

Applications for Broadcast

[**Reliable**, **Uniform**, **Causal**, and **Total-Order**]

Adi Seredinschi

Distributed Programming Laboratory

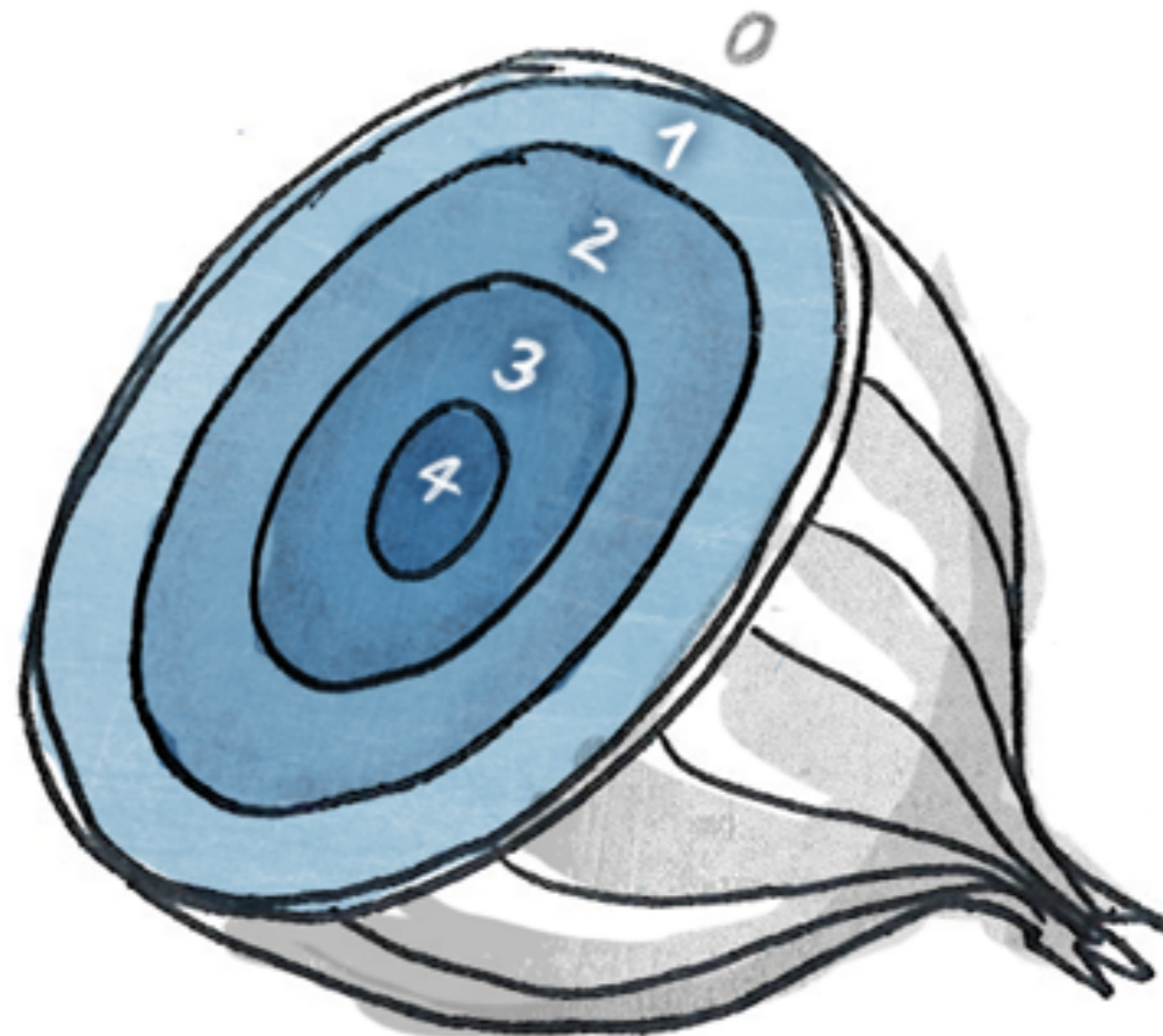
MIDTERM

23/11/2015

- ★ Lectures
 - including this one
- ★ Exercises
 - class (+ optional book)

Question:

What do Distributed Systems
have in common with Onions?

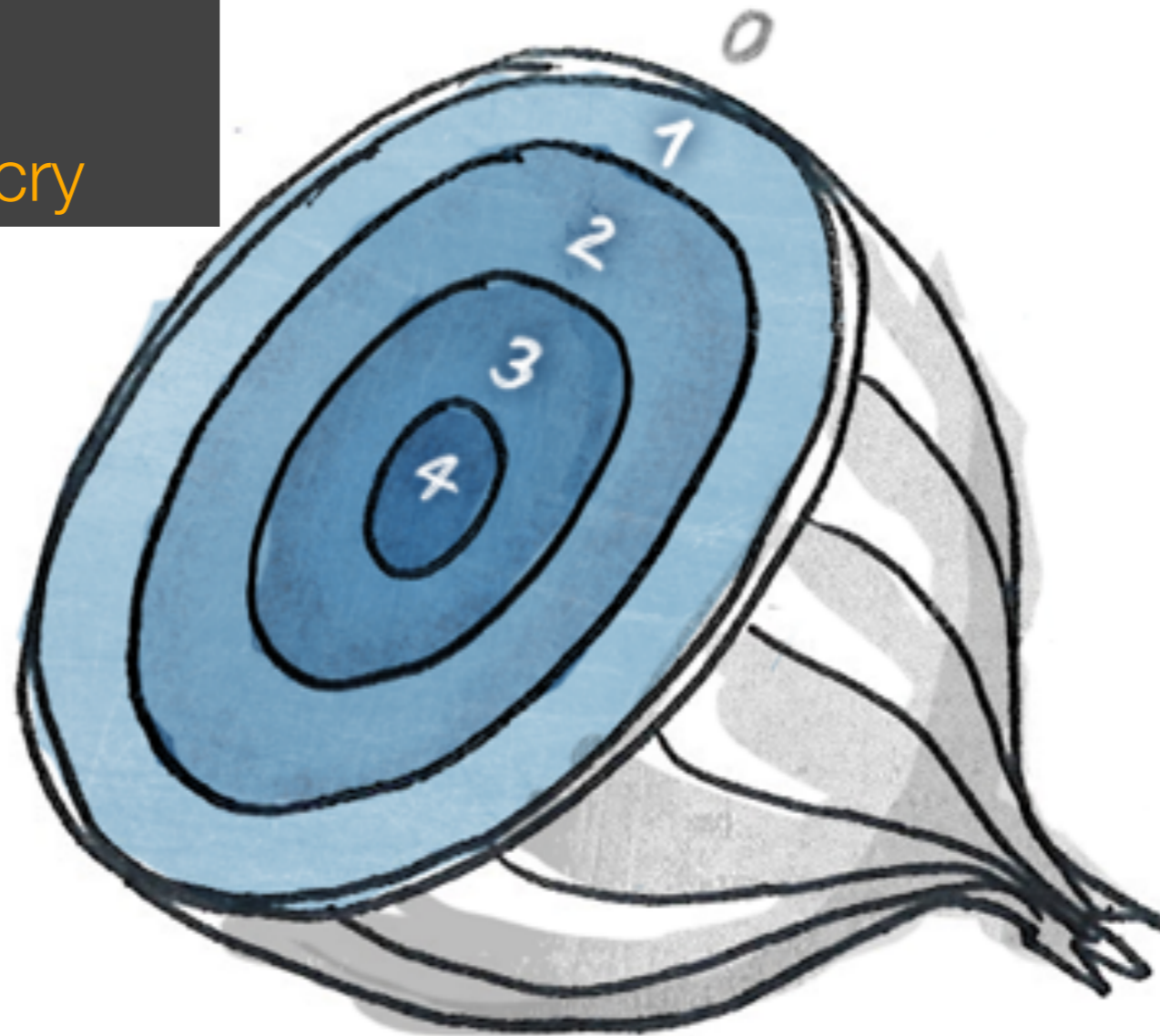


Question:

What do Distributed Systems have in common with Onions?

Answer:

1. Layering
2. Abstraction
3. They make you cry

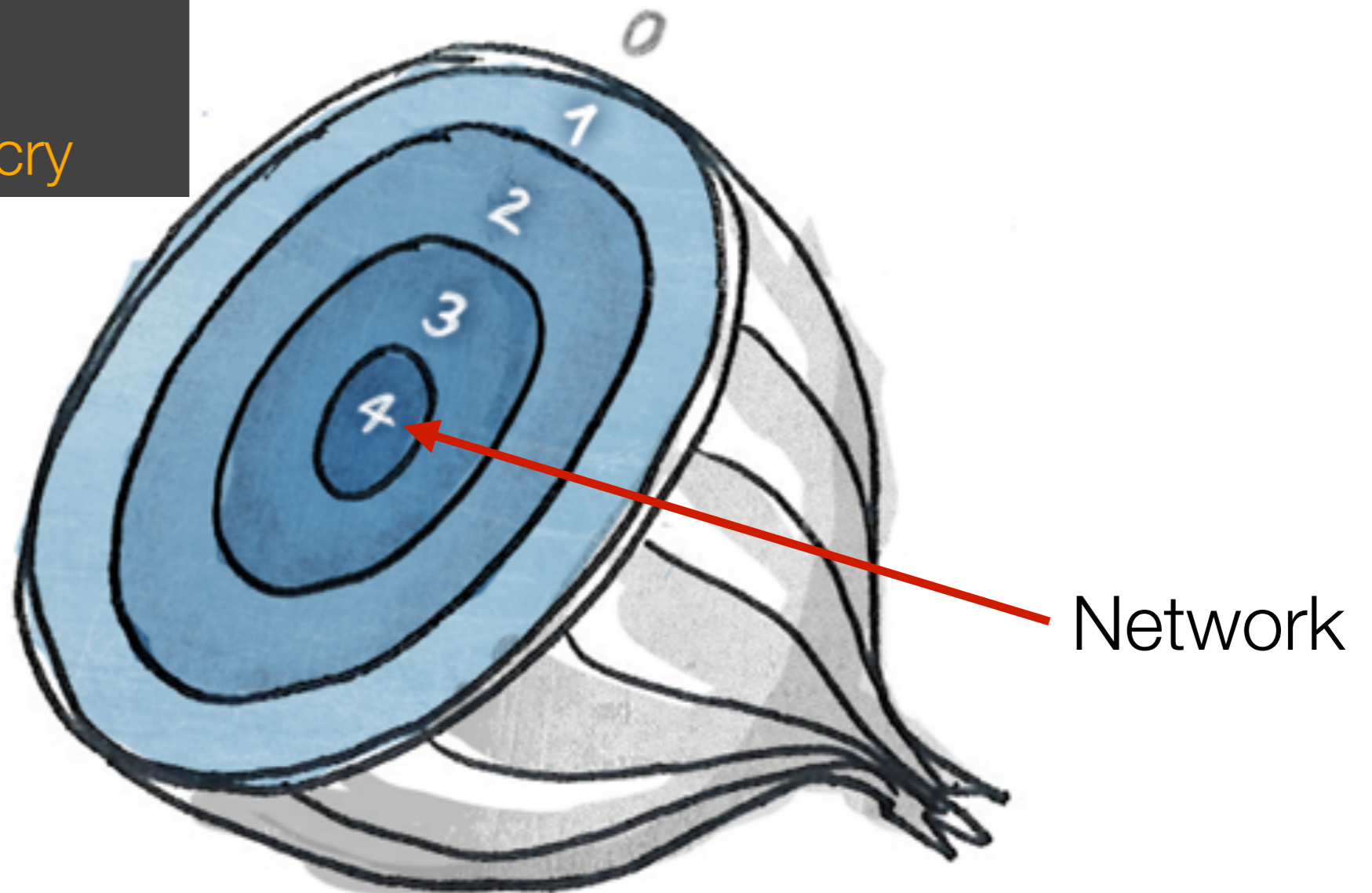


Question:

What do Distributed Systems have in common with Onions?

Answer:

1. Layering
2. Abstraction
3. They make you cry

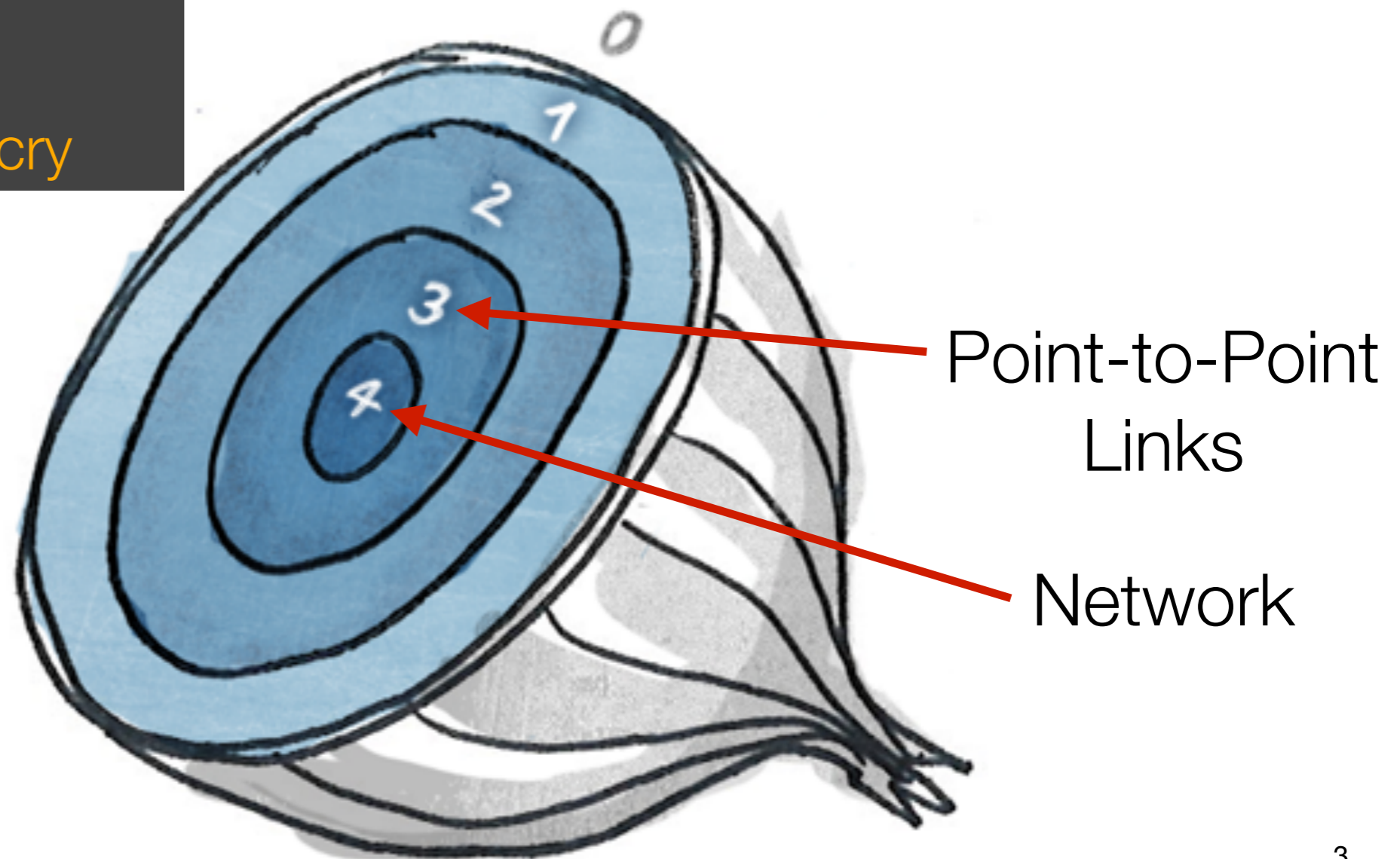


Question:

What do Distributed Systems have in common with Onions?

Answer:

1. Layering
2. Abstraction
3. They make you cry

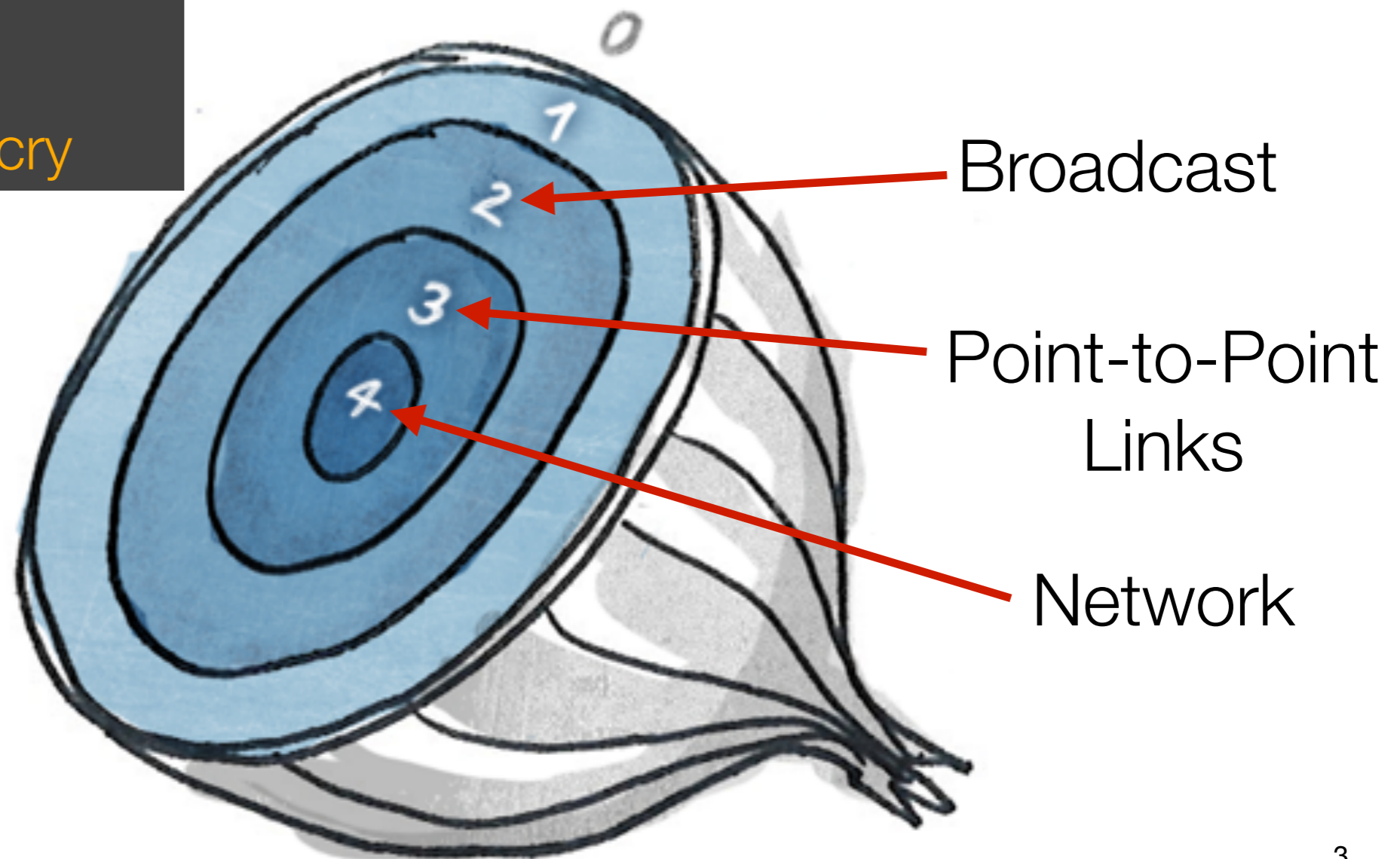


Question:

What do Distributed Systems have in common with Onions?

Answer:

1. Layering
2. Abstraction
3. They make you cry

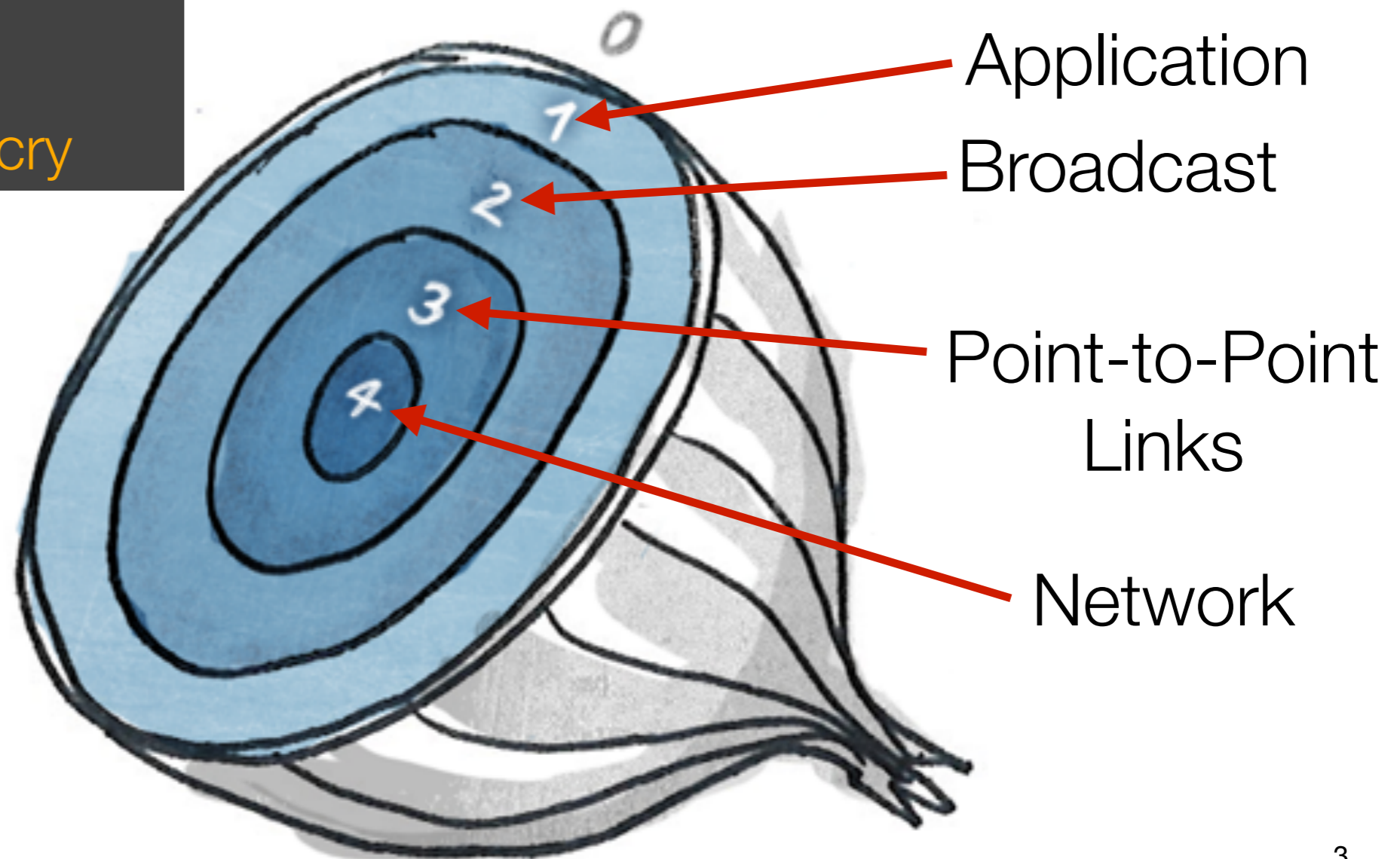


Question:

What do Distributed Systems have in common with Onions?

Answer:

1. Layering
2. Abstraction
3. They make you cry

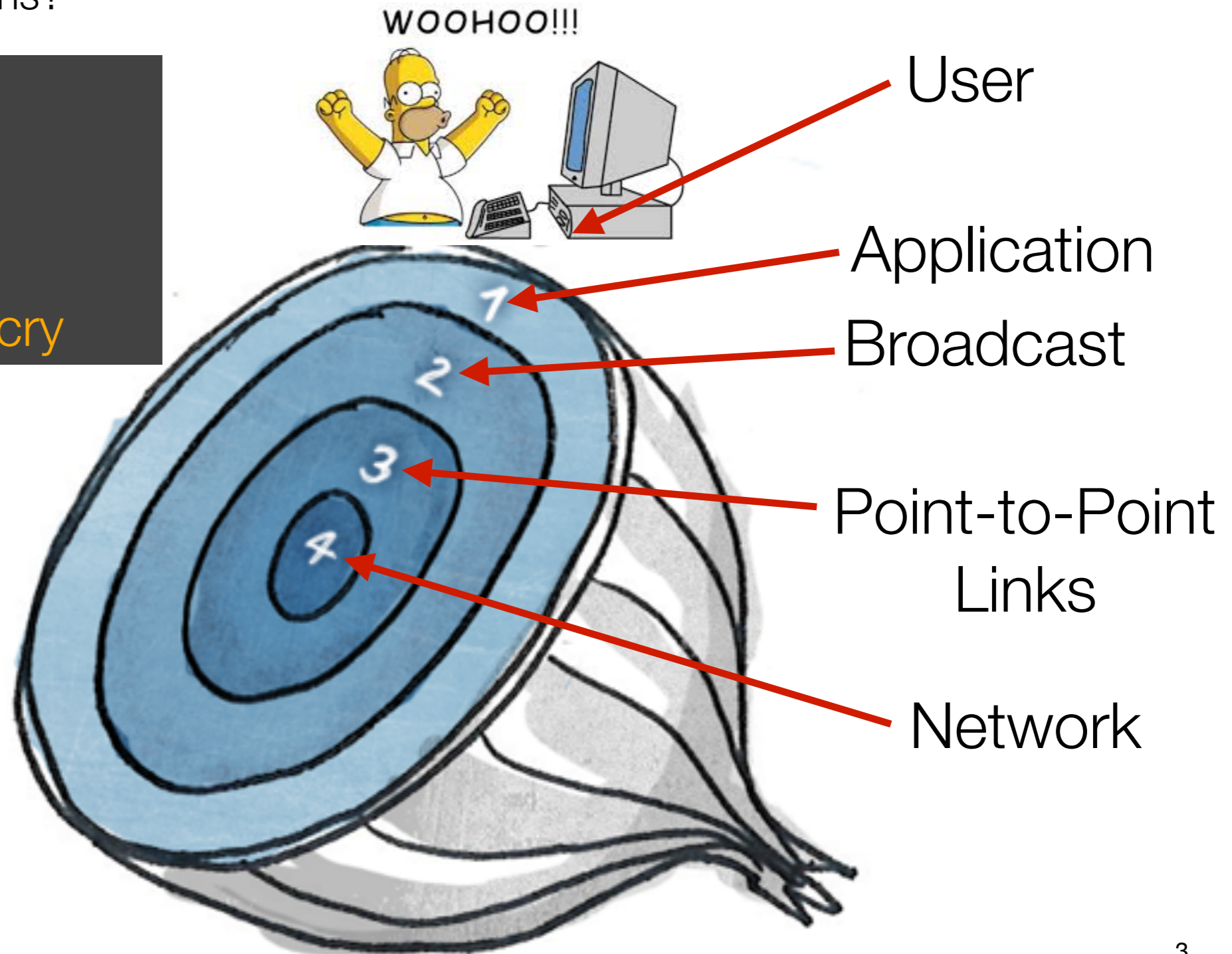


Question:

What do Distributed Systems have in common with Onions?

Answer:

1. Layering
2. Abstraction
3. They make you cry

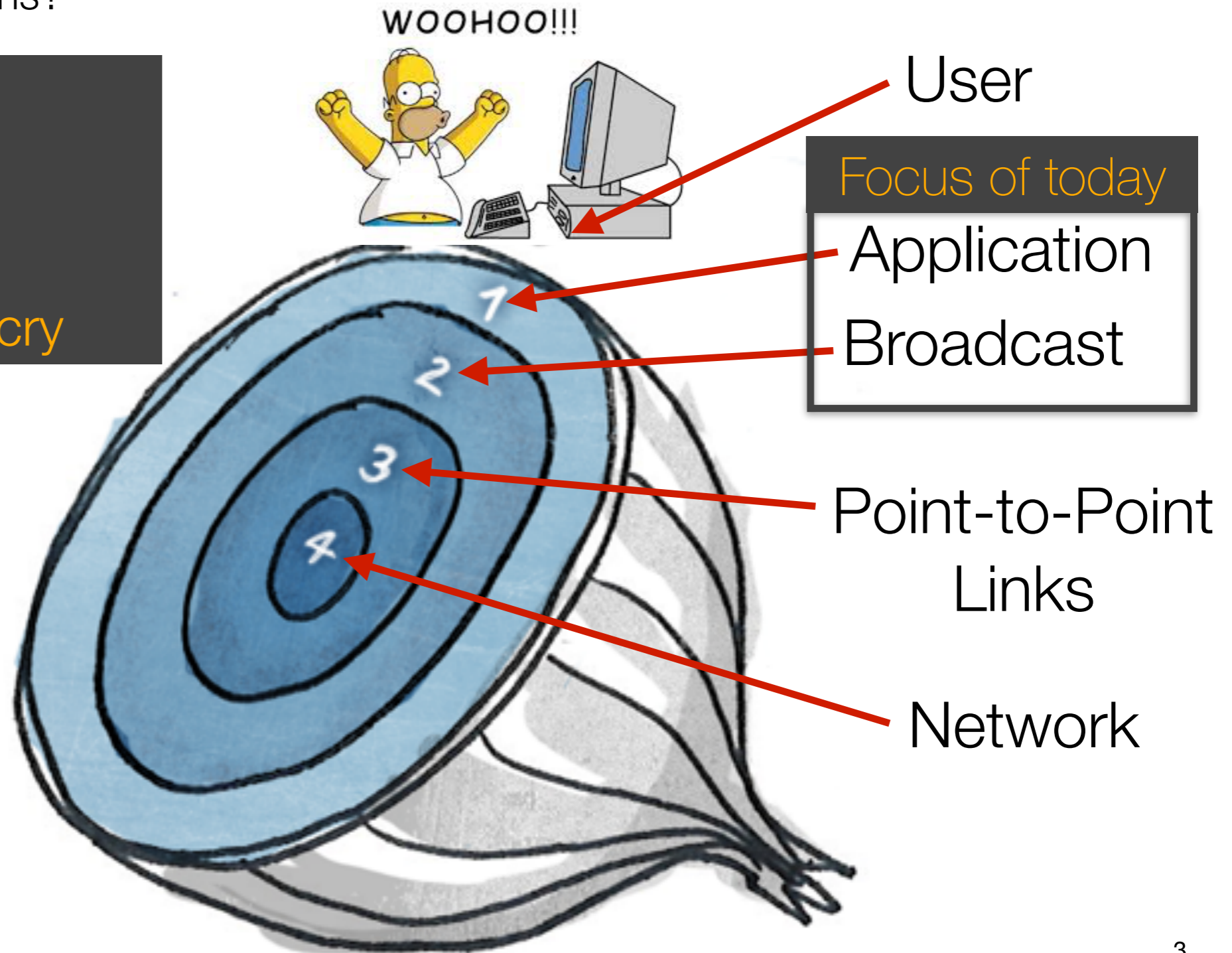


Question:

What do Distributed Systems have in common with Onions?

Answer:

1. Layering
2. Abstraction
3. They make you cry



Let's design some applications



Broadcast (we will investigate different properties..)

1. **CAMIPRO-Bitcoin**

- The canonical Bitcoin design
 - Uses gossip (best-effort broadcast)
 - Relies heavily on crypto — no time to cover that
- We will not discuss the canonical design
 - Instead, we will design our own version of Bitcoin
 - Optimized for CAMIPRO



Let's design some applications



Broadcast (we will investigate different properties..)

1. **CAMIPRO-Bitcoin**

- The canonical Bitcoin design
 - Uses gossip (best-effort broadcast)
 - Relies heavily on crypto — no time to cover that
- We will not discuss the canonical design
 - Instead, we will design our own version of Bitcoin
 - Optimized for CAMIPRO



Causal Broadcast

1. **S4: Storage for Social Networks**

- A simplified version of Twitter



[BONUS PROJECT]





Conceptual goals:

1. Replace traditional CAMIPRO
 - Based on CHF
2. Make banks obsolete



Conceptual goals:

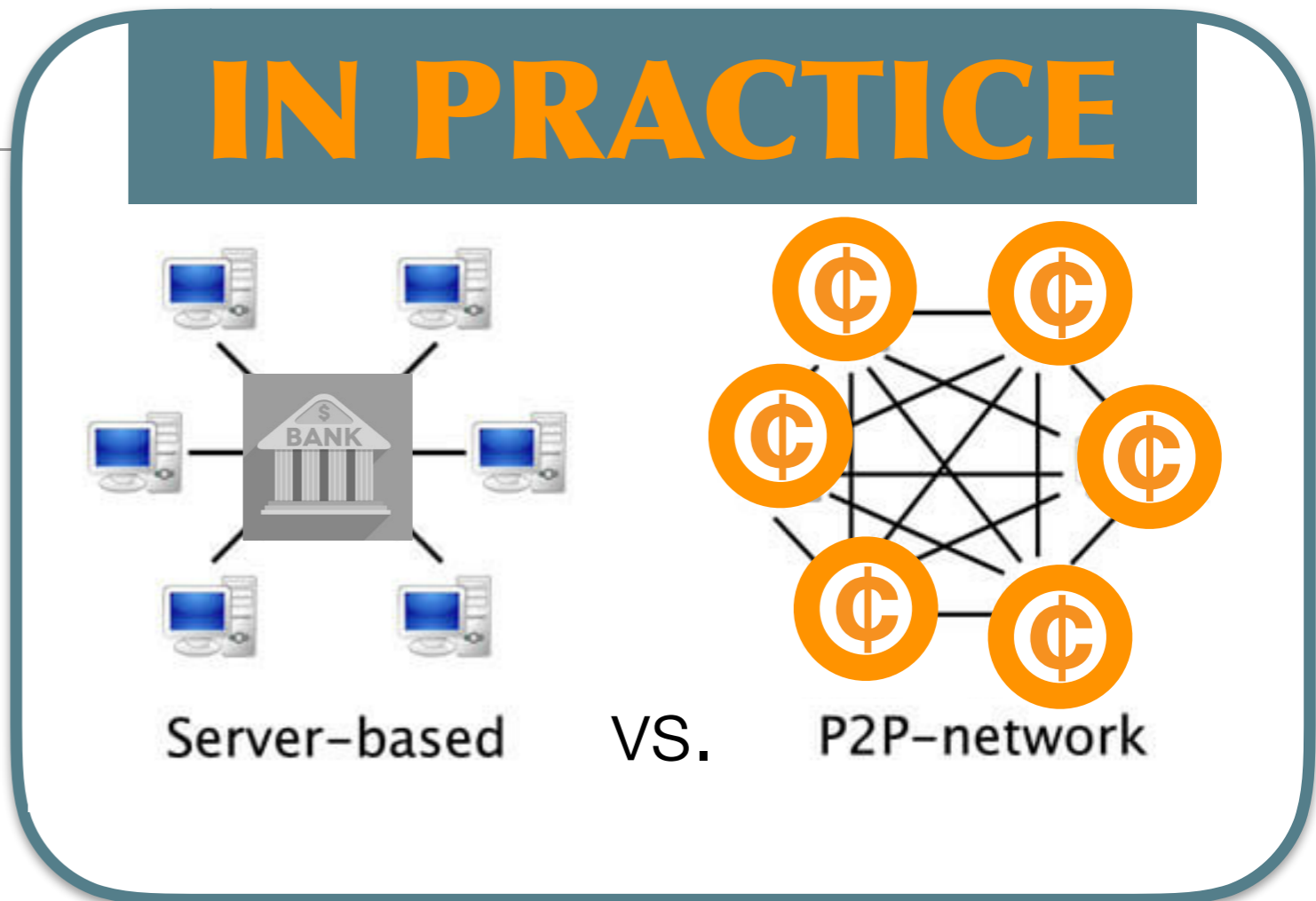
1. Replace traditional CAMIPRO
 - Based on CHF
2. Make banks obsolete





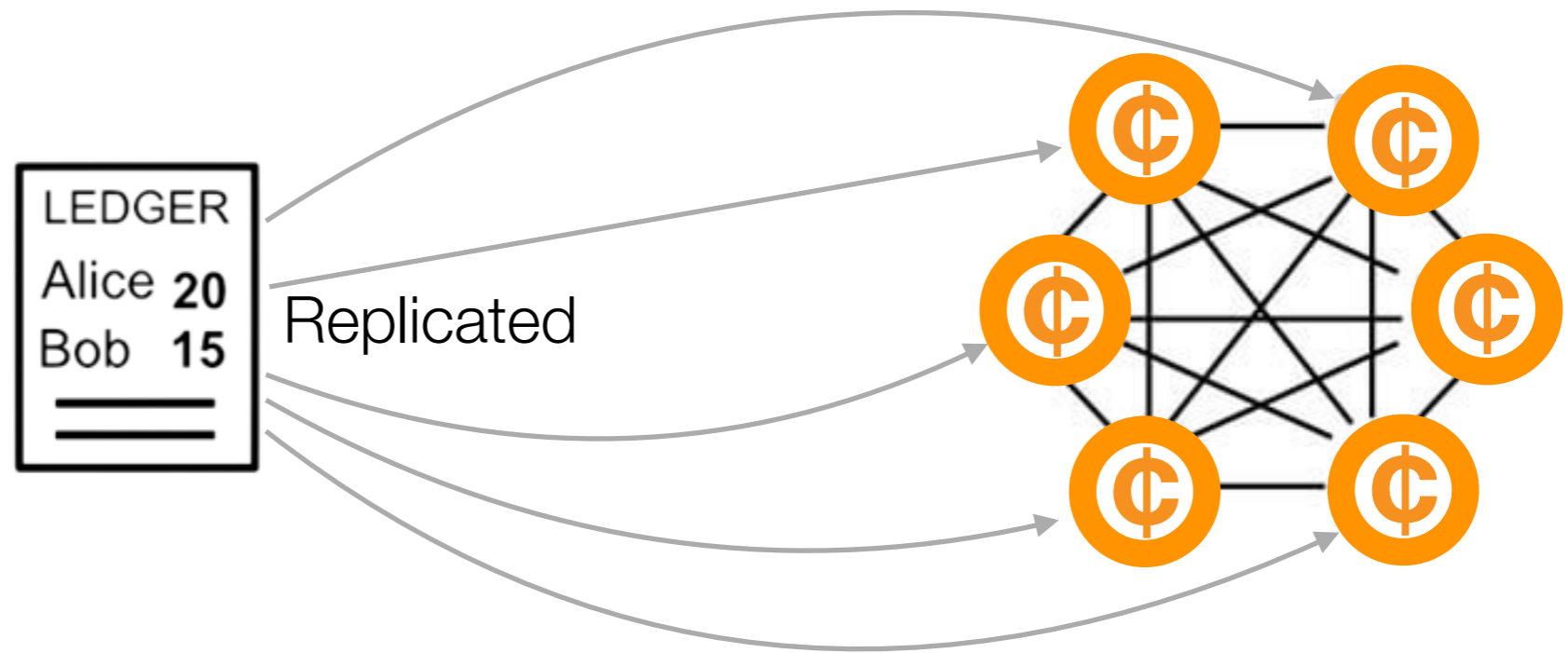
Conceptual goals:

1. Replace traditional CAMIPRO
 - Based on CHF
2. Make banks obsolete



Main concepts

- Ledger

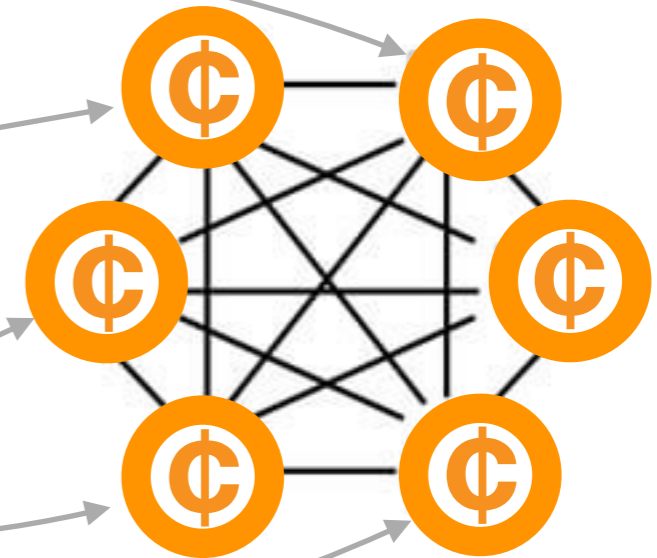


Main concepts

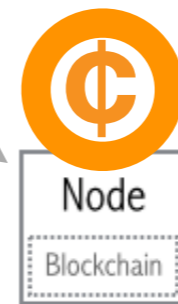
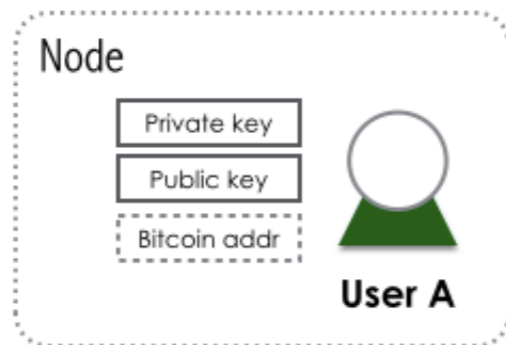
- Ledger

LEDGER	
Alice	20
Bob	15
<hr/> <hr/>	

Replicated



- Node

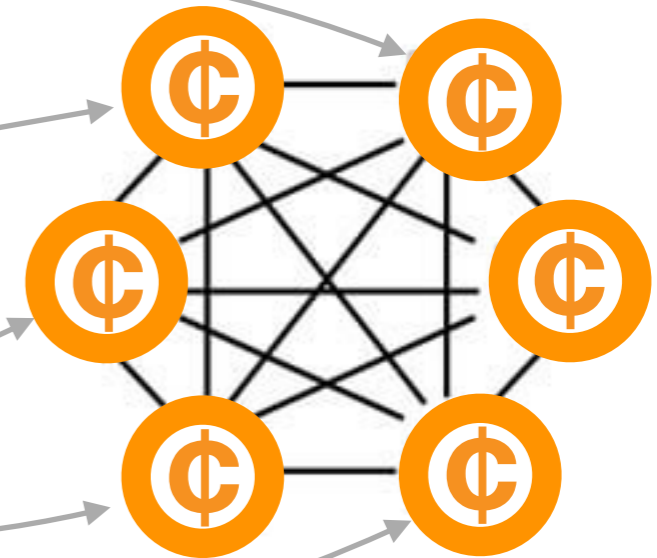


Main concepts

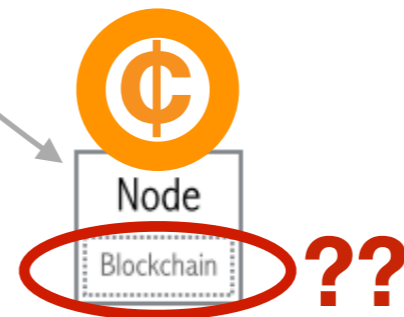
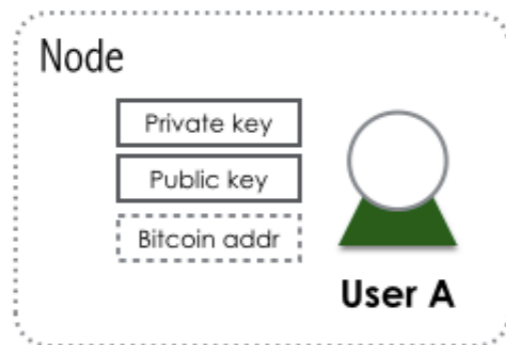
- Ledger

LEDGER	
Alice	20
Bob	15
=====	

Replicated

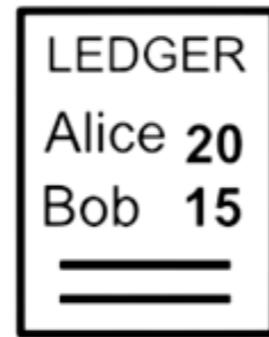


- Node

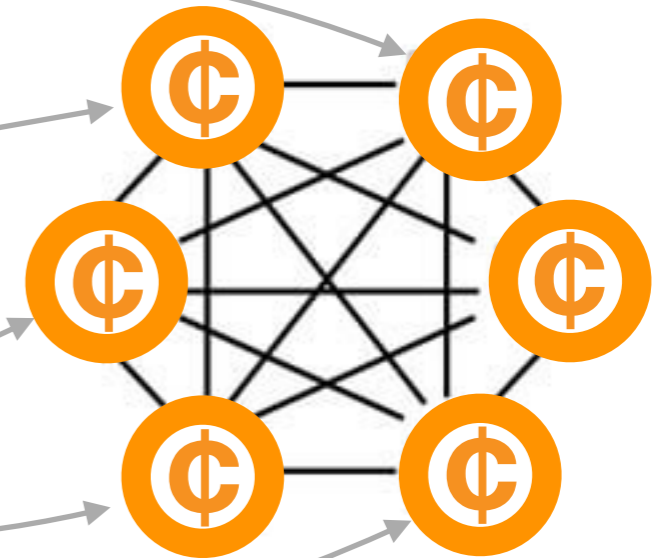


Main concepts

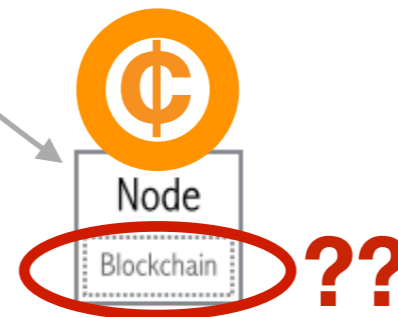
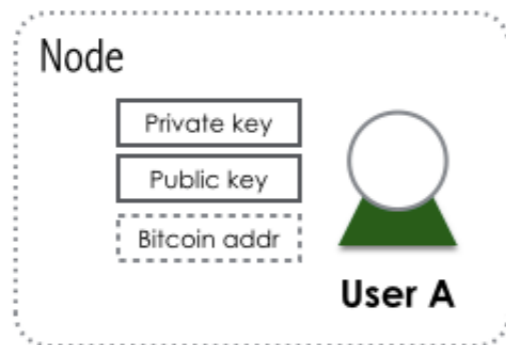
- Ledger



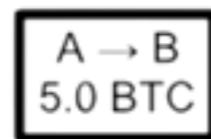
Replicated



- Node



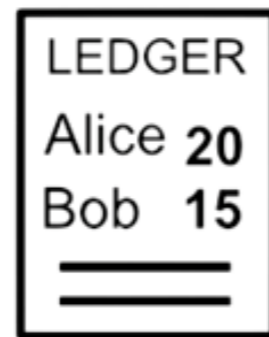
- Transaction



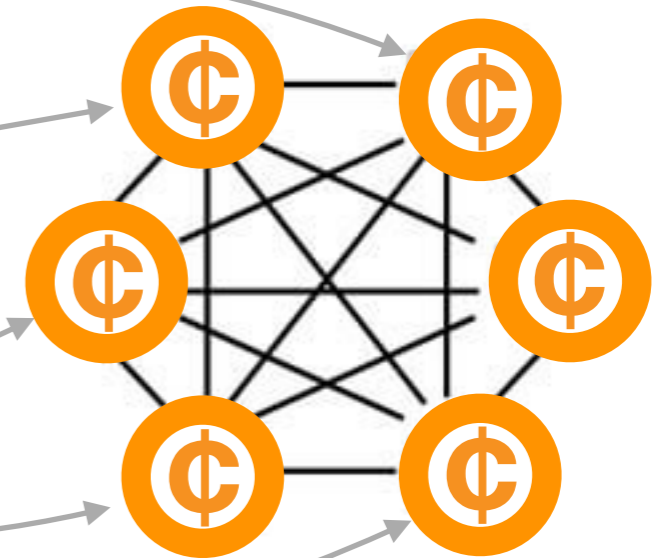
“Alice gives some money to Bob”

Main concepts

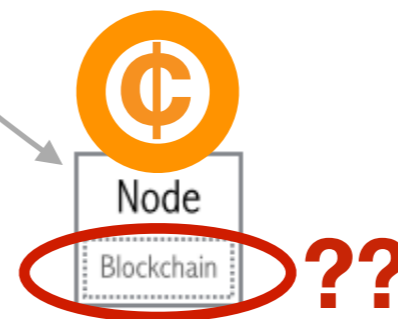
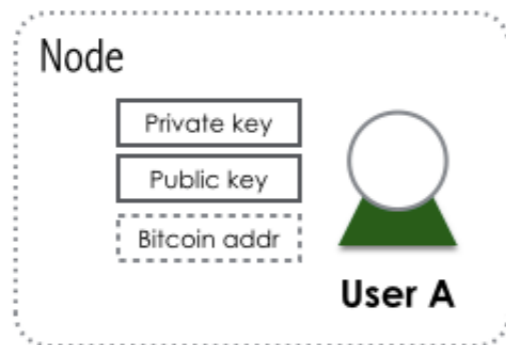
- Ledger



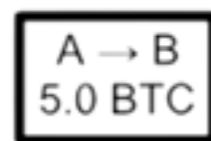
Replicated



- Node



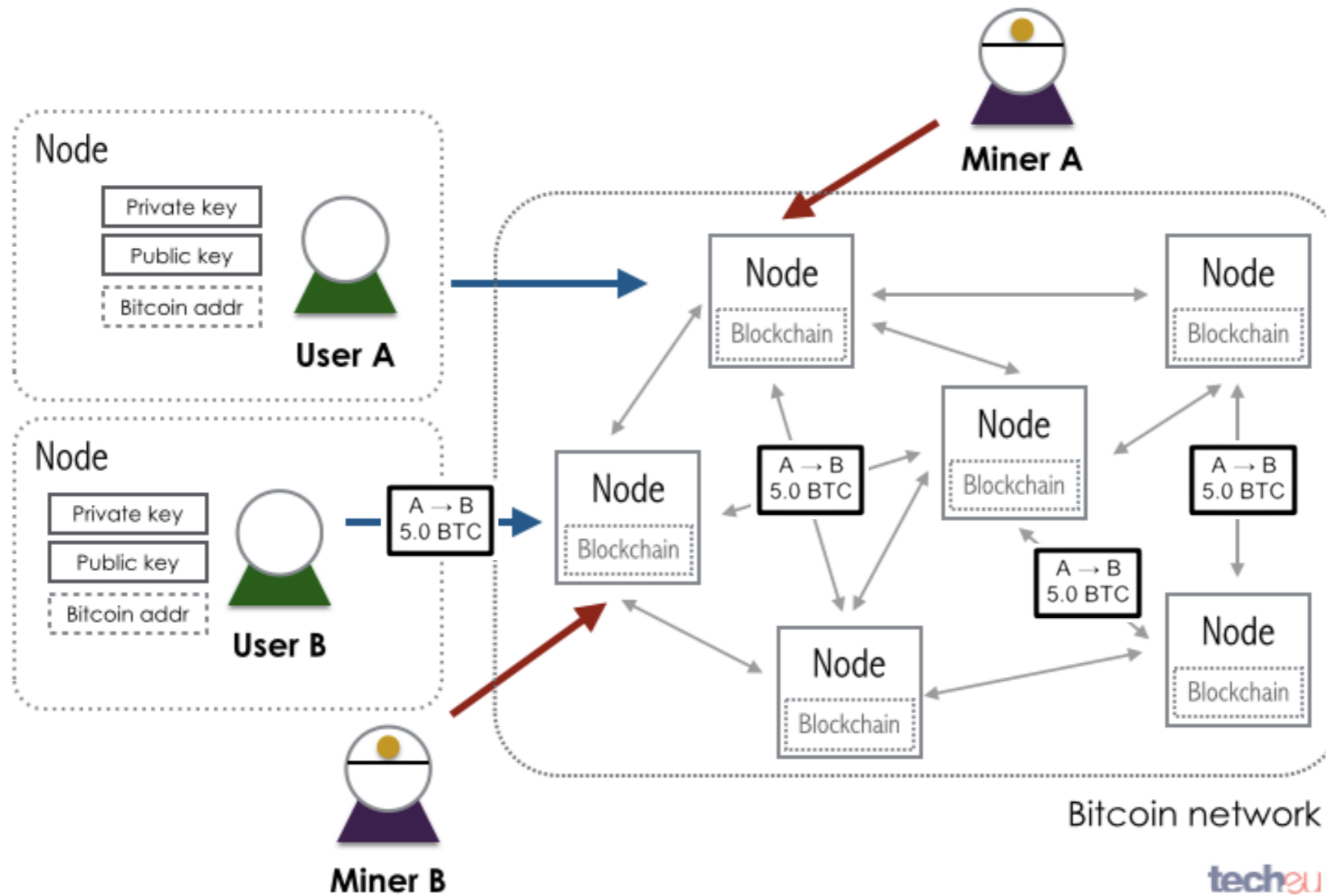
- Transaction



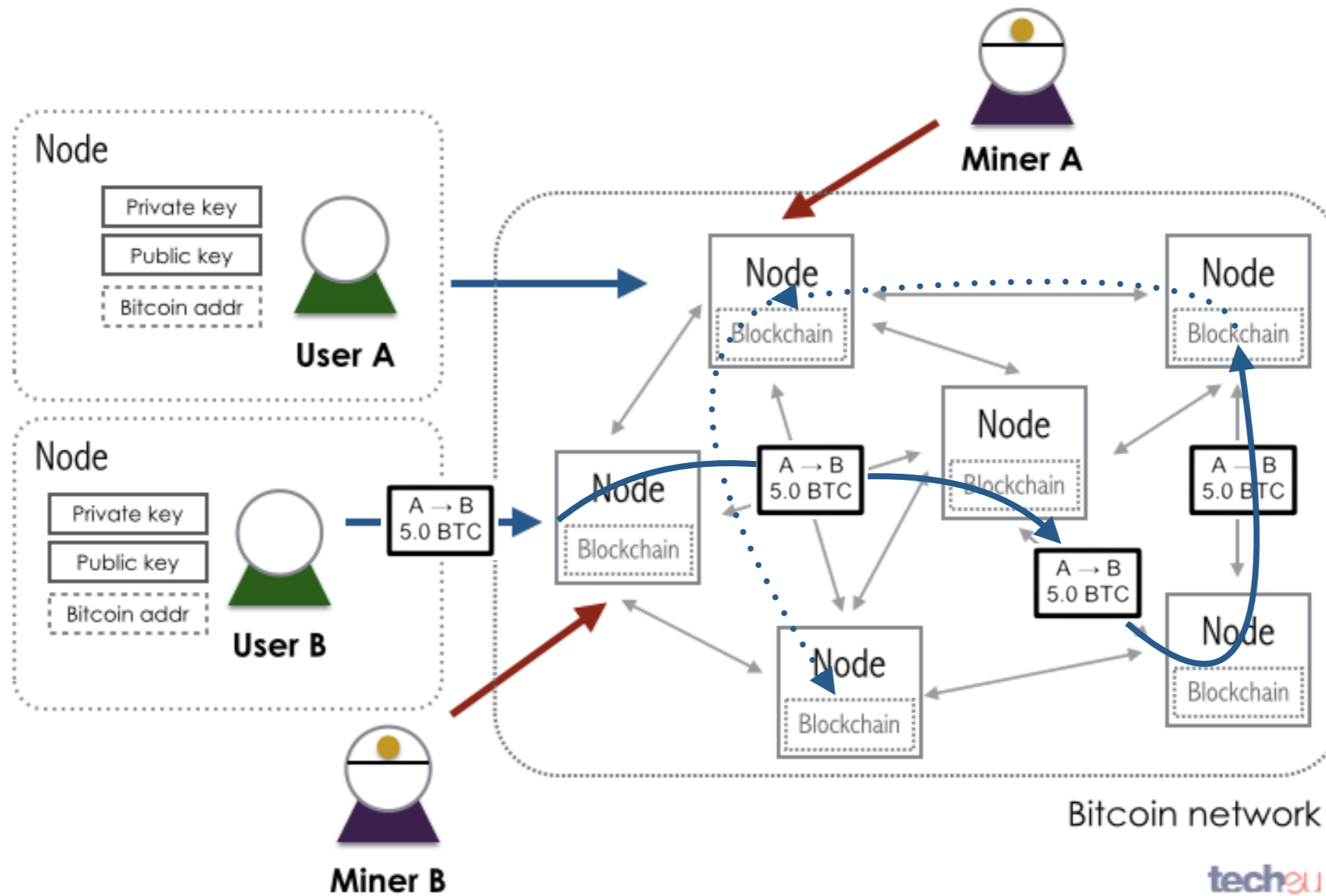
“Alice gives some money to Bob”

Let's see how they all fit together..

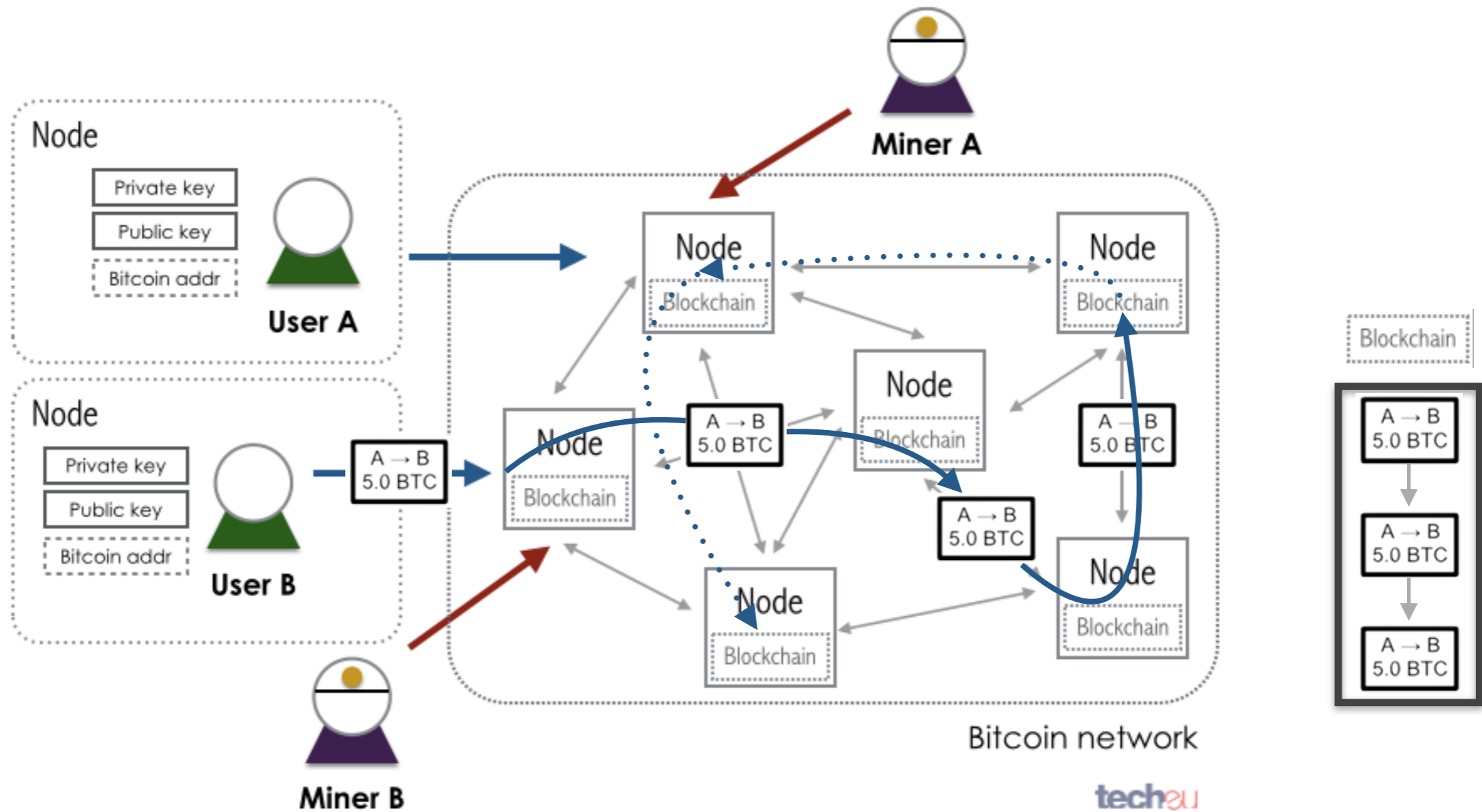
High-level interactions



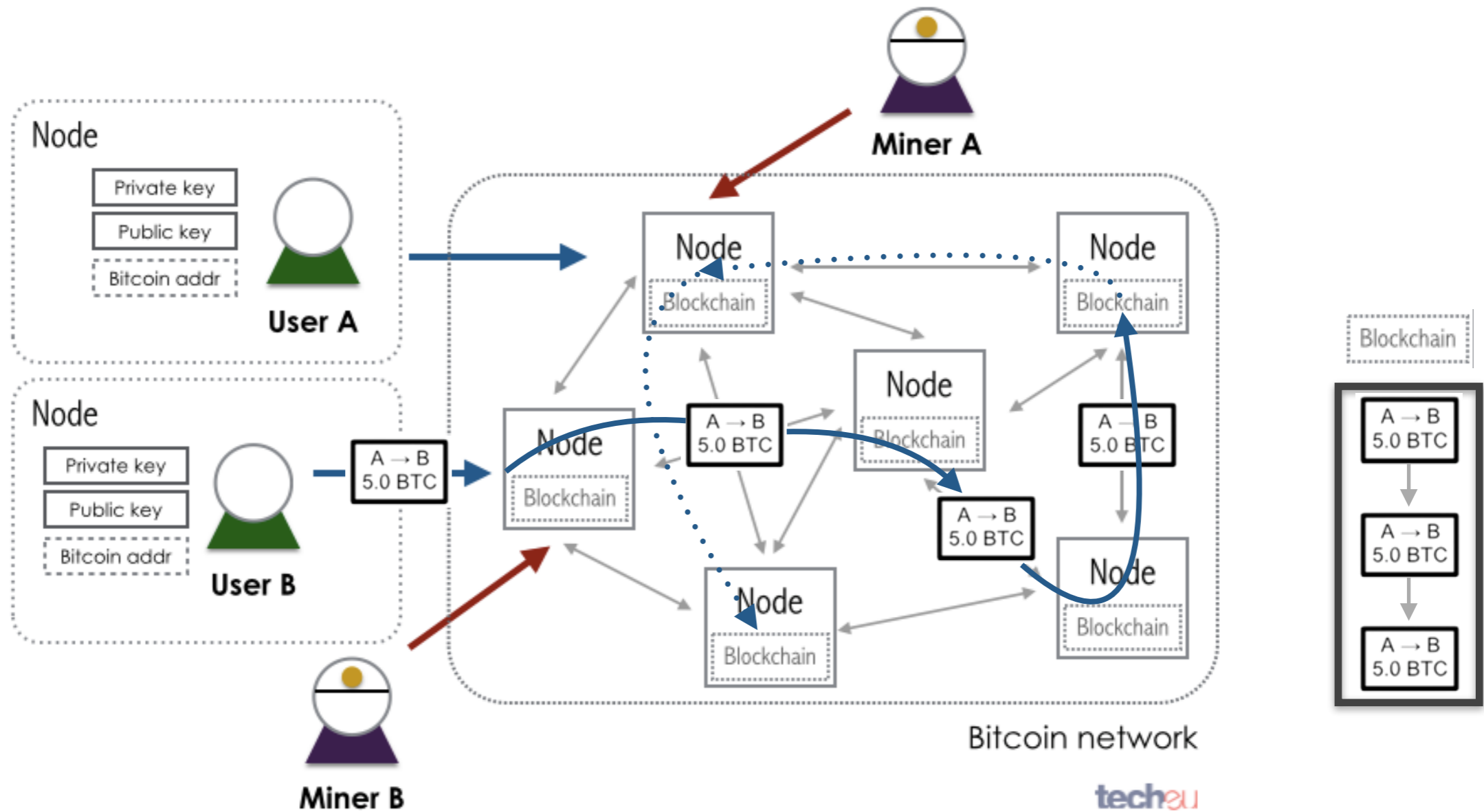
High-level interactions



High-level interactions

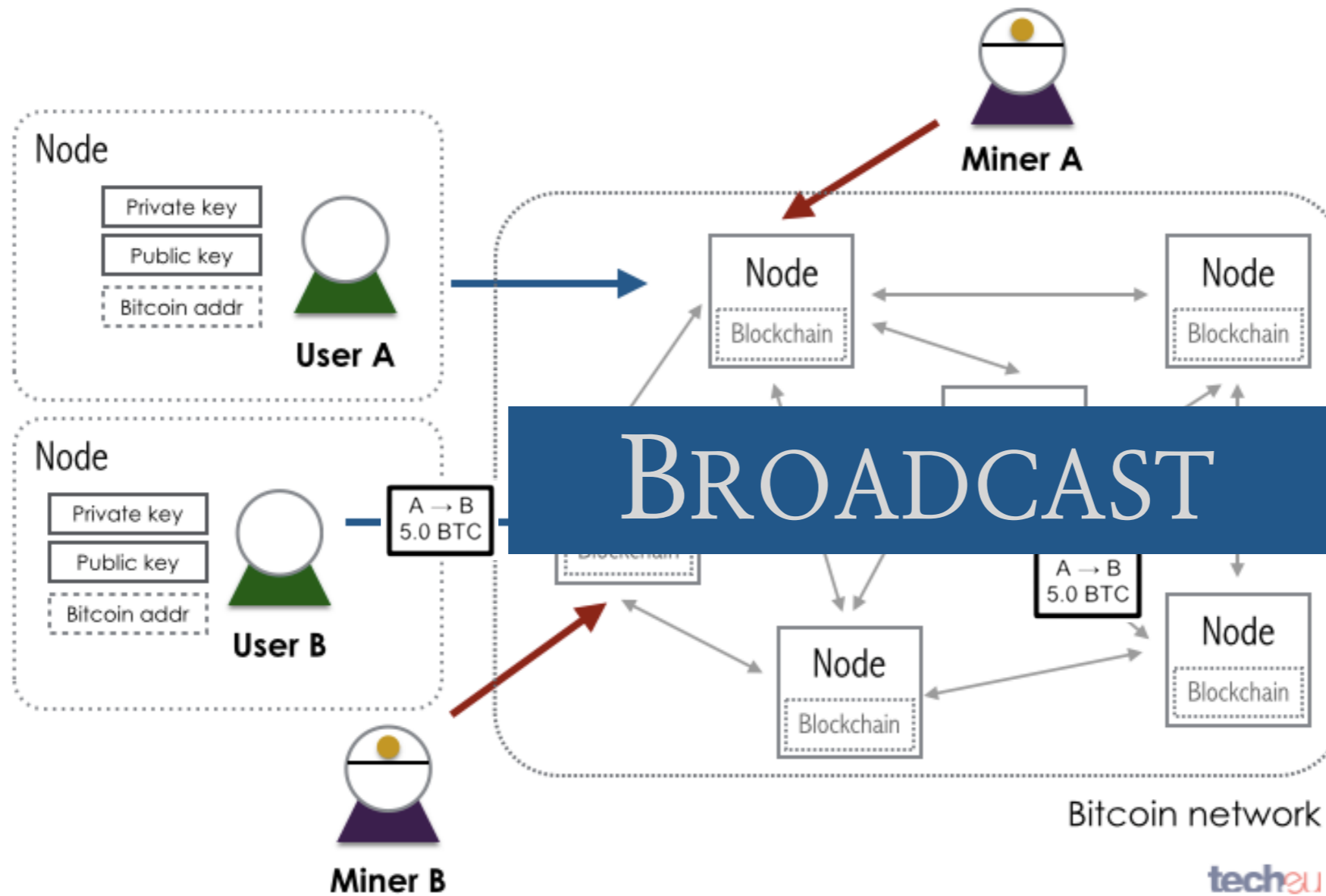


High-level interactions

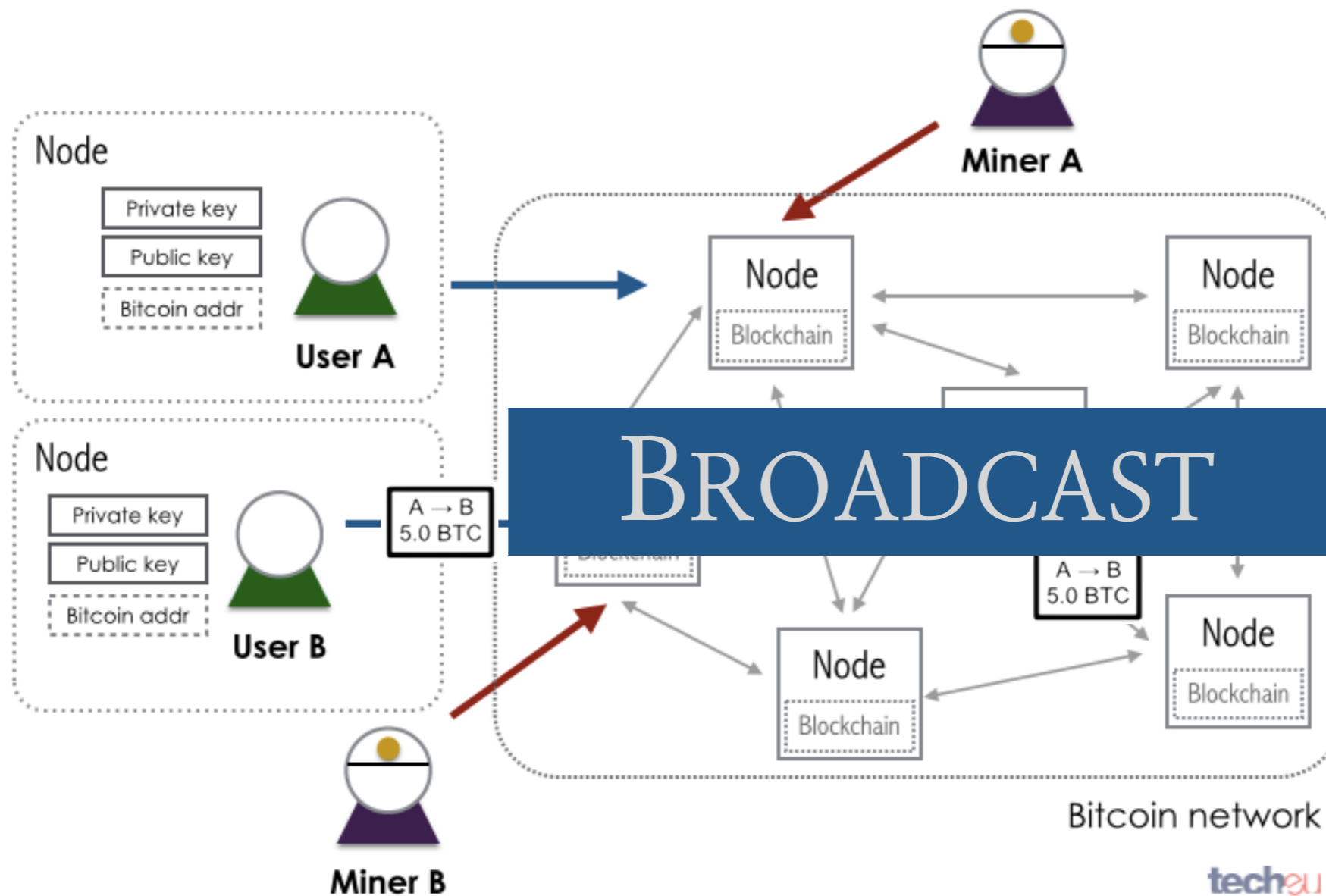


Recognise any abstractions?

High-level interactions



What kind of broadcast?



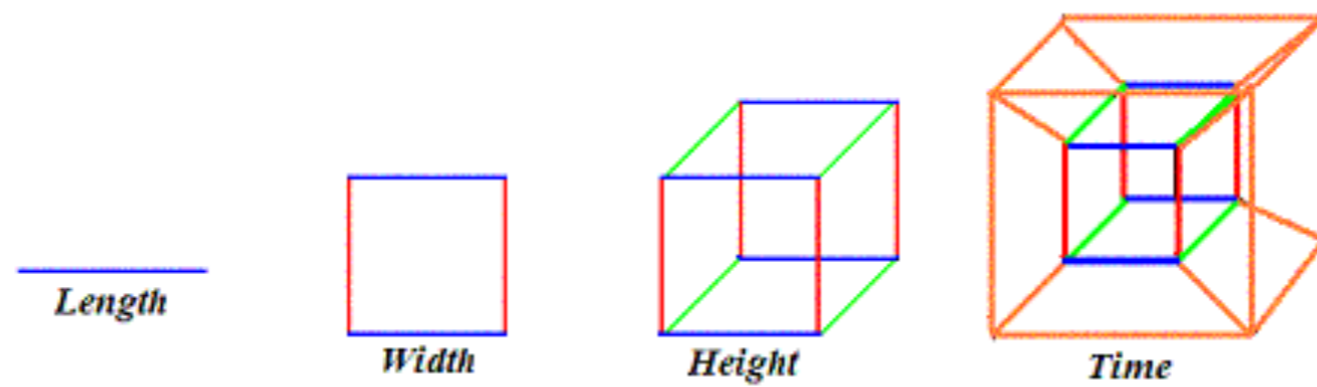
What kind of broadcast?

PROPERTIES
(guarantees)

PERFORMANCE

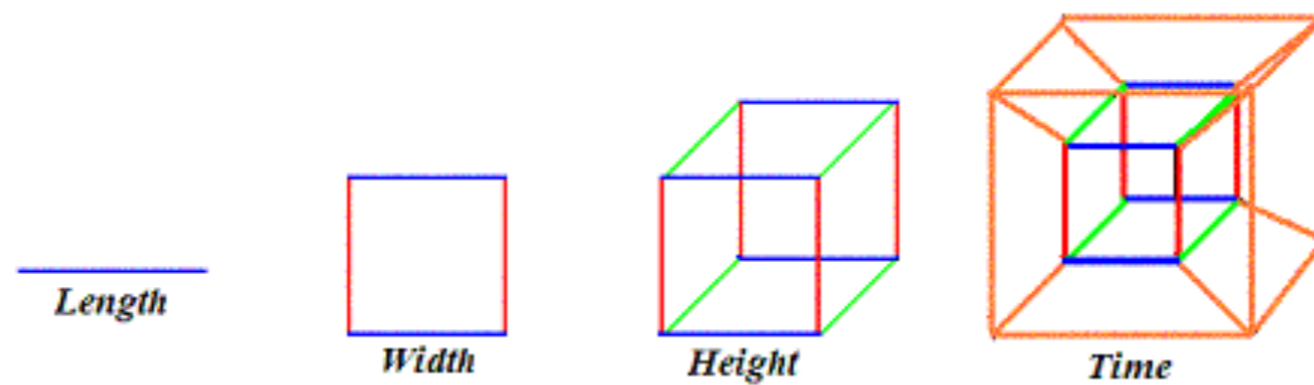
PROPERTIES

The Four Dimensions of a cube



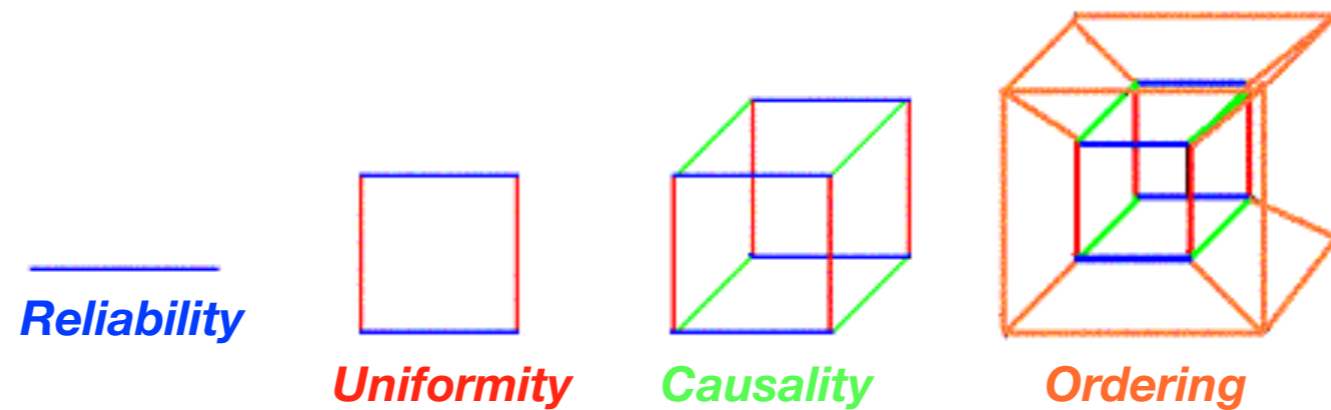
PROPERTIES

The Four Dimensions of a cube



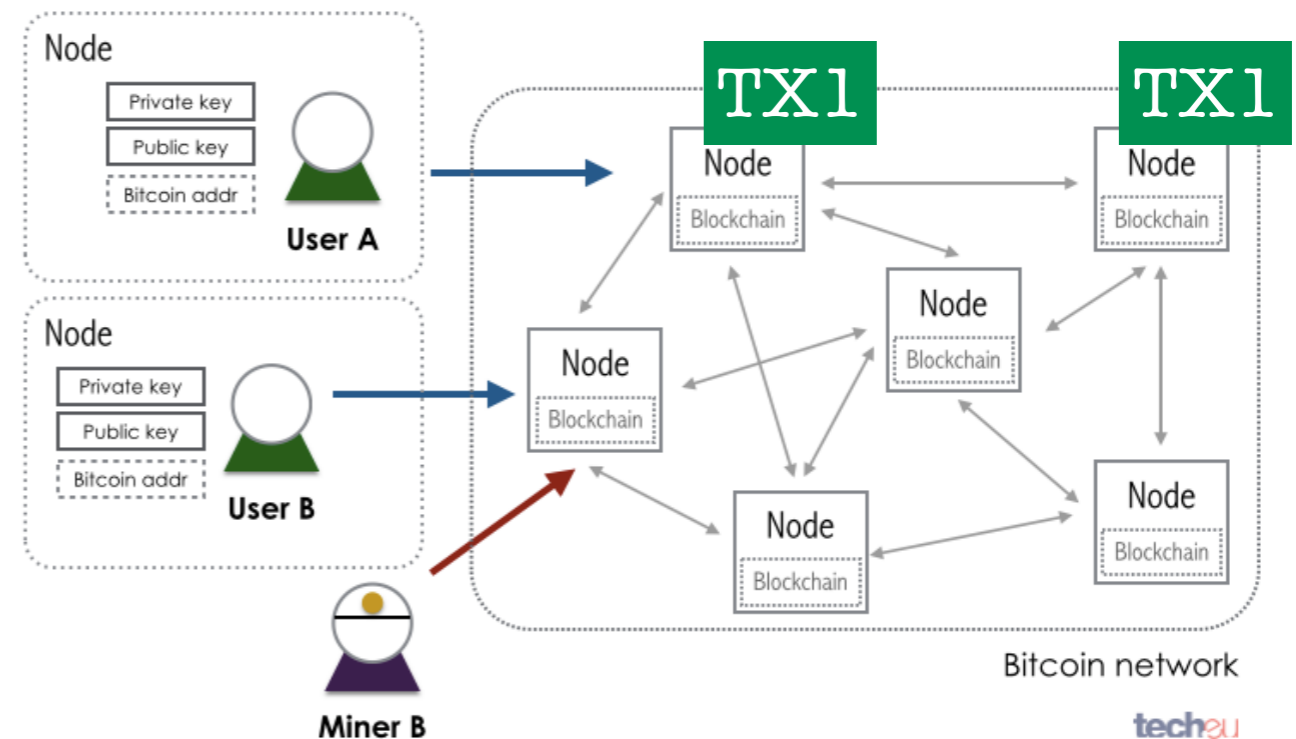
Analogy
(not formal)

The Four Dimensions of broadcast



PROPERTIES

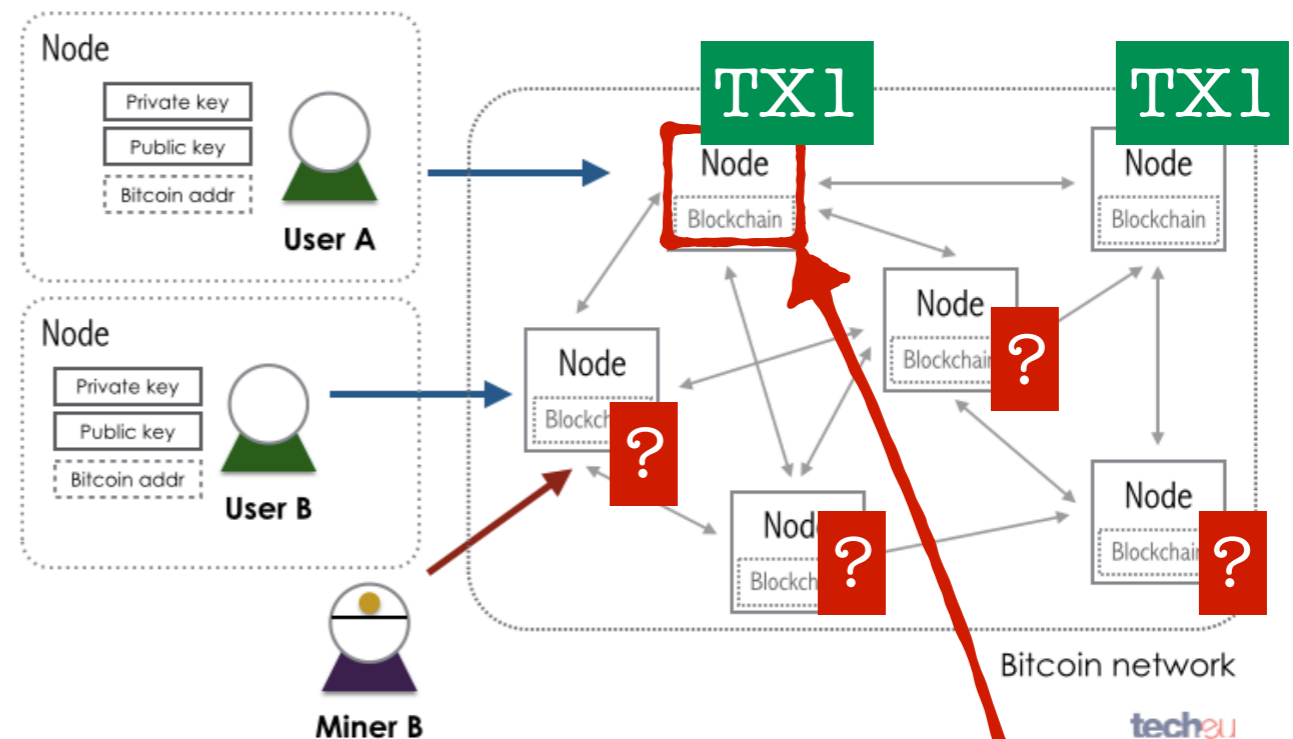
Do we need Reliability?



- Consider the following:
 - User A starts **TX1**
 - Use best-effort broadcast
 - Validity + !Duplication + !Creation
 - “the burden of reliability is on the sender”
 - Lacks Agreement \Rightarrow Nodes diverge

PROPERTIES

Do we need Reliability?

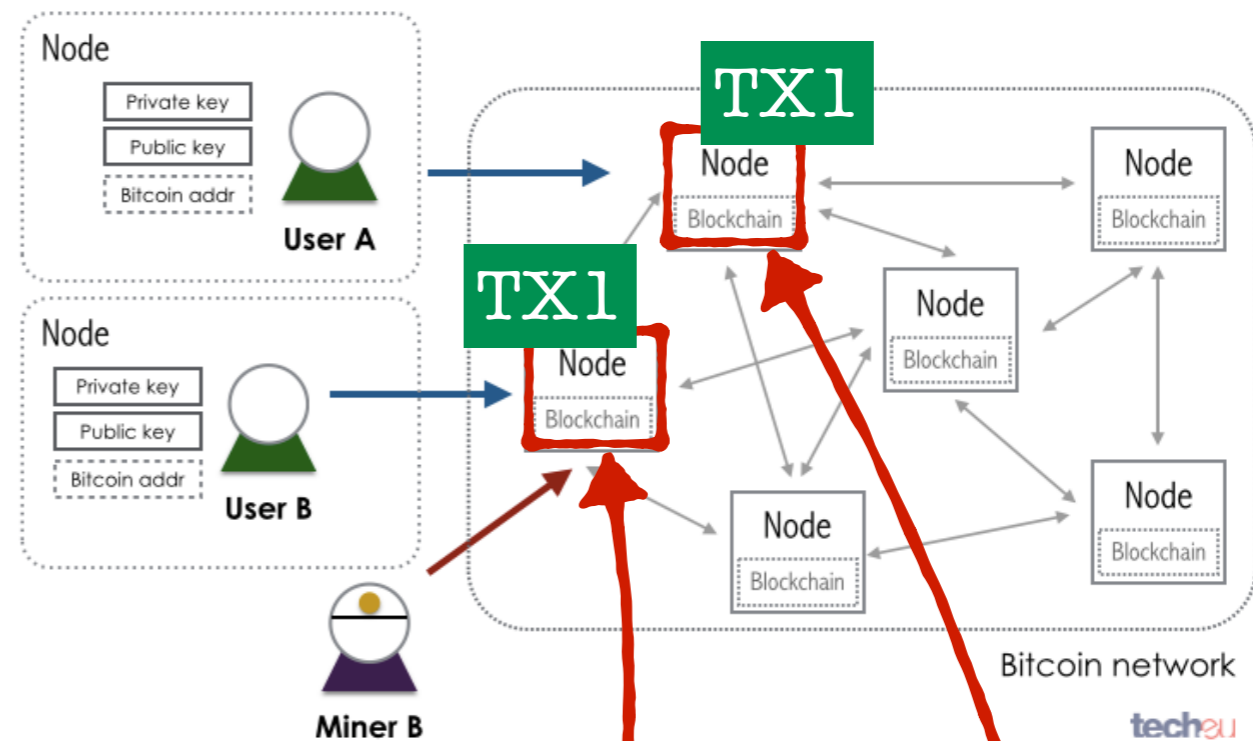


- Consider the following:
 - User A starts TX1
 - Use best-effort broadcast
 - Validity + !Duplication + !Creation
 - “the burden of reliability is on the sender”
 - Lacks Agreement \Rightarrow Nodes diverge

What if the sender crashes?

PROPERTIES

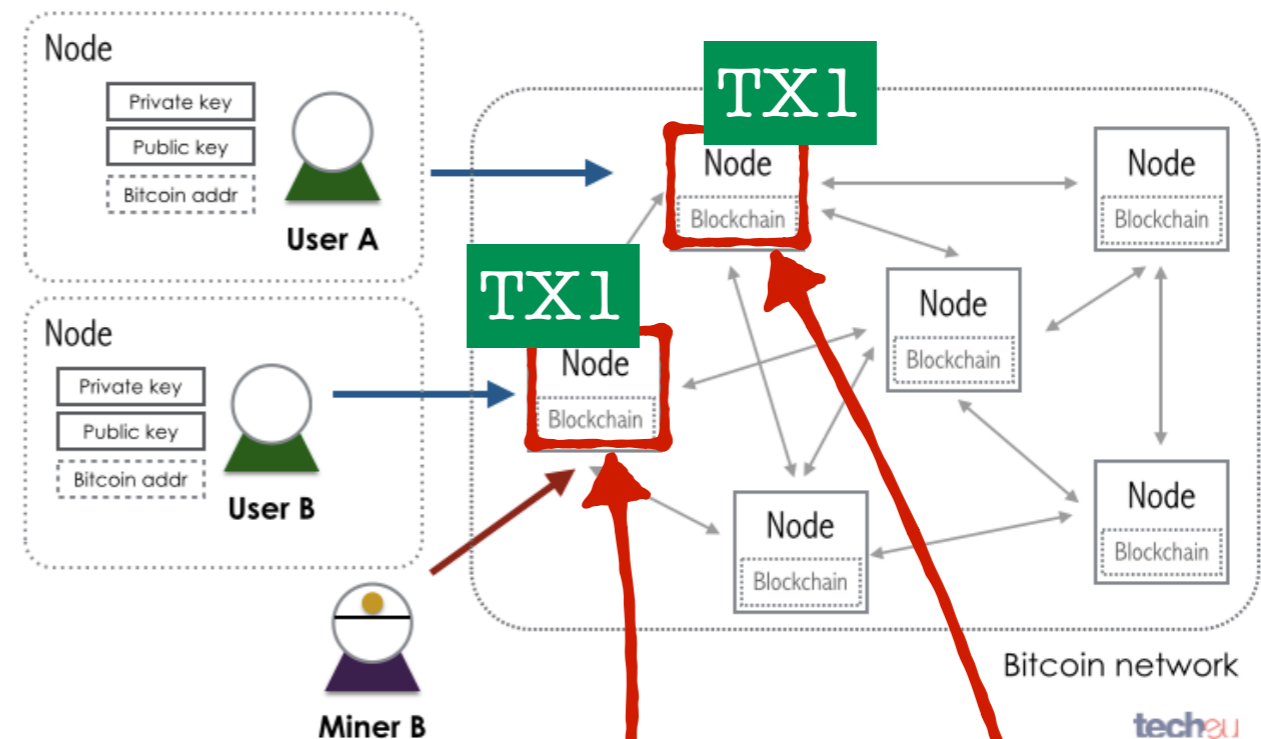
Do we need Uniformity?



- Consider the following:
 - User A starts TX1
 - Use regular reliable broadcast
 - Validity + !Duplication + !Creation + Agreement for correct nodes
 - Is it OK to deliver and crash?

PROPERTIES

Do we need Uniformity?

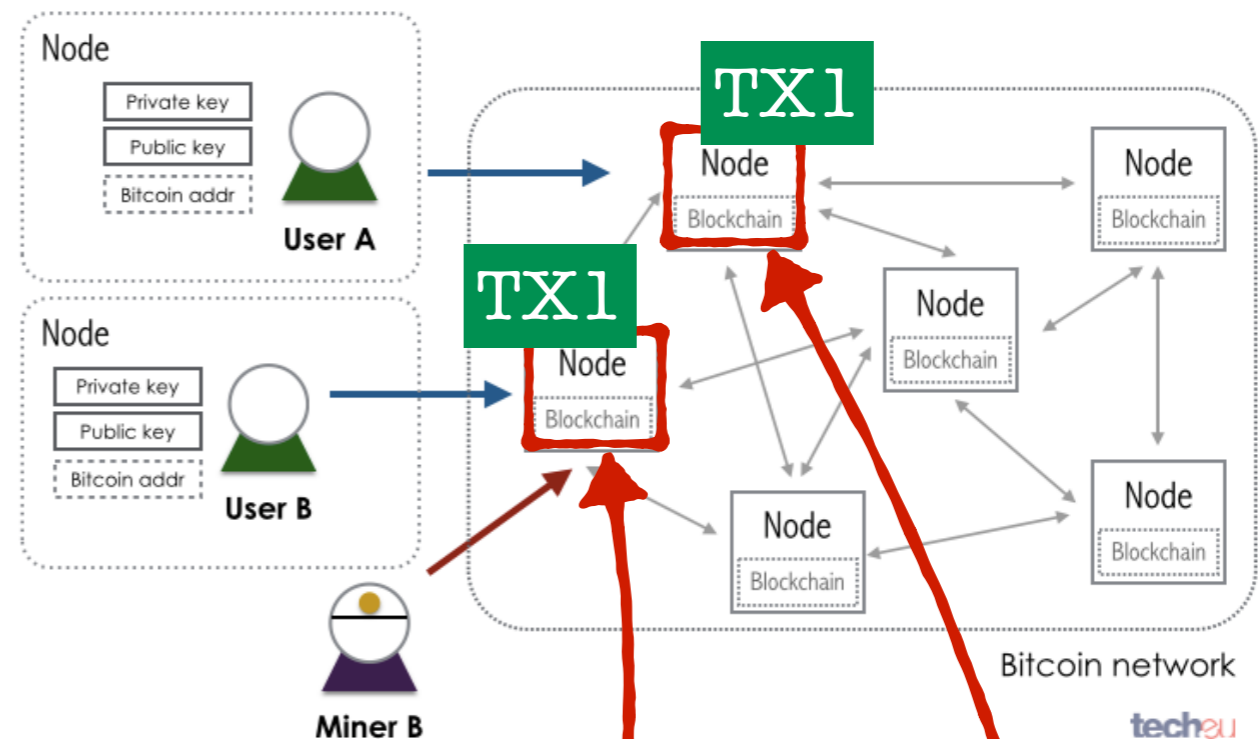


- Consider the following:
 - User A starts TX1
 - Use regular reliable broadcast
 - Validity + !Duplication + !Creation + Agreement for correct nodes
 - Is it OK to deliver and crash?

“Uniformity is important if nodes interact with the external world”

PROPERTIES

Do we need Uniformity?



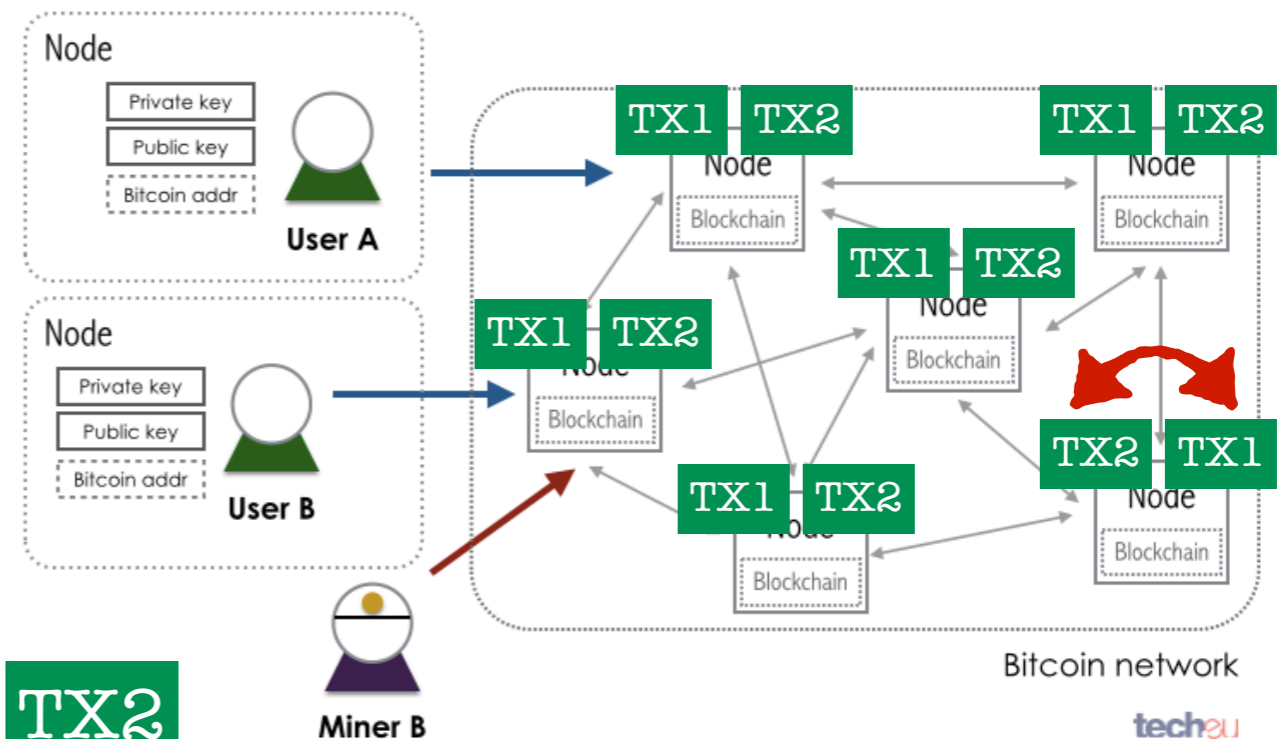
- Consider the following:
 - User A starts TX1
 - Use regular reliable broadcast
 - Validity + !Duplication + !Creation + Agreement for correct nodes
 - Is it OK to deliver and crash?

What if User B observes TX1 before the nodes crash?

“Uniformity is important if nodes interact with the external world”

PROPERTIES

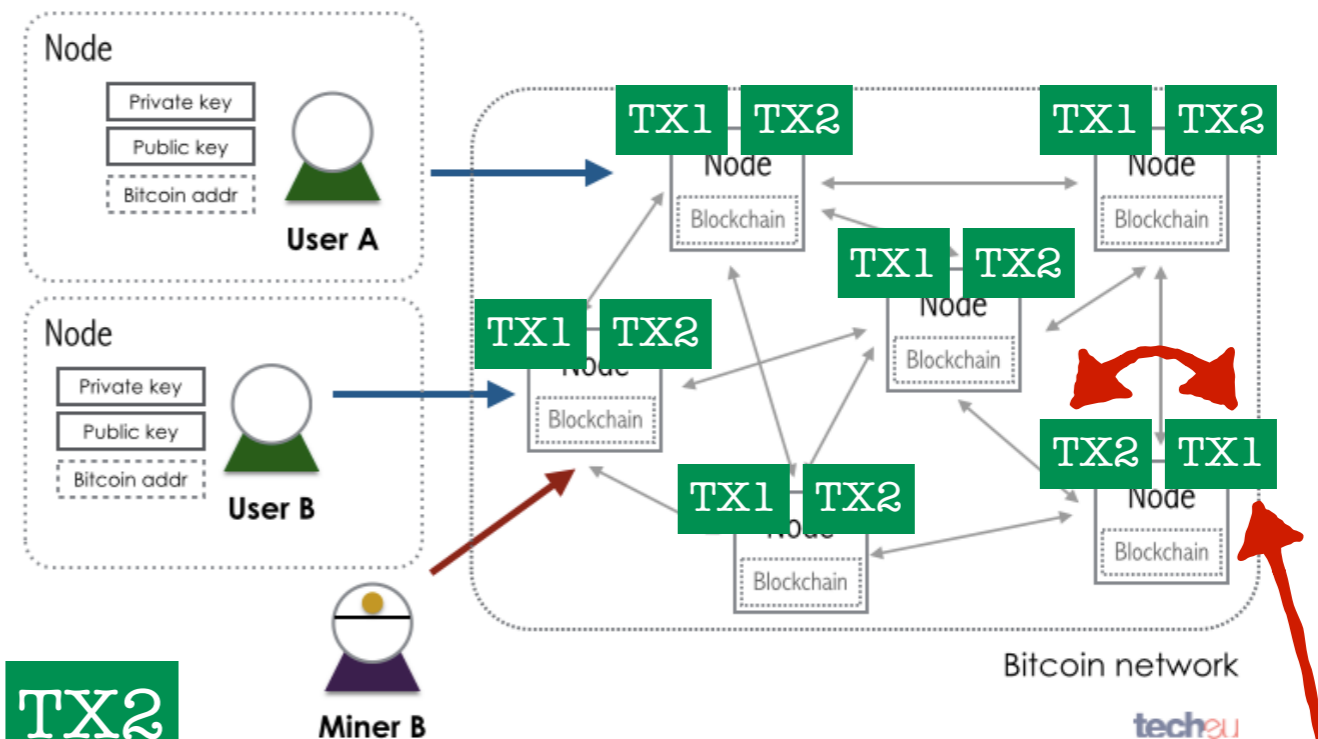
Do we need Causality? (partial ordering)



- Consider the following:
 - User A starts **TX1** and **TX2**
 - Use uniform reliable broadcast
 - Validity + !Duplication + !Creation + Uniform Agreement
⇒ Applies to all nodes
- All nodes deliver both TX, but the order may differ

PROPERTIES

Do we need Causality? (partial ordering)

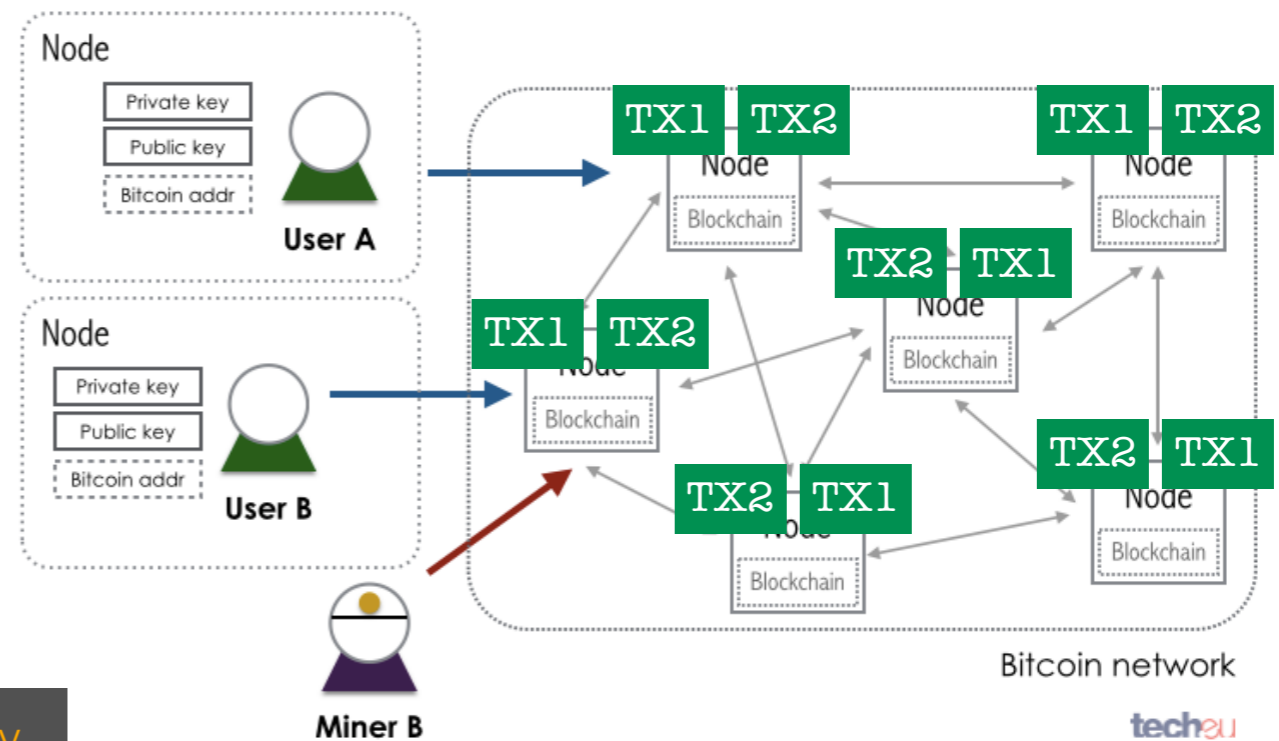


- Consider the following:
 - User A starts **TX1** and **TX2**
 - Use uniform reliable broadcast
 - Validity + !Duplication + !Creation + Uniform Agreement
⇒ Applies to all nodes
- All nodes deliver both TX, but the order may differ

What if **TX2** depends on
money from **TX1** ?

PROPERTIES

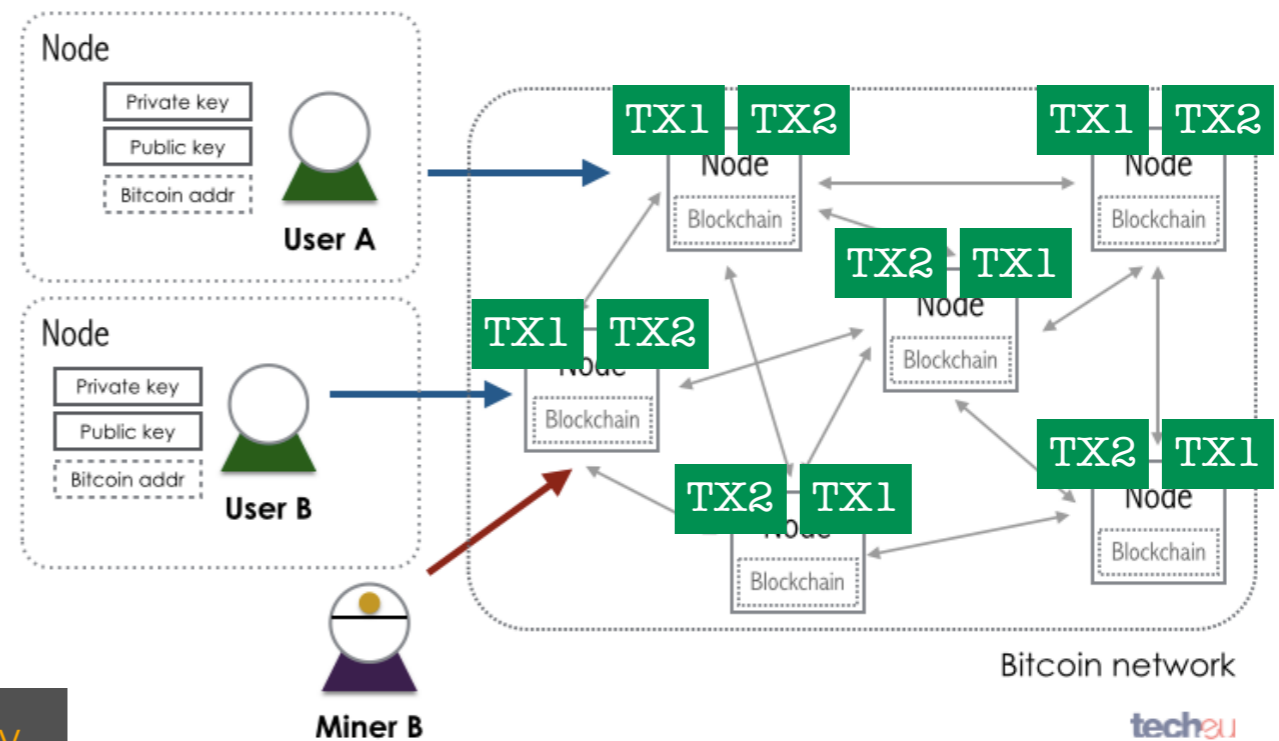
Do we need (Total) Ordering?



- Consider the following:
 - User A starts **TX1**
 - User B starts **TX2** No dependency among these two
 - Use causal-order uniform reliable broadcast
 - Validity + !Duplication + !Creation + Uniform Agreement + Causality
⇒ Respect causal dependencies among TXs
- All nodes deliver both TX, but the order may differ

PROPERTIES

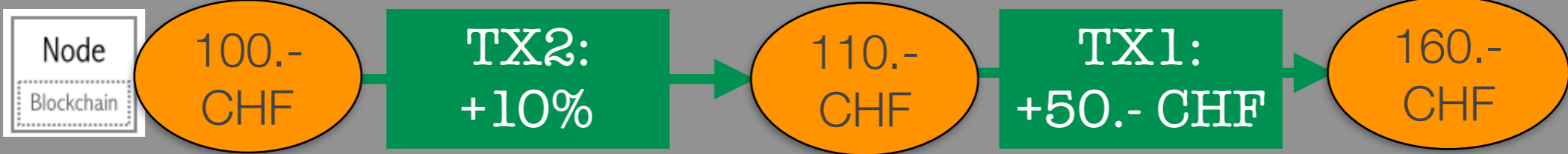
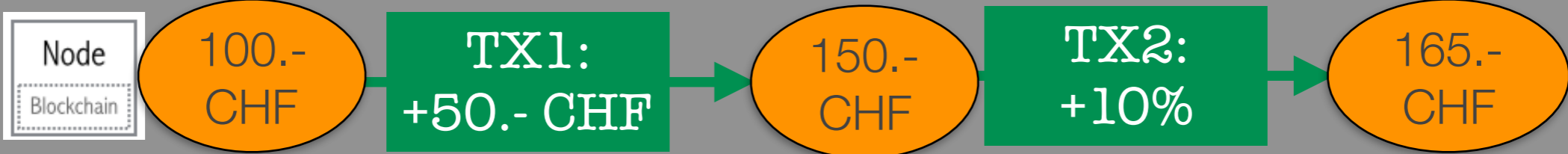
Do we need (Total) Ordering?



- Consider the following:
 - User A starts **TX1**
 - User B starts **TX2** No dependency among these two
 - Use causal-order uniform reliable broadcast
 - Validity + !Duplication + !Creation + Uniform Agreement + Causality
⇒ Respect causal dependencies among TXs
 - All nodes deliver both TX, but the order may differ

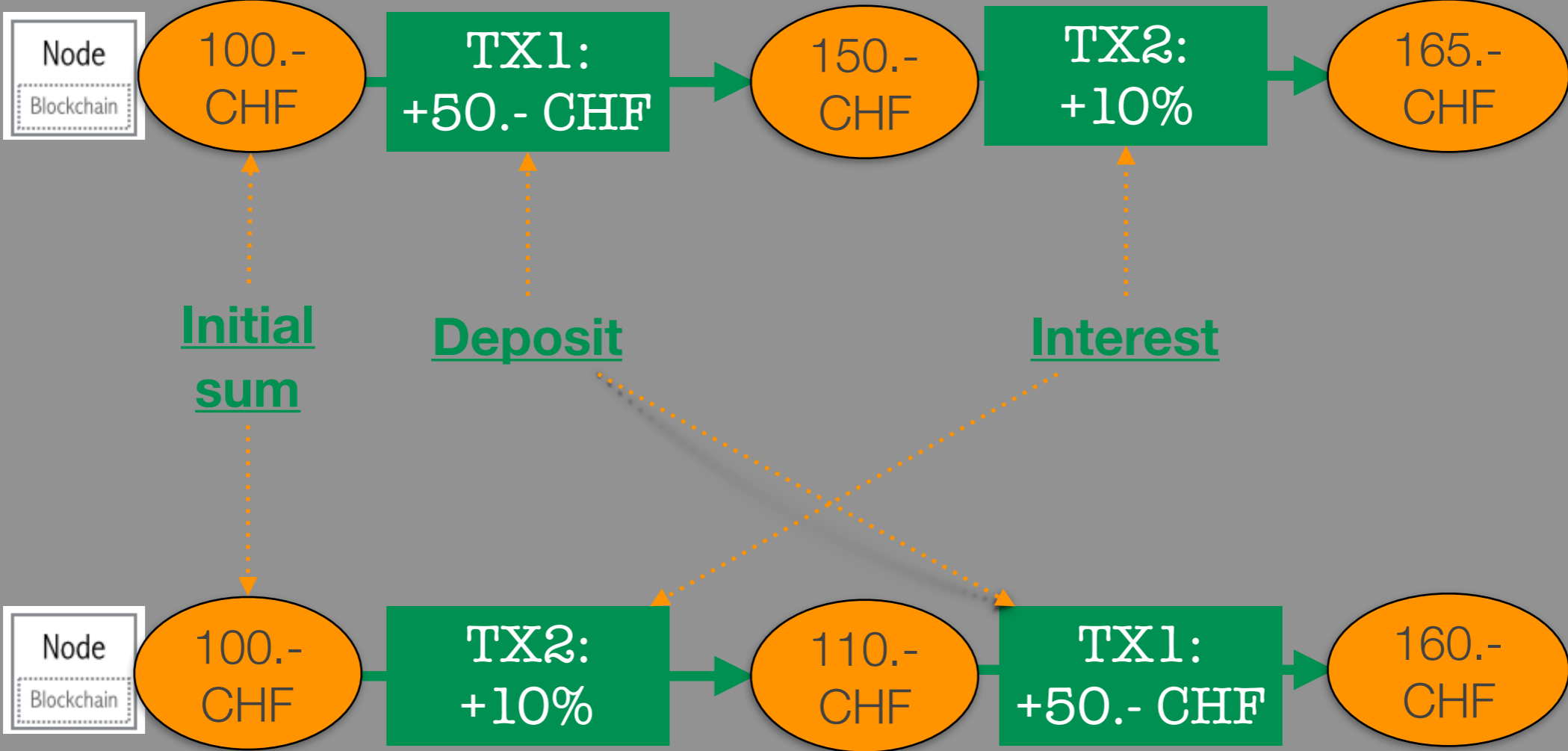
What if **TX1** and **TX2**
are not commutative?

Commutativity counter-example

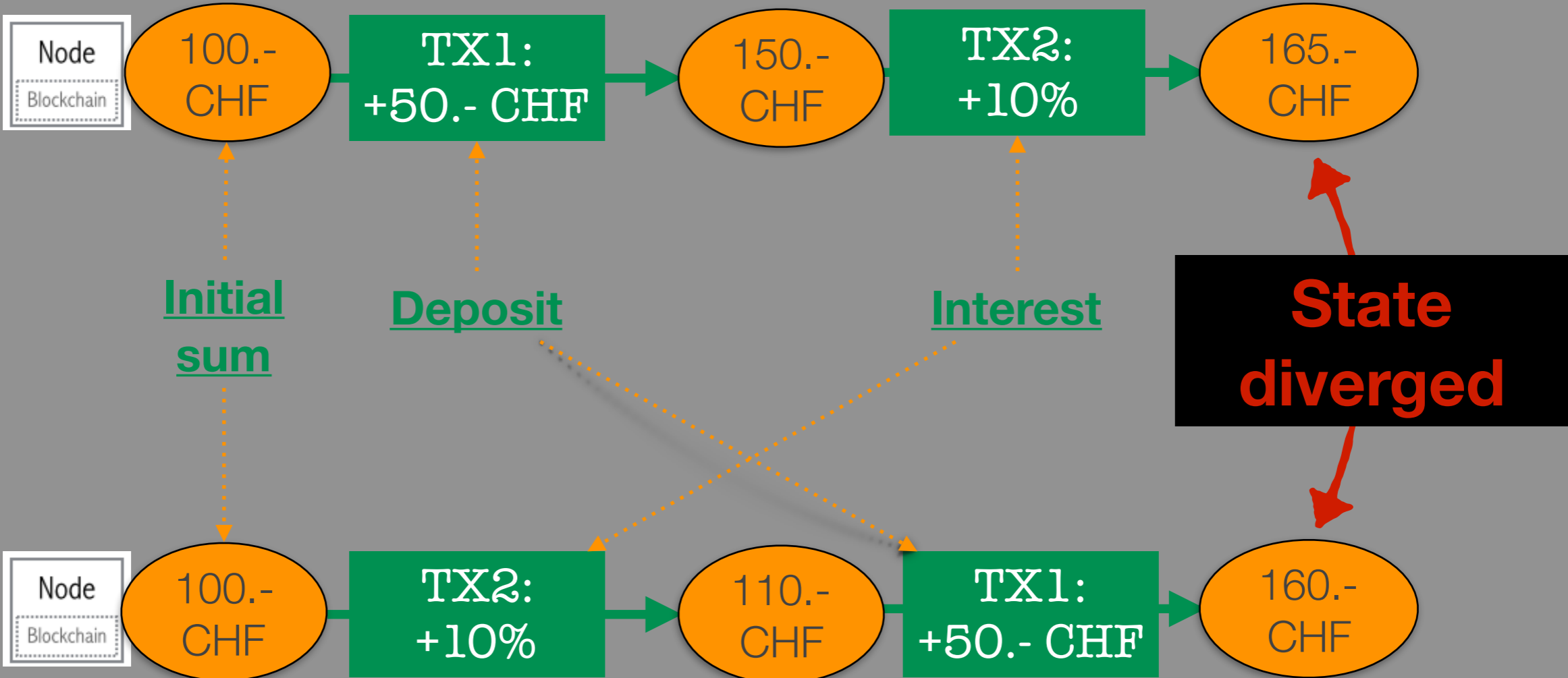


Commutativity

counter-example

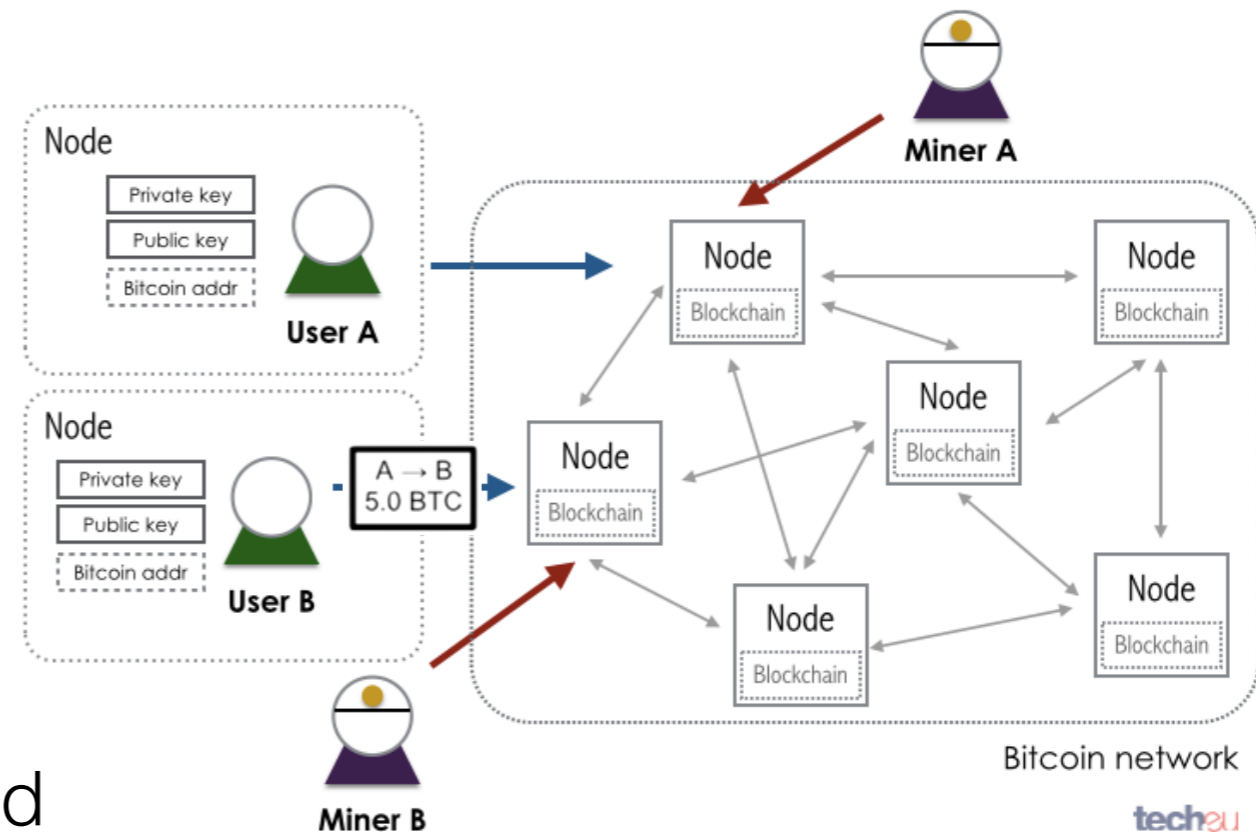


Commutativity counter-example



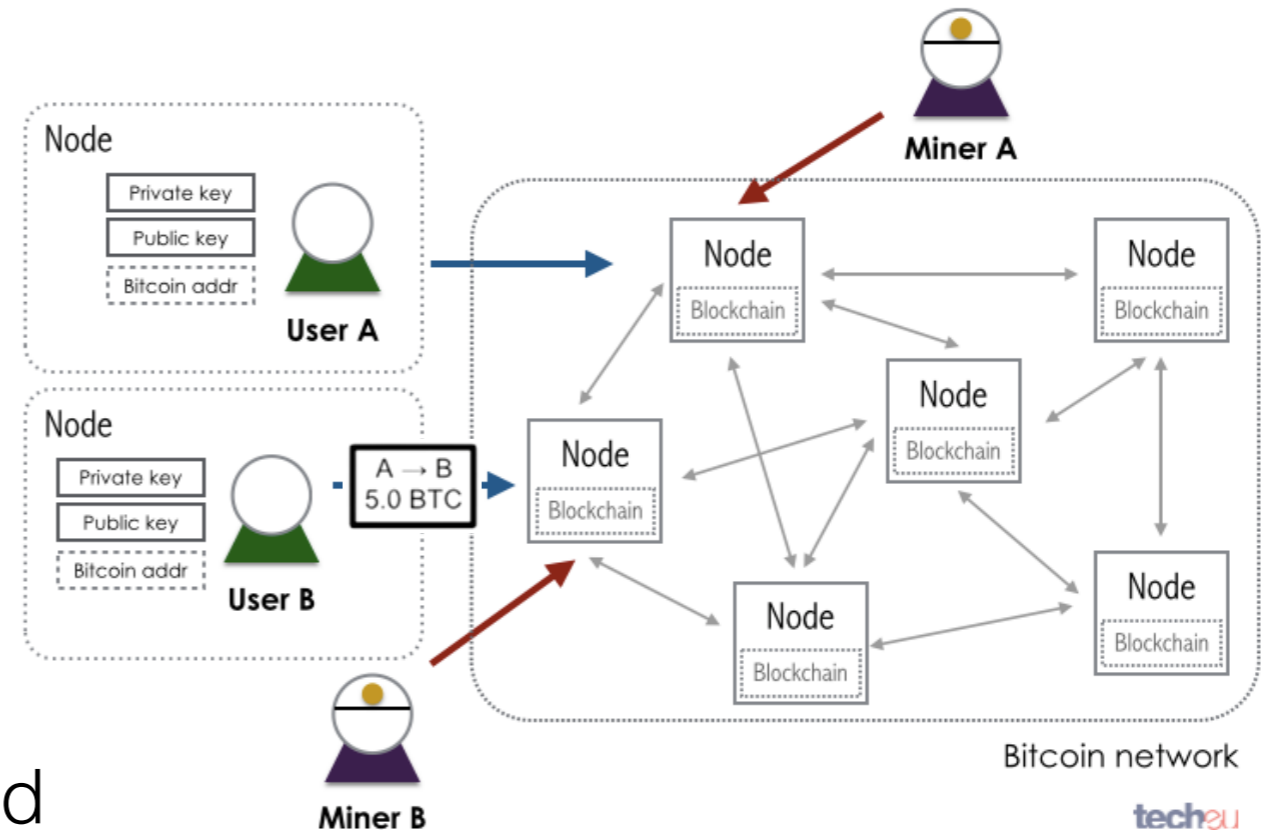
PROPERTIES

- Reliability
 - Sender crashes
 - Agreement
- Uniformity
 - Again, crashes
 - Interaction with outside world
- Causality
 - Partial order
 - Dependencies among TXs
- Total order
 - Commutativity



PROPERTIES

- Reliability
 - Sender crashes
 - Agreement
- Uniformity
 - Again, crashes
 - Interaction with outside world
- Causality
 - Partial order
 - Dependencies among TXs
- Total order
 - Commutativity

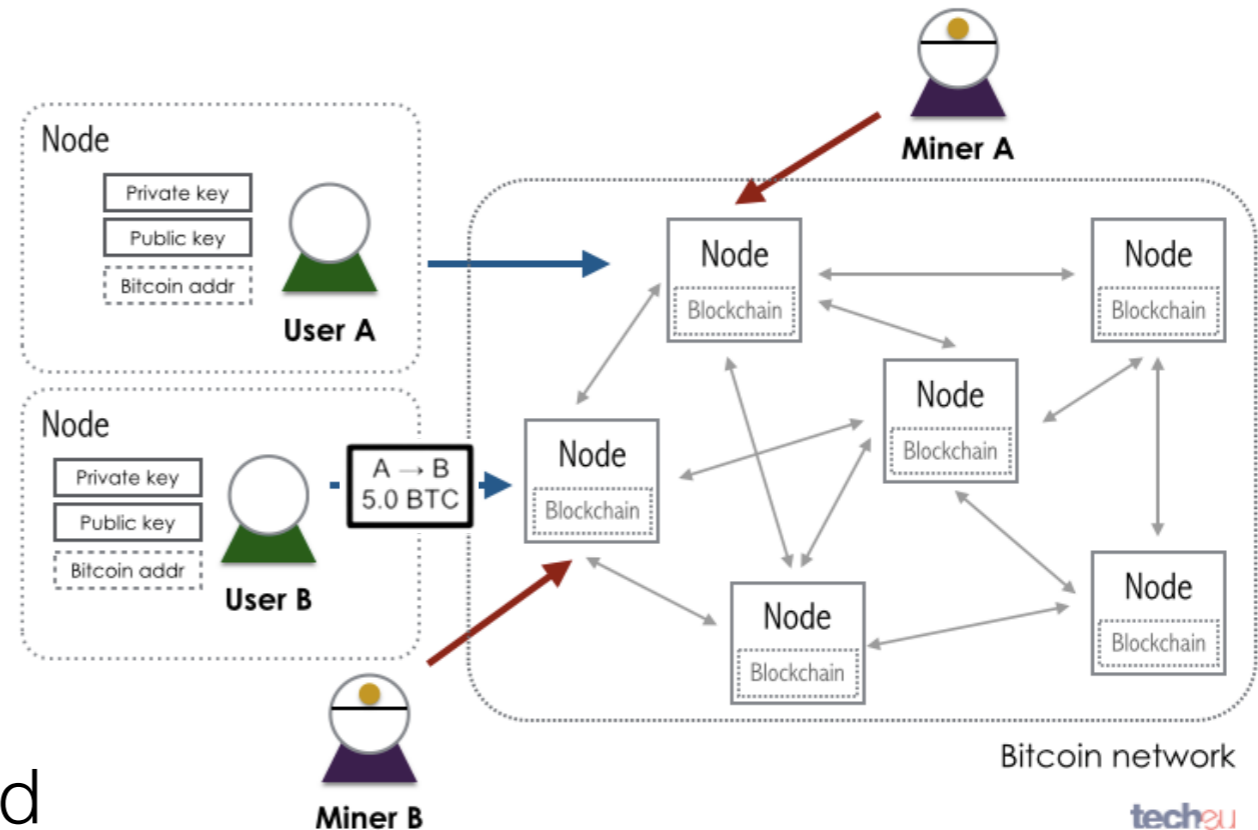


Contrast:

- Canonical Bitcoin requires best-effort broadcast
- Allows temp. inconsistencies
- Makes up with crypto

PROPERTIES

- Reliability
 - Sender crashes
 - Agreement
- Uniformity
 - Again, crashes
 - Interaction with outside world
- Causality
 - Partial order
 - Dependencies among TXs
- Total order
 - Commutativity



Contrast:

- Canonical Bitcoin requires best-effort broadcast
- Allows temp. inconsistencies
- Makes up with crypto

All properties are desirable

PERFORMANCE

Goals:

1. Replace traditional CAMIPRO
2. Make banks obsolete



AN INFORMAL COMPARISON...



PERFORMANCE

Goals:

1. Replace traditional CAMIPRO
2. Make banks obsolete



AN INFORMAL COMPARISON...



CANONICAL BITCOIN

- Best-effort broadcast
- Optimized for Internet (WAN)
- **TX~10 minutes**

How useful would that be?
(in the lunch queue)

CAMIPRO BITCOIN

- Bitcoin-Broadcast
- Optimized for small network (EPFL LAN)
- TX~1second (back of the envelope)
- Acceptable latency & throughput

Specification





Module:

Name: CAMIPRO-Bitcoin, instance *cbit*

Properties:

- RB1, RB2, RB3, RB4
- Causal Order (CO)
- Total order (TO)



What's left
?????

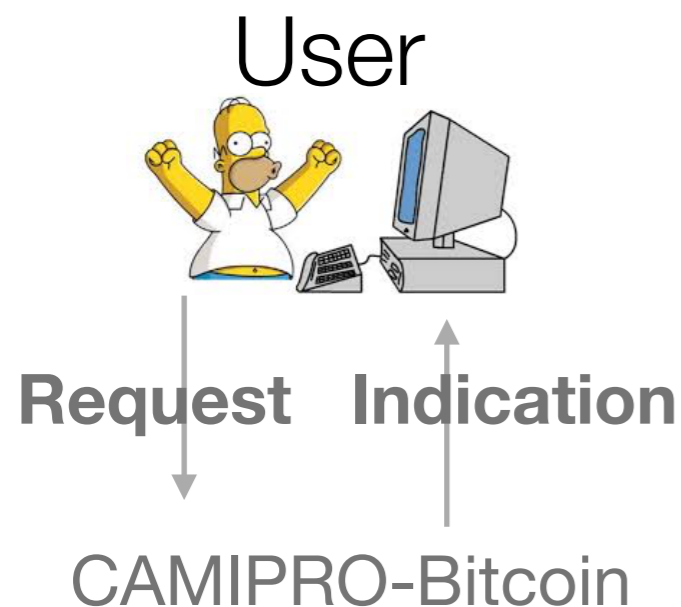
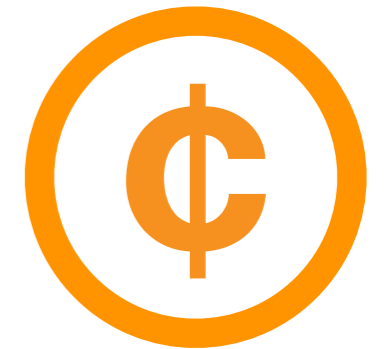


Module:

Name: CAMIPRO-Bitcoin, **instance** *cbit*

Properties:

- RB1, RB2, RB3, RB4
- Causal Order (CO)
- Total order (TO)



Events:

Request: $\langle cbit, Start \mid TX \rangle$:

Attempts to commit TX

Indication: $\langle cbit, Status \mid TX, s \rangle$:

Indicates the status $s \in \{ "Abort", "Commit" \}$ of TX

S4: (Student's) Simple Storage Service

Overview:

1. Setup & System model
2. Operations
3. Goals
4. Technicalities
5. One last look at causal consistency

S4: **(Student's) Simple Storage Service**

Overview:

1. Setup & System model
2. Operations
3. Goals
4. Technicalities
5. One last look at causal consistency



[BONUS PROJECT]

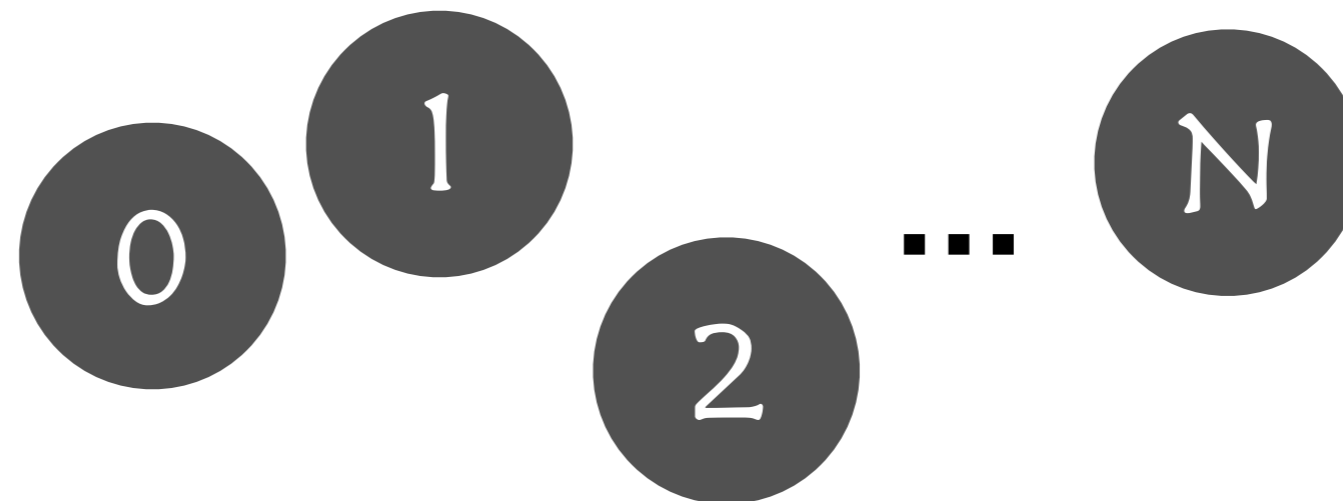
S4 — System model

Setup:

- Users



- Storage replicas (Nodes)



- Objects: stored on replicas

⟨key, value⟩

⟨A, 0⟩

⟨B, 1⟩

⟨C, 9⟩

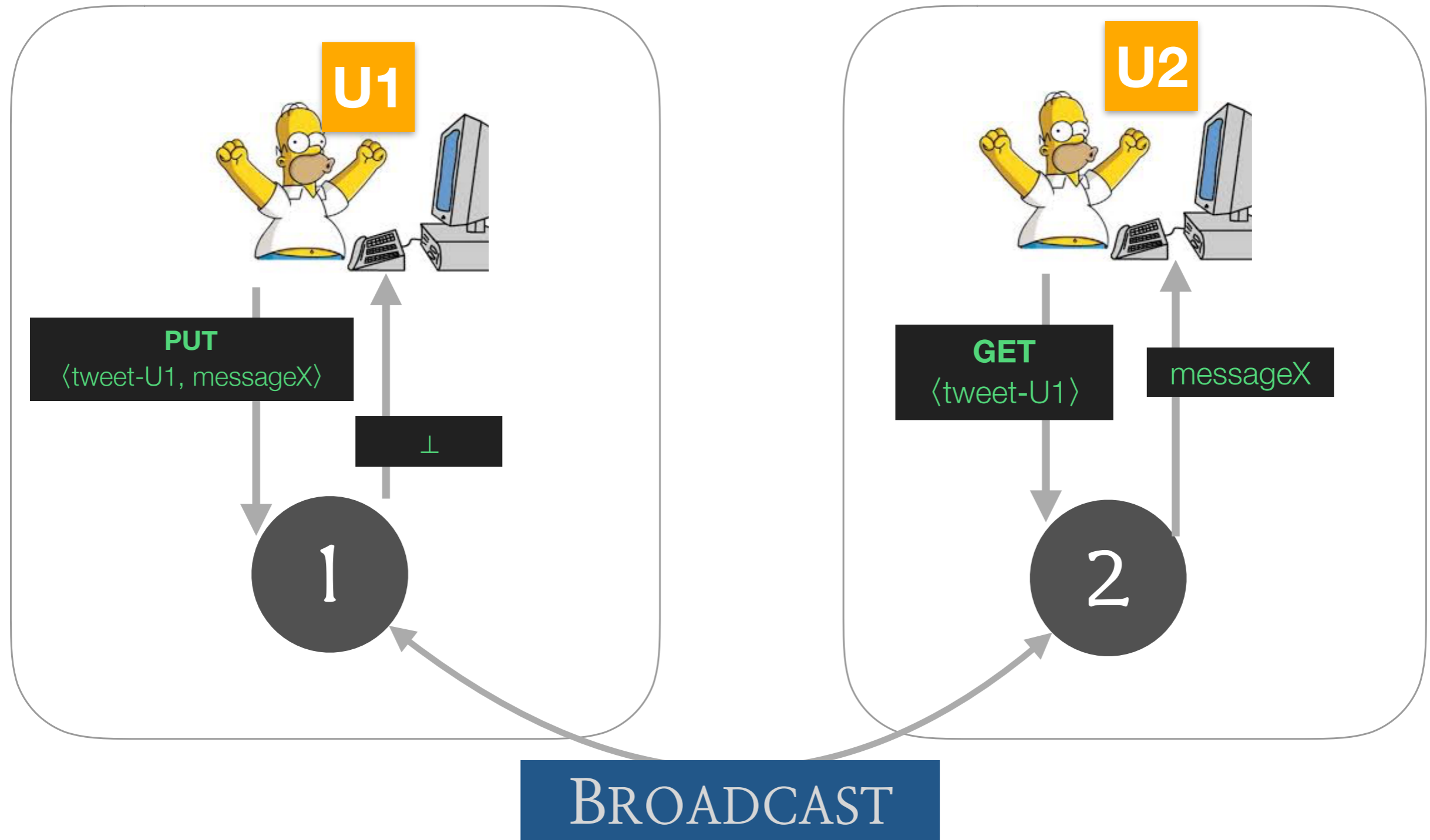
⟨MSGx, hello⟩

⟨tweetY, BYE⟩

S4 — Operations

PUT

GET



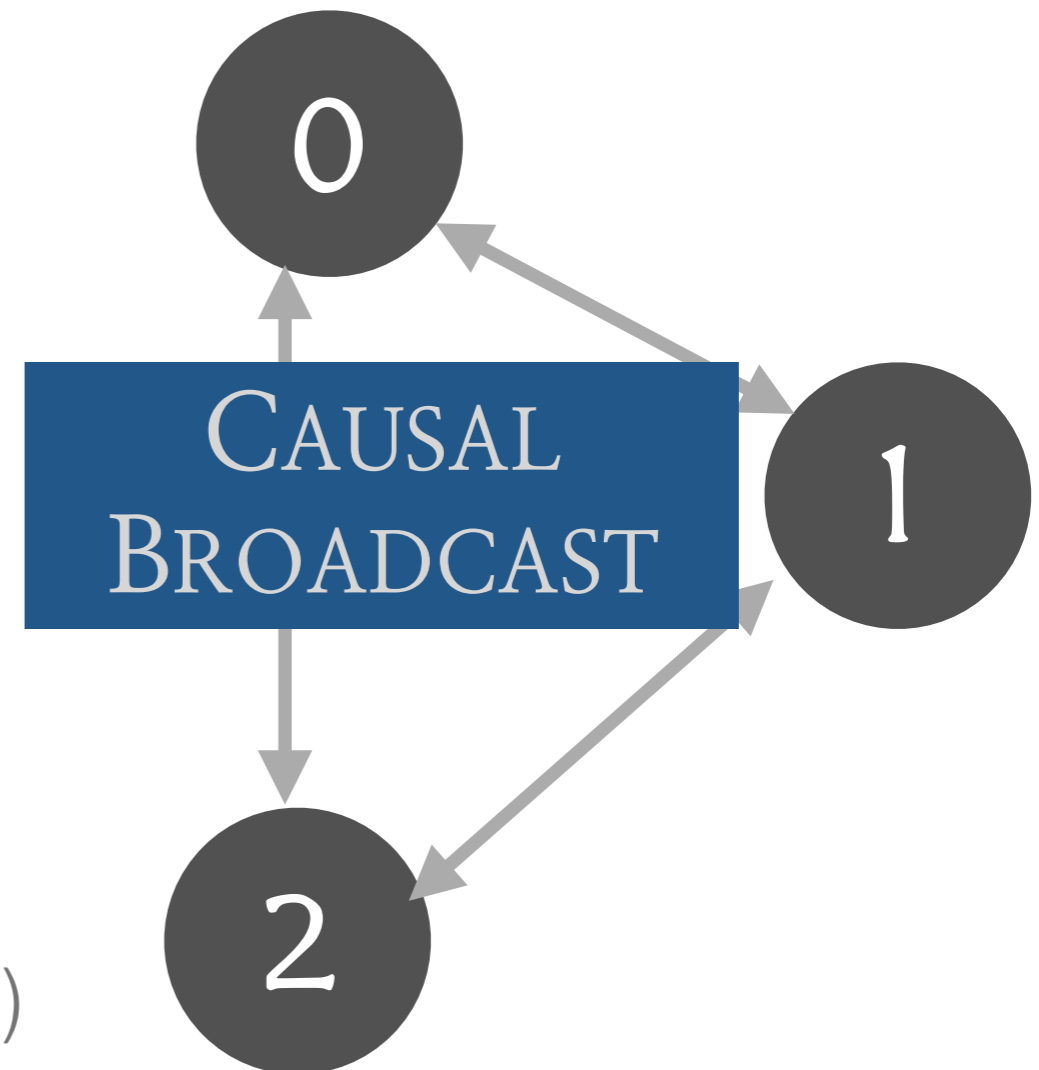
Goals:

1. Reliability

- Three replicas
- Communicate through message-passing

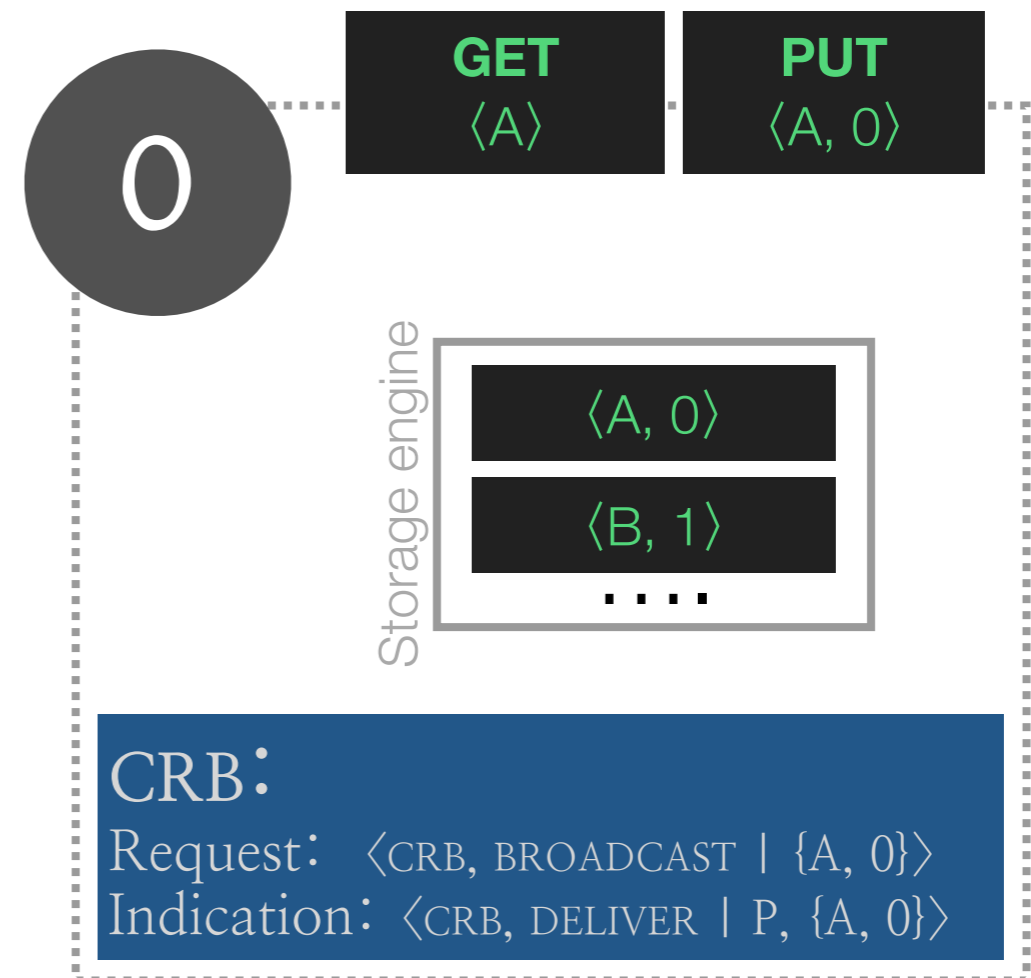
2. Consistency

- We want to use it in the Internet (WAN)
 - (Unlike CAMIPRO-Bitcoin, which was optimized for the EPFL network)
 - Can't afford total-order, too expensive!
- Causal consistency



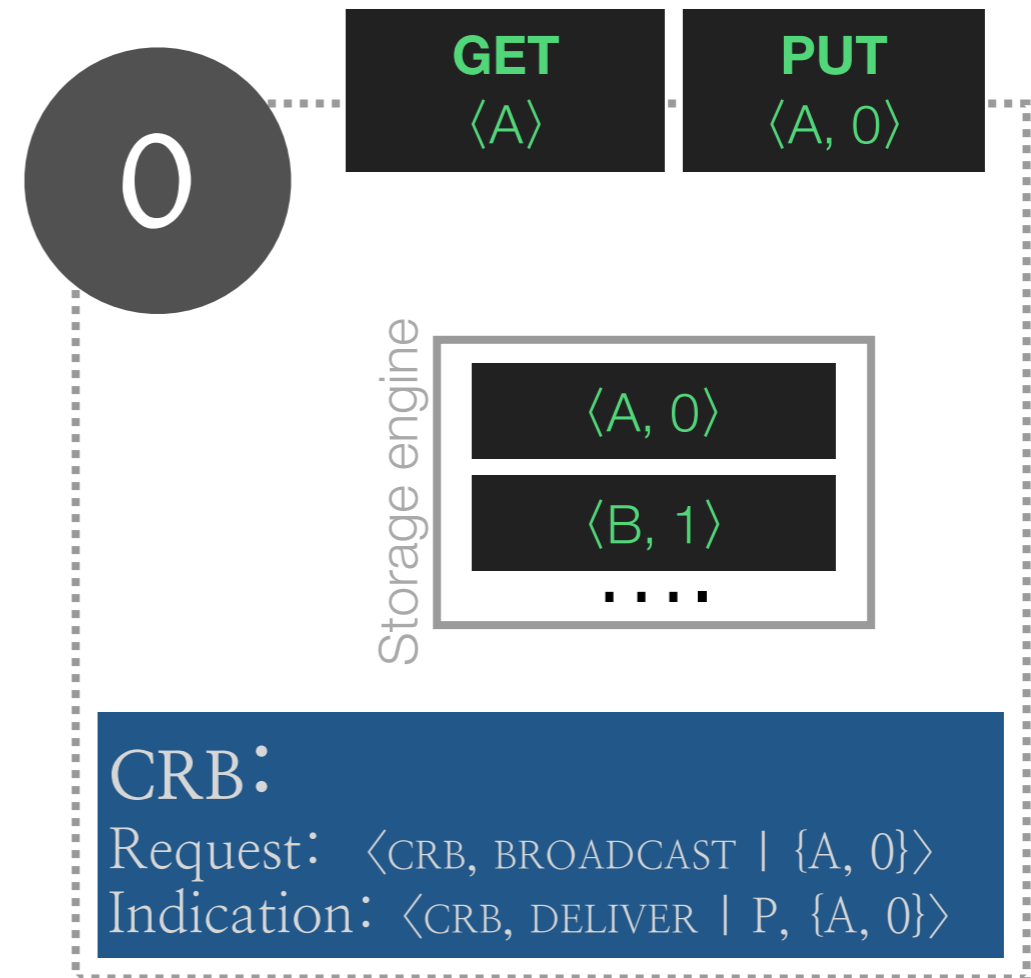
Nodes:

- Each node is a Linux process
- Contains:
 - A storage engine
 - A causal-broadcast implementation (*crb*)



Nodes:

- Each node is a Linux process
- Contains:
 - A storage engine
 - A causal-broadcast implementation (*crb*)



We start a process by calling:

```
s4 addr-0 port-0 addr-1 port-1 addr-2 port-2 n f input file
```

Red brackets group the arguments as follows:

- addr-0 port-0 → **Node 0: Address+Port**
- addr-1 port-1 → **Node 1: Address+Port**
- addr-2 port-2 → **Node 2: Address+Port**
- n f → **ID_∈{0,1,2}**

**To start all nodes,
we execute the commands:**

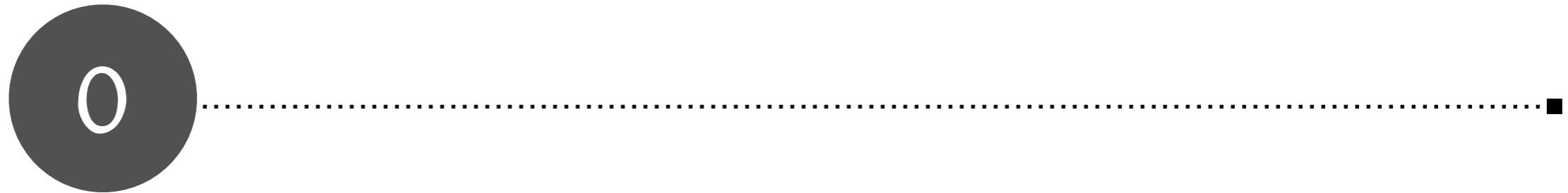
```
s4 127.0.0.1 8900 127.0.0.1 8901 127.0.0.1 8902 0 0.input  
s4 127.0.0.1 8900 127.0.0.1 8901 127.0.0.1 8902 1 1.input  
s4 127.0.0.1 8900 127.0.0.1 8901 127.0.0.1 8902 2 2.input
```

**To start all nodes,
we execute the commands:**

```
s4 127.0.0.1 8900 127.0.0.1 8901 127.0.0.1 8902 0 0.input  
s4 127.0.0.1 8900 127.0.0.1 8901 127.0.0.1 8902 1 1.input  
s4 127.0.0.1 8900 127.0.0.1 8901 127.0.0.1 8902 2 2.input
```

input files

INPUT FILES, OUTPUT FILES AND SIGNALS

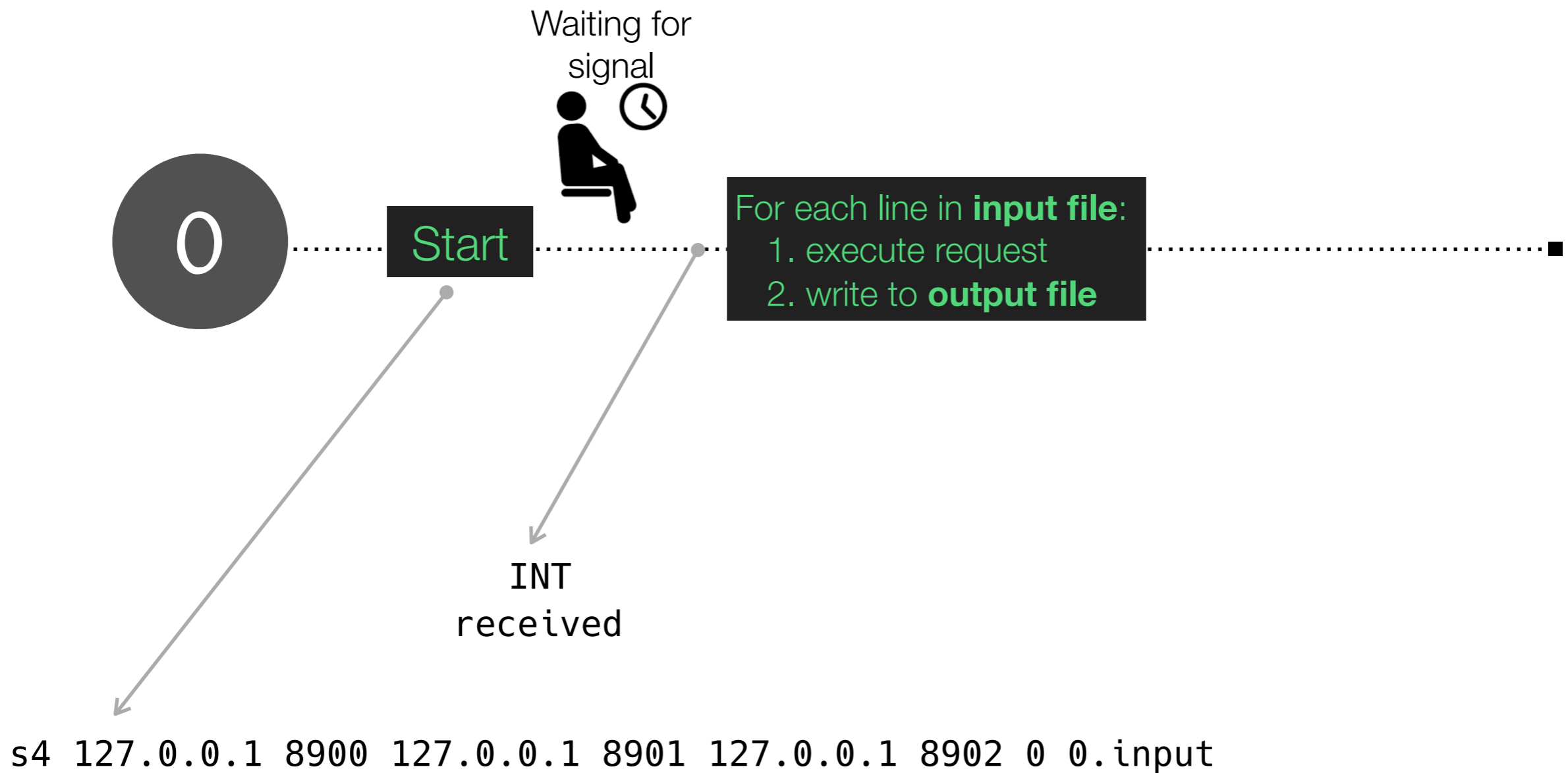


INPUT FILES, OUTPUT FILES AND SIGNALS

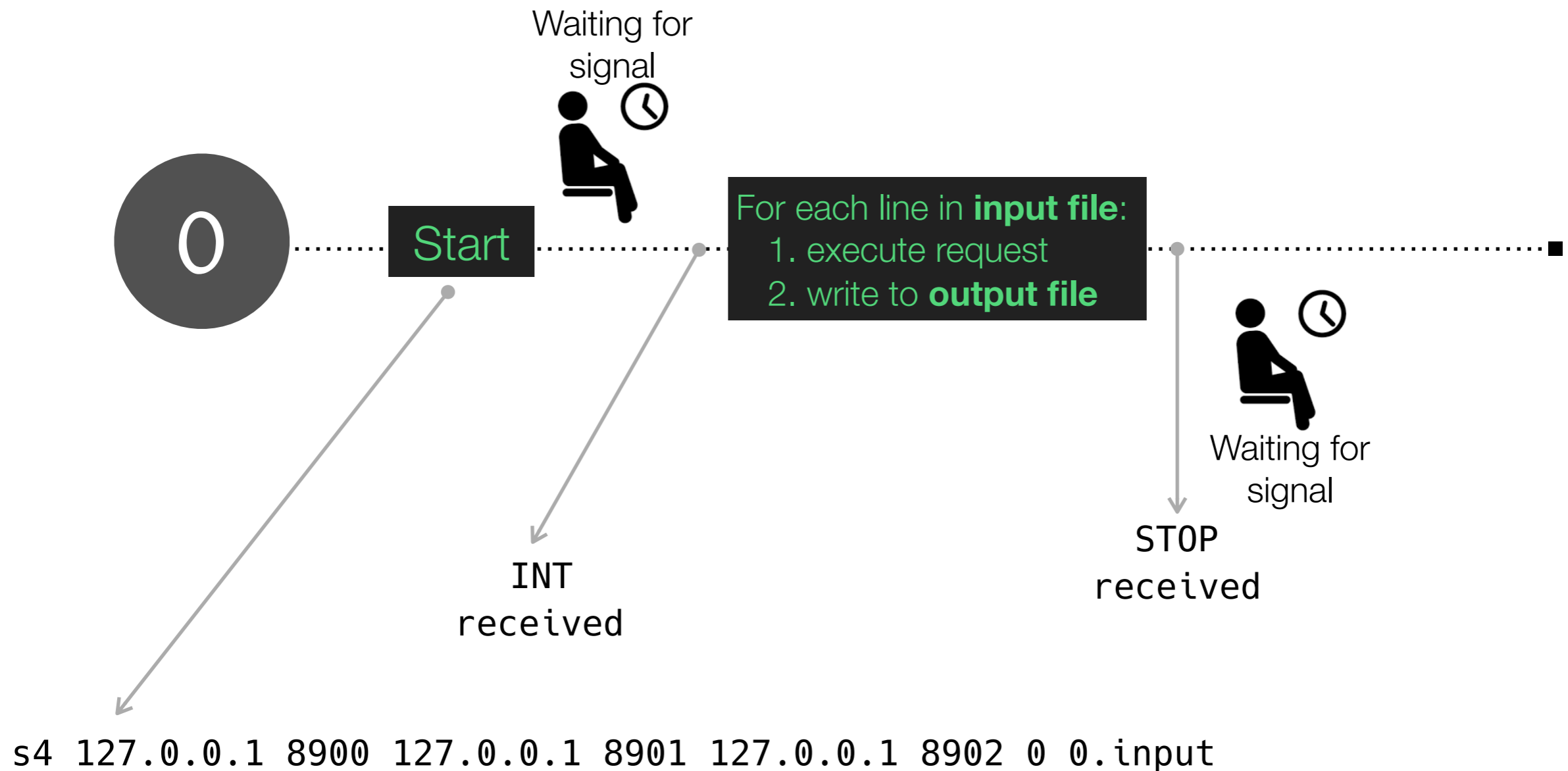


```
s4 127.0.0.1 8900 127.0.0.1 8901 127.0.0.1 8902 0 0.input
```

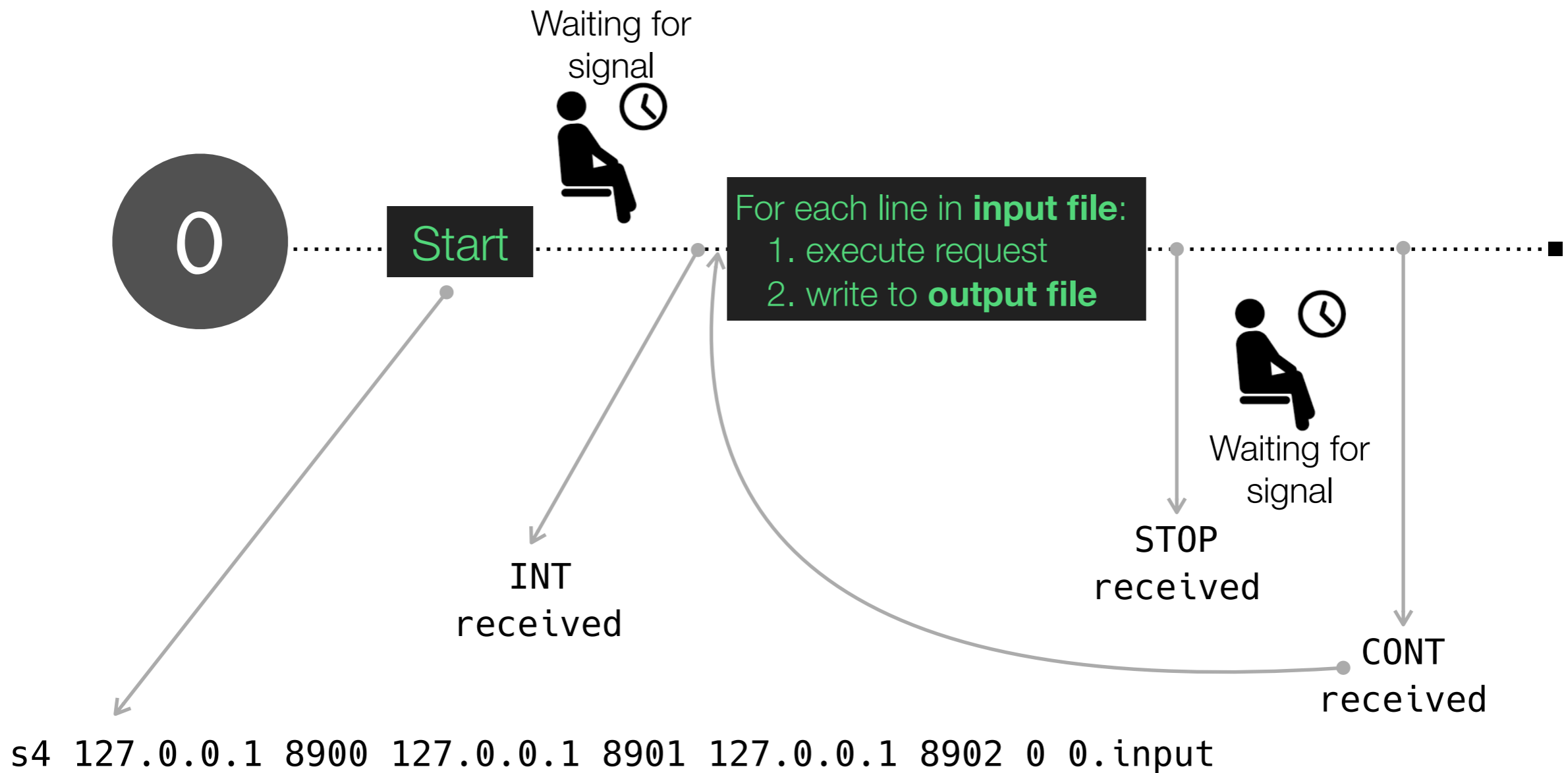
INPUT FILES, OUTPUT FILES AND SIGNALS



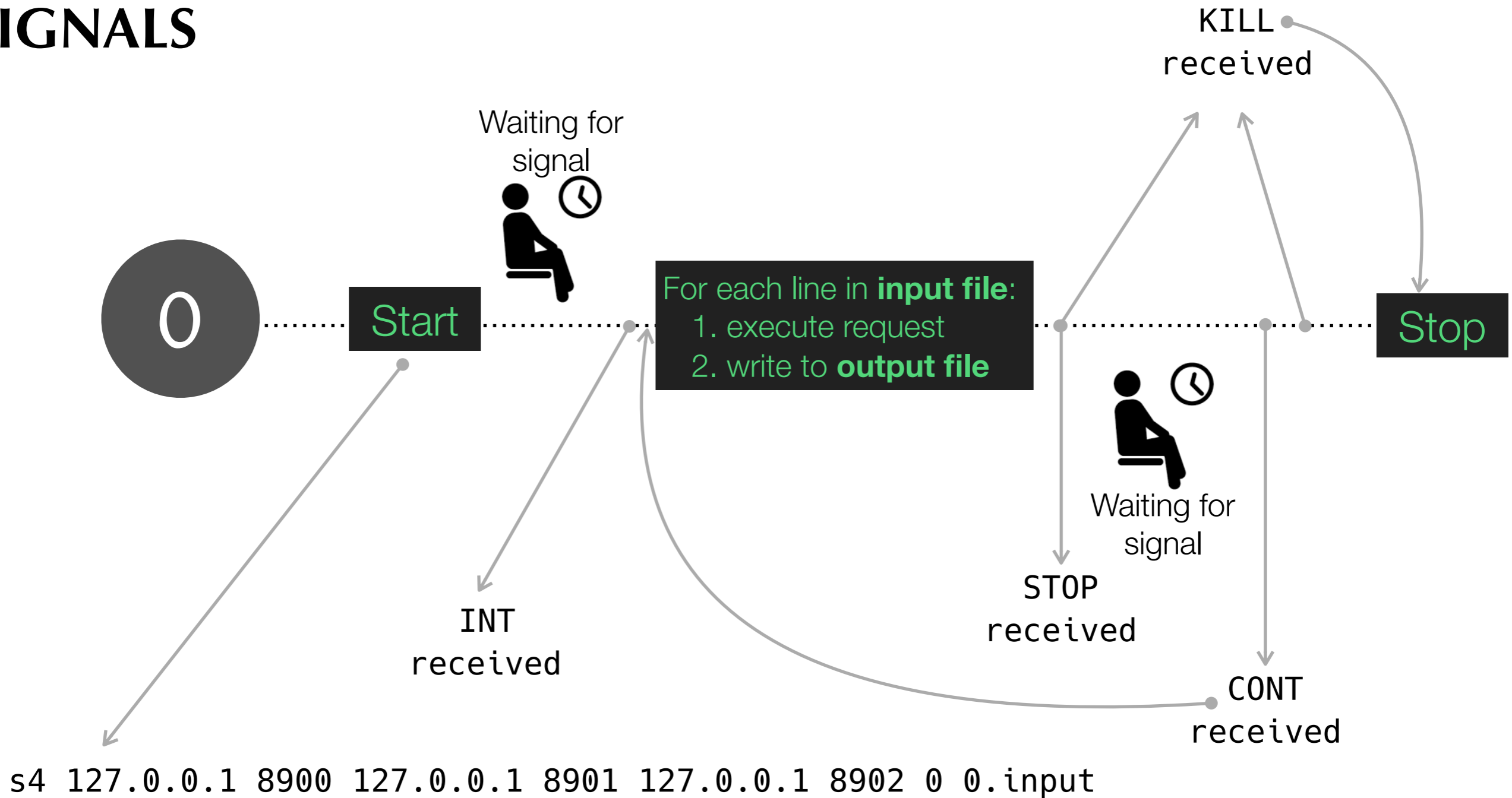
INPUT FILES, OUTPUT FILES AND SIGNALS



INPUT FILES, OUTPUT FILES AND SIGNALS

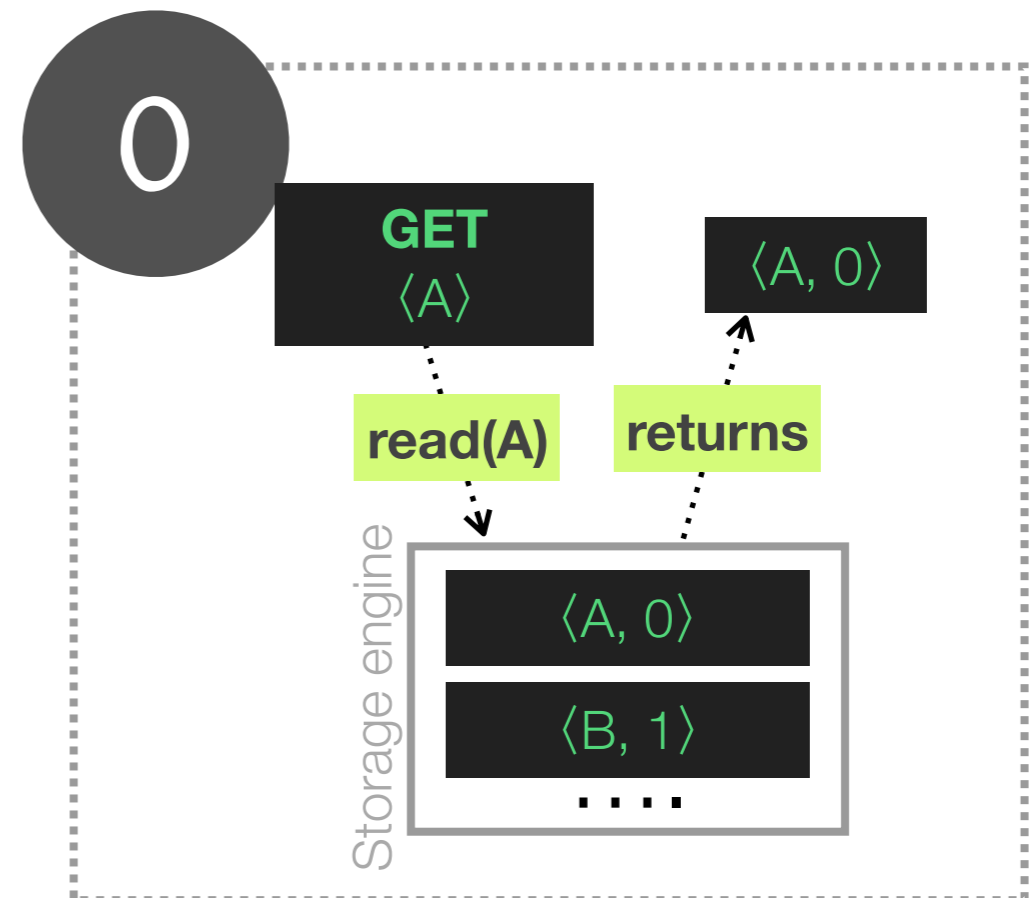


INPUT FILES, OUTPUT FILES AND SIGNALS



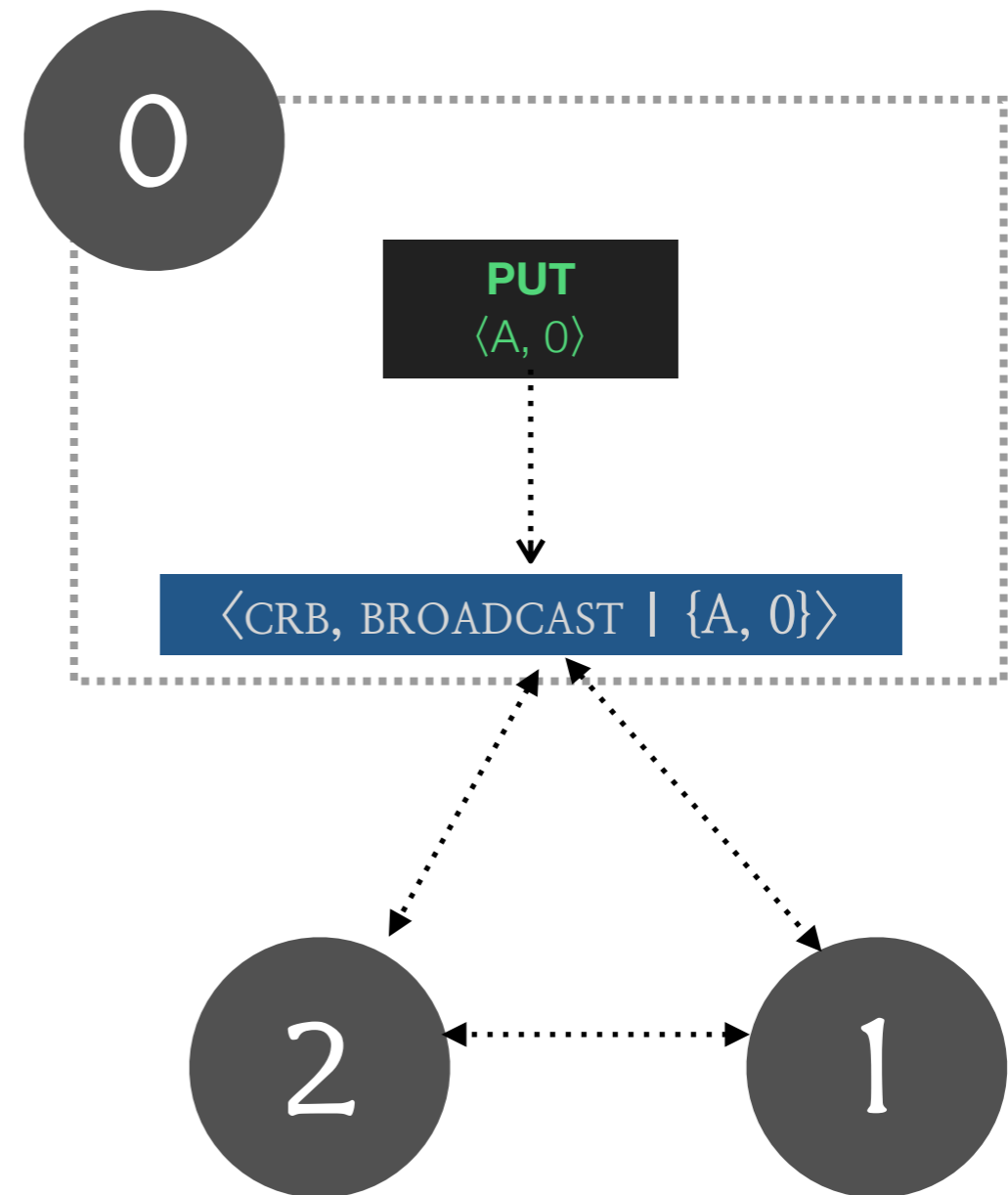
OPERATIONS — GET

- Translates to a local read operation
- Directly from the Storage engine
- Write the read value in the output file
- No need to contact the other nodes



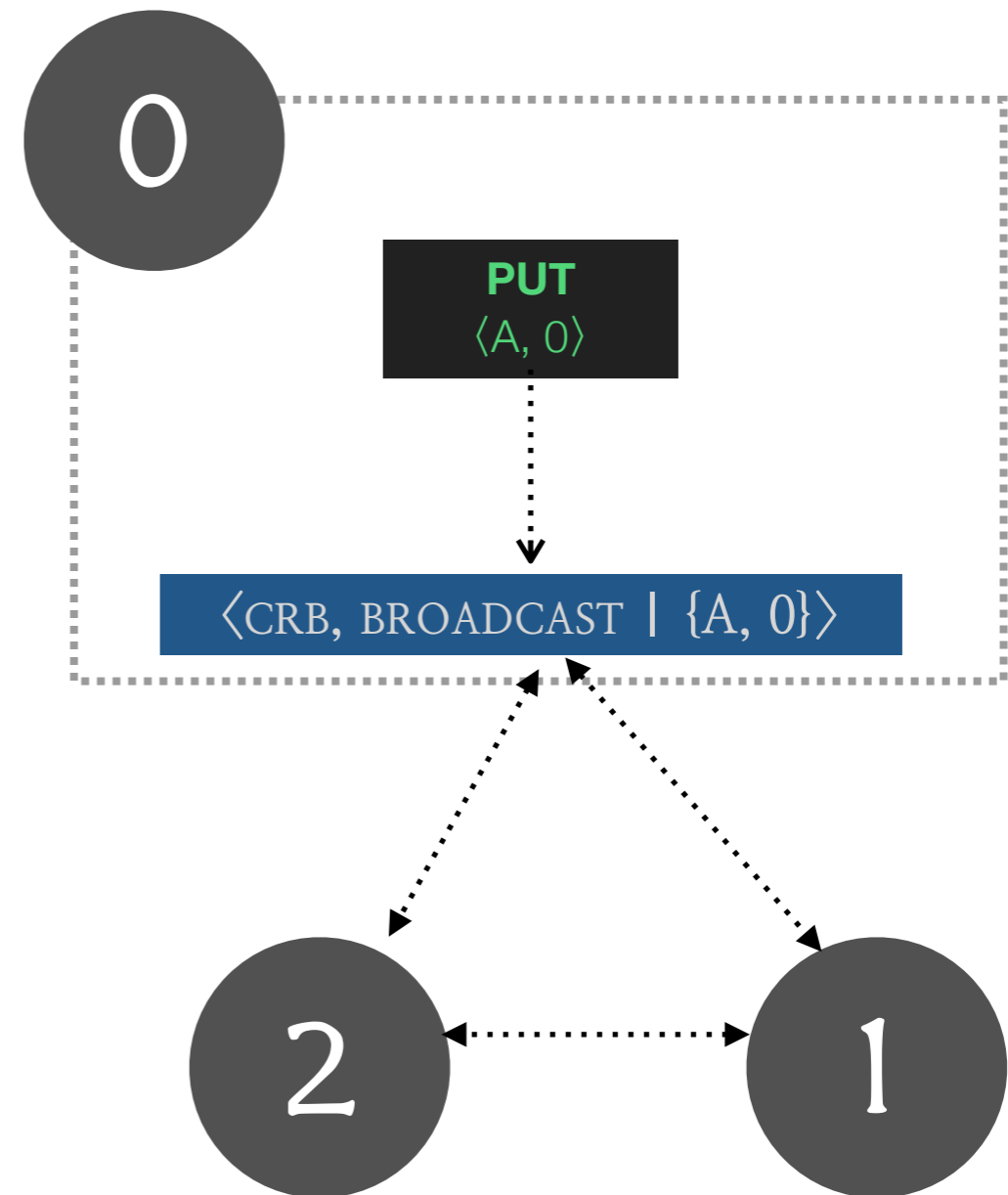
OPERATIONS — PUT

- Translates to a causal-order (crb) broadcast request
- Use the algorithm from the class
- No need to write anything (\perp) in the output file



OPERATIONS — PUT

- Translates to a causal-order (crb) broadcast request
- Use the algorithm from the class
- No need to write anything (\perp) in the output file



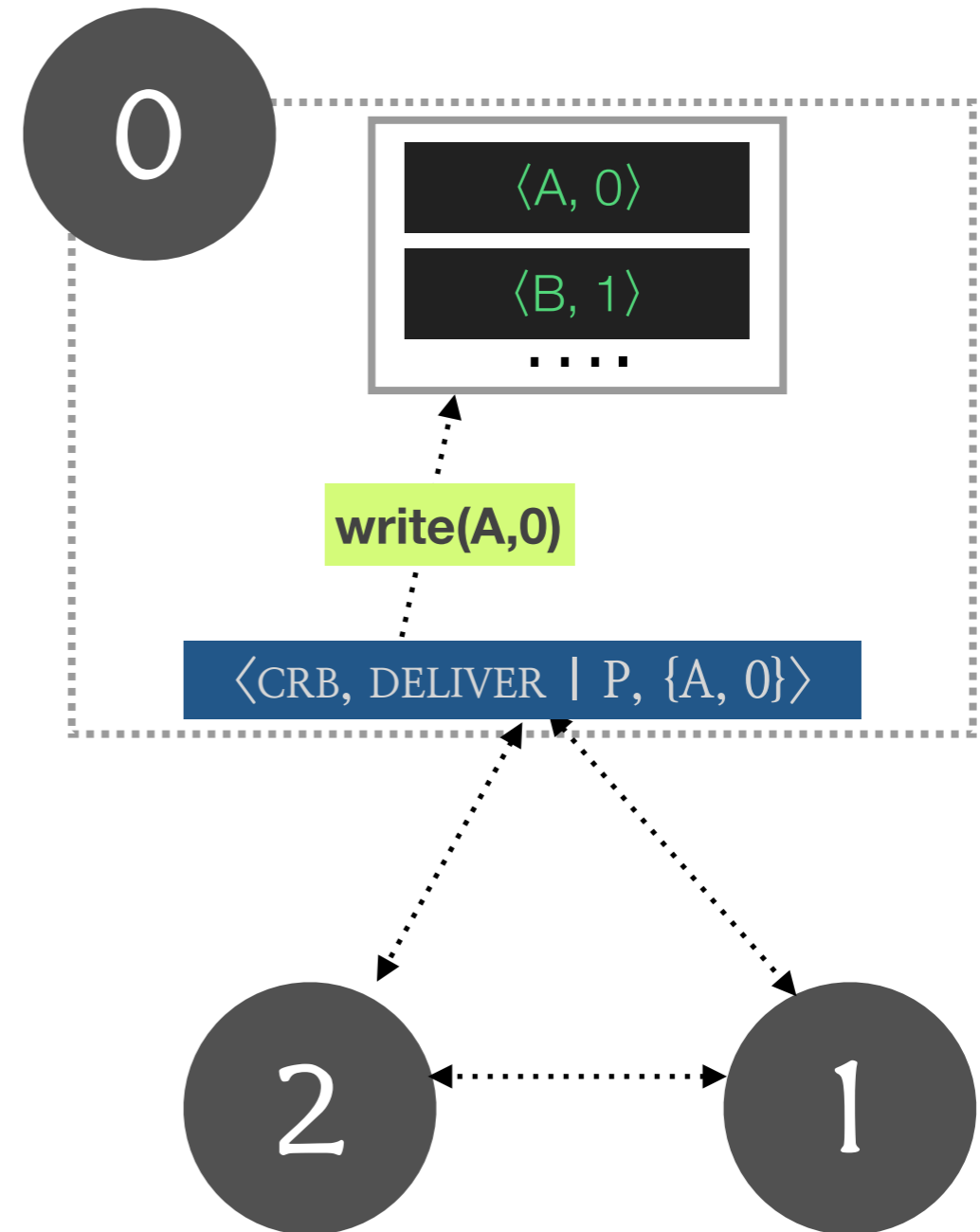
What happens when a node delivers a message?

<CRB, DELIVER | P, {A, 0}>

OPERATIONS

$\langle \text{CRB, DELIVER} \mid P, \{A, 0\} \rangle$

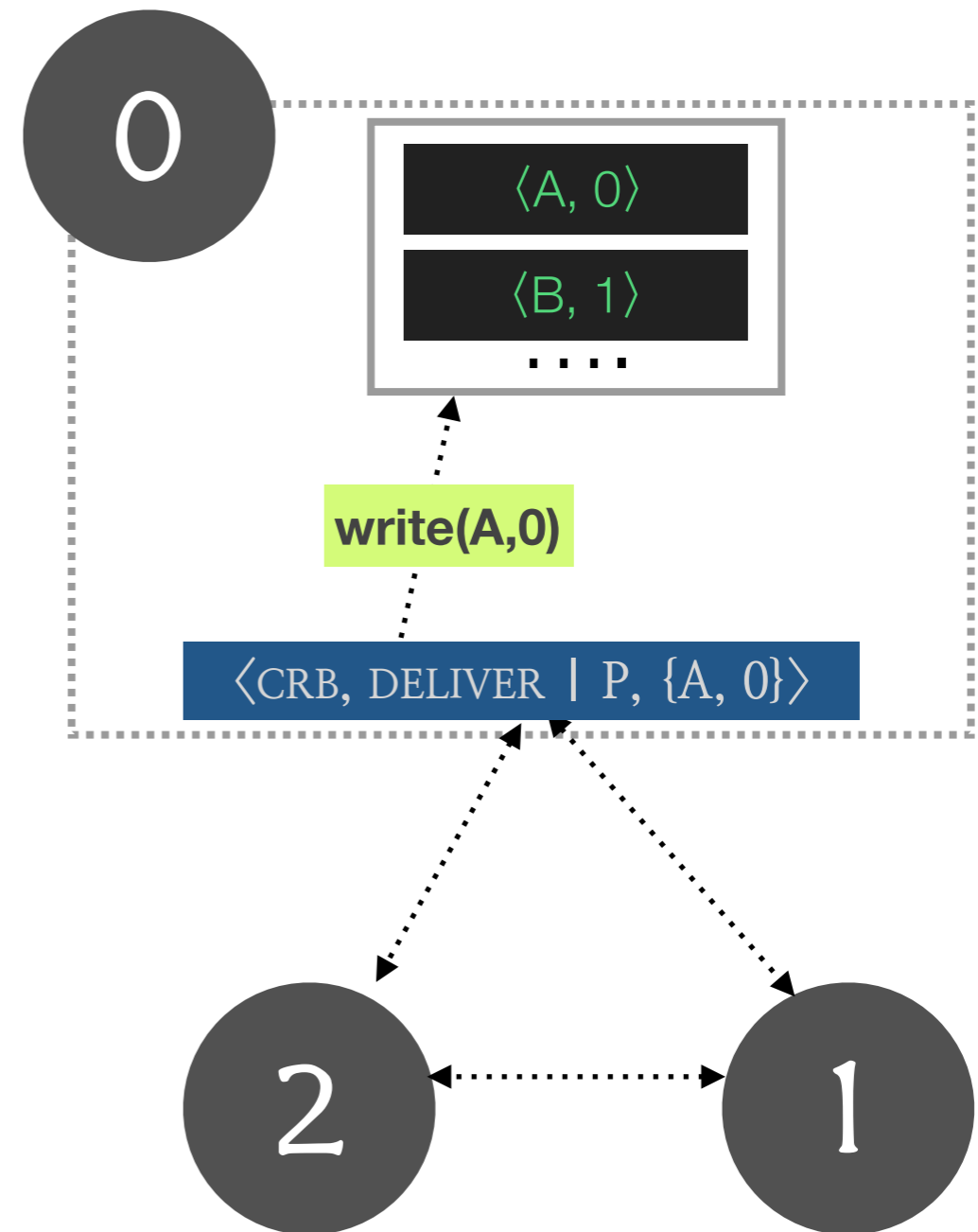
- Triggers an update in the storage engine



OPERATIONS

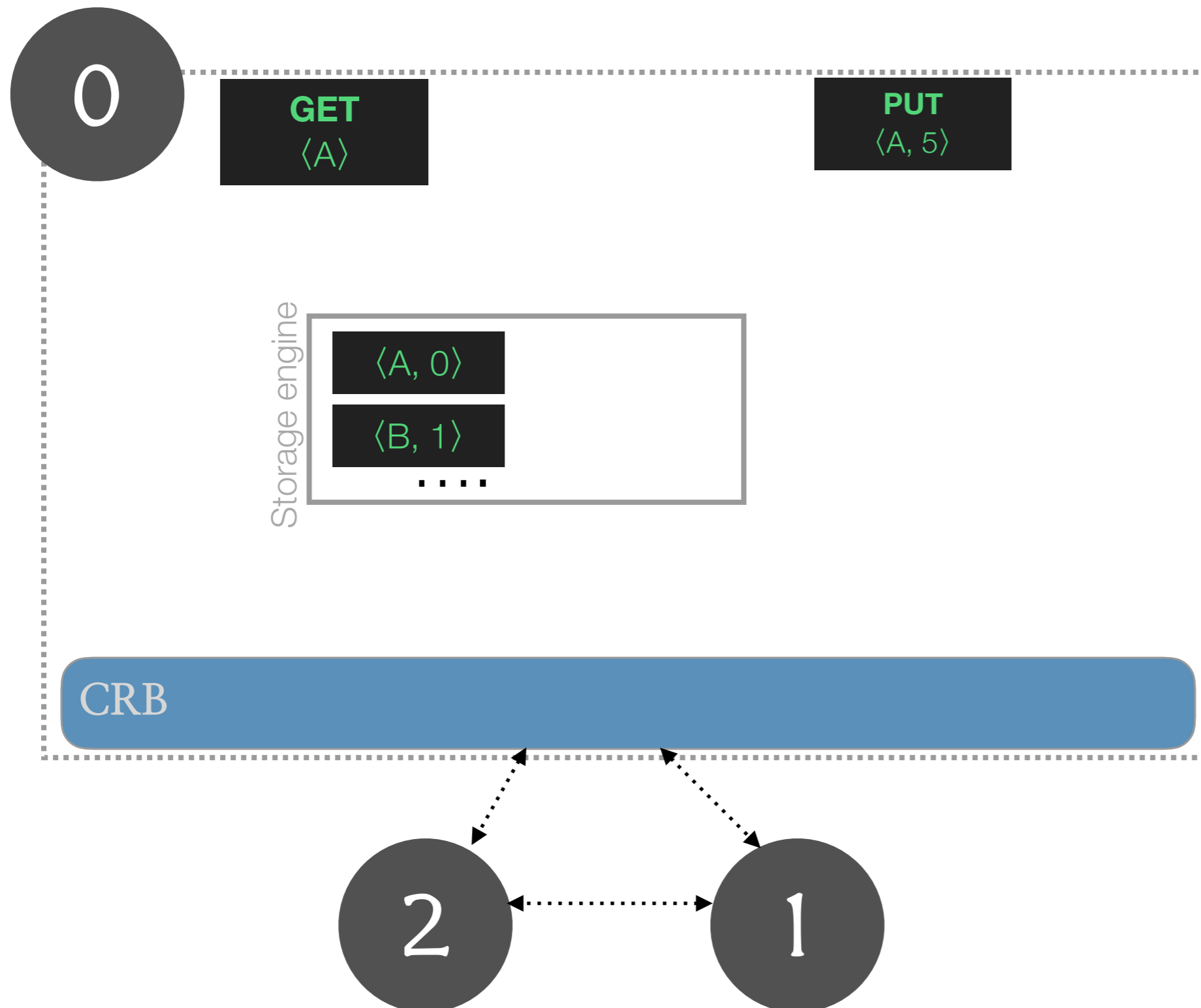
$\langle \text{CRB, DELIVER} \mid P, \{A, 0\} \rangle$

- Triggers an update in the storage engine

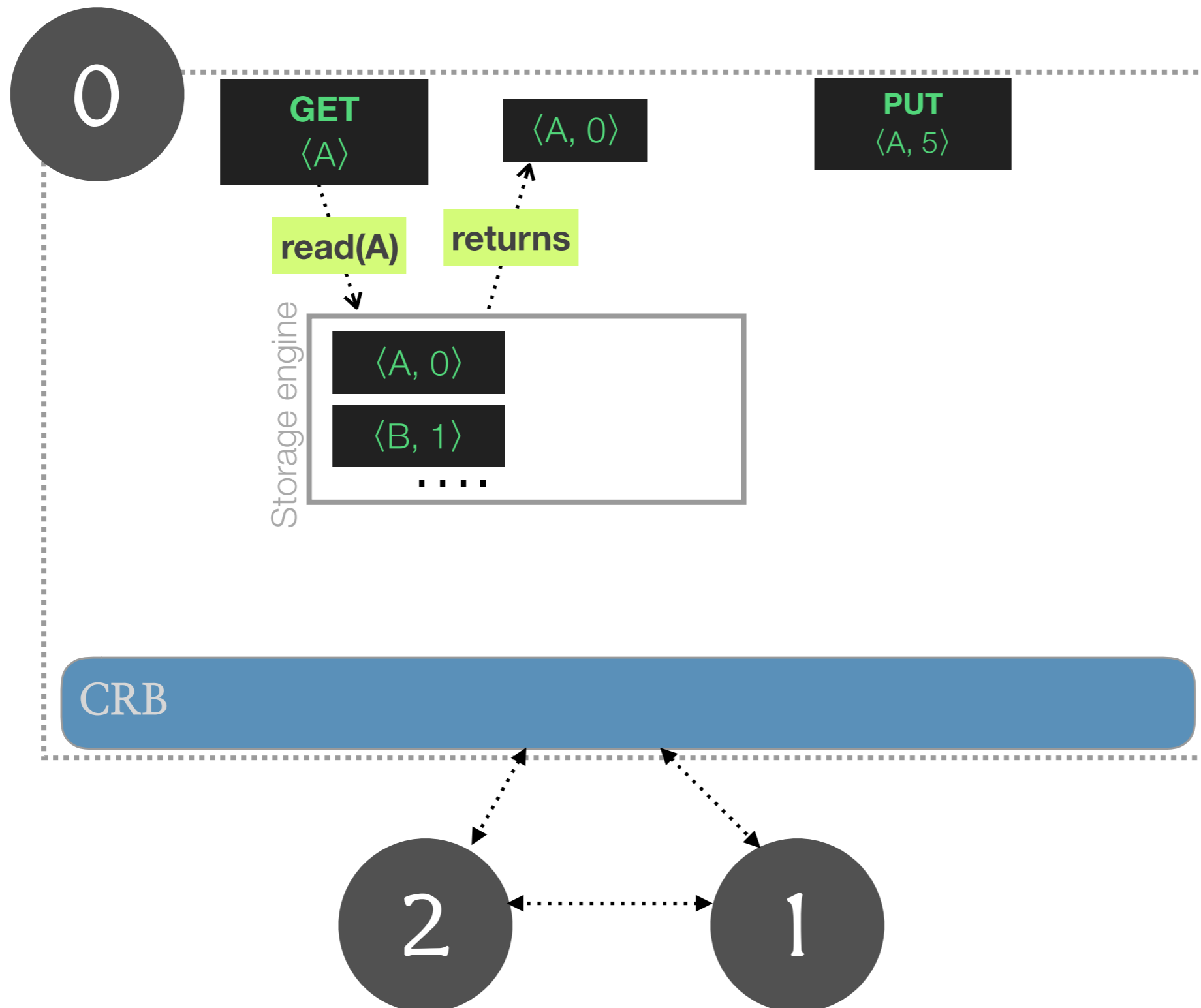


Now let's put all the pieces together...

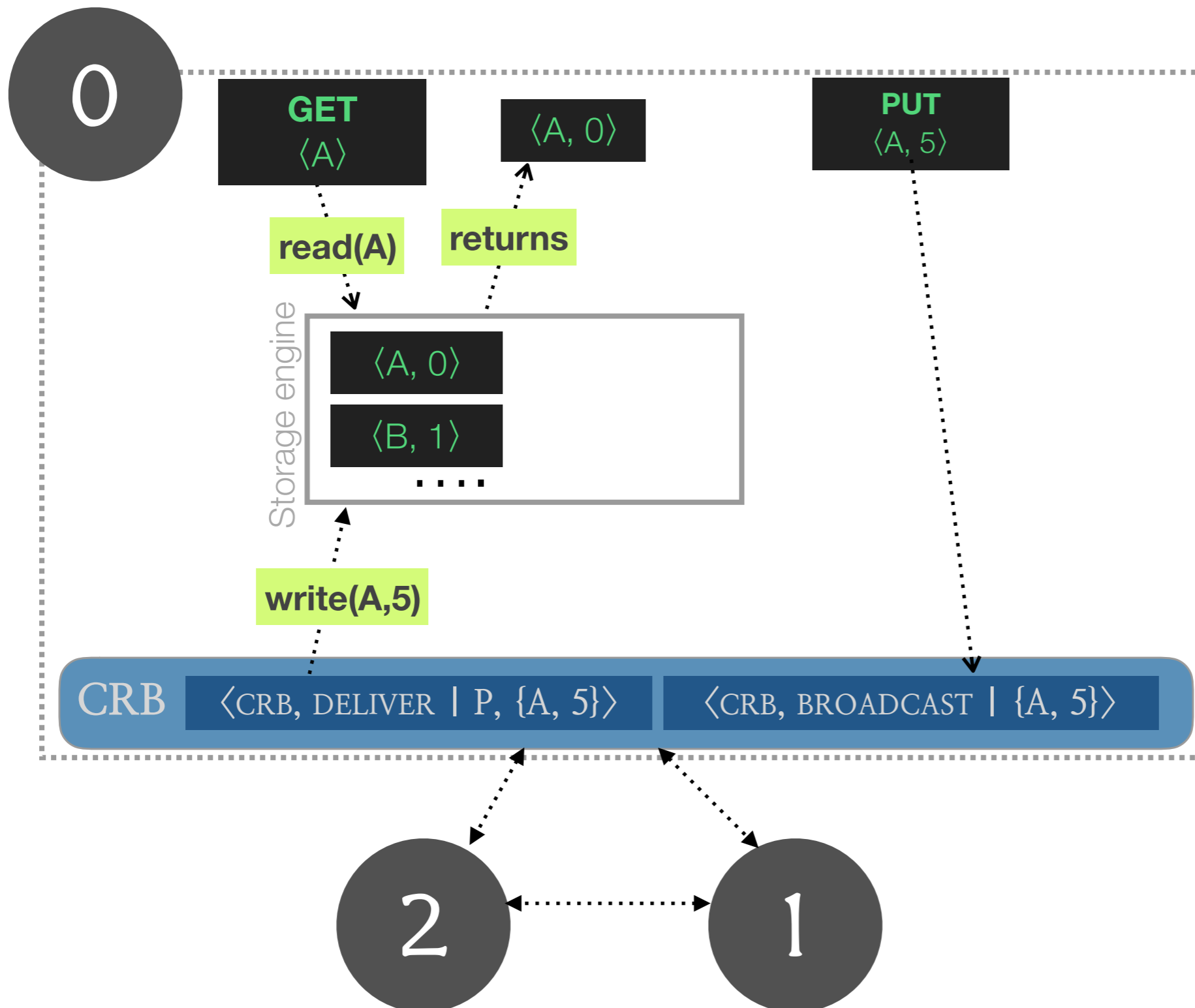
OVERVIEW OF A NODE



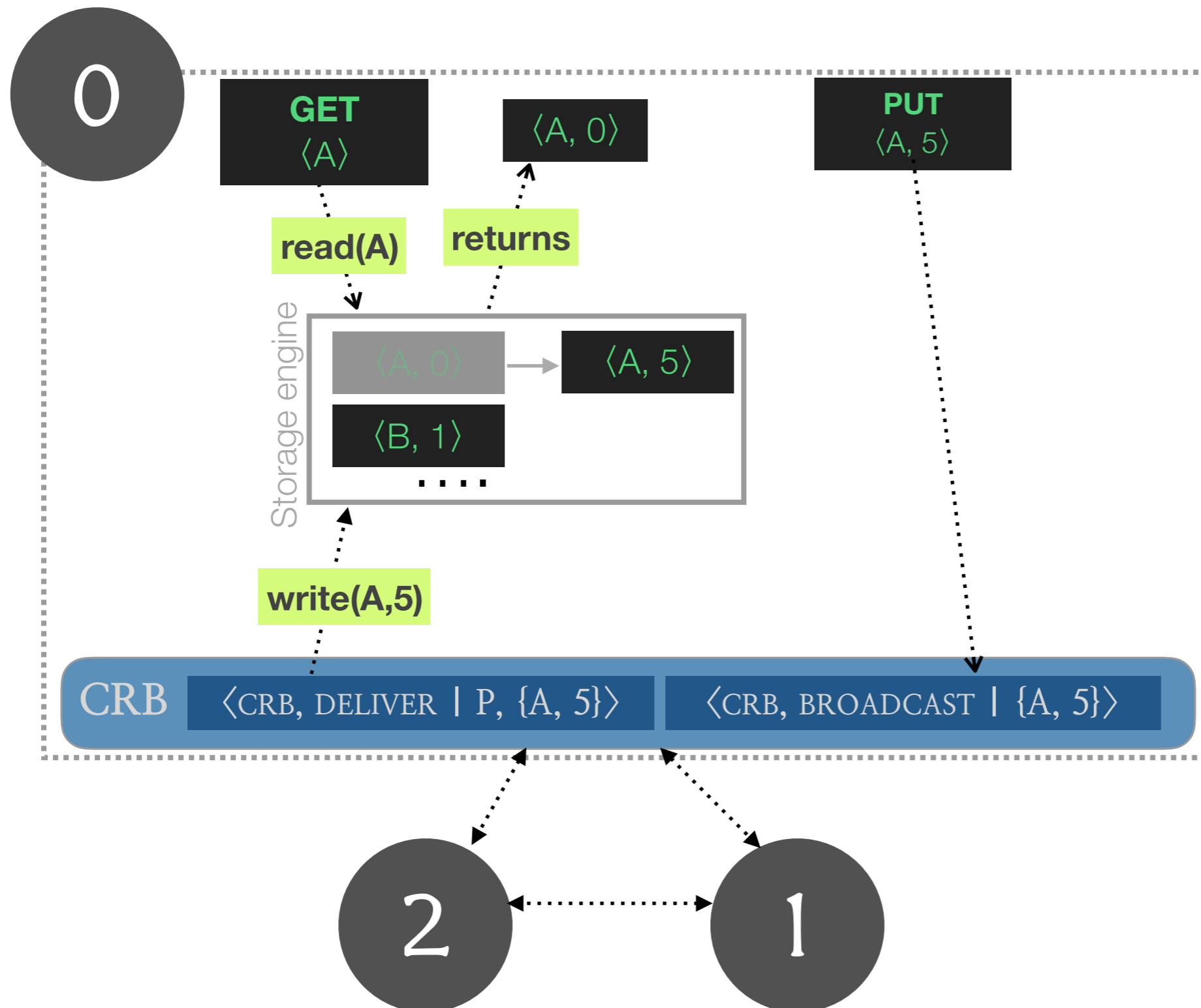
OVERVIEW OF A NODE



OVERVIEW OF A NODE



OVERVIEW OF A NODE



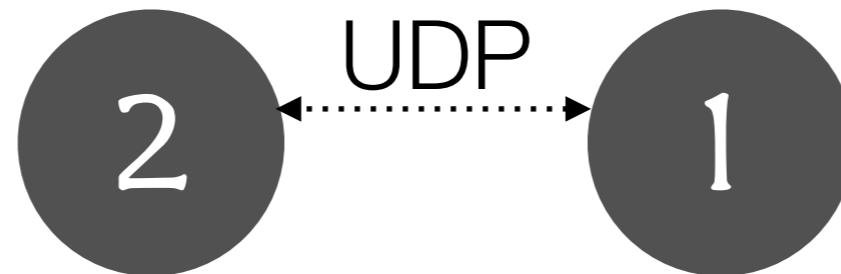
MESSAGES AND OBJECTS

- Use UDP

- Prohibited:

- TCP (or any other reliable communication protocol)

- Any other IPC



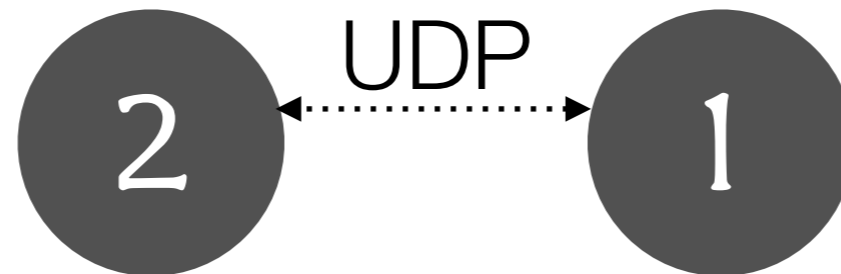
- For simplicity sake:

- Both keys and values are short ASCII strings

- At most 10 chars each

MESSAGES AND OBJECTS

- Use UDP
- Prohibited:
 - TCP (or any other reliable communication protocol)
 - Any other IPC



⟨Key, Value⟩

- For simplicity sake:
 - Both keys and values are short ASCII strings
 - At most 10 chars each

TESTING AND GRADING

```
put,A,1      /* write to object A value 1 */
put,B,2      /* write to object B value 2 */
get,A        /* read the value of object A */
get,B        /* read the value of object B */
```

<u>.input</u>	<u>.output</u>
put,A,1	⊥ (ignore in the output file)
put,B,2	⊥ (ignore in the output file)
get,A	A,1
get,B	B,2

The results of all the
GET operations

TESTING AND GRADING

For each line in **input file**:

1. execute request
2. write to **output file**

INPUT

```
put,A,1      /* write to object A value 1 */
put,B,2      /* write to object B value 2 */
get,A        /* read the value of object A */
get,B        /* read the value of object B */
```

INPUT + OUTPUT

.input	.output
put,A,1	⊥ (ignore in the output file)
put,B,2	⊥ (ignore in the output file)
get,A	A,1
get,B	B,2

The results of all the
GET operations

S4

Overview:

1. Setup & System model

2. Operations

3. Goals

4. Technicalities

- Java template
- Compilation
- Testing
 - 1. Correctness
 - 2. Performance

5. One last look at causal consistency



S4

Overview:

1. Setup & System model

2. Operations

3. Goals

4. Technicalities

- Java template
- Compilation
- Testing
 - 1. Correctness
 - 2. Performance



**NOT COVERED
SEE PROJECT DOCUMENTATION**

5. One last look at causal consistency

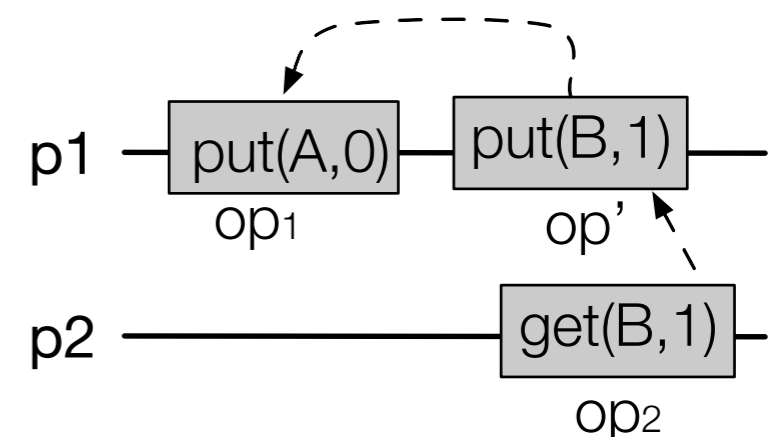
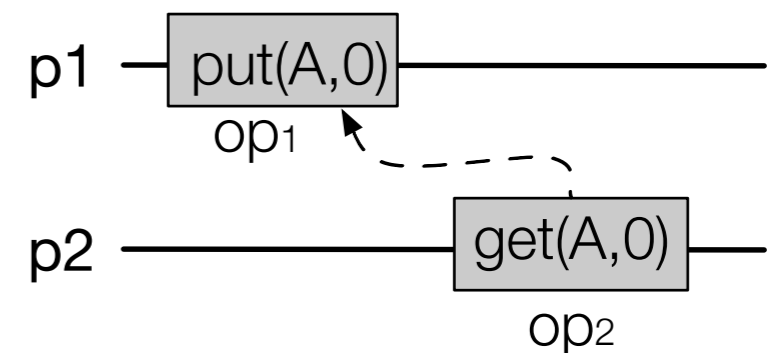
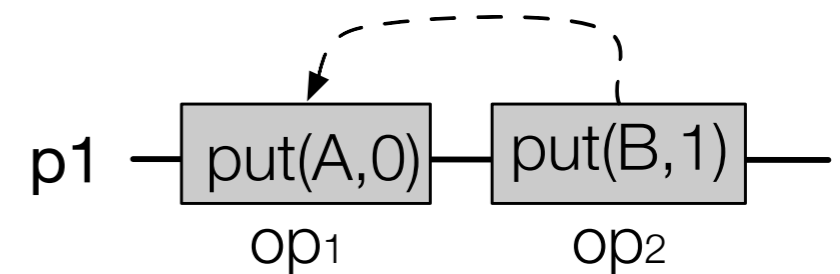
Causal consistency — one last look

- Given two operations, op_1 and op_2
- $op_1 \rightarrow op_2$ (causally precedes)
if any of the following three cases applies:

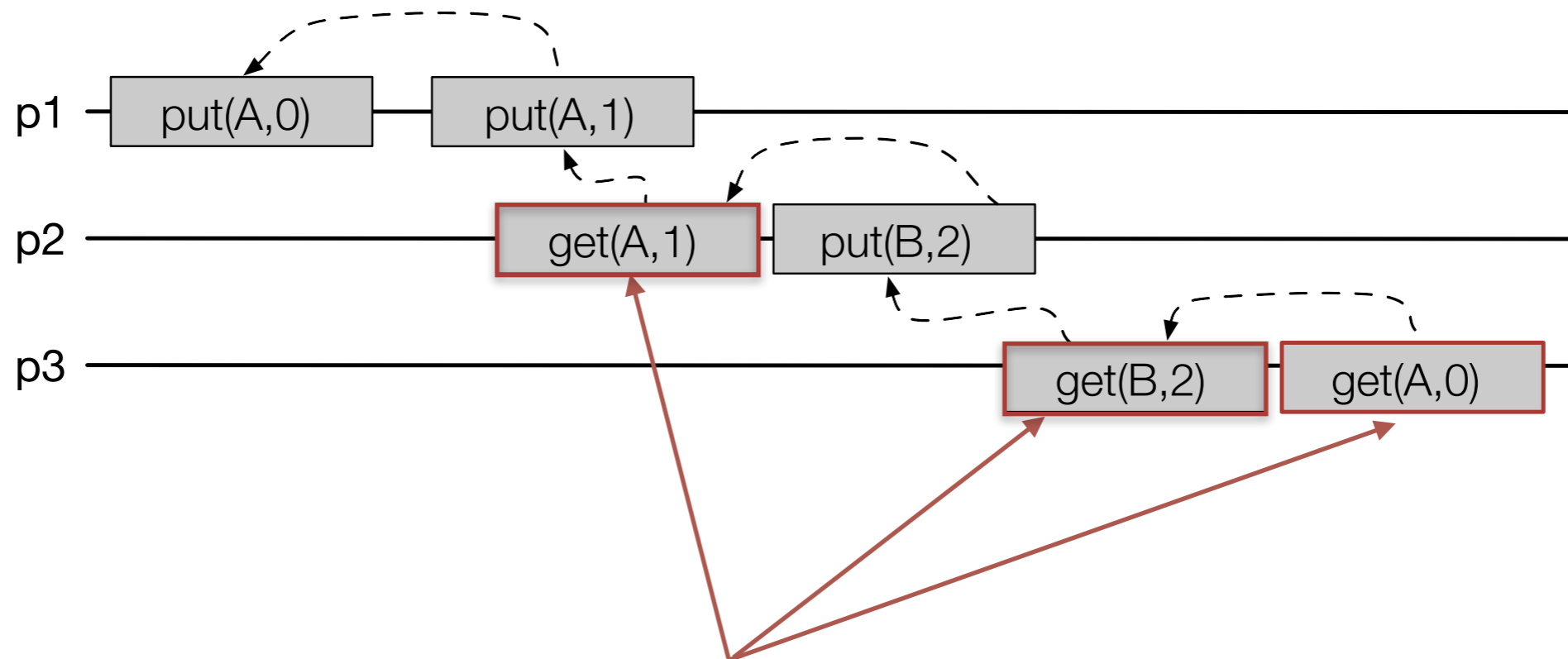
(a) FIFO: A process invokes op_1 and then invokes op_2

(b) Local: A process invokes the PUT operation op_1 and another process invokes the GET operation op_2 , where op_2 observes the written value of op_1

(c) Transitivity: There exists an intermediate operation op' such that $op_1 \rightarrow op'$ and $op' \rightarrow op_2$.

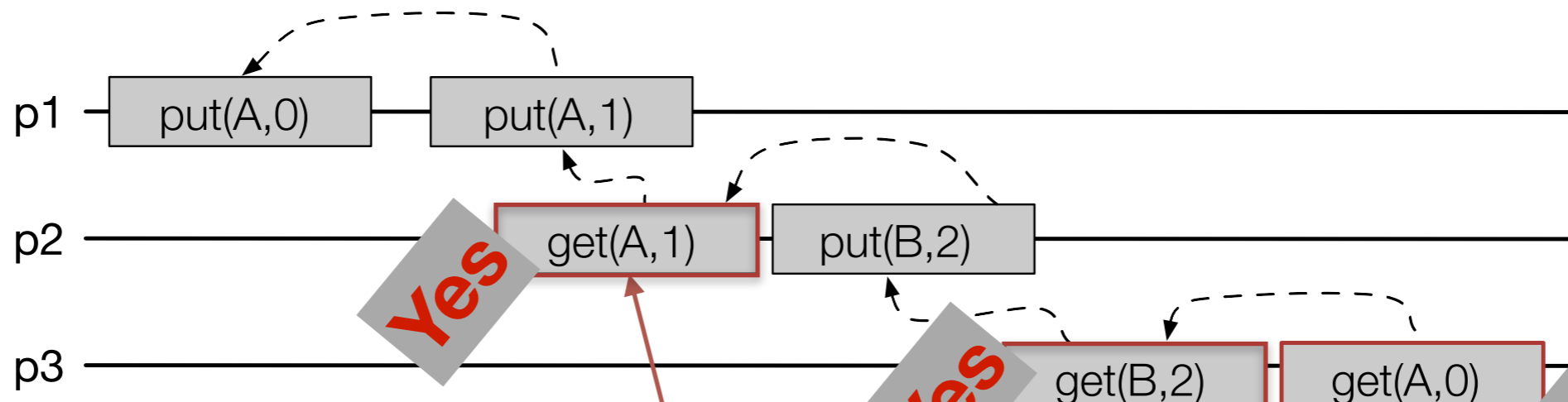


Causal consistency — one last look



Does it respect causality?

Causal consistency — one last look



Does it respect causality?

Yes

Yes

NO:
It should return (A,1)

1. **Bitcoin** – not needed for exam

- Nakamoto, Satoshi (24 May 2009).
"Bitcoin: A Peer-to-Peer Electronic Cash System".
<https://bitcoin.org/bitcoin.pdf>
(canonical Bitcoin)



Causal Broadcast

2. **Social networking (= bonus project)**

- CS-451 Bonus Project (2015).
"A Storage System for Social Networks".
<https://github.com/LPD-EPFL/da15-s4>
- Ahamad, M., Neiger, G., Burns, J. E., Kohli, P., & Hutto, P. W. (1995). **Causal memory: definitions, implementation, and programming. (§5)**. Distributed Computing, 9(1).

