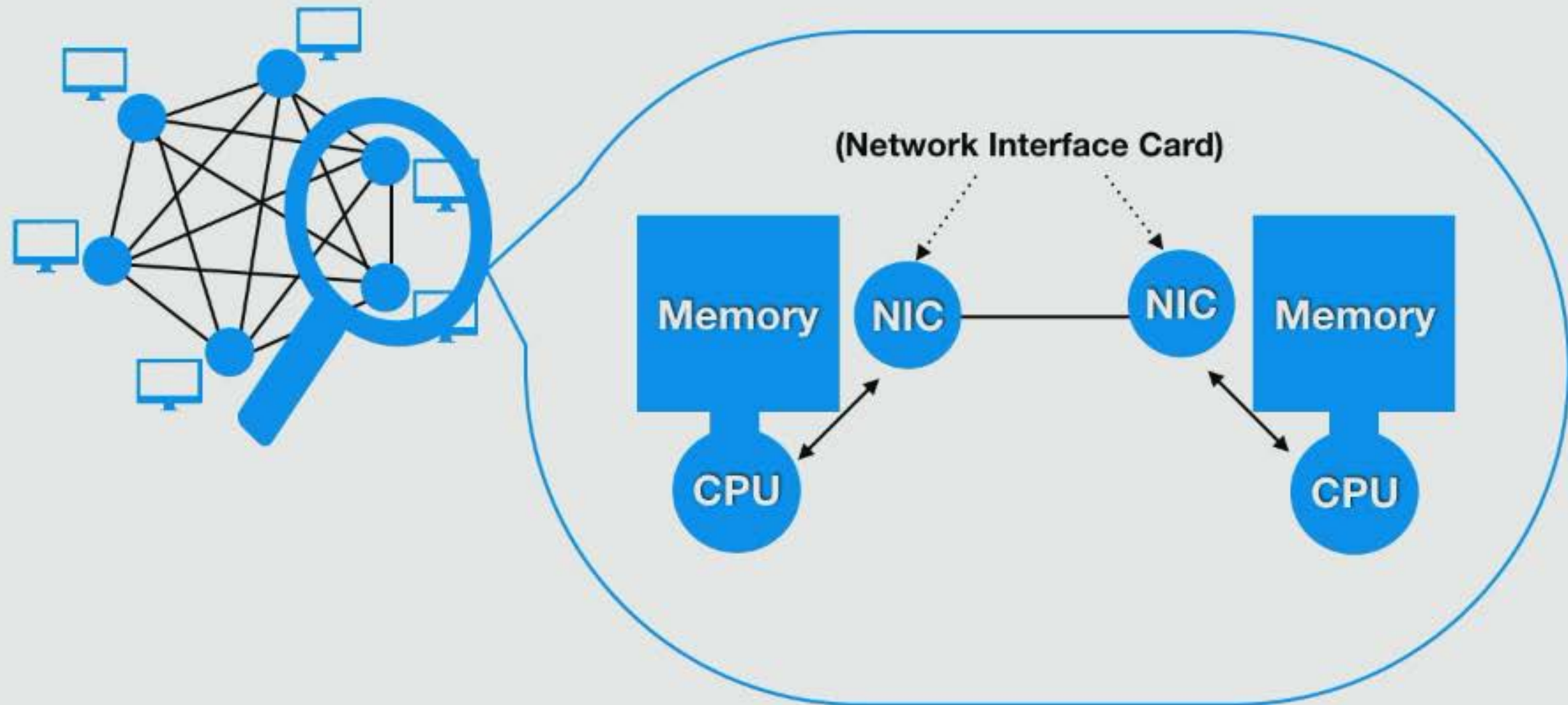Erez Petrank, Sam Toueg, Igor Zablotchi
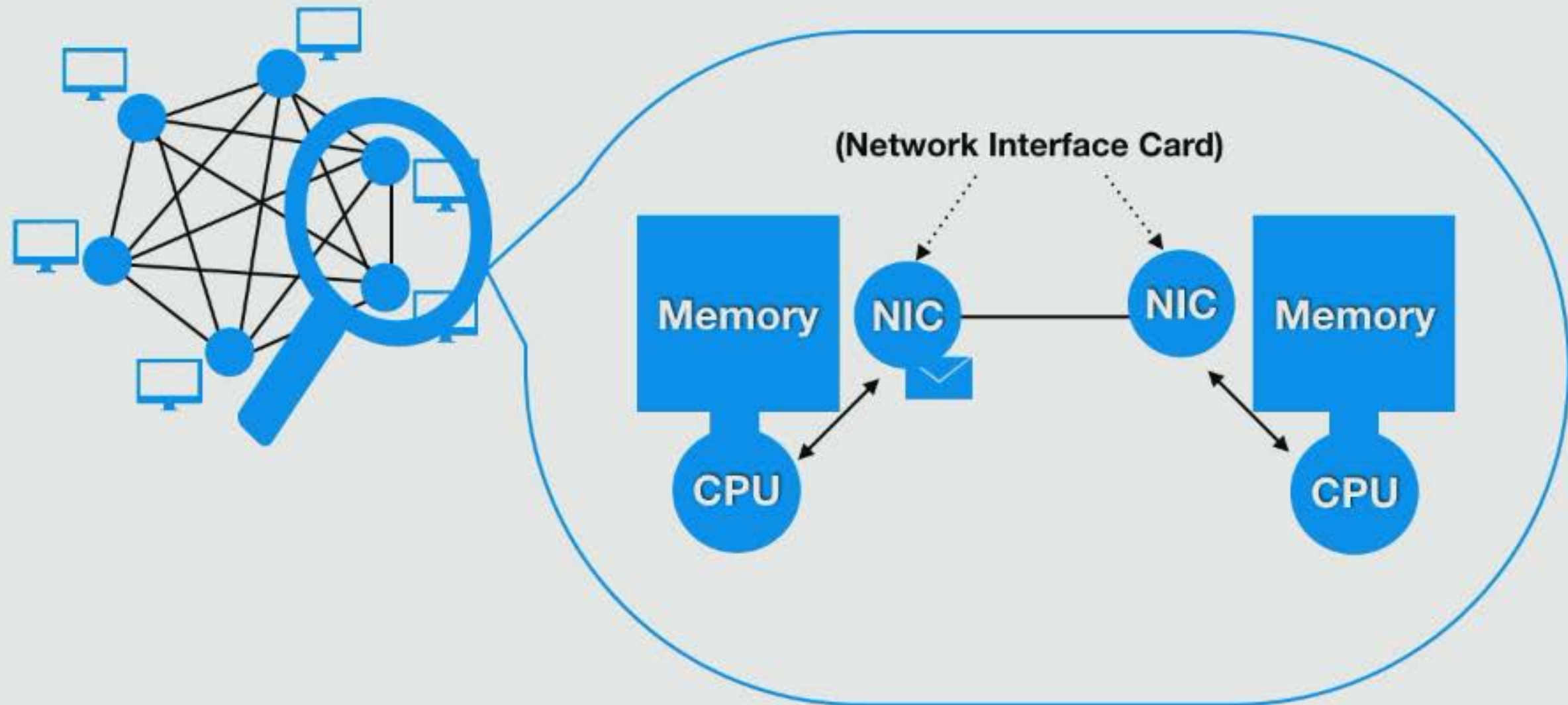
1

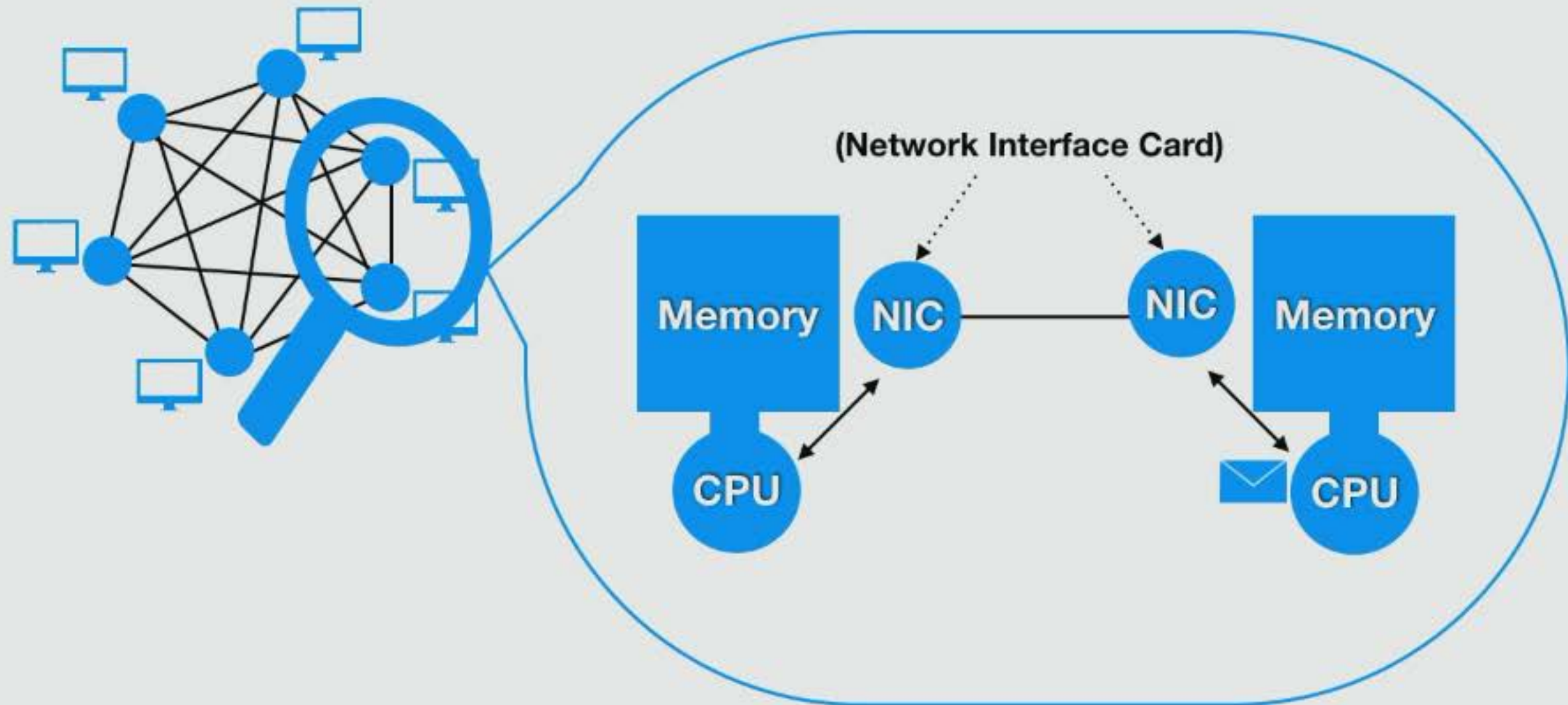# Data Center Technology: RDMA

# Data Center Technology: RDMA

# Data Center Technology: RDMA

# Data Center Technology: RDMA

(Network Interface Card)
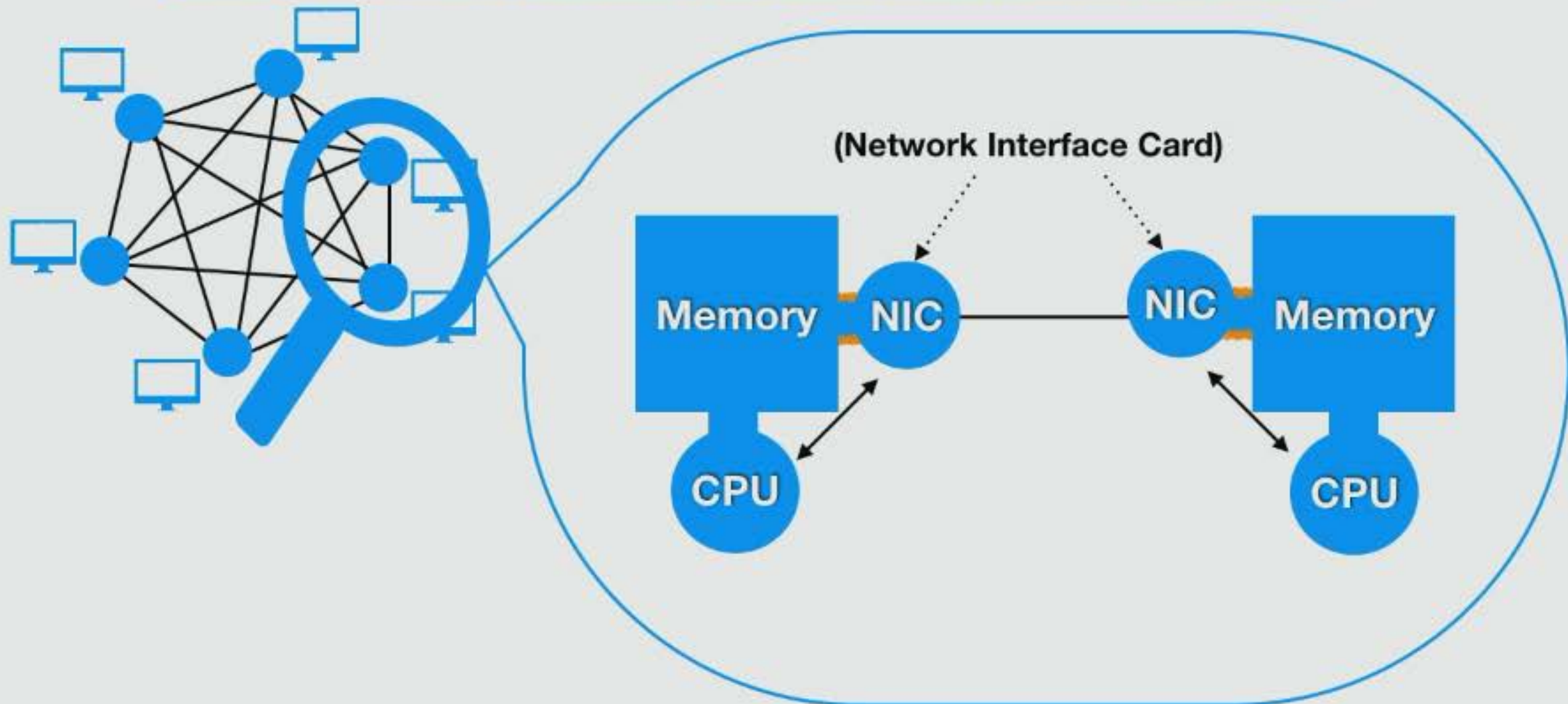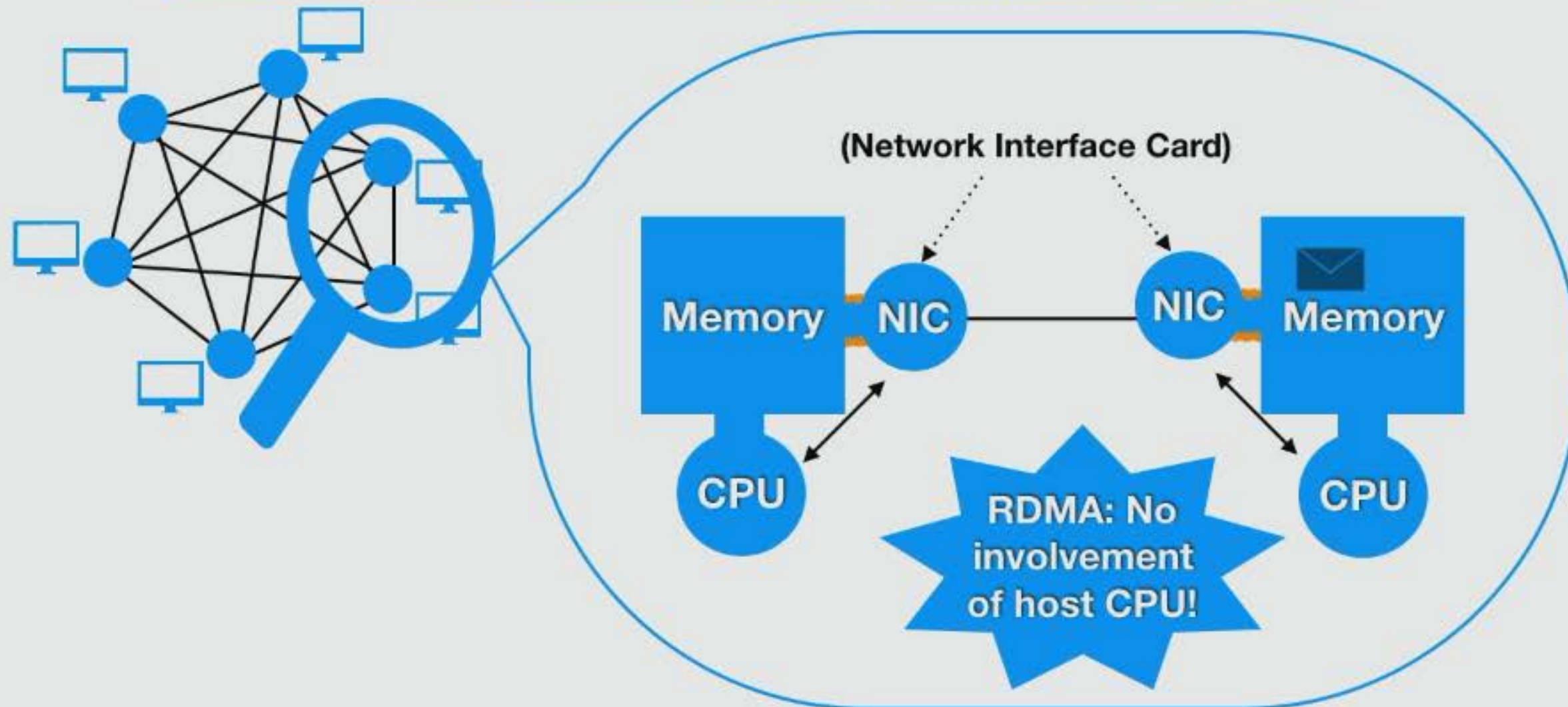
Memory — NIC — NIC — Memory

CPU          CPU

# Data Center Technology: RDMA

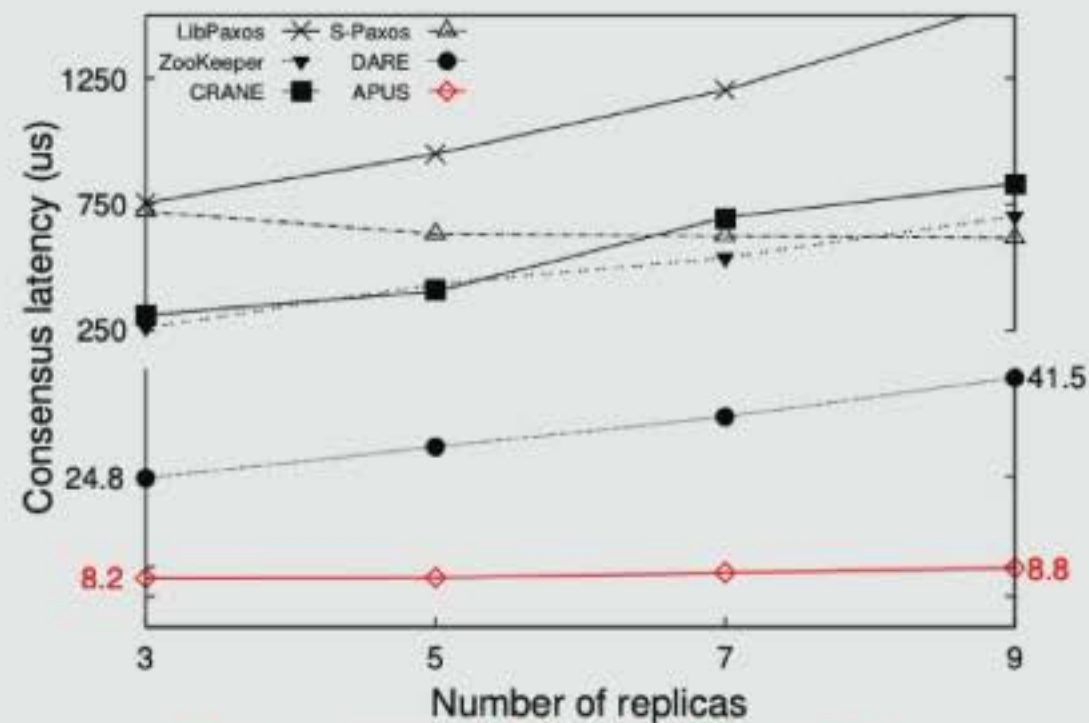# Data Center Technology: RDMA

# RDMA in Practice

- Huge speedups with RDMA

# RDMA in Practice

- Huge speedups with RDMA



[WangJaingChenYiCui'17]

# RDMA in Practice
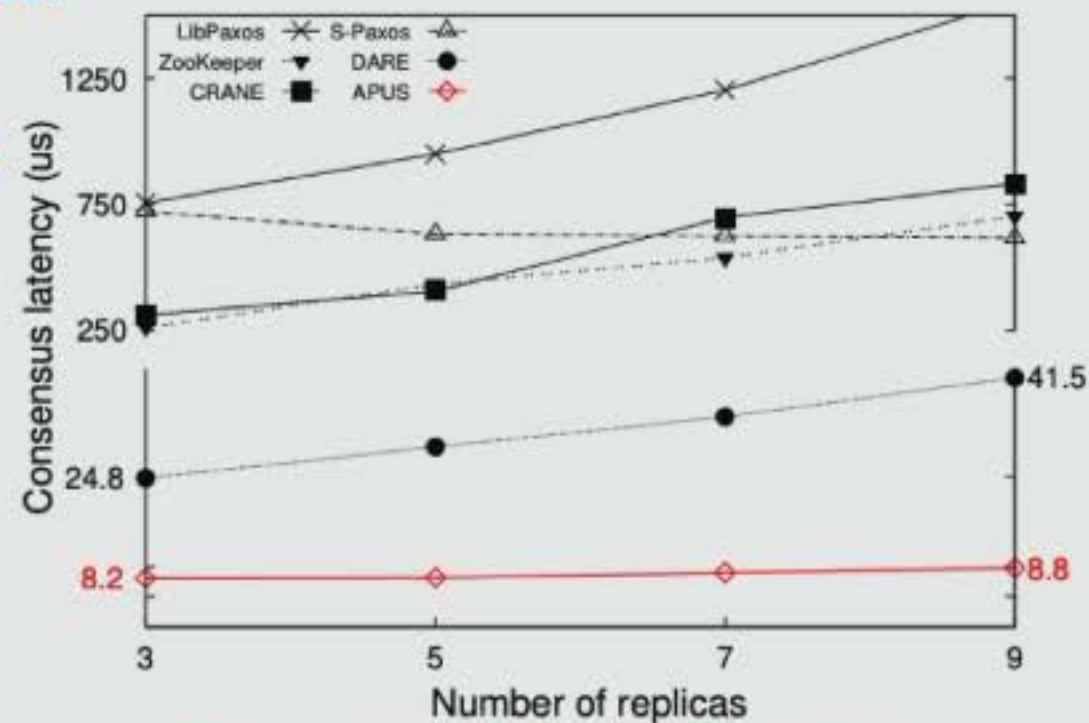
- Huge speedups with RDMA



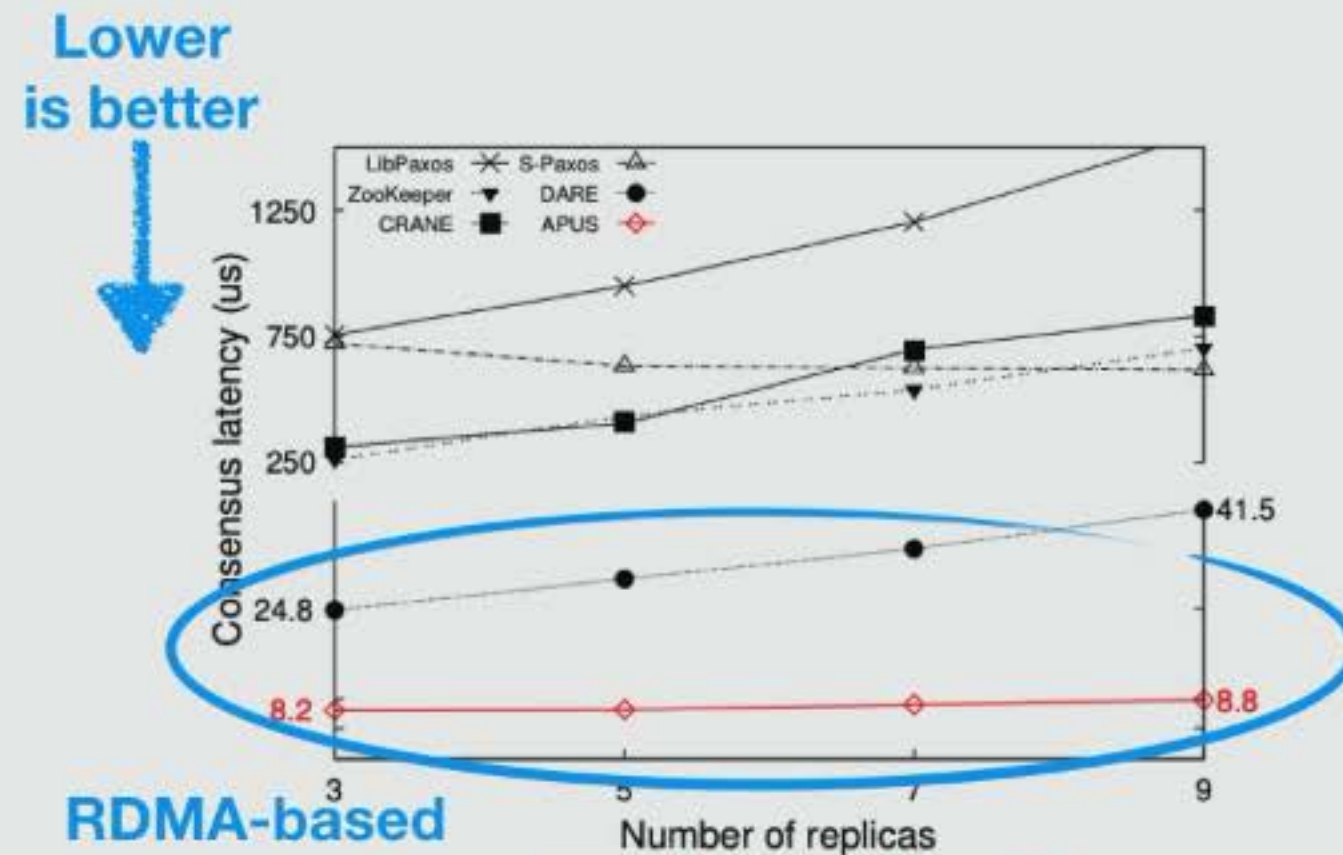[WangJaingChenYiCui'17]

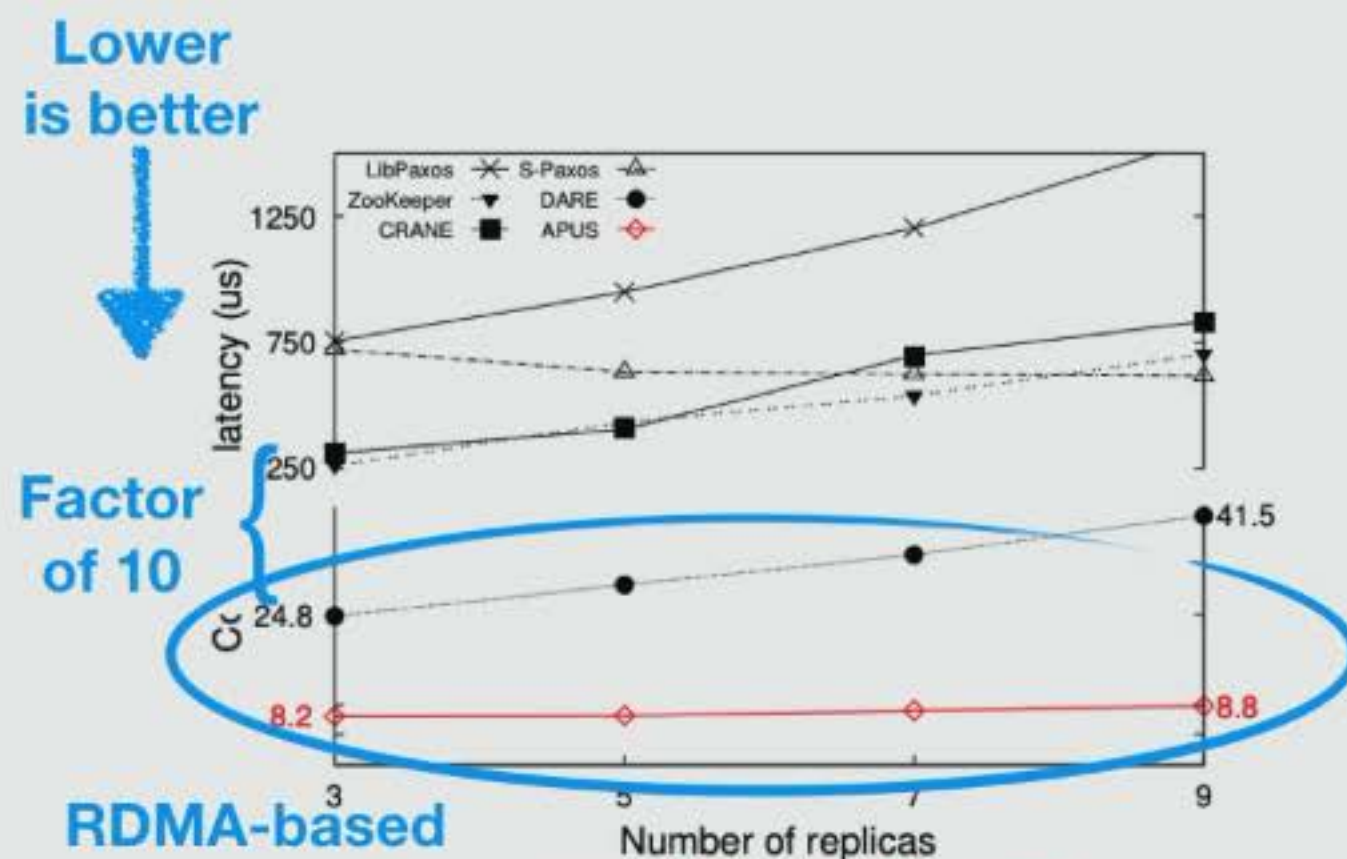# RDMA in Practice

- Huge speedups with RDMA



Lower is better

RDMA-based

[WangJaingChenYiCui'17]
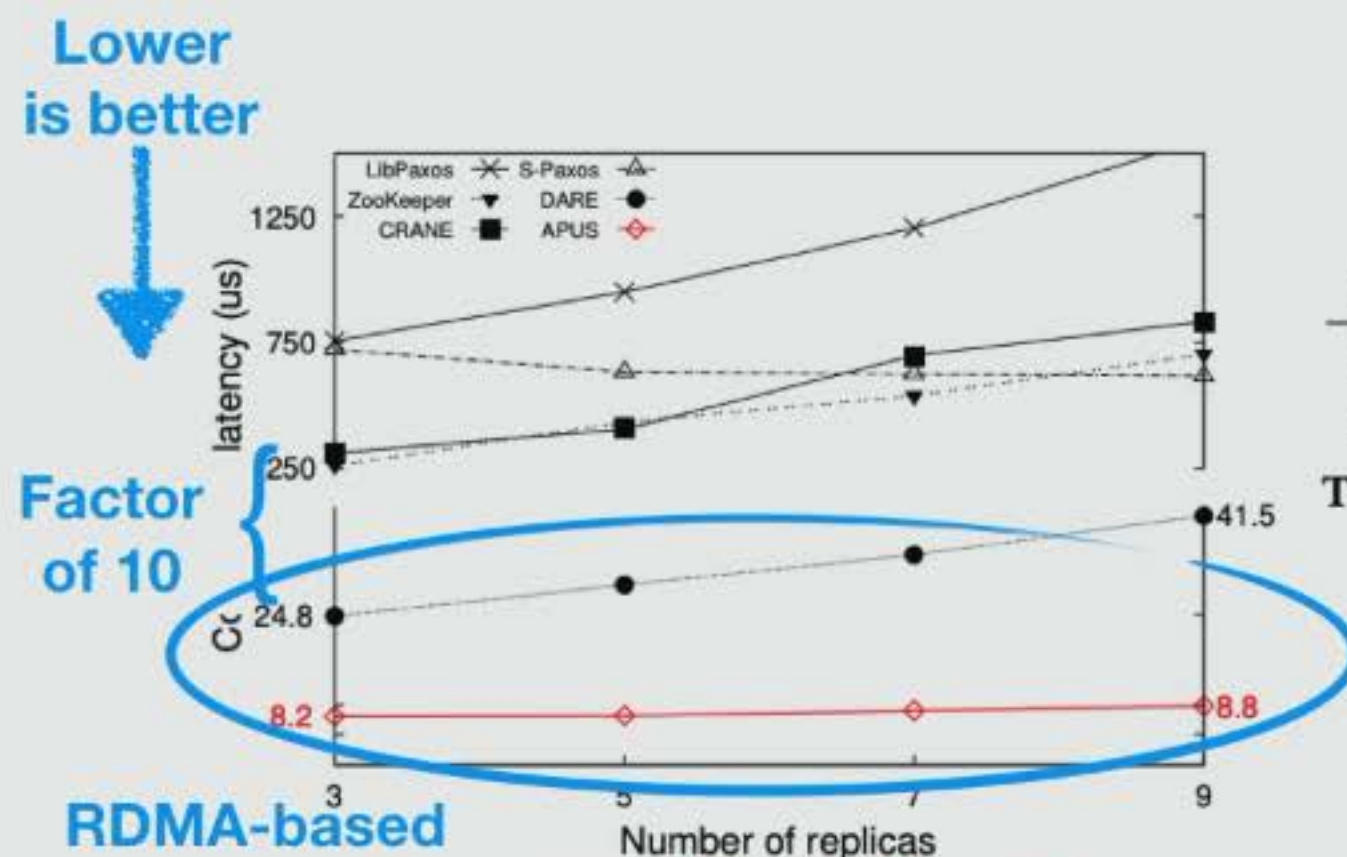
# RDMA in Practice

- Huge speedups with RDMA

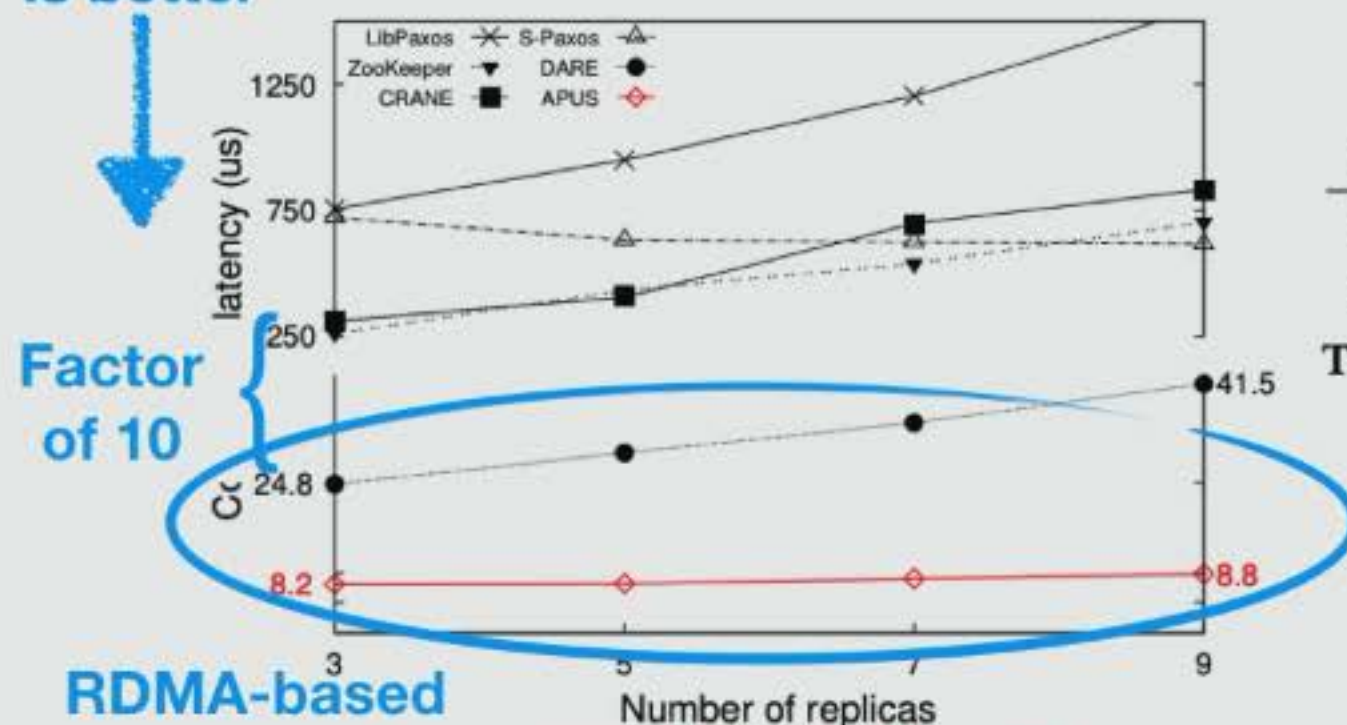- But is this improvement fundamental to RDMA?



[WangJaingChenYiCui'17]

# RDMA in Practice

- Huge speedups with RDMA

- But is this improvement fundamental to RDMA?



**Lower is better**

**Factor of 10**

**RDMA-based**

latency (us)

LibPaxos, S-Paxos, ZooKeeper, DARE, CRANE, APUS

1250, 750, 250

41.5, 24.8, 8.2, 8.8

Number of replicas: 3, 5, 7, 9

[WangJaingChenYiCui'17]

| Cluster | CX3 (InfiniBand) | CX4 (Eth) | CX5 (Eth) |
|---------|------------------|-----------|-----------|
| RDMA read | 1.7 µs | 2.9 µs | 2.0 µs |
| eRPC | 2.1 µs | 3.7 µs | 2.3 µs |

Table 2: Comparison of median latency with eRPC and RDMA

3

# RDMA in Practice

- Huge speedups with RDMA

- But is this improvement fundamental to RDMA?



**Lower is better**

**Factor of 10**

**RDMA-based**
[WangJaingChenYiCui'17]

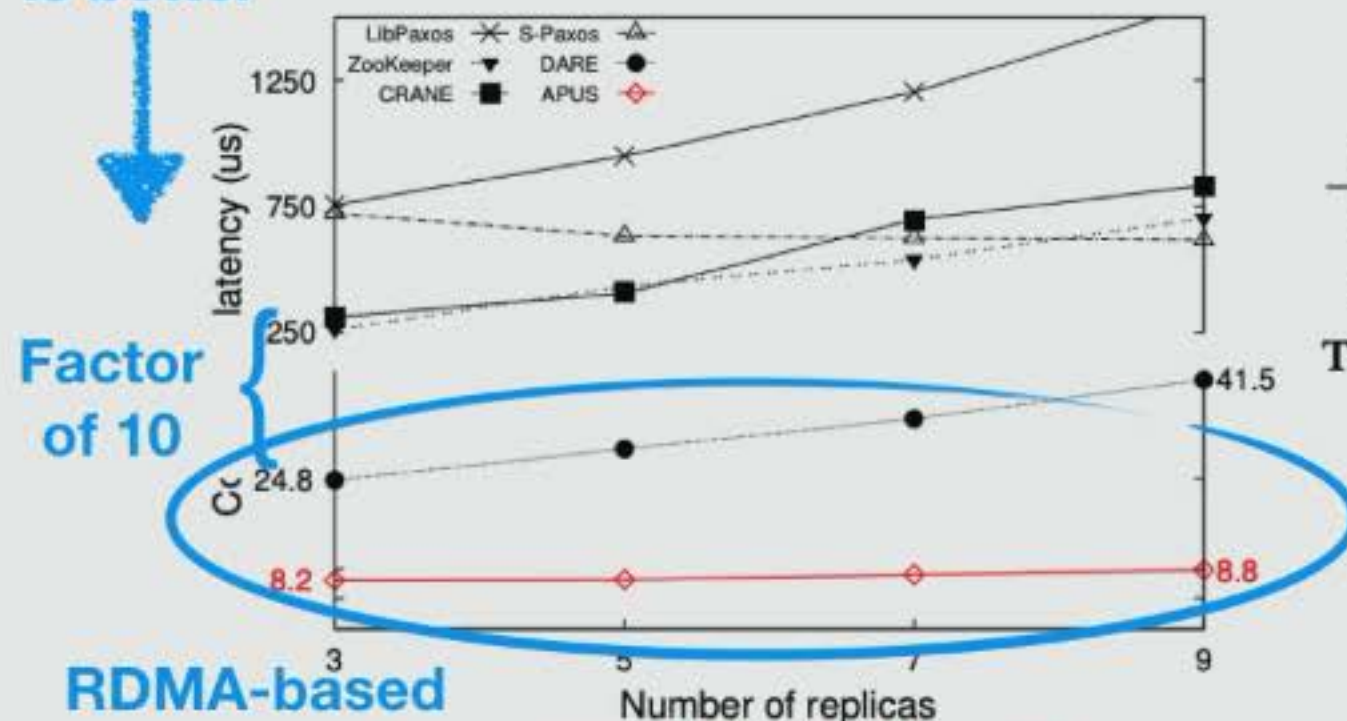| Cluster | CX3 (InfiniBand) | CX4 (Eth) | CX5 (Eth) |
|---------|------------------|-----------|-----------|
| RDMA read | 1.7 μs | 2.9 μs | 2.0 μs |
| eRPC | 2.1 μs | 3.7 μs | 2.3 μs |

Table. Comparison of median latency with eRPC and RDMA

**Old technology, optimized software**
[KaliaKaminskiAndersen'19]

3

# RDMA in Practice

- Huge speedups with RDMA

- But is this improvement fundamental to RDMA?



**Lower is better**

**Factor of 10**

**RDMA-based**
[WangJaingChenYiCui'17]

**Performance comparable with RDMA**

| Cluster | CX3 (InfiniBand) | CX4 (Eth) | CX5 (Eth) |
|---------|------------------|-----------|-----------|
| RDMA read | 1.7 μs | 2.9 μs | 2.0 μs |
| eRPC | 2.1 μs | 3.7 μs | 2.3 μs |

Table : Comparison of median latency with eRPC and RDMA

**Old technology, optimized software**
[KaliaKaminskiAndersen'19]

3

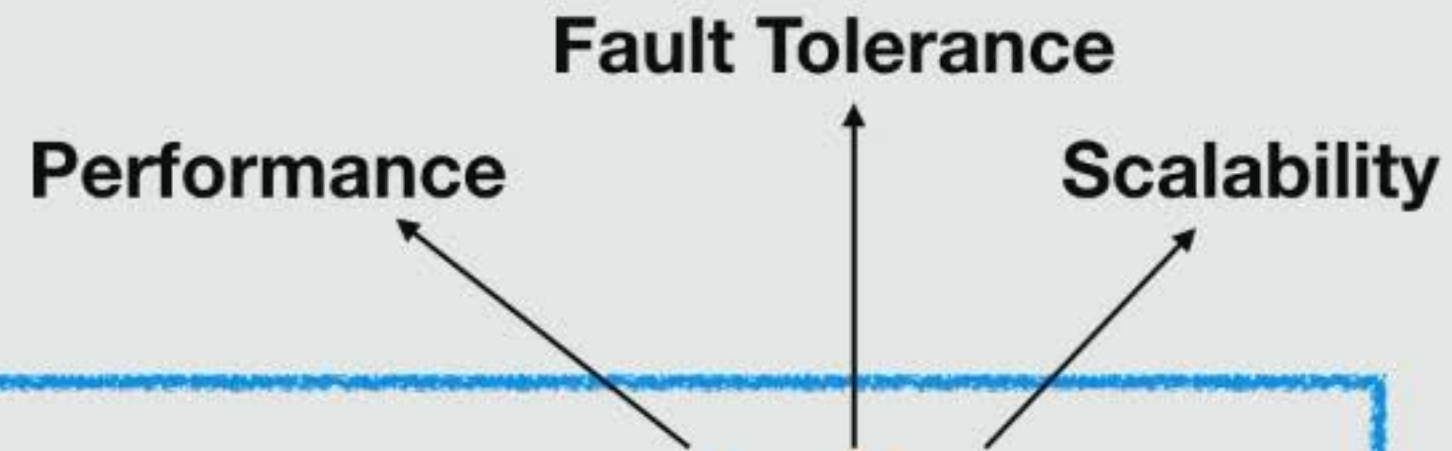**Is RDMA fundamentally *better* than other communication mechanisms?**

Performance

**Is RDMA fundamentally *better* than other communication mechanisms?**

Performance Fault Tolerance Scalability

**Is RDMA fundamentally *better* than other communication mechanisms?**
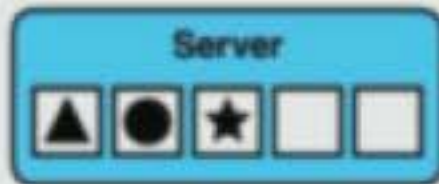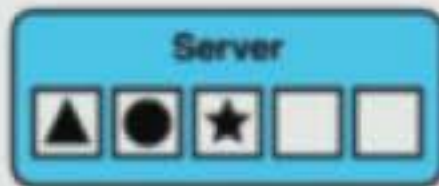
# Consensus as a Lens

Consensus: **Agreement**

# Consensus as a Lens

Consensus: **Agreement**

Replication, shared data structures, blockchains

# Consensus as a Lens

Consensus: **Agreement**

Replication, shared data structures, blockchains

# Consensus as a Lens

Consensus: **Agreement**

Replication, shared data structures, blockchains

# Consensus: Definition

- **Input**: every process gets input
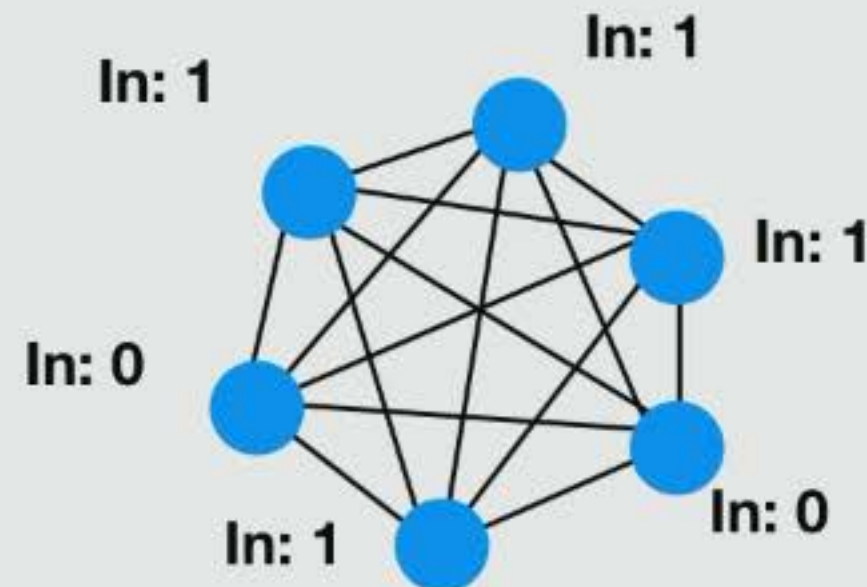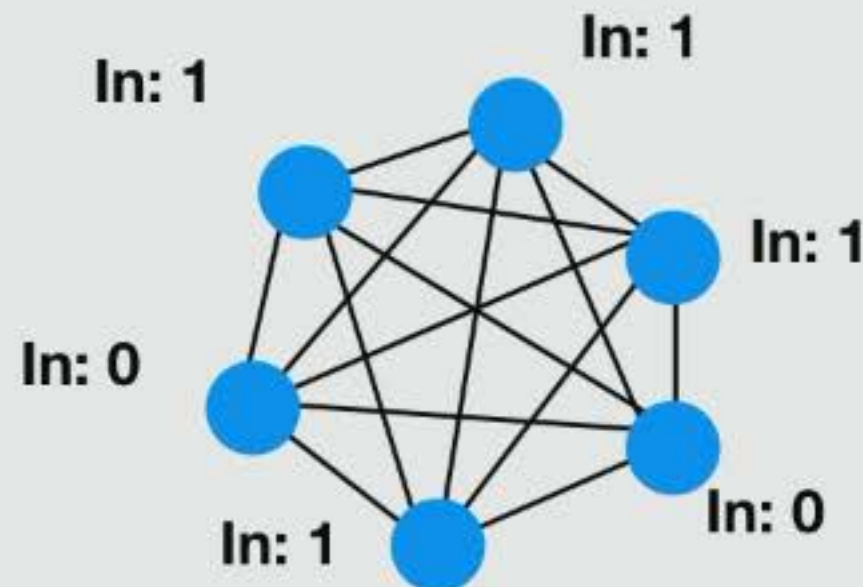
- **Output**: Every process outputs something

# Consensus: Definition

- **Input**: every process gets input
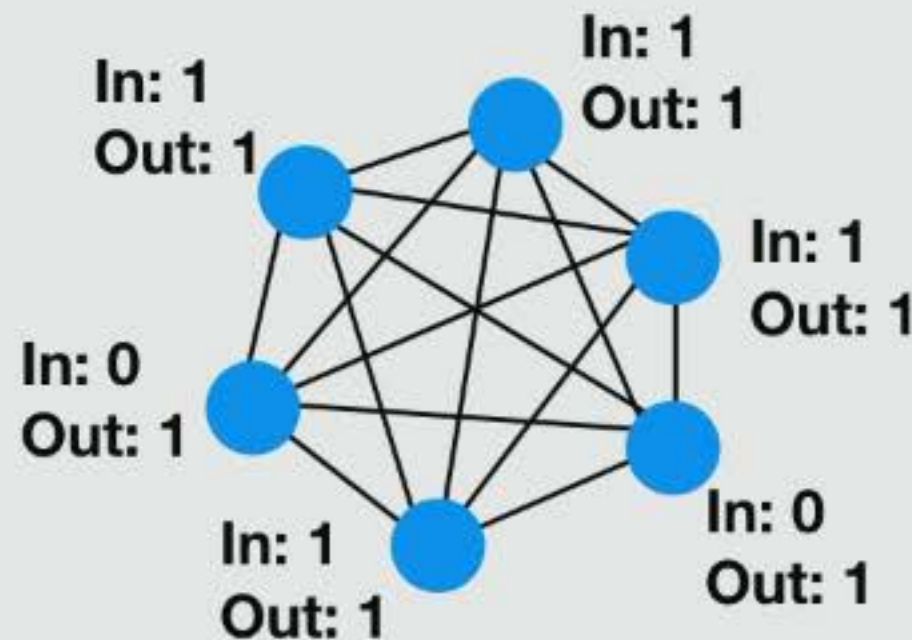
- **Output**: Every process outputs something

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something
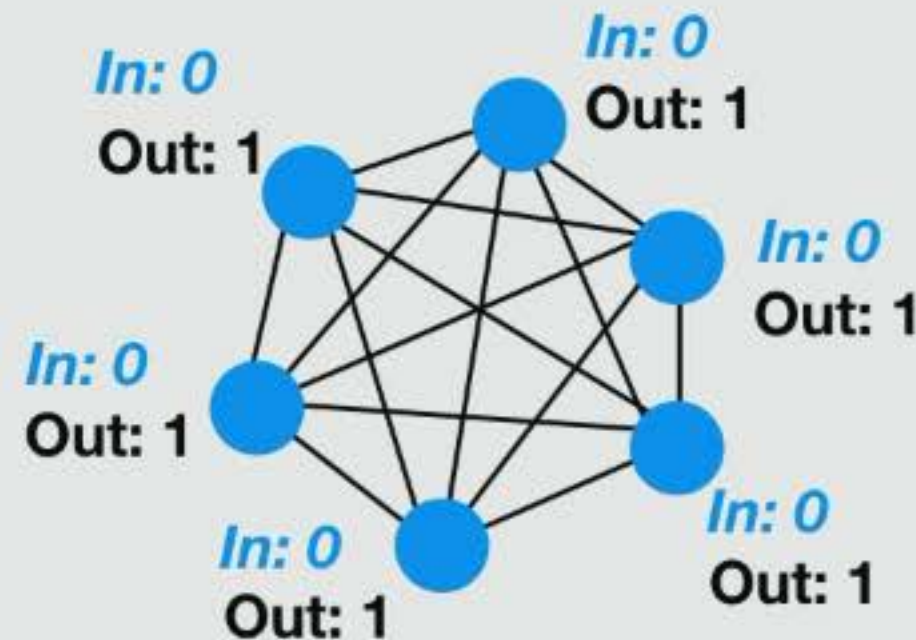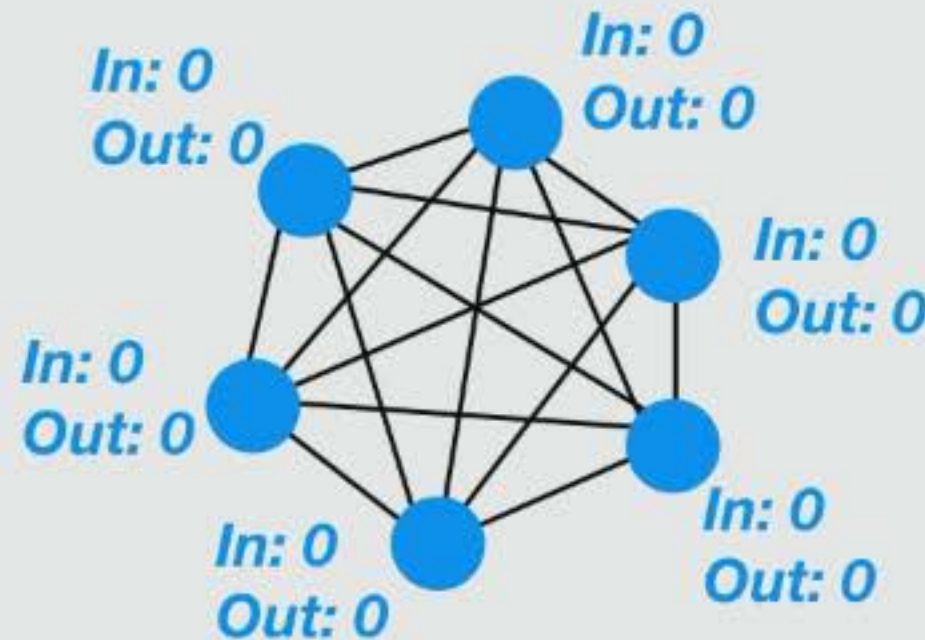
  - Agreement: Every process outputs the same value

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something

  - Agreement: Every process outputs the same value

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something

  - Agreement: Every process outputs the same value

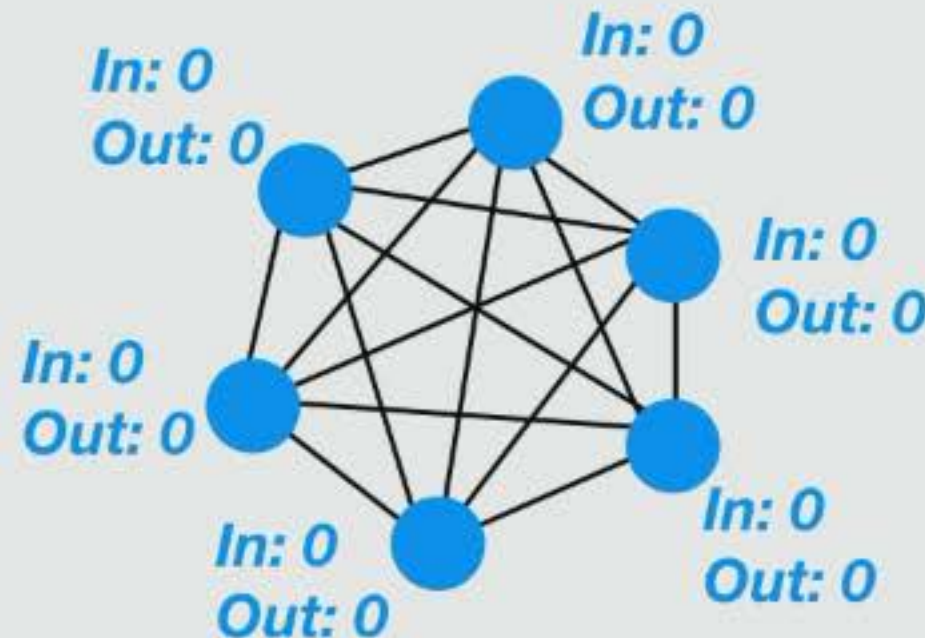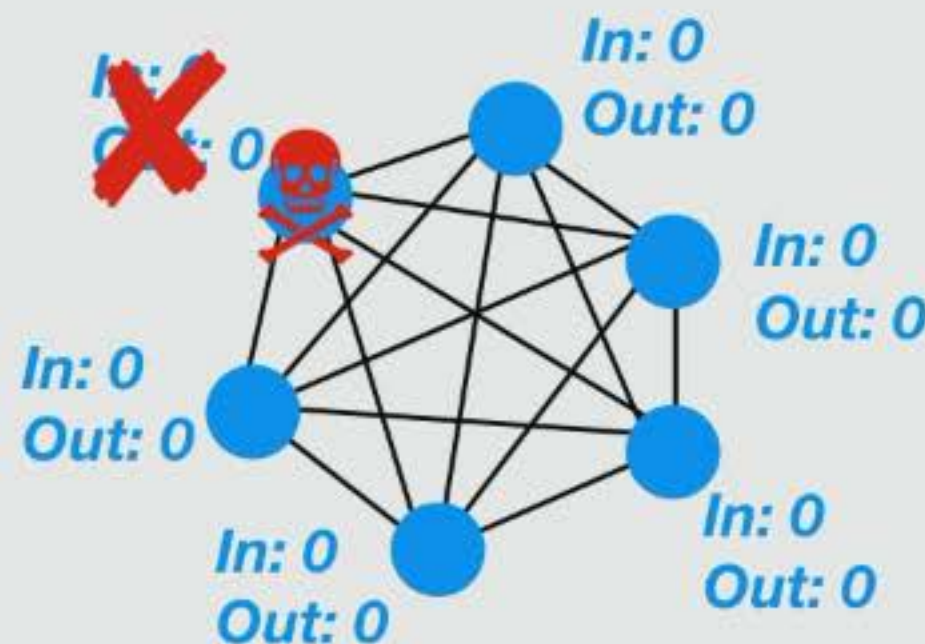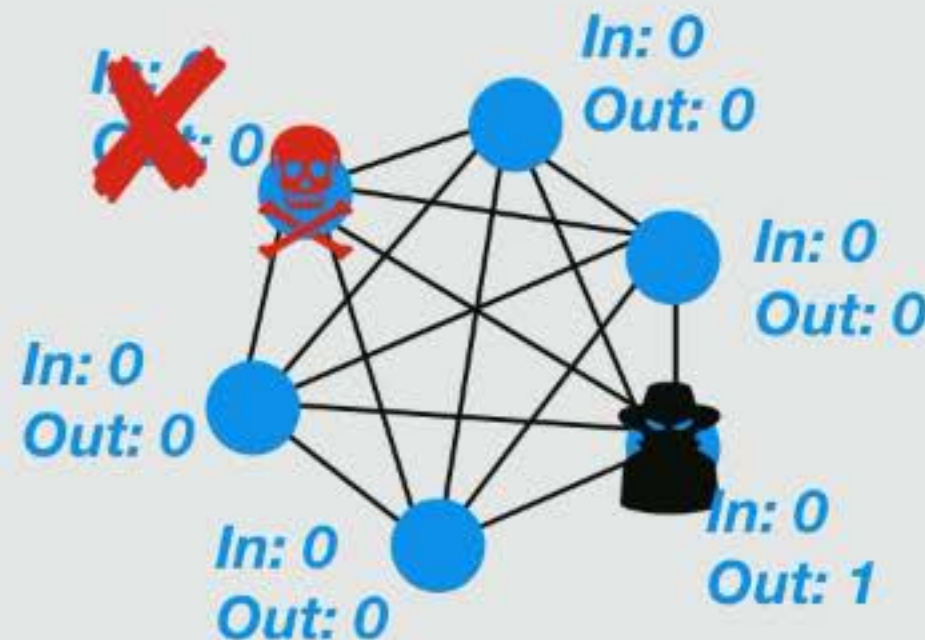  - Validity: output value must be input of some process

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something
    - Agreement: Every process outputs the same value
    - Validity: output value must be input of some process
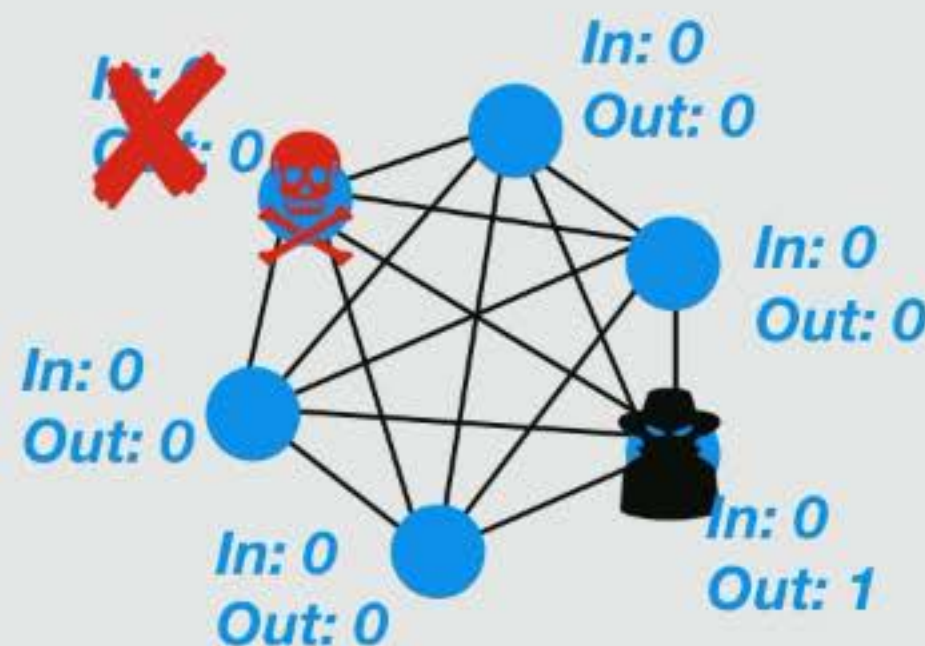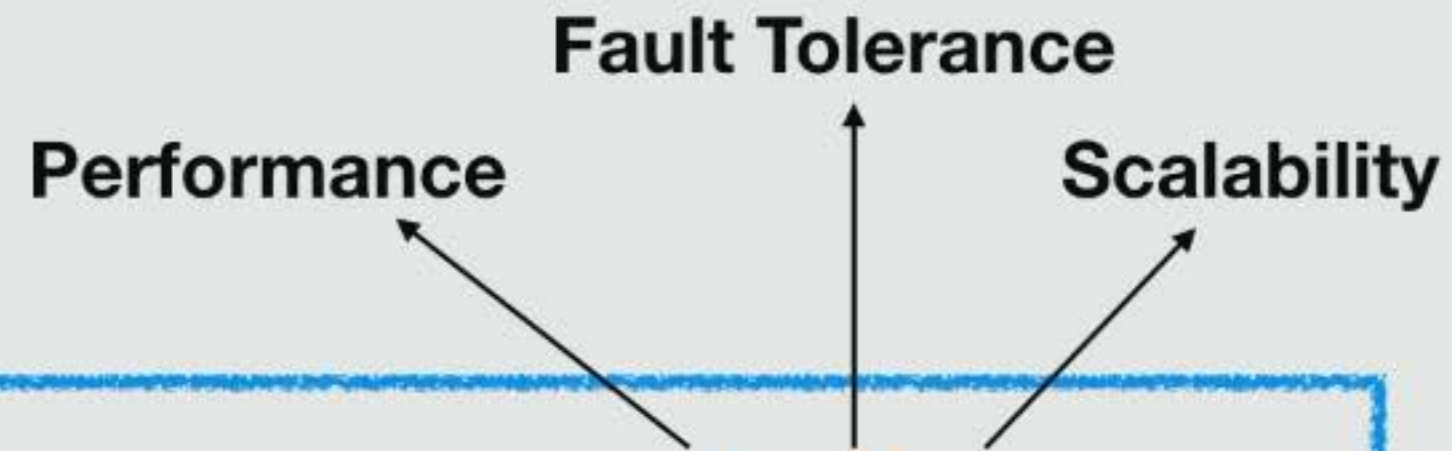
# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something
    - Agreement: Every process outputs the same value

    - Validity: output value must be input of some process

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something

  - Agreement: Every process outputs the same value

  - Validity: output value must be input of some process

- **Challenges**: **Asynchrony**, processes **crash** or are **Byzantine**
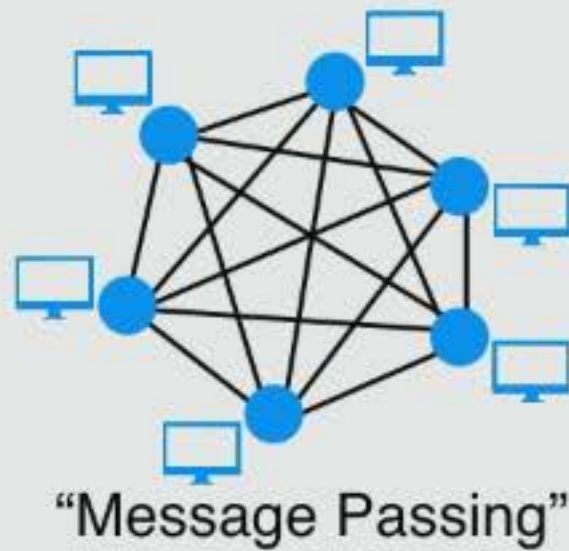
# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something

  - Agreement: Every process outputs the same value

  - Validity: output value must be input of some process

- **Challenges**: **Asynchrony**, processes **crash** or are **Byzantine**

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every process outputs something

    - Agreement: Every process outputs the same value

    - Validity: output value must be input of some process

- **Challenges**: **Asynchrony**, processes **crash** or are **Byzantine**

# Consensus: Definition

- **Input**: every process gets input

- **Output**: Every *correct* process outputs something

  - Agreement: Every process outputs the same value

  - Validity: output value must be input of some process

- **Challenges**: **Asynchrony**, processes **crash** or are **Byzantine**

# Communication Mechanisms

# Communication Mechanisms



"Message Passing"

# Communication Mechanisms



"Message Passing"



- Data centers, Internet

# Communication Mechanisms

"Message Passing"

"Shared Memory"

- Data centers, Internet

# Communication Mechanisms



"Message Passing"

"Shared Memory"

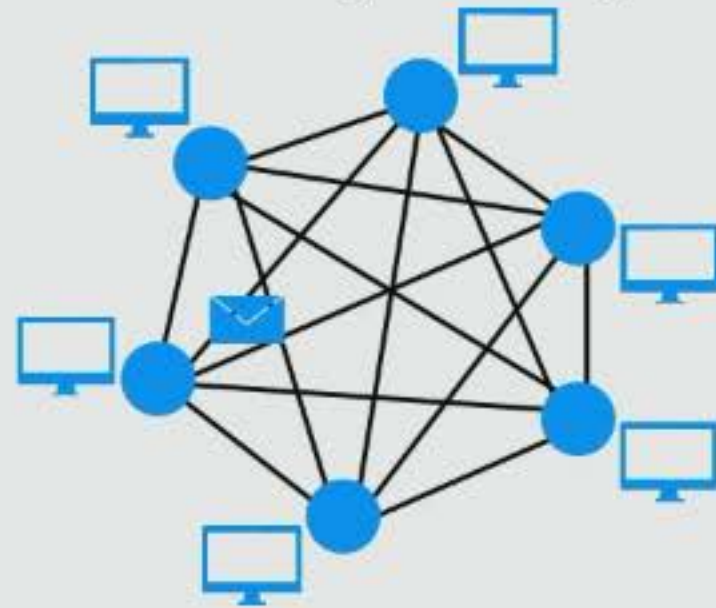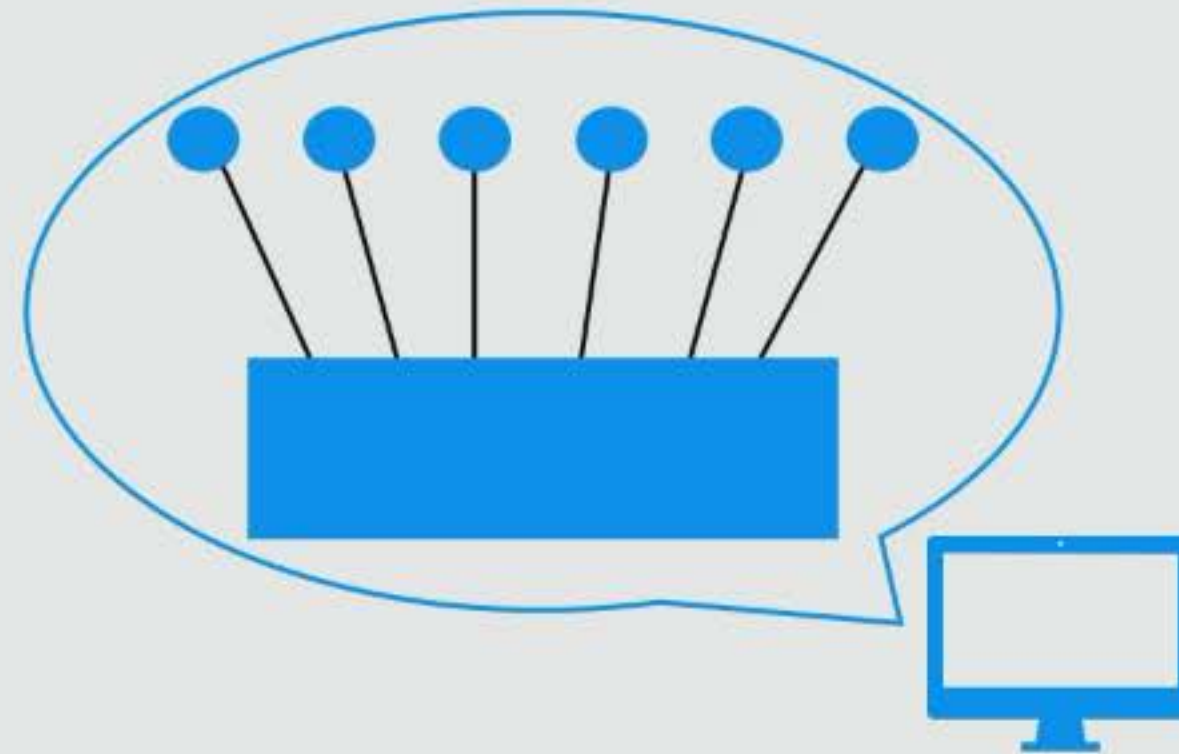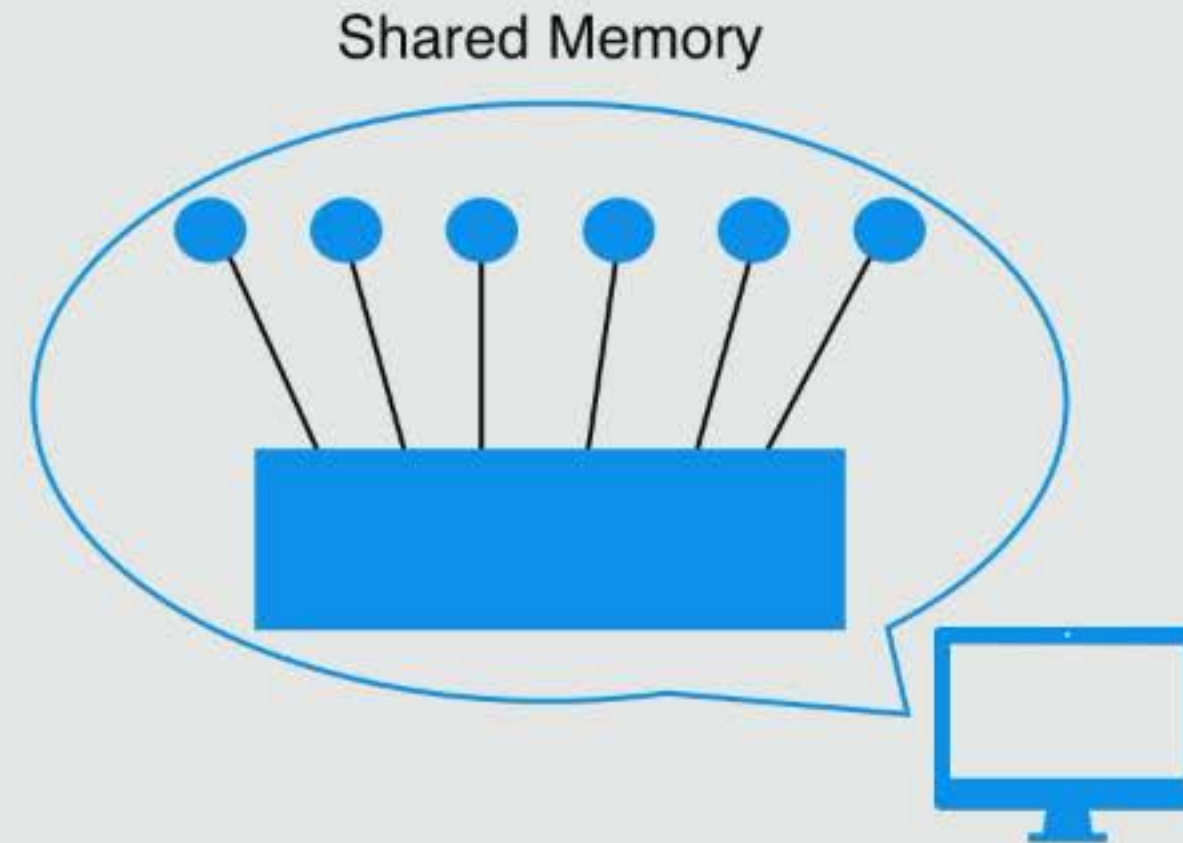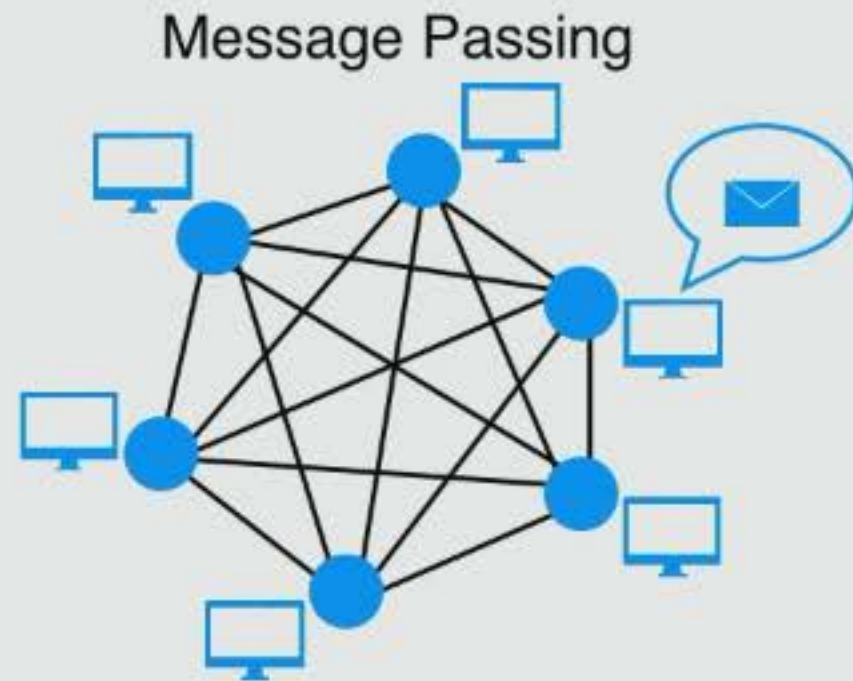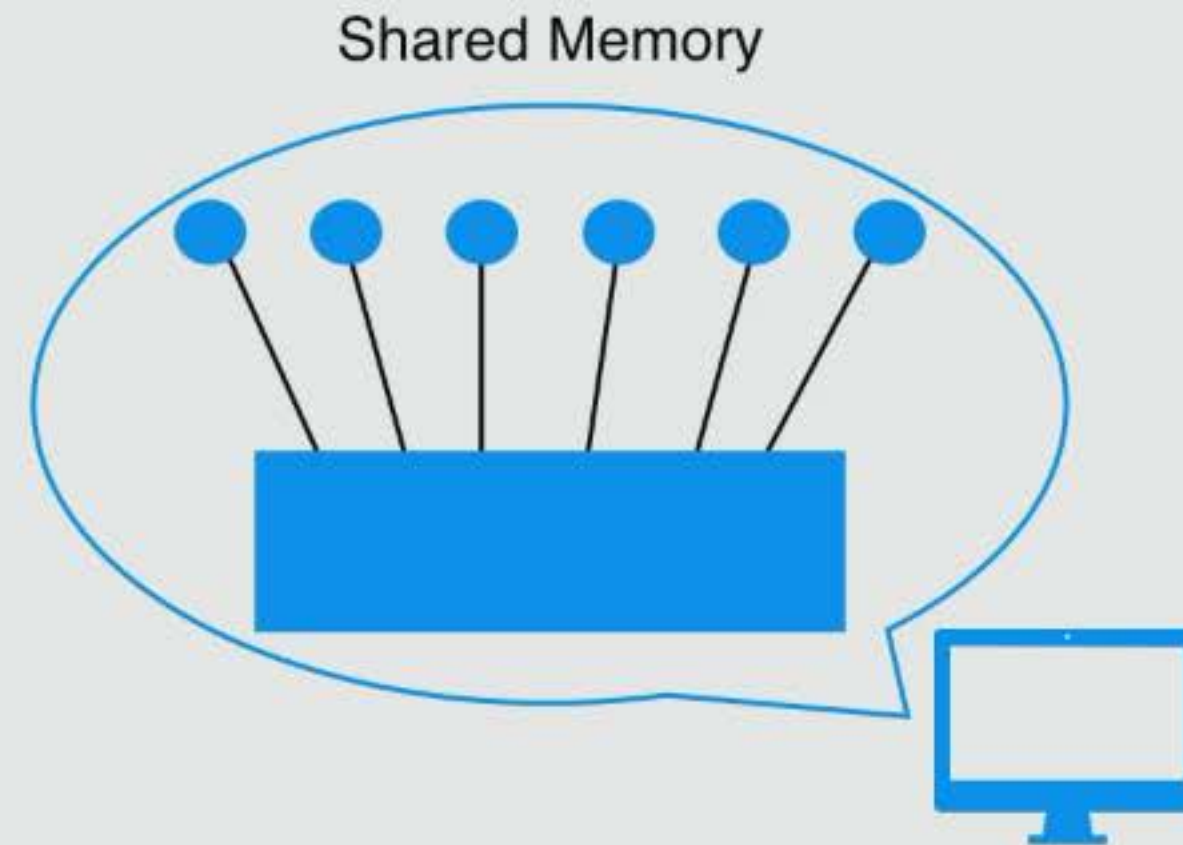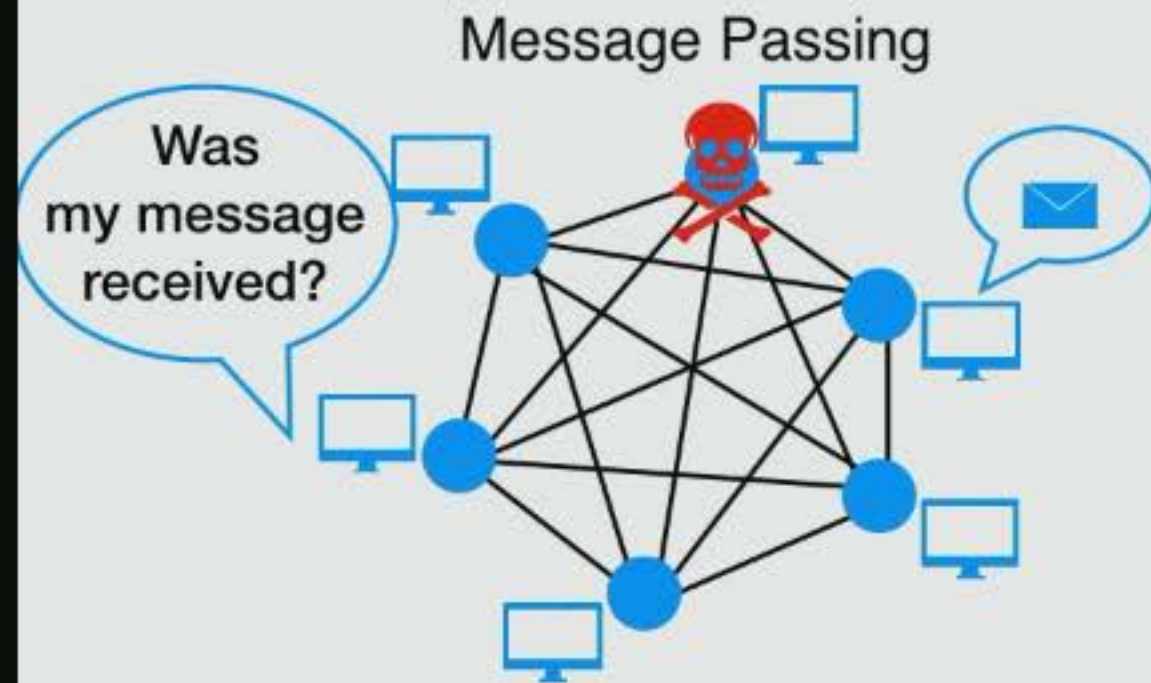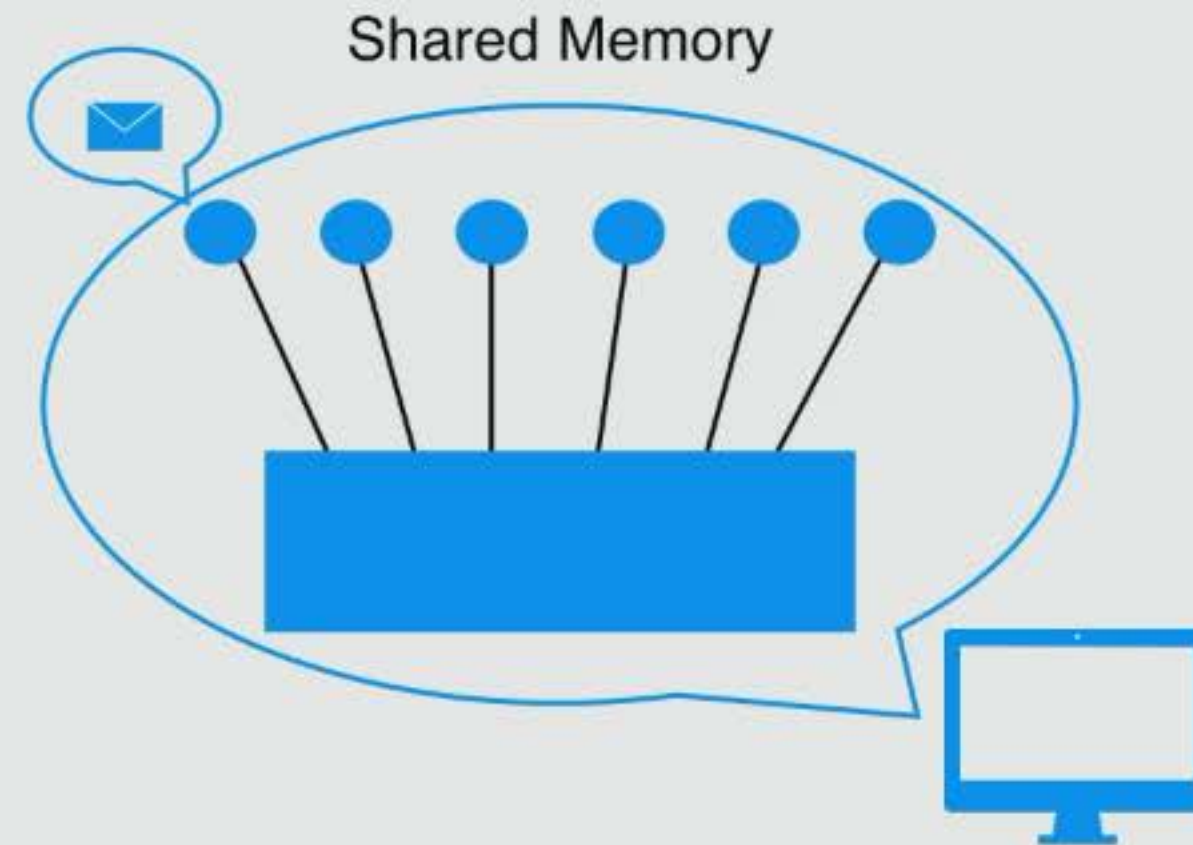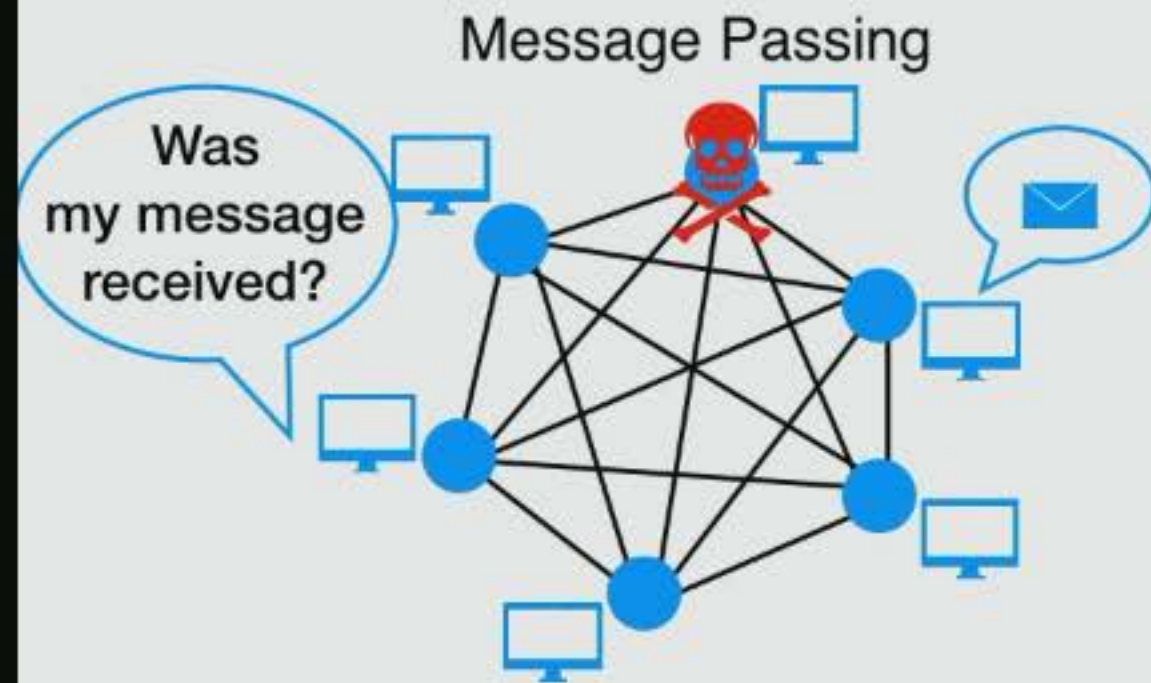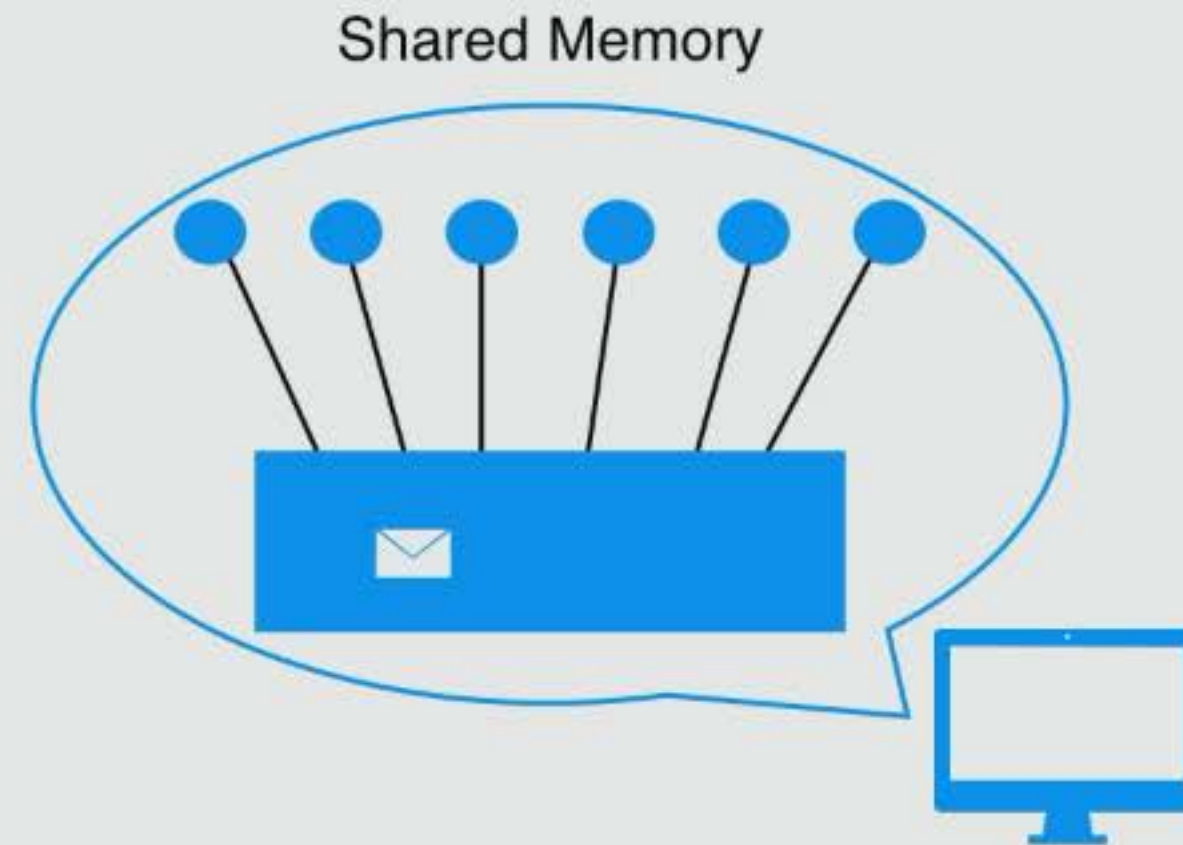- Data centers, Internet
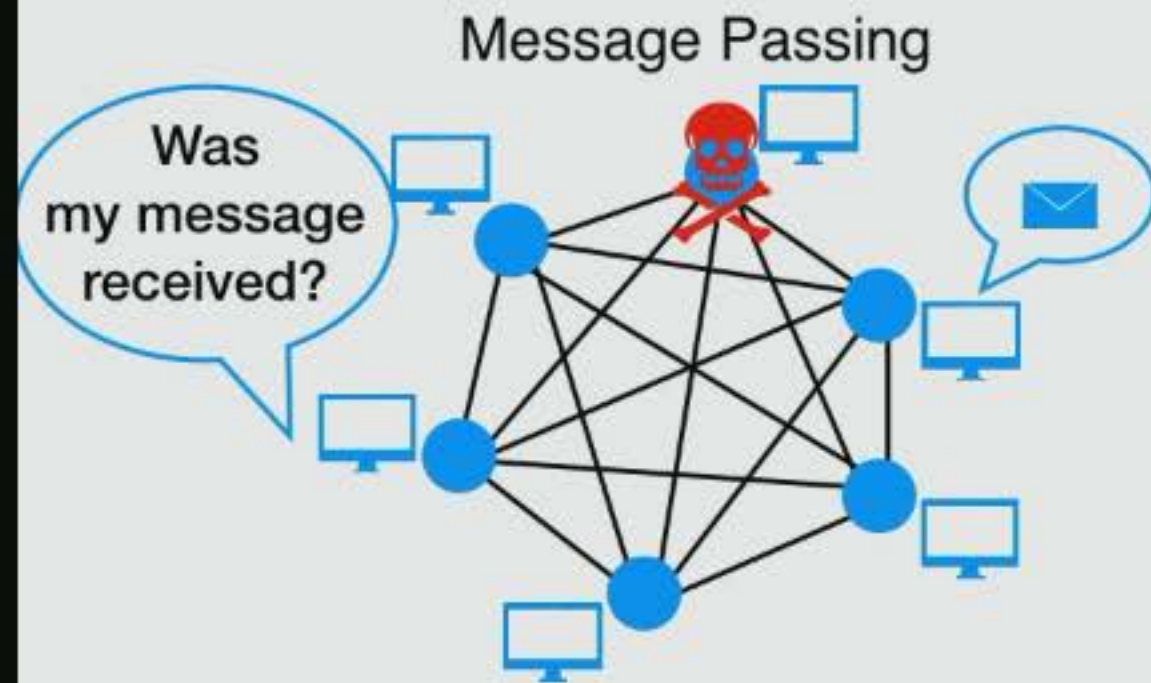
- Multicore machines

# Message Passing vs Shared Memory

# Message Passing vs Shared Memory

# Message Passing vs Shared Memory

# Message Passing vs Shared Memory

# Message Passing vs Shared Memory

# Message Passing vs Shared Memory

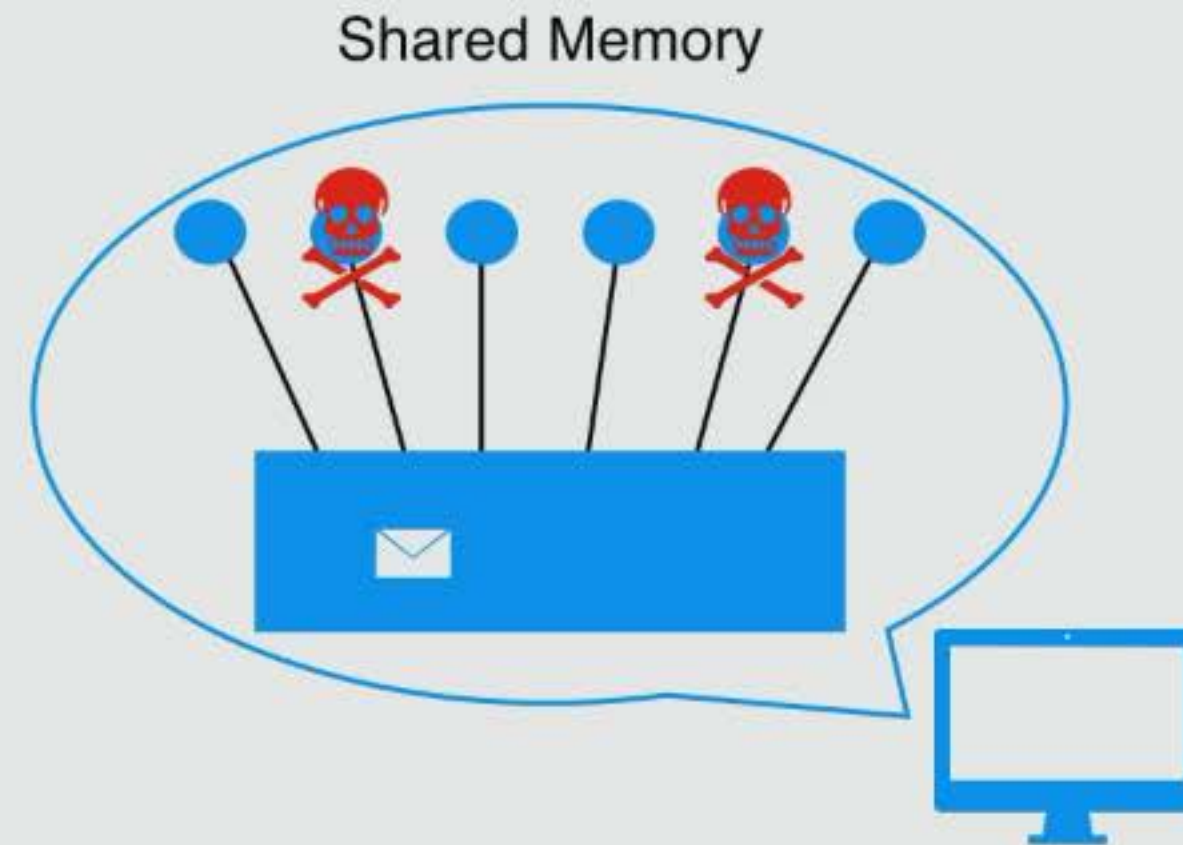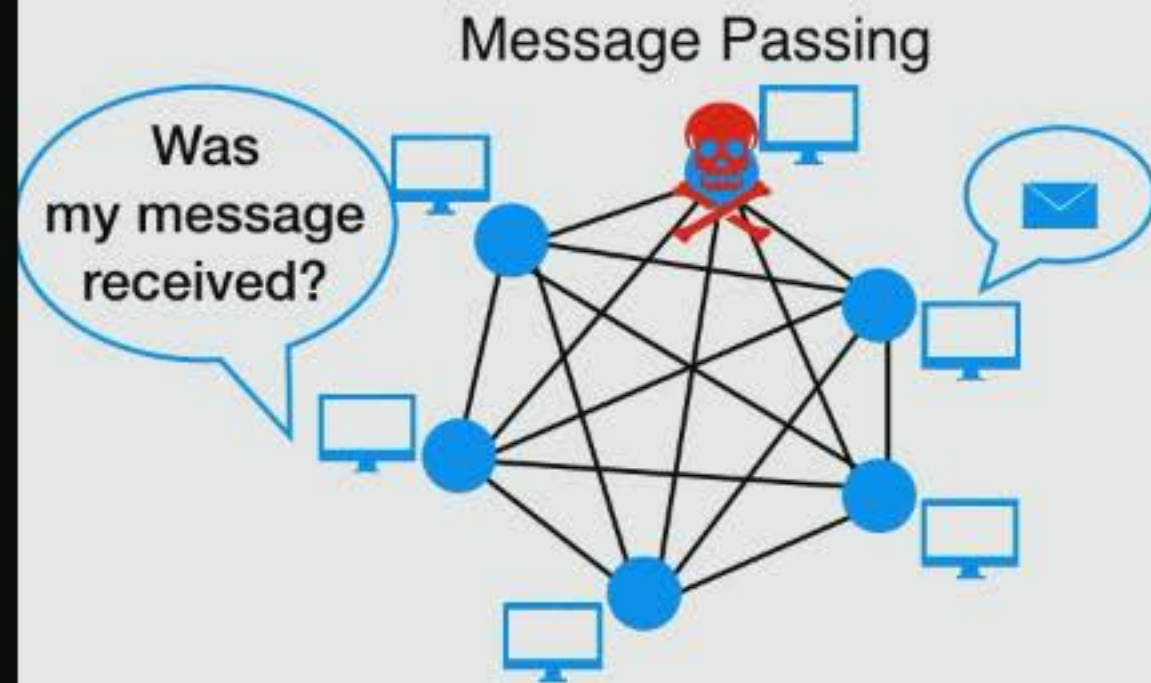# Message Passing vs Shared Memory

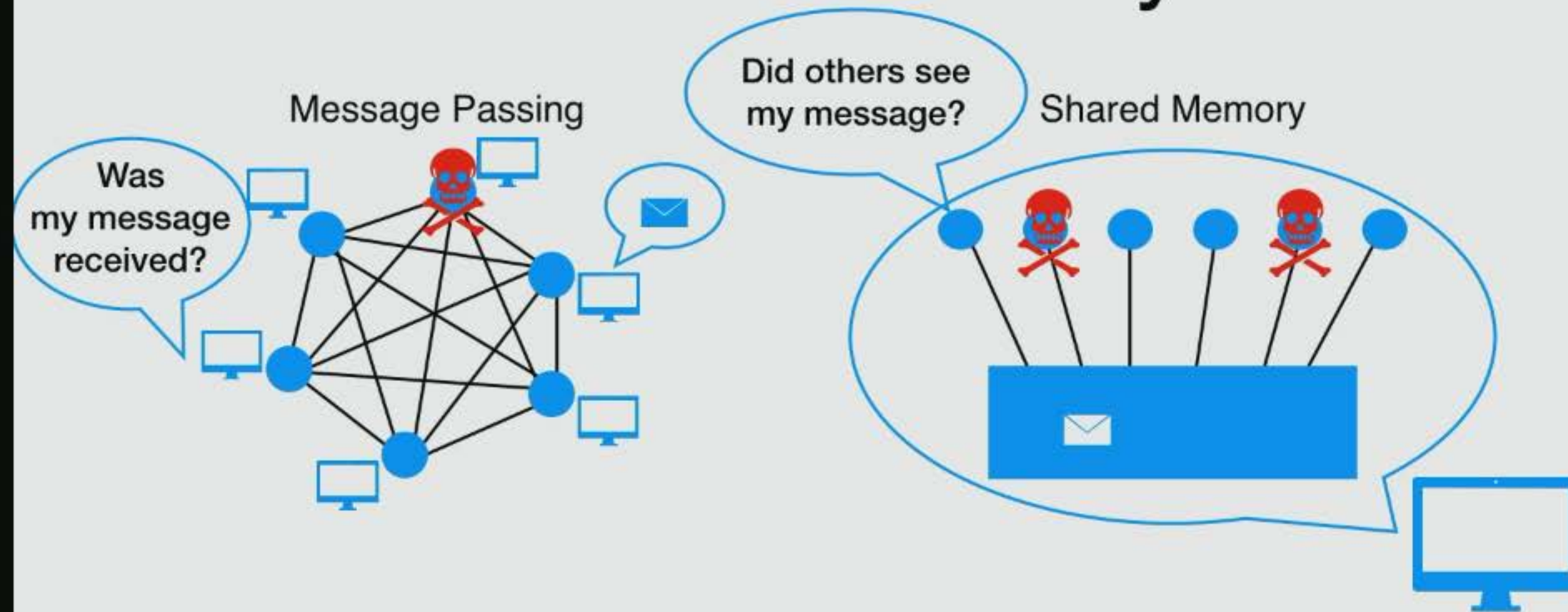# Message Passing vs Shared Memory

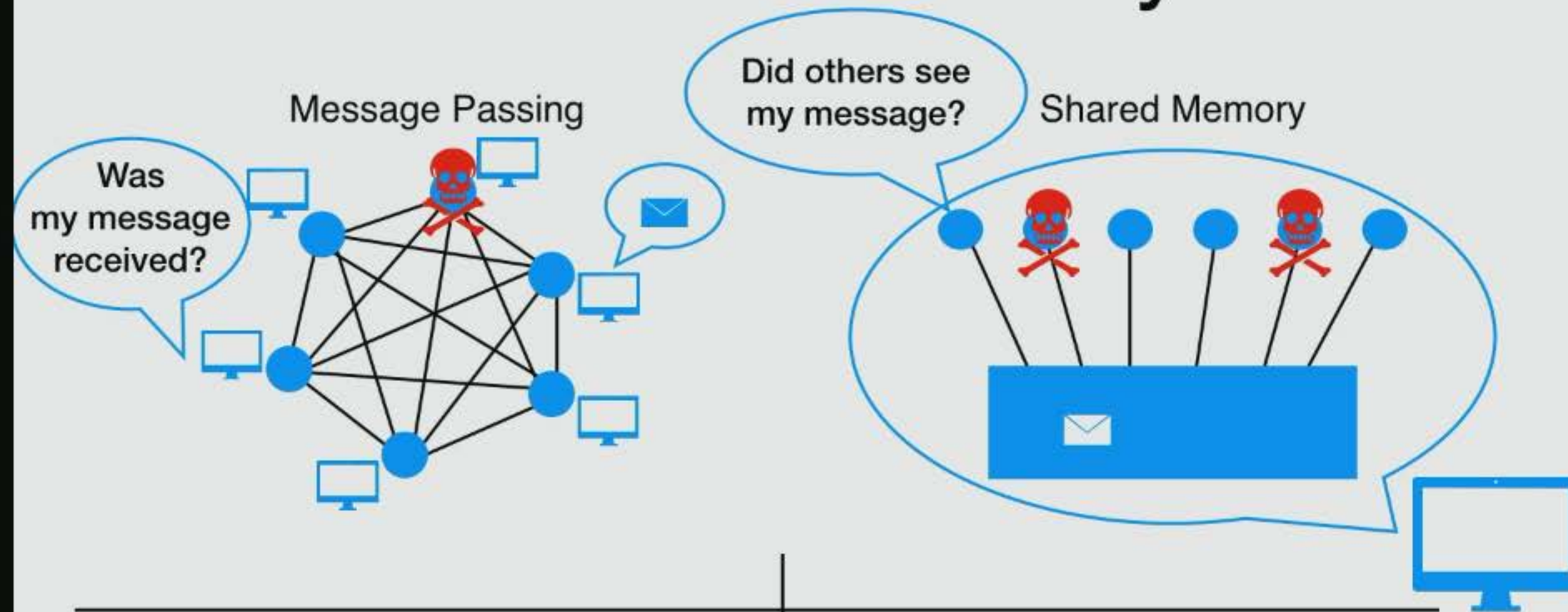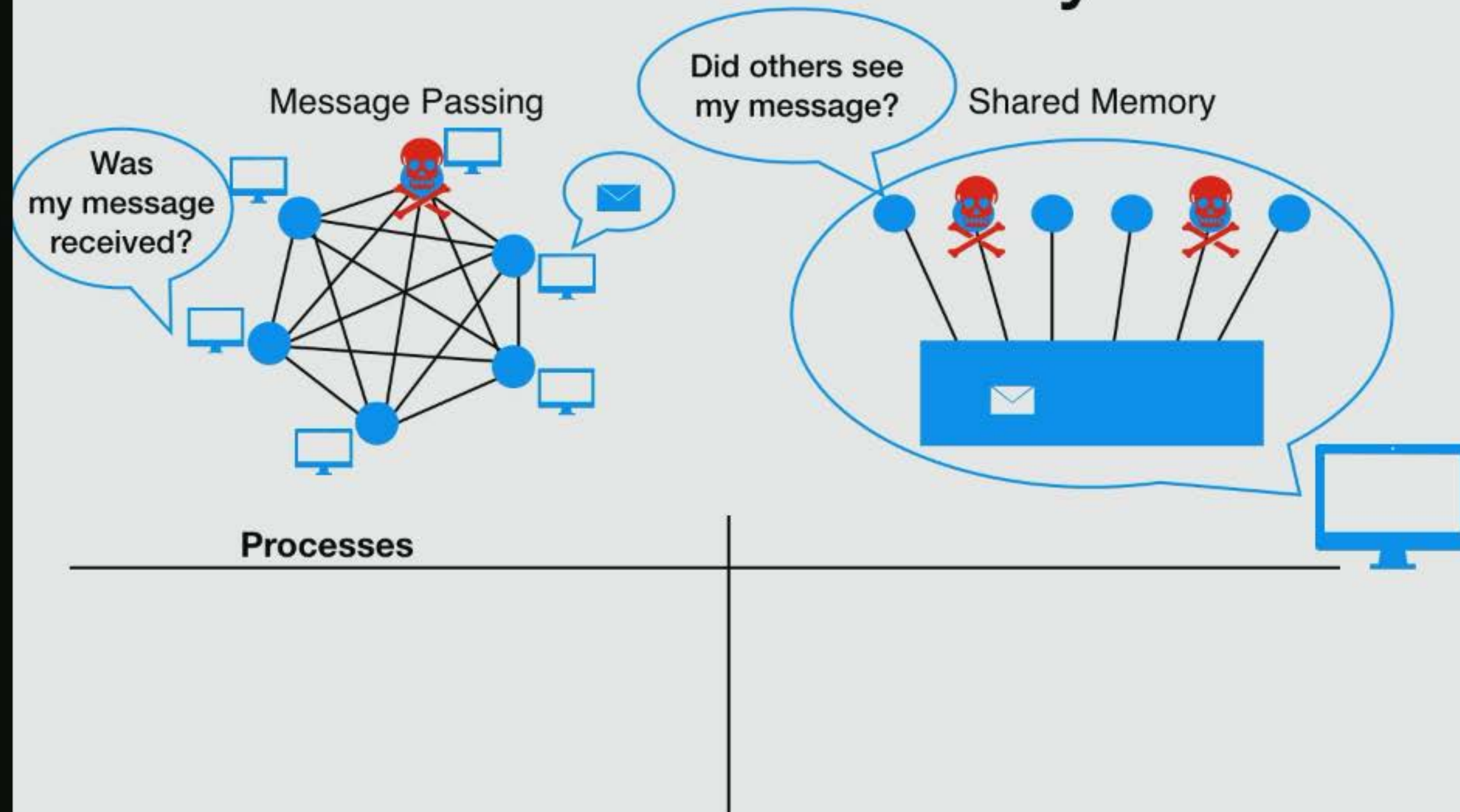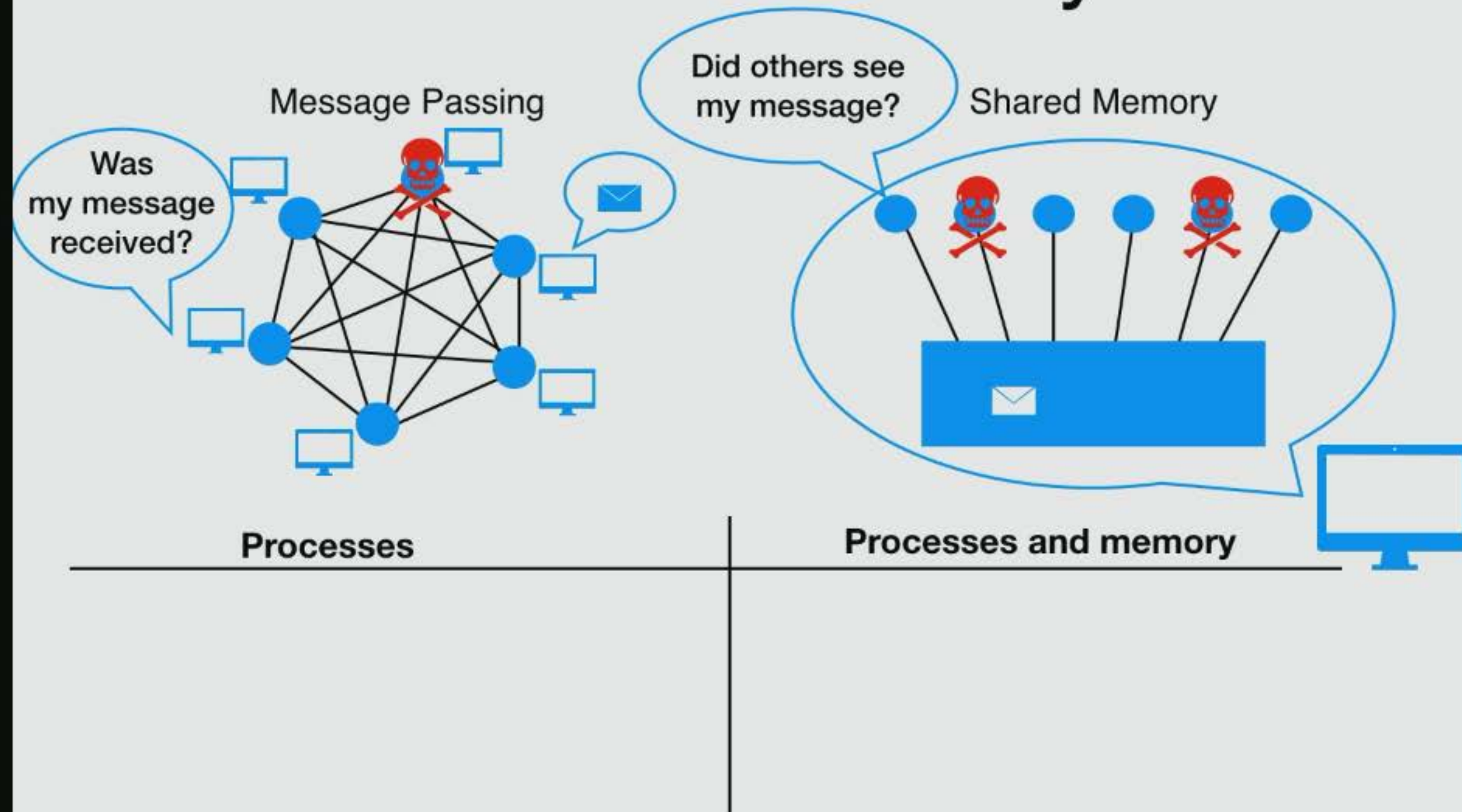# Message Passing vs Shared Memory

Message Passing vs Shared Memory

# Message Passing vs Shared Memory

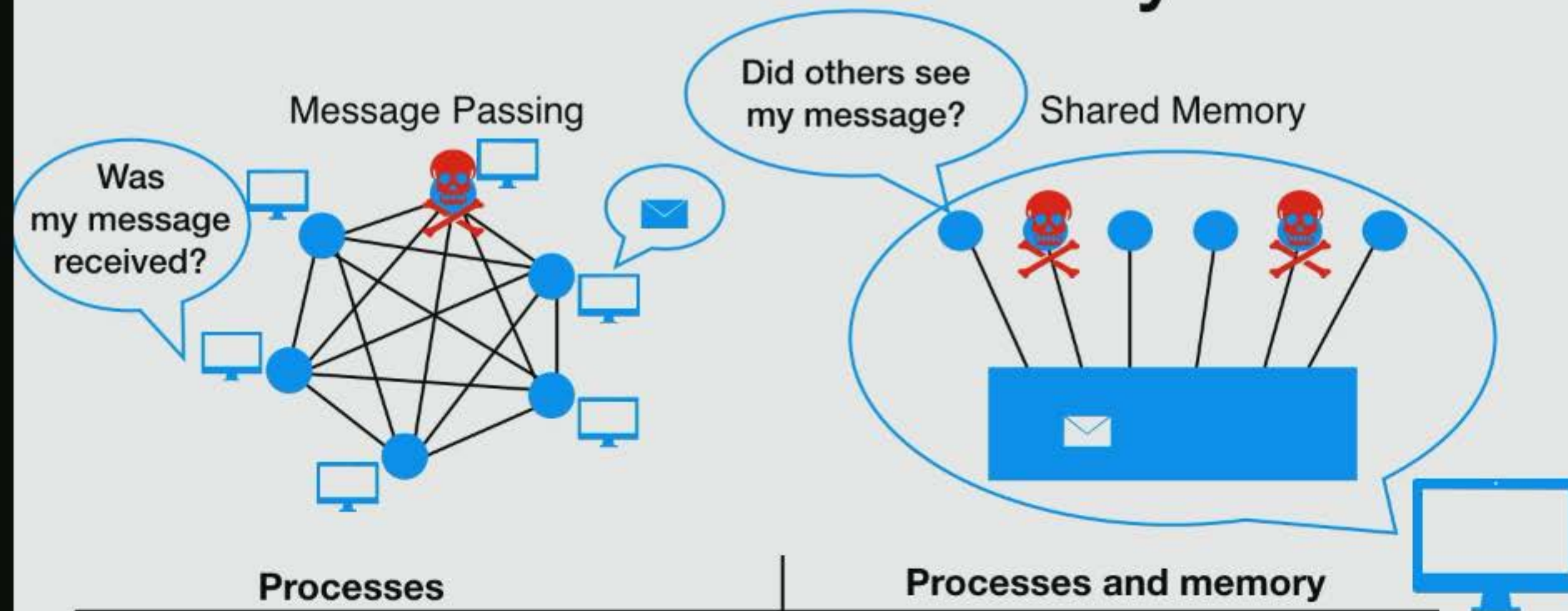# RDMA: Messages and Memory

Remote Direct Memory Access (RDMA)



(Network Interface Card)

# RDMA: Messages and Memory

**Remote Direct Memory Access (RDMA)**



(Network Interface Card)

Memory — NIC — NIC — Memory

CPU          CPU

RDMA: No involvement of host CPU!

# Main Take-Away

RDMA improves tradeoff between fault tolerance, performance and scalability

# Main Take-Away

Byzantine or crash failures of processes

RDMA improves tradeoff between fault tolerance, performance and scalability

# Main Take-Away

Byzantine or crash failures of processes

RDMA improves tradeoff between fault tolerance, performance and scalability

Common-case running time

# Main Take-Away

Byzantine or crash failures of processes

RDMA improves tradeoff between fault tolerance, performance and scalability

Common-case running time

Common case:
- Synchronous
- No Failures

Running time (agreement):
- Time until *first process* decides

# Main Take-Away

Byzantine or crash failures of processes

RDMA improves tradeoff between fault tolerance, performance and scalability

Common-case running time

Best case performance
Worst case resilience

Common case:
- Synchronous
- No Failures

Running time (agreement):
- Time until *first process* decides

# Outline

- **RDMA details and previous results**

- **Improving fault tolerance and performance**

- **Improving scalability**

# Previous Results

# Previous Results

| n = num processes f = num failures | Shared Memory | Message Passing |
| --- | --- | --- |

# Previous Results

| n = num processes<br>f = num failures | | Shared<br>Memory | Message<br>Passing |
|---|---|---|---|
| **Fault Tolerance** | Crash | n>f | n>2f |
| | Byzantine | N/A | n>3f |

# Previous Results

| n = num processes<br>f = num failures | | Shared<br>Memory | Message<br>Passing |
|---|---|---|---|
| Fault<br>Tolerance | Crash | n>f | n>2f |
| | Byzantine | N/A | n>3f |

# Previous Results

| n = num processes f = num failures | | Shared Memory | Message Passing |
|---|---|---|---|
| **Fault Tolerance** | Crash | n>f | n>2f |
| | Byzantine | N/A | n>3f |

# Previous Results

| n = num processes<br>f = num failures | | Shared Memory | Message Passing |
|---|---|---|---|
| Fault Tolerance | Crash | n>f | n>2f |
| | Byzantine | N/A | n>3f |

# Previous Results

| n = num processes<br>f = num failures | | Shared<br>Memory | Message<br>Passing |
|---|---|---|---|
| Fault<br>Tolerance | Crash | n>f | n>2f |
| | Byzantine | N/A | n>3f |
| Complexity*<br>(Best Case Round<br>Trips) | | 2 | 1 |

# Previous Results

| n = num processes<br>f = num failures | | Shared Memory | Message Passing |
|---|---|---|---|
| Fault Tolerance | Crash | n>f | n>2f |
| | Byzantine | N/A | n>3f |
| Complexity*<br>(Best Case Round Trips) | | 2 | 1 |
| Scalability<br>(processes in network) | | 10-100 | 10,000 - 100,000 |

# Previous Results

| n = num processes<br>f = num failures | | Shared Memory | Message Passing | RDMA Full<br>[ABGMZ'19] |
|---|---|---|---|---|
| **Fault Tolerance** | Crash | n>f | n>2f | n>f |
| | Byzantine | N/A | n>3f | n>2f |
| Complexity*<br>(Best Case Round Trips) | | 2 | 1 | 1 |
| Scalability<br>(processes in network) | | 10-100 | 10,000 - 100,000 | 10-100 |

# Previous Results

| n = num processes<br>f = num failures | | Shared Memory | Message Passing | RDMA Full [ABGMZ'19] | RDMA Scale [ABCGPT'18] |
|---|---|---|---|---|---|
| **Fault Tolerance** | Crash | n>f | n>2f | n>f | n>f+x (x∈[0,f]) |
| | Byzantine | N/A | n>3f | n>2f | - |
| **Complexity\* (Best Case Round Trips)** | | 2 | 1 | 1 | - |
| **Scalability (processes in network)** | | 10-100 | 10,000 - 100,000 | 10-100 | 10-100,000 |

# Data Center Technology: RDMA

# Data Center Technology: RDMA

- Can choose RDMA **connections**

# RDMA Scalability



Figure 5: Impact of connection multiplexing

[DragojevicNarayananHodsonCastro'14]

# RDMA Scalability



Figure 5: Impact of connection multiplexing

[DragojevicNarayananHodsonCastro'14]

# RDMA Scalability



Several connections per server

Figure 5: Impact of connection multiplexing

# RDMA Scalability



**Several connections per server**

As number of connections per server increases, throughput decreases.

Approximate number of connections per machine

Figure 5: Impact of connection multiplexing

[DragojevicNarayananHodsonCastro'14]

# Modeling Scalability

# Modeling Scalability

# Modeling Scalability

# Modeling Scalability



| RDMA Full [A**B**GMZ'19] | RDMA Scale [A**B**CGPT'18] |
|---|---|
| n>f | n>f+x (xε[0,f]) |
| n>2f | - |
| 1 | - |
| **10-100** | **10-100,000** |

Shared Memory Graph

Fault tolerance will depend on topology of SM graph

# Outline

✅ **RDMA details**

- Setting 1: **RDMA's full power** (complete graph)

  - **Crash-only** algorithm: $n>f$ tolerant, 1 round-trip

  - **Byzantine** algorithm: $n>2f$ tolerant, 1 round-trip

- Setting 2: **Scalability: Using RDMA sparingly** (incomplete graph)

  - Crash-only Algorithm: tolerance vs topology

# Data Center Technology: RDMA



- Can choose RDMA **connections** and **permissions**

- Can give different permissions for different **memory regions**

| | |
|---|---|
| p1: read | R1 |
| p3: write | R1& R2 |
| p6: read & write | R2 |
| p2, p5: none | — |

# Data Center Technology: RDMA

- Can choose RDMA **connections** and **permissions**

- Can give different permissions for different **memory regions**

| p1: read | R1 |
| p3: write | R1& R2 |
| p6: read & write | R2 |
| p2, p5: none | — |

# Data Center Technology: RDMA

- Can choose RDMA **connections** and **permissions**

- Can give different permissions for different **memory regions**

| p1: read | R1 |
| --- | --- |
| p3: write | R1 & R2 |
| p6: read & write | R2 |
| p2, p5: none | — |

| p1: read | R1 |
| --- | --- |
| p3: write | R1 & R2 |
| p6: read & write | R2 |
| p2, p4: none | — |

# Representing an RDMA Network

# Representing an RDMA Network

# Representing an RDMA Network

# Representing an RDMA Network

# Representing an RDMA Network



Kernel-level code executed on "memories"

Memories

All-to-all Connections

Processes

$p_1$  $p_2$  $p_3$  $p_4$  $p_5$  $p_6$

User-level code executed on processes

# RDMA Model

- **Asynchronous** network of ***n* processes** and ***m* memories**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

- **Memories** respond to changePermission requests with **acceptChange** policy

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

- **Memories** respond to changePermission requests with **acceptChange** policy

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

- **Memories** respond to changePermission requests with **acceptChange** policy

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

- **Memories** respond to changePermission requests with **acceptChange** policy

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

- **Memories** respond to changePermission requests with **acceptChange** policy

# RDMA Model

- **Asynchronous** network of *n* **processes** and *m* **memories**

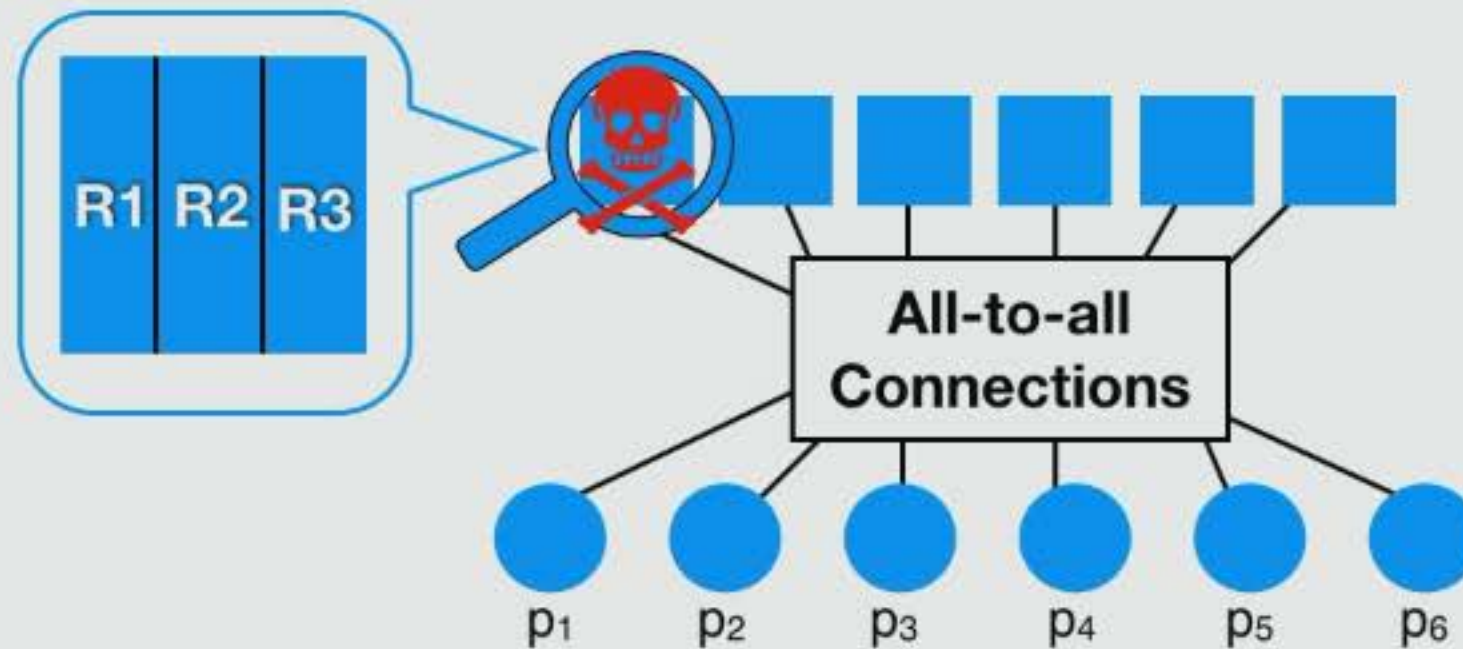- **Memories** fail by **crashing**, **processes** fail by **crashing** or being **Byzantine**

- Processes access memories through **read**, **write**, and **changePermission**

- **Memories** respond to changePermission requests with **acceptChange** policy

# Handling Memory Failures

Replication: Treat all memories the same

# Handling Memory Failures

Replication: Treat all memories the same

Send all write/read requests to all memories, wait to hear acknowledgement from majority

# Handling Memory Failures

Replication: Treat all memories the same

Send all write/read requests to all memories, wait to hear acknowledgement from majority

# Handling Memory Failures

Replication: Treat all memories the same

Send all write/read requests to all memories, wait to hear acknowledgement from majority

# Handling Memory Failures

Replication: Treat all memories the same

Send all write/read requests to all memories, wait to hear acknowledgement from majority

# Handling Memory Failures

Replication: Treat all memories the same

Send all write/read requests to all memories, wait to hear acknowledgement from majority

# Handling Memory Failures

Replication: Treat all memories the same

Send all write/read requests to all memories, wait to hear acknowledgement from majority

Instead of many faulty memories, we can now think of *one non-faulty memory!*

Acks: 4 ✓

All-to-all Connections

$p_1$ $p_2$ $p_3$ $p_4$ $p_5$ $p_6$

# Representing an RDMA Network

# Representing an RDMA Network

# Representing an RDMA Network

# Representing an RDMA Network

# Outline

✅ **RDMA details**

✅ Setting 1: **RDMA's full power** (complete graph)

- **Crash-only** algorithm: $n>f$ tolerant, 1 round-trip

- **Byzantine** algorithm: $n>2f$ tolerant, 1 round-trip

- Setting 2: **Scalability: Using RDMA sparingly** (incomplete graph)

  - Crash-only Algorithm: tolerance vs topology

# Paxos

Message passing consensus

# Representing an RDMA Network

# Paxos

Message passing consensus

# Paxos

Message passing consensus

- Two rounds:

# Paxos

Message passing consensus

- Two rounds:

  1. **prepare:** I want to propose a value!

# Paxos

Message passing consensus

- Two rounds:

    1. **prepare:** I want to propose a value!



1. I want to propose

# Paxos

Message passing consensus

- Two rounds:

  1. **prepare:** I want to propose a value!

# Paxos

[Lamport'98]

Message passing consensus

- Two rounds:

  1. **prepare:** I want to propose a value!

# Paxos

Message passing consensus

- Two rounds:

  1. **prepare:** I want to propose a value!

# Paxos

Message passing consensus

- Two rounds:

  1. **prepare:** I want to propose a value!

  2. **accept:** Here is my value!

# Paxos

Message passing consensus

- Two rounds:

    1. **prepare:** I want to propose a value!

    2. **accept:** Here is my value!

# Paxos

[Lamport'98]

Message passing consensus

- Two rounds:

    1. **prepare:** I want to propose a value!

    2. **accept:** Here is my value!

# Paxos

Message passing consensus

- Two rounds:

  1. **prepare:** I want to propose a value!

  2. **accept:** Here is my value!

# Paxos

Message passing consensus

- Two rounds:

    1. **prepare:** I want to propose a value!

    2. **accept:** Here is my value!

- Complications: If multiple processes try, **accept only last prepared**

**In the best case, only need one round!**



**Idea:**
- Choose **leader** a priori, let it **skip prepare phase**
- If leader is slow, others start executing prepare

# Disk Paxos

Idea: run Paxos on shared memory

**To send:** write your message in **your slot in disk**
**To receive:** read **others' slots in disk**

# Disk Paxos

Idea: run Paxos on shared memory

**Fault tolerance
n>f instead of n>2f**

**To send:** write your message in your slot in disk
**To receive:** read others' slots in disk

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

# Disk Paxos

Idea: run Paxos on shared memory

**Fault tolerance
n>f instead of n>2f**

**To send:** write your message in **your slot in disk**
**To receive:** read **others' slots in disk**

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-
case execution?

# Disk Paxos

Idea: run Paxos on shared memory

**Fault tolerance
n>f instead of n>2f**

**To send:** write your message in **your slot in disk**
**To receive:** read **others' slots in disk**

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-case execution?

- p1 proposes by writing

# Disk Paxos

Idea: run Paxos on shared memory

**Fault tolerance**
**n>f instead of n>2f**

**To send:** write your message in your slot in disk
**To receive:** read others' slots in disk

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-case execution?

- p1 proposes by writing

# Disk Paxos

Idea: run Paxos on shared memory

**Fault tolerance**
**n>f instead of n>2f**

**To send:** write your message in **your slot in disk**
**To receive:** read **others' slots in disk**

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-case execution?

- p1 proposes by writing

[GafniLamport'02]

# Disk Paxos

**Fault tolerance n>f instead of n>2f**

Idea: run Paxos on shared memory

**To send:** write your message in **your slot in disk**
**To receive:** read **others' slots in disk**

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-case execution?

- p1 proposes by writing

- p1 must read all slots to ensure it's the only proposer

# Disk Paxos

[GafniLamport'02]

Idea: run Paxos on shared memory

**Fault tolerance**
**n>f instead of n>2f**

**To send:** write your message in **your slot in disk**
**To receive:** read **others' slots in disk**

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-case execution?

- p1 proposes by writing

- p1 must read all slots to ensure it's the only proposer

**P1 doesn't know it's a good execution!**

# Disk Paxos

Idea: run Paxos on shared memory

**Fault tolerance
n>f instead of n>2f**

**To send:** write your message in **your slot in disk**
**To receive:** read **others' slots in disk**

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-case execution?

- p1 proposes by writing

- p1 must read all slots to ensure it's the only proposer

**P1 doesn't know it's a good execution!**

# Disk Paxos

[GafniLamport'02]

Idea: run Paxos on shared memory

**Fault tolerance
n>f instead of n>2f**

**To send:** write your message in your slot in disk
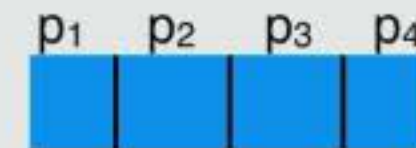**To receive:** read others' slots in disk

**Paxos:**
- **Leader proposes its value by sending it to everyone.**
- **Everyone tells leader whether they heard from anyone else**

What happens in common-case execution?

- p1 proposes by writing

- p1 must read all slots to ensure it's the only proposer

**P1 doesn't know it's a good execution!**

Time (in round trips):

1

1

total: 2

# Disk Paxos with Permissions

In Disk, proposer must **read every value from disk** to know whether someone is competing with it.

Idea: leverage RDMA dynamic permissions to get rid of this step.



Memory

p₁   p₂   p₃   p₄

# Disk Paxos with Permissions

[ABGMZ'19]

In Disk, proposer must **read every value from disk** to know whether someone is competing with it.

Idea: leverage RDMA dynamic permissions to get rid of this step.

Memory

I will give *write permission* only to the last person who requested it.

p₁    p₂    p₃    p₄

Request permission from memory

# Disk Paxos with Permissions

[ABGMZ'19]

In Disk, proposer must **read every value from disk** to know whether someone is competing with it.

Idea: leverage RDMA dynamic permissions to get rid of this step.

I will give *write permission* only to the last person who requested it.

Memory

$p_1$    $p_2$    $p_3$    $p_4$

Request permission from memory

# Outline

✓ **RDMA details**

- Column 1: **RDMA's full power** (complete graph)

  ✓ **Crash-only** algorithm: $n>f$ tolerant, 1 round-trip

  - **Byzantine** algorithm: $n>2f$ tolerant, 1 round-trip

- Column 2: **Scalability: Using RDMA sparingly** (incomplete graph)

  - Crash-only Algorithm: tolerance vs topology

# Byzantine Algorithm Breakdown

Two pieces:

# Byzantine Algorithm Breakdown

Two pieces:

- **CheapQuorum:** Fast (1 round trip) algorithm that **aborts** at first sign of trouble

# Byzantine Algorithm Breakdown

Two pieces:

- **CheapQuorum:** Fast (1 round trip) algorithm that **aborts** at first sign of trouble

CheapQuorum

29

# Byzantine Algorithm Breakdown

Two pieces:

- **CheapQuorum:** Fast (1 round trip) algorithm that **aborts** at first sign of trouble

  CheapQuorum

  - Use permissions to get speed

# Byzantine Algorithm Breakdown

Two pieces:

- **CheapQuorum:** Fast (1 round trip) algorithm that **aborts** at first sign of trouble

  - Use permissions to get speed

- **Robust Backup:** Slow algorithm that is tolerant to $n > 2f$ Byzantine failures

CheapQuorum

# Byzantine Algorithm Breakdown

Two pieces:

- **CheapQuorum:** Fast (1 round trip) algorithm that **aborts** at first sign of trouble

  - Use permissions to get speed

- **Robust Backup:** Slow algorithm that is tolerant to $n > 2f$ Byzantine failures

  - Shared memory algorithm

CheapQuorum

Robust Backup

# [A**B**GMZ'19] Robust Algorithm

- High level idea: run Paxos, but replace messaging primitives (send/receive) with special *non-equivocating broadcast/ deliver*

# [ABGMZ'19] Robust Algorithm

- High level idea: run Paxos, but replace messaging primitives (send/receive) with special *non-equivocating broadcast/ deliver*

  **Equivocation:**

# [ABGMZ'19] Robust Algorithm

- High level idea: run Paxos, but replace messaging primitives (send/receive) with special **non-equivocating broadcast/deliver**

  **Equivocation:**

# [A**B**GMZ'19] Robust Algorithm

- High level idea: run Paxos, but replace messaging primitives (send/receive) with special *non-equivocating broadcast/ deliver*

  **Equivocation:**

# Robust Algorithm

- High level idea: run Paxos, but replace messaging primitives (send/receive) with special *non-equivocating broadcast/ deliver*

    **Equivocation:**

# [A**B**GMZ'19] Robust Algorithm

- High level idea: run Paxos, but replace messaging primitives (send/receive) with special *non-equivocating broadcast/ deliver*

  **Equivocation:**

  **Broadcast primitive prevents this behavior**

# [A**B**GMZ'19] Robust Algorithm

- High level idea: run Paxos, but replace messaging primitives (send/receive) with special *non-equivocating broadcast/ deliver*

  **Equivocation:**

  **Broadcast primitive prevents this behavior**

- If we can prevent equivocation, then we can solve Byzantine agreement with n > 2f [ClementJunqueiraKateRodrigues'12]

# Preventing Equivocation

**Single Writer Multi Reader region per process**



| R1 | R2 | R3 | R4 | R5 | R6 |

Each process gets its own SWMR region

# Preventing Equivocation

**Single Writer Multi Reader region per process**

| R1 | R2 | R3 | R4 | R5 | R6 |
|----|----|----|----|----|----|
|    |    |    |    |    |    |

$p_1$ $p_2$ $p_3$ $p_4$ $p_5$ $p_6$

Each process gets its own SWMR region

Protocol: Sign and copy over everything that you see

# Preventing Equivocation

**Single Writer Multi Reader region per process**



Each process gets its own SWMR region

Protocol: Sign and copy over everything that you see

# Preventing Equivocation

**Single Writer Multi Reader region per process**

| R1 | R2 | R3 | R4 | R5 | R6 |
|----|----|----|----|----|----|
| (v, p1) | | ((v, p1), p3) | | | |

$p_1$  $p_2$  $p_3$  $p_4$  $p_5$  $p_6$

Each process gets its own SWMR region

Protocol: Sign and copy over everything that you see

Can now verify that others read the same value

# Preventing Equivocation

**Single Writer Multi Reader region per process**

| R1 | R2 | R3 | R4 | R5 | R6 |
|----|----|----|----|----|----|
| $(v, p1)$ | | $((v, p1), p3)$ | | | |

$p_1$  $p_2$  $p_3$  $p_4$  $p_5$  $p_6$

Each process gets its own SWMR region

Protocol: Sign and copy over everything that you see

Can now verify that others read the same value

# Preventing Equivocation

**Single Writer Multi Reader region per process**



| R1 | R2 | R3 | R4 | R5 | R6 |
|----|----|----|----|----|----|
| $(v, p1)$ | | $((v, p1), p3)$ | | | |

$p_1$  $p_2$  $p_3$  $p_4$  $p_5$  $p_6$

Each process gets its own SWMR region

Protocol: Sign and copy over everything that you see

Can now verify that others read the same value

# Preventing Equivocation

**Single Writer Multi Reader region per process**

| R1 | R2 | R3 | R4 | R5 | R6 |
|----|----|----|----|----|----|
| (v, p1) | | ((v, p1), p3) | | | |

p₁ p₂ p₃ p₄ p₅ p₆

Each process gets its own SWMR region

**p3 read the same value I did from p1**

Protocol: Sign and copy over everything that you see

Can now verify that others read the same value

# RDMA vs Previous Results

| n = num processes<br>f = num failures | | Shared Memory | Message Passing | RDMA Full<br>[ABGMZ'19] | RDMA Scale<br>[ABCGPT'18] |
|---|---|---|---|---|---|
| **Fault Tolerance** | Crash | n>f | n>2f | n>f | n>f+x<br>(x∈[0,f]) |
| | Byzantine | N/A | n>3f | n>2f | - |
| **Complexity*** (Best Case Round Trips) | | 2 | 1 | 1 | - |
| **Scalability** (processes in network) | | 10-100 | 10,000 - 100,000 | 10-100 | 10-100,000 |

32

# RDMA vs Previous Results

| n = num processes f = num failures | | Shared Memory | Message Passing | RDMA Full* [ABGMZ'19] | RDMA Scale [ABCGPT'18] |
|---|---|---|---|---|---|
| **Fault Tolerance** | Crash | n>f | n>2f | n>f | n>f+x (x∈[0,f]) |
| | Byzantine | N/A | n>3f | n>2f | - |
| **Complexity\* (Best Case Round Trips)** | | 2 | 1 | 1 | - |
| **Scalability (processes in network)** | | 10-100 | 10,000 – 100,000 | 10-100 | 10-100,000 |

*With up to half of the memories crashing

Performance

Fault Tolerance

Scalability

**Is RDMA fundamentally *better* than other communication mechanisms?**

**Yes!**
**RDMA gives us the power of shared memory without compromising performance**

**Plus better Byzantine algorithms**

✓ ✓ Fault Tolerance ✗ Scalability
Performance

**Is RDMA fundamentally *better* than other communication mechanisms?**

**Yes!**

**RDMA gives us the power of shared memory without compromising performance**

**Plus better Byzantine algorithms**

Can we scale better and still retain some of RDMA's advantages?

# Scalability

What prevented our algorithms from scaling?

# Scalability

What prevented our algorithms from scaling?

**Many open connections**

# Scalability

What prevented our algorithms from scaling?

**Many open connections**

**NIC experiences frequent cache misses**

# Scalability

What prevented our algorithms from scaling?

**Many open connections**

↓

**NIC experiences frequent cache misses**

↓

# Scalability

What prevented our algorithms from scaling?

**Many open connections**

↓

**NIC experiences frequent cache misses**

↓

**Slower communication**

# Scalability

What prevented our algorithms from scaling?

**Many open connections**

↓

**NIC experiences frequent cache misses**

↓

**Slower communication**

Solution: don't open all connections.

# Scalability



What prevented our algorithms from scaling?

**Many open connections**

↓

**NIC experiences frequent cache misses**

↓

**Slower communication**

Solution: don't open all connections.

# Scalability



What prevented our algorithms from scaling?

**Many open connections *per machine***

⬇

**NIC experiences frequent cache misses**

⬇

**Slower communication**

Solution: don't open all connections.

# Scalability

What prevented our algorithms from scaling?

**Many open connections _per machine_**

**NIC experiences frequent cache misses**

**Slower communication**

Solution: don't open all connections.

**Shared Memory Graph**

# Scalability



What prevented our algorithms from scaling?

**Many open connections *per machine***

↓

**NIC experiences frequent cache misses**

↓

**Slower communication**

**Shared Memory Graph**



Solution: don't open all connections.

Goal: Keep **degree** of shared memory graph low

# Simplifying the Model

- Can choose RDMA **connections** ~~and~~ **permissions**
- ~~Can give different permissions for different~~ **memory regions**

$p_1$ $p_2$ $p_3$ $p_4$ $p_5$ $p_6$

Memory NIC CPU

# Simplifying the Model

- Can choose RDMA **connections** and **permissions**
- ~~Can give different permissions for different **memory regions**~~

**No Byzantine failures**
**No memory failures**

# The M&M model

- **Asynchronous** network of $n$ processes with up to $f$ **crash failures**

# The M&M model

- **Asynchronous** network of $n$ processes with up to $f$ **crash failures**

# The M&M model

- **Asynchronous** network of *n* processes with up to *f crash failures*

- Fully-connected message passing network: **nodes=procs, edges=links**

# The M&M model

- *Asynchronous* network of $n$ processes with up to $f$ *crash failures*

- Fully-connected message passing network: *nodes=procs, edges=links*

- *Each node owns a piece of memory*

# The M&M model

- *Asynchronous* network of $n$ processes with up to $f$ *crash failures*

- Fully-connected message passing network: *nodes=procs, edges=links*

- *Each node owns a piece of memory*

- Shared memory graph, $G_{SM} = (V, E)$

# The M&M model

- *Asynchronous* network of *n* processes with up to *f crash failures*

- Fully-connected message passing network: *nodes=procs, edges=links*

- *Each node owns a piece of memory*

- Shared memory graph, $G_{SM} = (V, E)$

- Nodes *u* and *v* can access each other's memory iff $(u,v) \in E$

# The M&M model

- **Asynchronous** network of *n* processes with up to *f crash failures*

- Fully-connected message passing network: **nodes=procs, edges=links**

- **Each node owns a piece of memory**

- **Shared memory graph, $G_{SM} = (V, E)$**

- **Nodes *u* and *v* can access each other's memory iff $(u,v) \epsilon$ E**

# The M&M model

- *Asynchronous* network of *n* processes with up to *f crash failures*

- Fully-connected message passing network: *nodes=procs, edges=links*

- *Each node owns a piece of memory*

- Shared memory graph, $G_{SM} = (V, E)$

- Nodes *u* and *v* can access each other's memory iff $(u,v) \in E$

- Processes may crash, but their *memory remains accessible*

# The M&M model

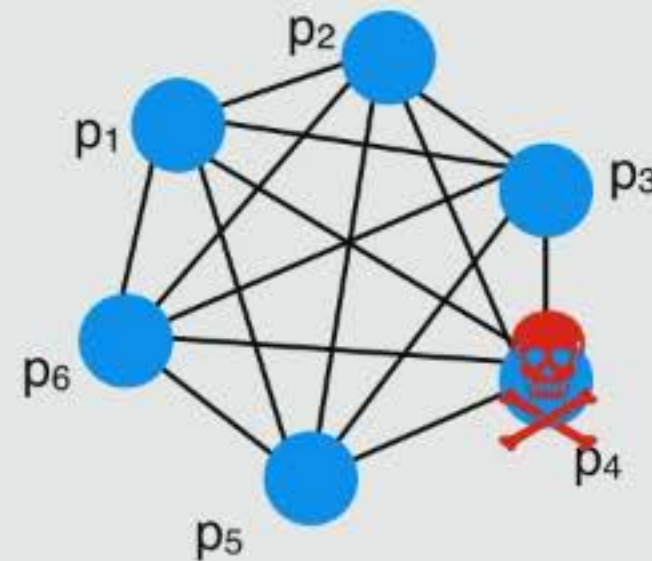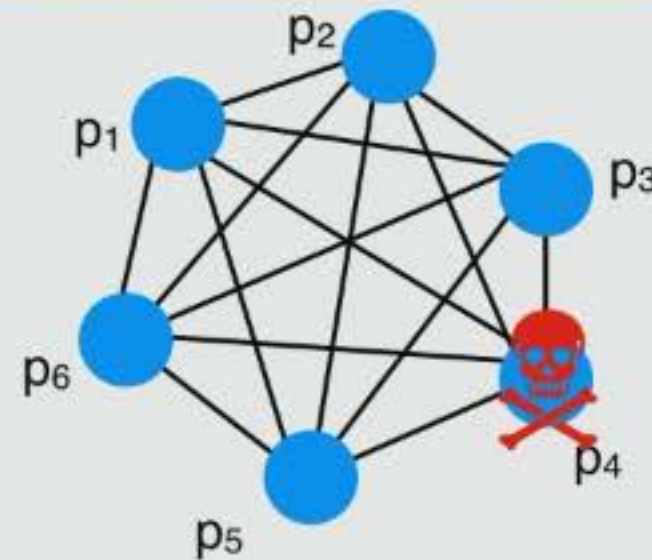- *Asynchronous* network of $n$ processes with up to $f$ *crash failures*

- Fully-connected message passing network: *nodes=procs, edges=links*

- *Each node owns a piece of memory*

- Shared memory graph, $G_{SM} = (V, E)$

- Nodes $u$ and $v$ can access each other's memory iff $(u,v) \in E$

- Processes may crash, but their *memory remains accessible*



p₂ was here

# The M&M model

- *Asynchronous* network of $n$ processes with up to $f$ *crash failures*

- Fully-connected message passing network: *nodes=procs, edges=links*
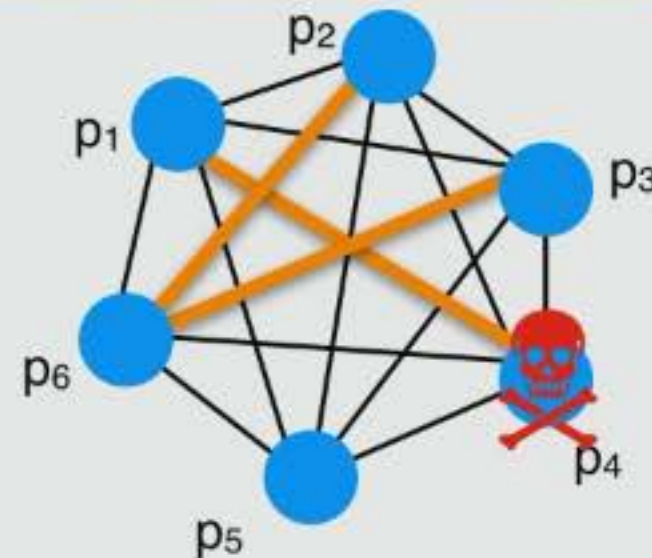
- *Each node owns a piece of memory*
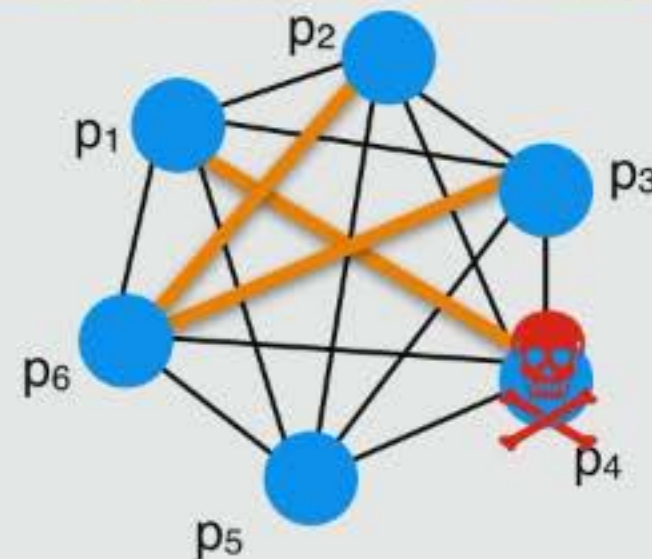
- Shared memory graph, $G_{SM} = (V, E)$

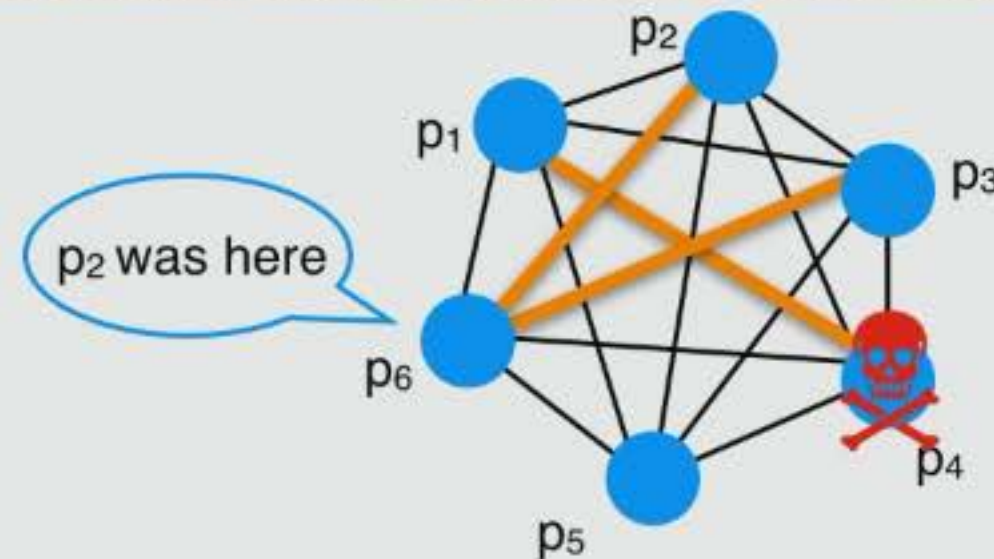- Nodes $u$ and $v$ can access each other's memory iff $(u,v) \in E$

- Processes may crash, but their *memory remains accessible*

$p_2$ was here

# The M&M model

- *Asynchronous* network of *n* processes with up to *f crash failures*

- Fully-connected message passing network: *nodes=procs, edges=links*

- *Each node owns a piece of memory*

- Shared memory graph, $G_{SM} = (V, E)$

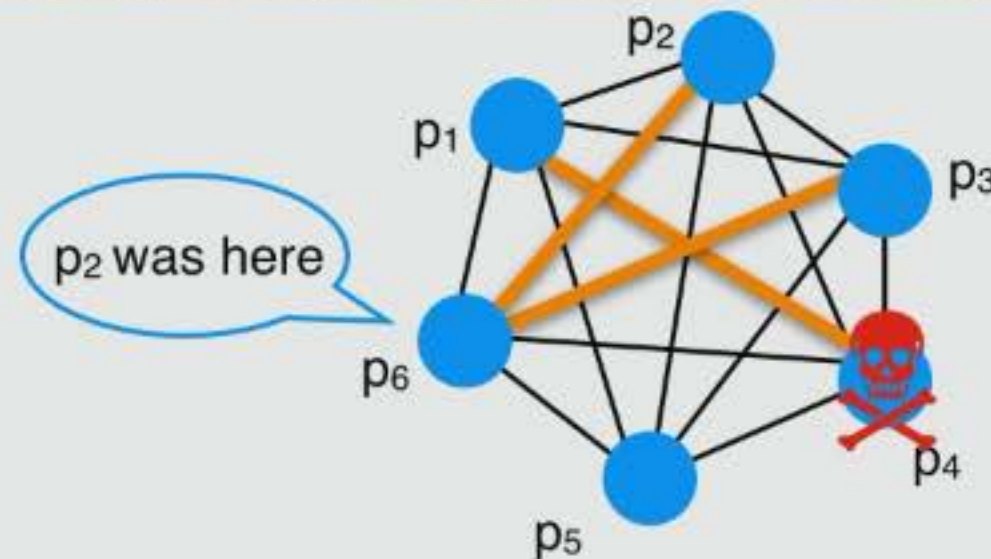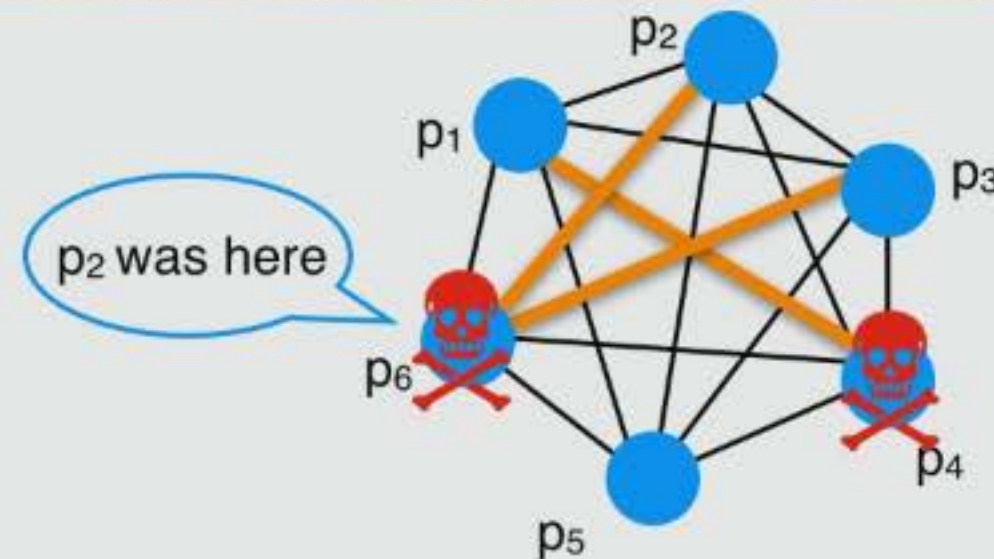- Nodes *u* and *v* can access each other's memory iff $(u,v) \in E$

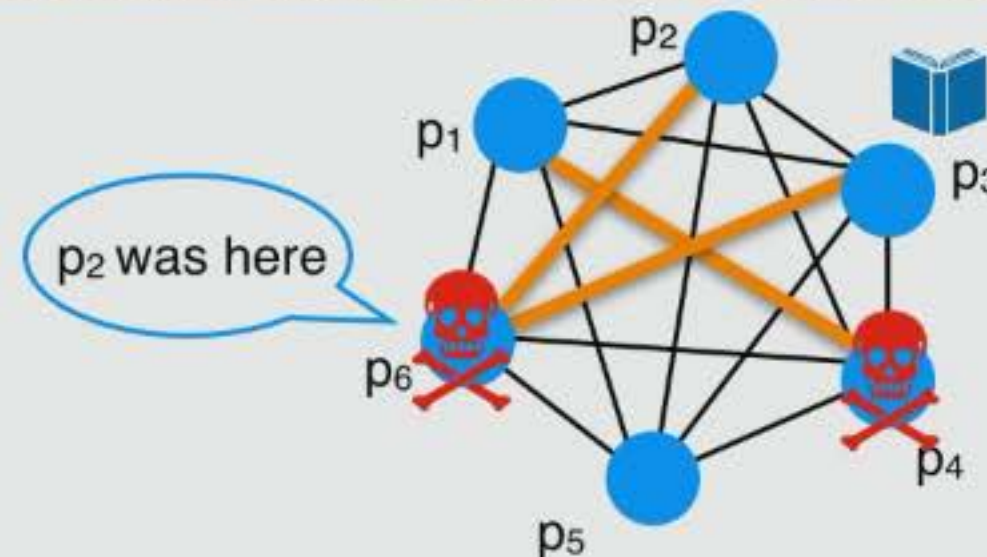- Processes may crash, but their *memory remains accessible*

p₂ was here

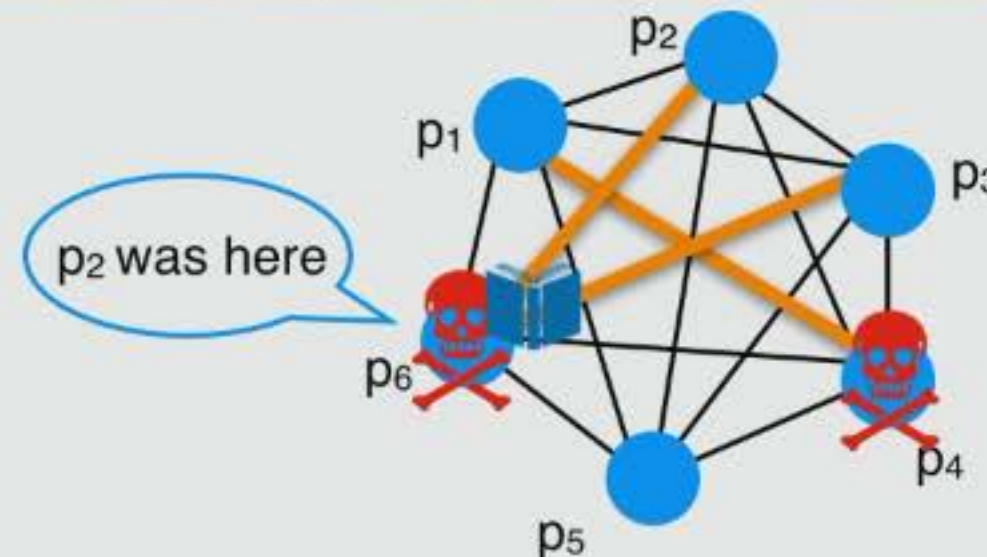# The M&M model

- **_Asynchronous_** network of _n_ processes with up to _f_ **_crash failures_**

- Fully-connected message passing network: **_nodes=procs, edges=links_**

- **_Each node owns a piece of memory_**

- **Shared memory graph, $G_{SM} = (V, E)$**

- **Nodes _u_ and _v_ can access each other's memory iff $(u,v) \in E$**
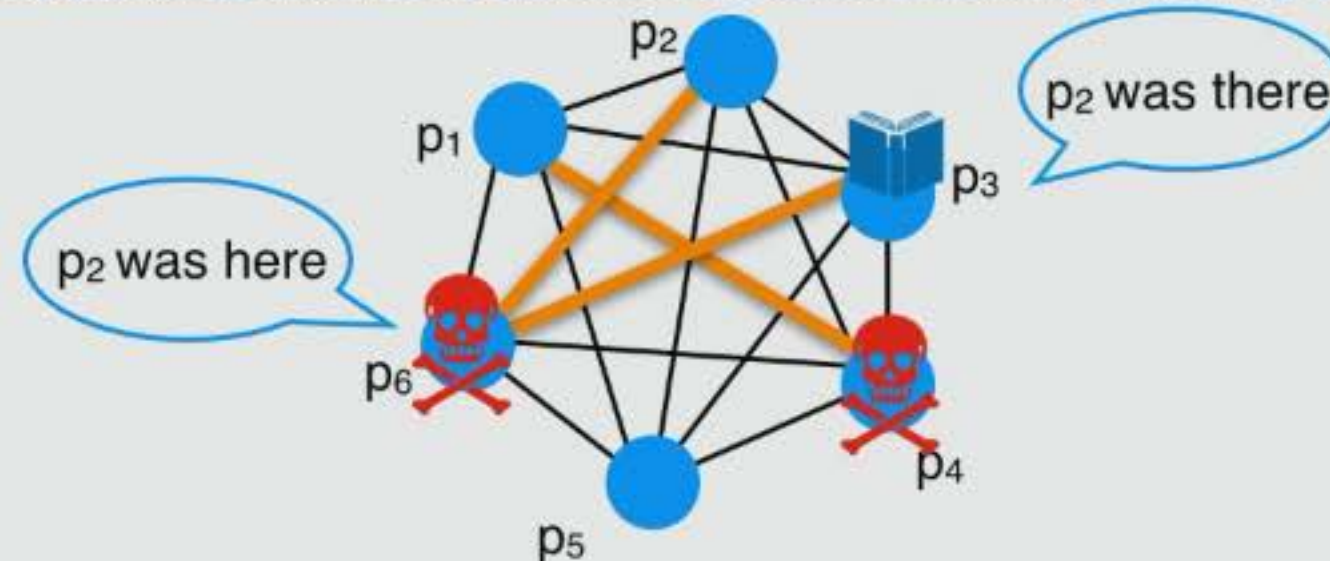
- **Processes may crash, but their _memory remains accessible_**

# RDMA vs Previous Results

| n = num processes<br>f = num failures | | Shared Memory | Message Passing | RDMA Full*<br>[ABGMZ'19] | RDMA Scale<br>[ABCGPT'18] |
|---|---|---|---|---|---|
| **Fault Tolerance** | **Crash** | n>f | n>2f | n>f | n>f+x<br>(x∈[0,f]) |
| | Byzantine | N/A | n>3f | n>2f | - |
| | Complexity*<br>(Best Case Round Trips) | 2 | 1 | 1 | - |
| | Scalability<br>(processes in network) | 10-100 | 10,000 - 100,000 | 10-100 | 10-100,000 |

# [ABGMZ'19] M&M Consensus

"Pretend" more processes are alive by sending their messages too

# M&M Consensus

"Pretend" more processes are alive by sending their messages too

Simulate a *message passing algorithm* by replacing messages with
*list of messages* representing your *shared memory neighbors*

# [ABGMZ'19] M&M Consensus

"Pretend" more processes are alive by sending their messages too

n>2f: Tolerates n/2 failures

Simulate a *message passing algorithm* by replacing messages with *list of messages* representing your *shared memory neighbors*

# [A**B**GMZ'19] M&M Consensus

"Pretend" more processes are alive by sending their messages too

**n>2f: Tolerates n/2 failures**

**n>f: SM consensus needs 1 process**

Simulate a *message passing algorithm* by replacing messages with *list of messages* representing your *shared memory neighbors*

# M&M Consensus
[ABGMZ'19]

"Pretend" more processes are alive by sending their messages too

n>2f: Tolerates n/2 failures

n>f: SM consensus needs 1 process

Simulate a *message passing algorithm* by replacing messages with *list of messages* representing your *shared memory neighbors*
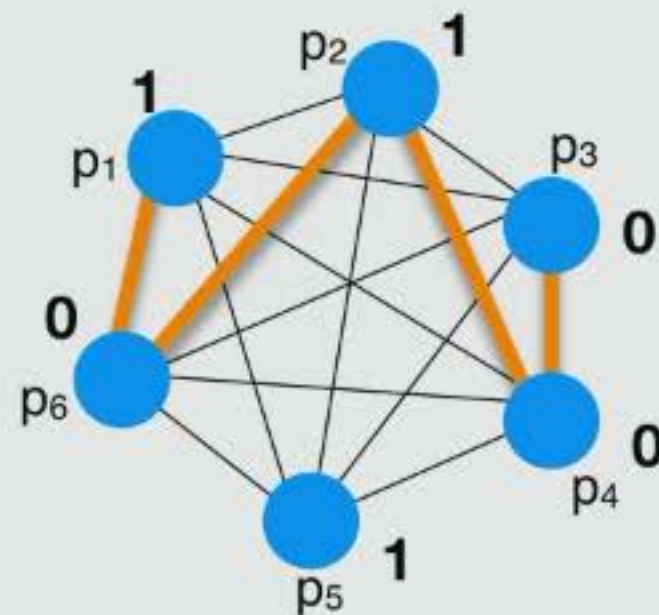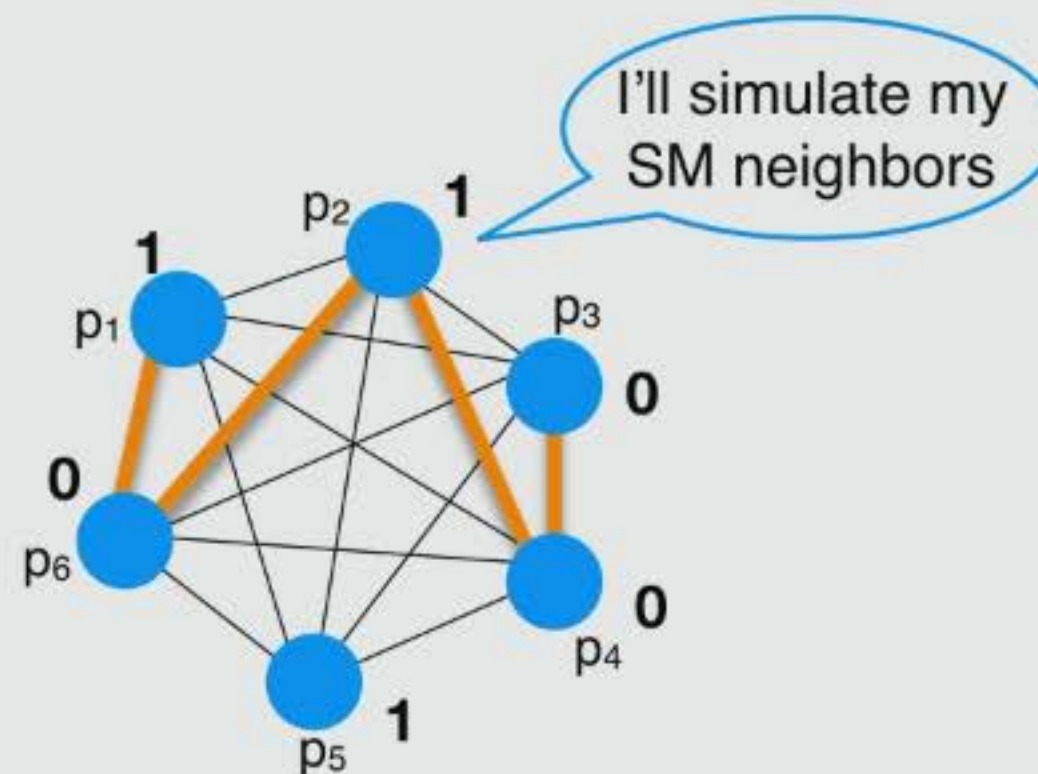
# [ABGMZ'19] M&M Consensus

"Pretend" more processes are alive by simulating neighbors

# M&M Consensus

"Pretend" more processes are alive by simulating neighbors

Original Algorithm

# M&M Consensus

[ABGMZ'19]

"Pretend" more processes are alive by simulating neighbors

Original Algorithm

**Message passing can only tolerate n>2f**

# M&M Consensus

"Pretend" more processes are alive by simulating neighbors

Original Algorithm

**Algorithm is doomed**

M&M Algorithm

**Message passing can only tolerate n>2f**

# [A**B**GMZ'19] M&M Consensus

"Pretend" more processes are alive by simulating neighbors

Original Algorithm

**Algorithm is doomed**

M&M Algorithm

**More than half -> Success!**

**Message passing can only tolerate n>2f**

# M&M Consensus

[ABGMZ'19]

"Pretend" more processes are alive by simulating neighbors

Original Algorithm

**Algorithm is doomed**

M&M Algorithm

**Message passing can only tolerate n>2f**

40

# M&M Consensus

"Pretend" more processes are alive by simulating neighbors

M&M Algorithm

Exactly half -> Failure!

We care about the number of neighbors of correct processes

Message passing can only tolerate n>2f

# Expander Graphs

Fault tolerance depends on shared memory graph:

*Number of neighbors of correct processes*

# Expander Graphs

Fault tolerance depends on shared memory graph:

**_Number of neighbors of correct processes_**

Expander graphs to the rescue!

# Expander Graphs



Fault tolerance depends on shared memory graph:

**_Number of neighbors of correct processes_**

Expander graphs to the rescue!

"G has high expansion"

↔

"Every subset of the vertices has many neighbors"

# Expander Graphs

Fault tolerance depends on shared memory graph:

**_Number of neighbors of correct processes_**

Expander graphs to the rescue!

"G has high expansion"

↔

"Every subset of the vertices has many neighbors"

$$h(G) = \min_{S \text{ s.t. } |S| \leq |V|/2} |\delta S|/|S|$$

# Expander Graphs

Fault tolerance depends on shared memory graph:

*Number of neighbors of correct processes*

Expander graphs to the rescue!

"G has high expansion"

↔

"Every subset of the vertices has many neighbors"

Neighbors of set S

$$h(G)=\min_{S \text{ s.t. } |S|\leq |V|/2} |\delta S|/|S|$$

Subset S of vertices

# Expander Graphs

Fault tolerance depends on shared memory graph:

***Number of neighbors of correct processes***

Expander graphs to the rescue!

"G has high expansion"

$\leftrightarrow$

"Every subset of the vertices has many neighbors"

Neighbors of set S

$$h(G) = \min_{S \text{ s.t. } |S| \leq |V|/2} |\delta S| / |S|$$

Subset S of vertices

S with worst ratio defines graph's expansion

# Putting it Together

- Think of **set of live processes** as S

- Adversary will pick S to be the set with the least expansion

*Graph with high expansion can tolerate more failures*

# [ABGMZ'19] Putting it Together

- Think of *set of live processes* as S

- Adversary will pick S to be the set with the least expansion

*Graph with high expansion can tolerate more failures*

**Theorem** [AguileraBCalciuGuerraouiPetrankToueg'18]:
If shared memory graph has vertex expansion ratio $h$,

then we can tolerate $f < \left(1 - \dfrac{1}{2 \cdot h}\right) \cdot n$ failures

# Putting it Together

- Think of **set of live processes** as S

- Adversary will pick S to be the set with the least expansion

*Graph with high expansion can tolerate more failures*

> **Theorem** [AguileraBCalciuGuerraouiPetrankToueg'18]:
> If shared memory graph has vertex expansion ratio *h*,
>
> then we can tolerate $f < \left(1 - \dfrac{1}{2 \cdot h}\right) \cdot n$ failures

> **Also show impossibility result:**
> relation to expansion is *inherent*.

# Outline

✓ **RDMA details**

- Setting 1: **RDMA's full power** (complete graph)

  ✓ **Crash-only** algorithm: $n>f$ tolerant, 1 round-trip

  ✓ **Byzantine** algorithm: $n>2f$ tolerant, 1 round-trip

- Setting 2: **Scalability: Using RDMA sparingly** (incomplete graph)

  ✓ Crash-only Algorithm: tolerance vs topology

# [ABGMZ'19] Putting it Together

- Think of **set of live processes** as S

- Adversary will pick S to be the set with the least expansion

*Graph with high expansion can tolerate more failures*

# [ABGMZ'19] Putting it Together

- Think of **set of live processes** as S

- Adversary will pick S to be the set with the least expansion

**Graph with high expansion can tolerate more failures**

**Theorem** [AguileraBCalciuGuerraouiPetrankToueg'18]:
If shared memory graph has vertex expansion ratio $h$,

then we can tolerate $f < \left(1 - \frac{1}{2 \cdot h}\right) \cdot n$ failures

**Also show impossibility result:**
relation to expansion is *inherent*.

# Outline

✓ **RDMA details**

- Setting 1: **RDMA's full power** (complete graph)

  ✓ **Crash-only** algorithm: $n>f$ tolerant, 1 round-trip

  ✓ **Byzantine** algorithm: $n>2f$ tolerant, 1 round-trip

- Setting 2: **Scalability: Using RDMA sparingly** (incomplete graph)

  ✓ Crash-only Algorithm: tolerance vs topology

# RDMA vs Previous Results

| n = num processes<br>f = num failures | | Shared Memory | Message Passing | RDMA Full*<br>[ABGMZ'19] | RDMA Scale<br>[ABCGPT'18] |
|---|---|---|---|---|---|
| **Fault Tolerance** | Crash | n>f | n>2f | n>f | n>f+x<br>(x∈[0,f]) |
| | Byzantine | N/A | n>3f | n>2f | - |
| **Complexity***<br>(Best Case Round Trips) | | 2 | 1 | 1 | - |
| **Scalability**<br>(processes in network) | | 10-100 | 10,000 - 100,000 | 10-100 | 10-100,000 |

*With up to half of the memories crashing

**Can we scale better and still retain some of RDMA's advantages?**

**Can we scale better and still retain some of RDMA's advantages?**

**Yes!**
**We can tolerate more crash failures**
**than message passing**

**Can we scale better and still retain some of RDMA's *advantages*?**

Yes!
*We can tolerate more crash failures than message passing*

**Crash Tolerance**

**Can we scale better and still retain some of RDMA's *advantages*?**

**Yes!**
**We can tolerate more crash failures than message passing**

47

**Crash Tolerance**

**Byzantine Tolerance**

*Can we scale better and still retain some of RDMA's advantages?*

Yes!
*We can tolerate more crash failures than message passing*

Performance

Crash Tolerance

Byzantine Tolerance

## Can we scale better and still retain some of RDMA's *advantages*?

## Yes!
### We can tolerate more crash failures than message passing

✔ Crash Tolerance

? Performance  ? Byzantine Tolerance

**Can we scale better and still retain some of RDMA's *advantages*?**

Yes!
*We can tolerate more crash failures than message passing*

# Summary



- **Consensus** as a lens to study **RDMA**

  - RDMA improves tradeoff between **fault tolerance** and **performance**

  - RDMA could **scale** to large networks

**Future Directions**

- Strengthening scalability model

- Implementing these solutions

- Problems beyond consensus

# Summary

- **Consensus** as a lens to study **RDMA**

  - RDMA improves tradeoff between **fault tolerance** and **performance**

  - RDMA could **scale** to large networks

**Future Directions**

- Strengthening scalability model

- Implementing these solutions

- Problems beyond consensus

**Thank you!**