

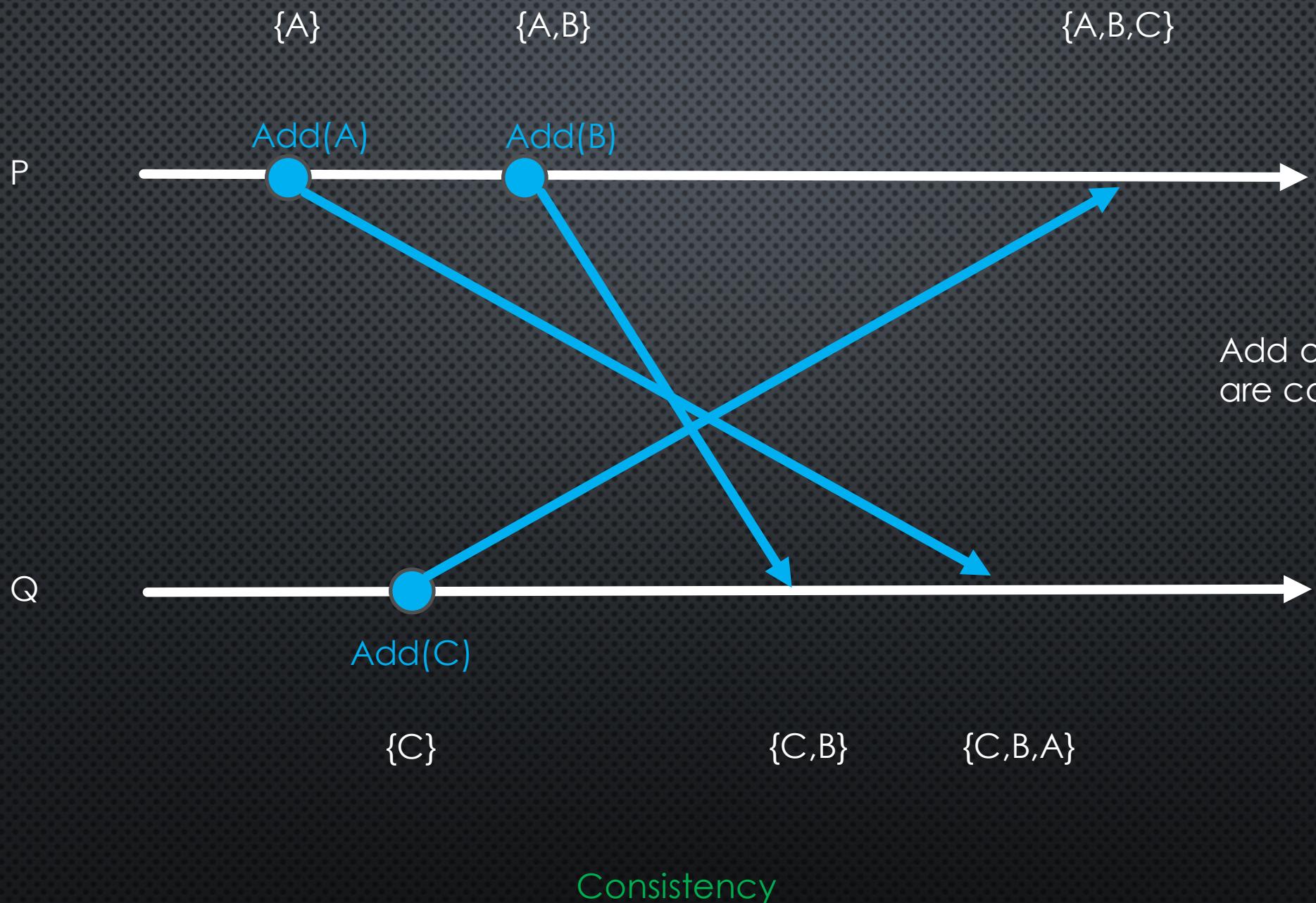
THE PITFALLS OF ACHIEVING TAGGED CAUSAL DELIVERY

GEORGES YOUNES

PAPOC'18
APRIL 23, 2018

USE CASE 1: REPLICATED GSET

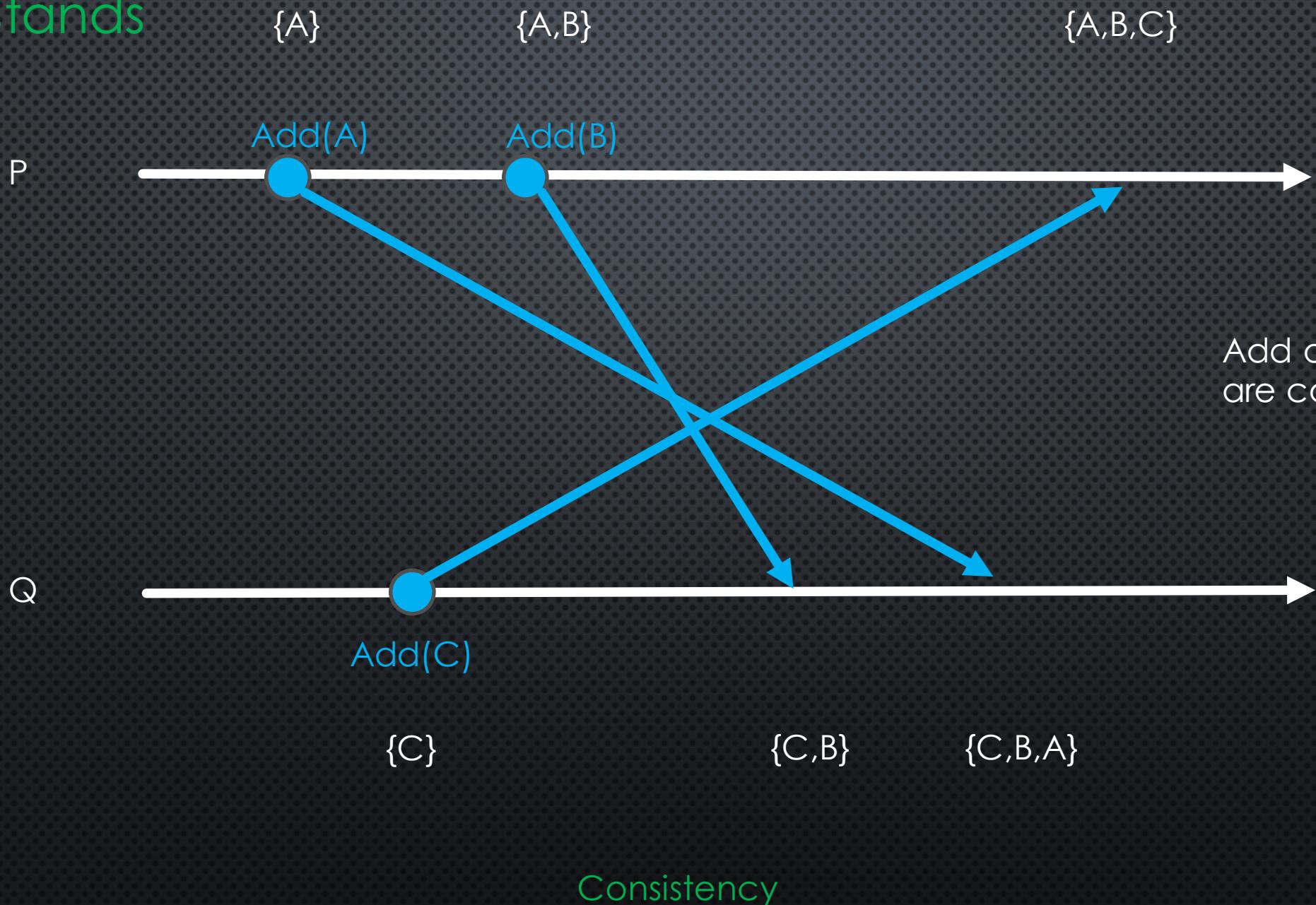
- A GROW-ONLY SET
- API:
 - ADD(ELEMENT): ADDS THE ELEMENT TO THE SET
 - QUERY(): RETURNS ALL THE ELEMENTS IN THE SET

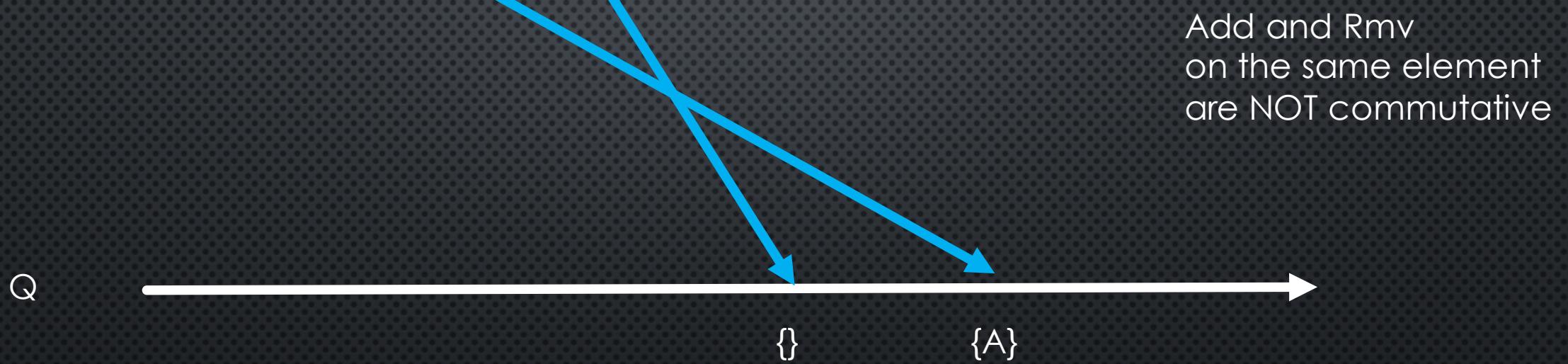


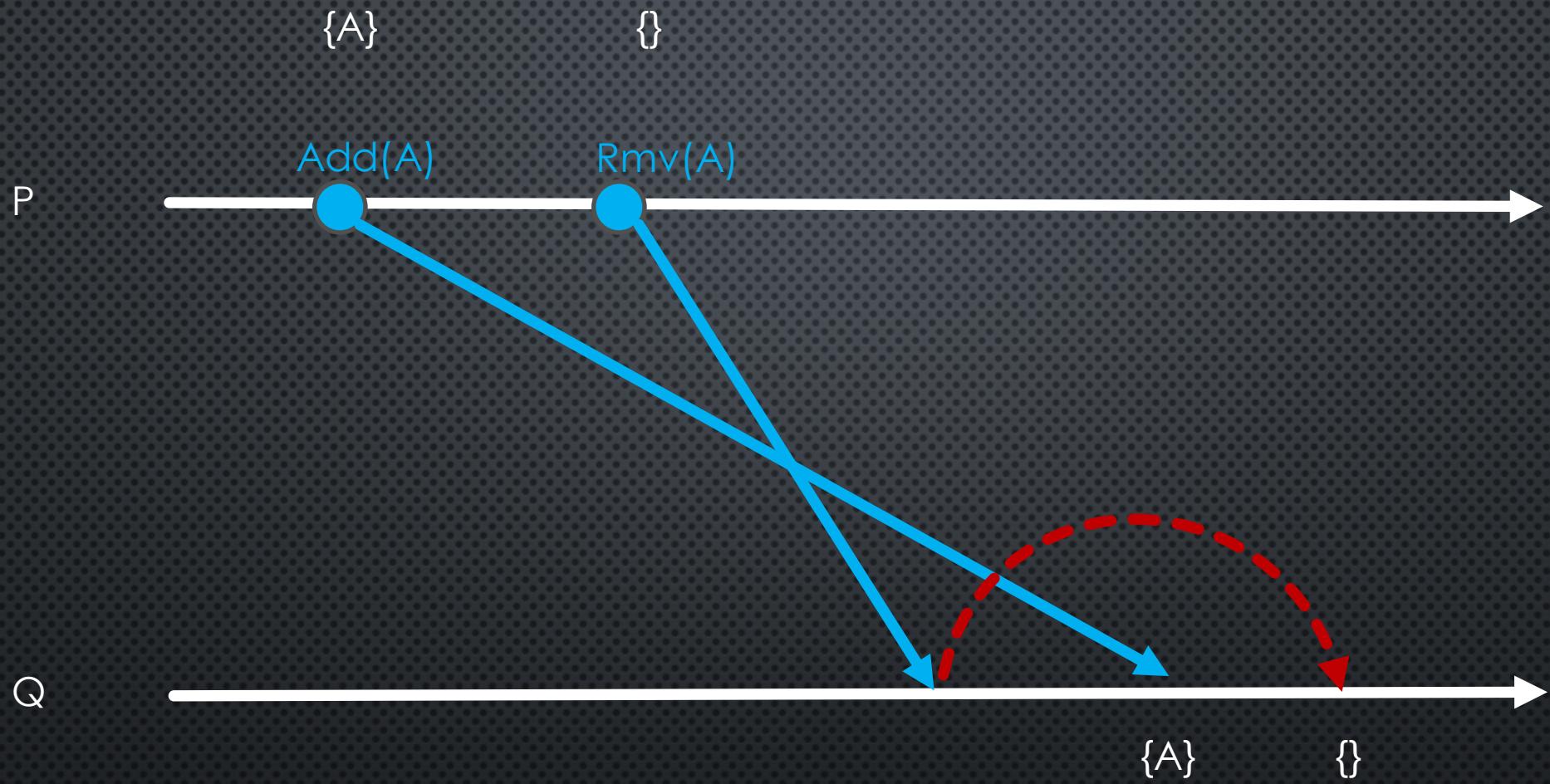
USE CASE 2: REPLICATED SET

- **Not** GROW-ONLY SET
- API:
 - ADD(ELEMENT): ADDS THE ELEMENT TO THE SET
 - RMV(ELEMENT): REMOVES THE ELEMENT FROM THE SET
 - QUERY(): RETURNS ALL THE ELEMENTS IN THE SET

Still Stands

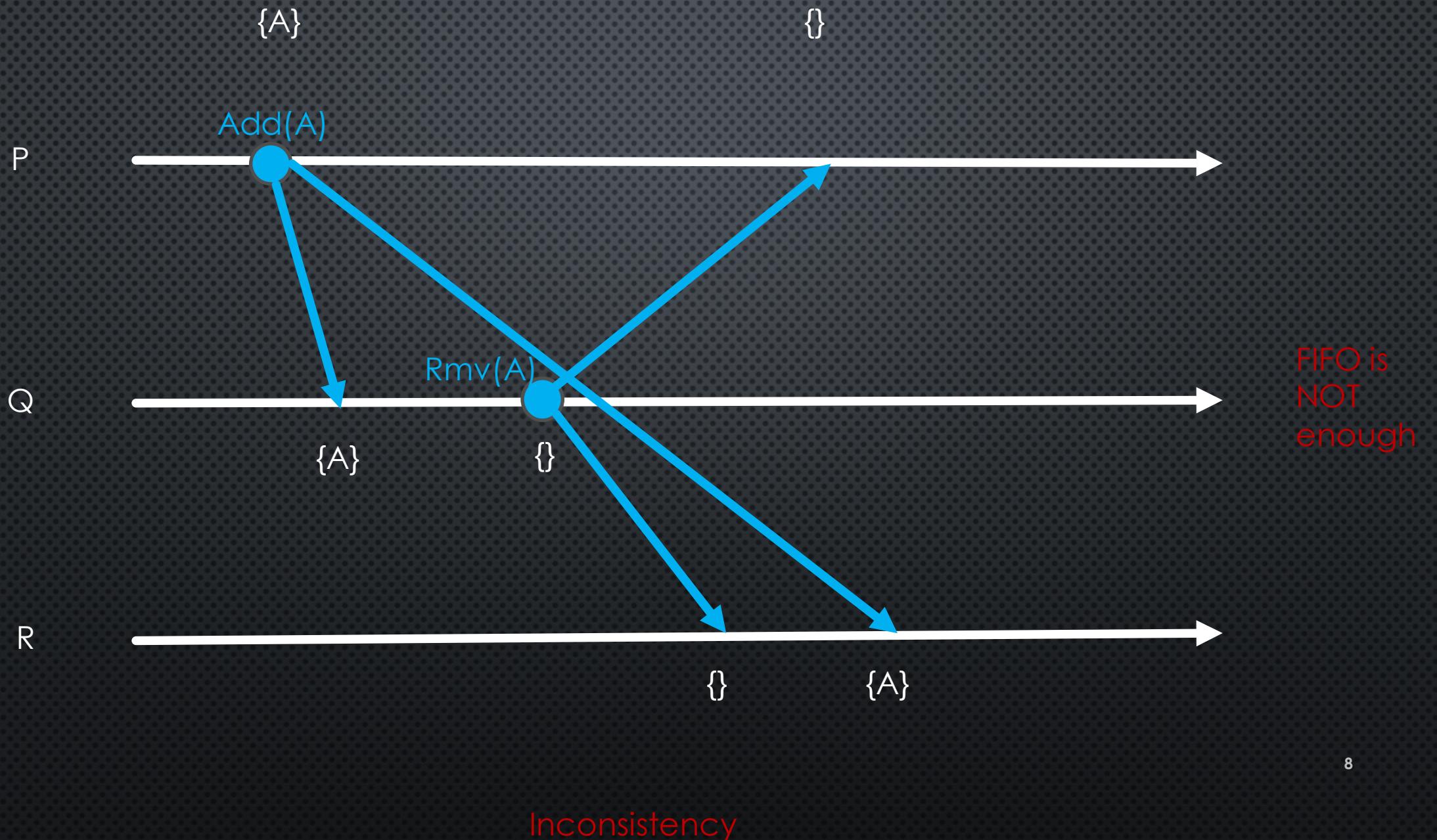


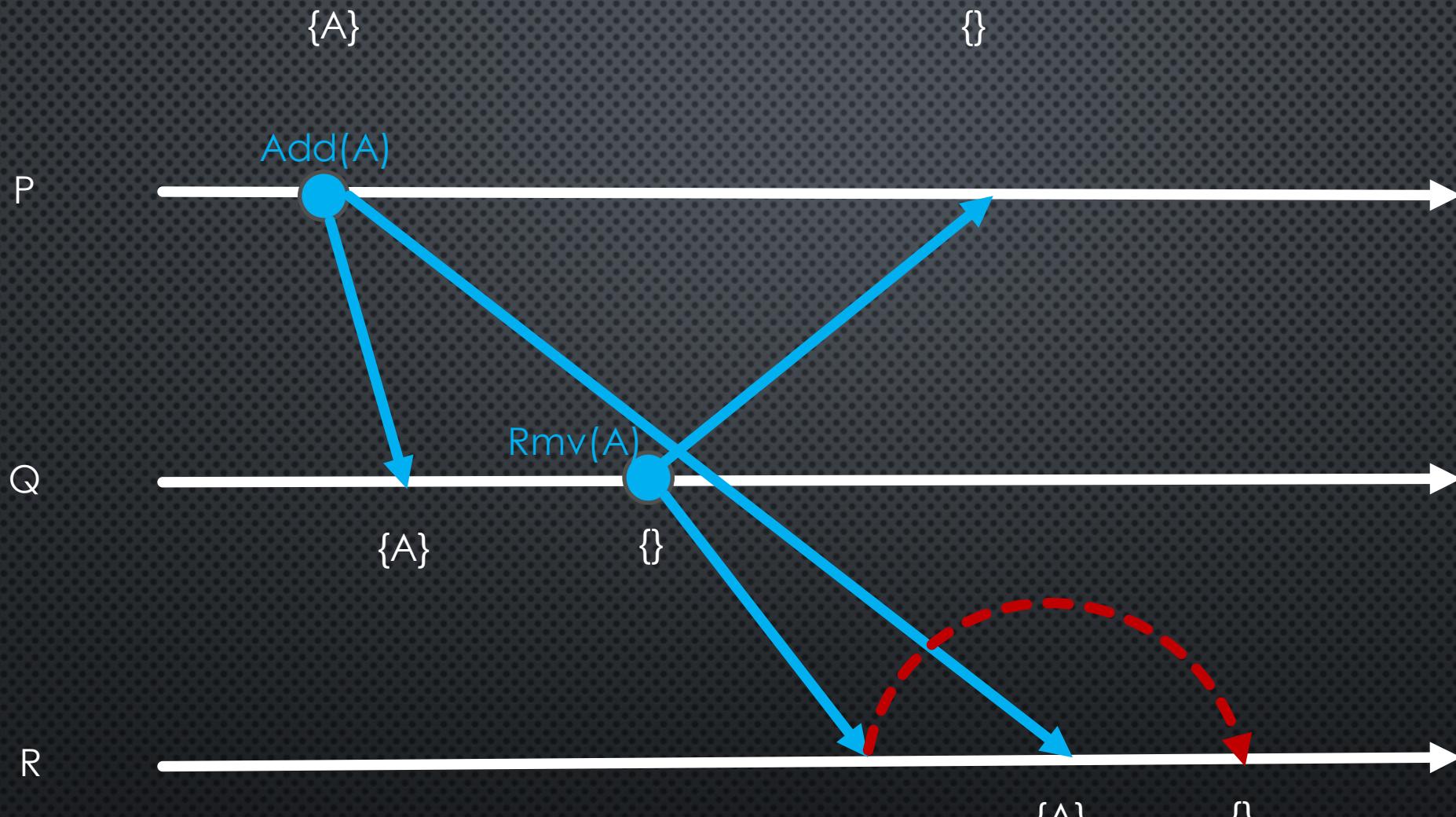


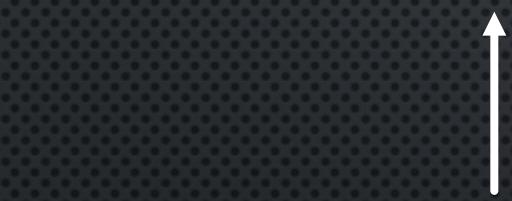
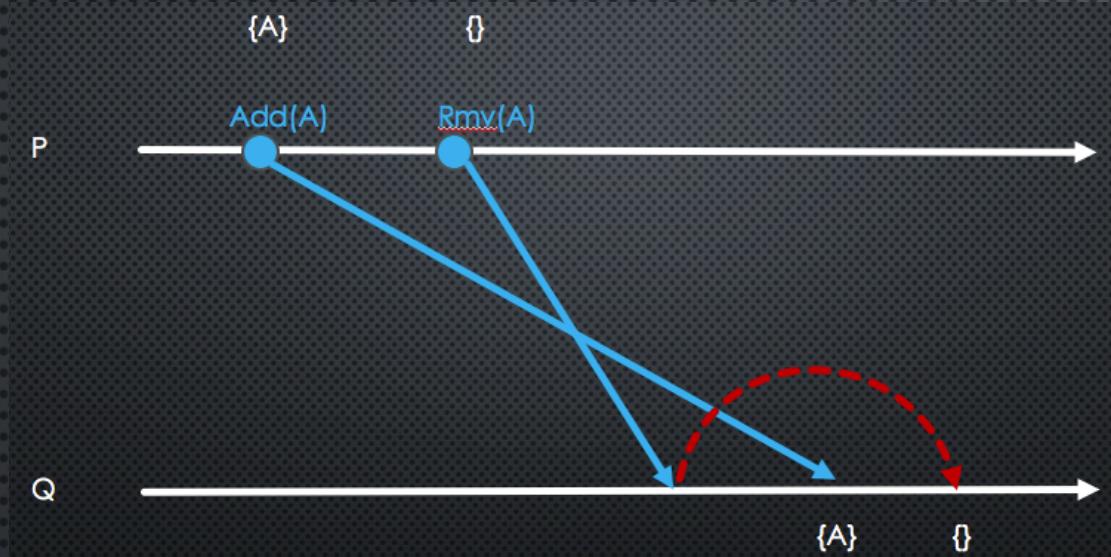


FIFO ordering abstraction

Consistency

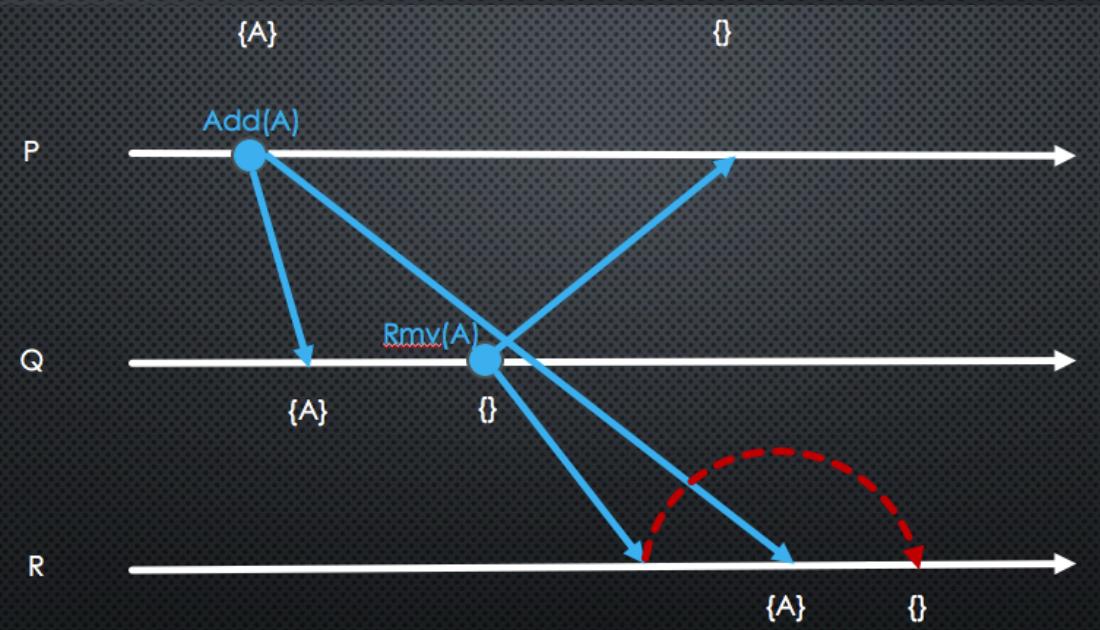






Forcing delivery of operations
based on the order of
their emission at the sender

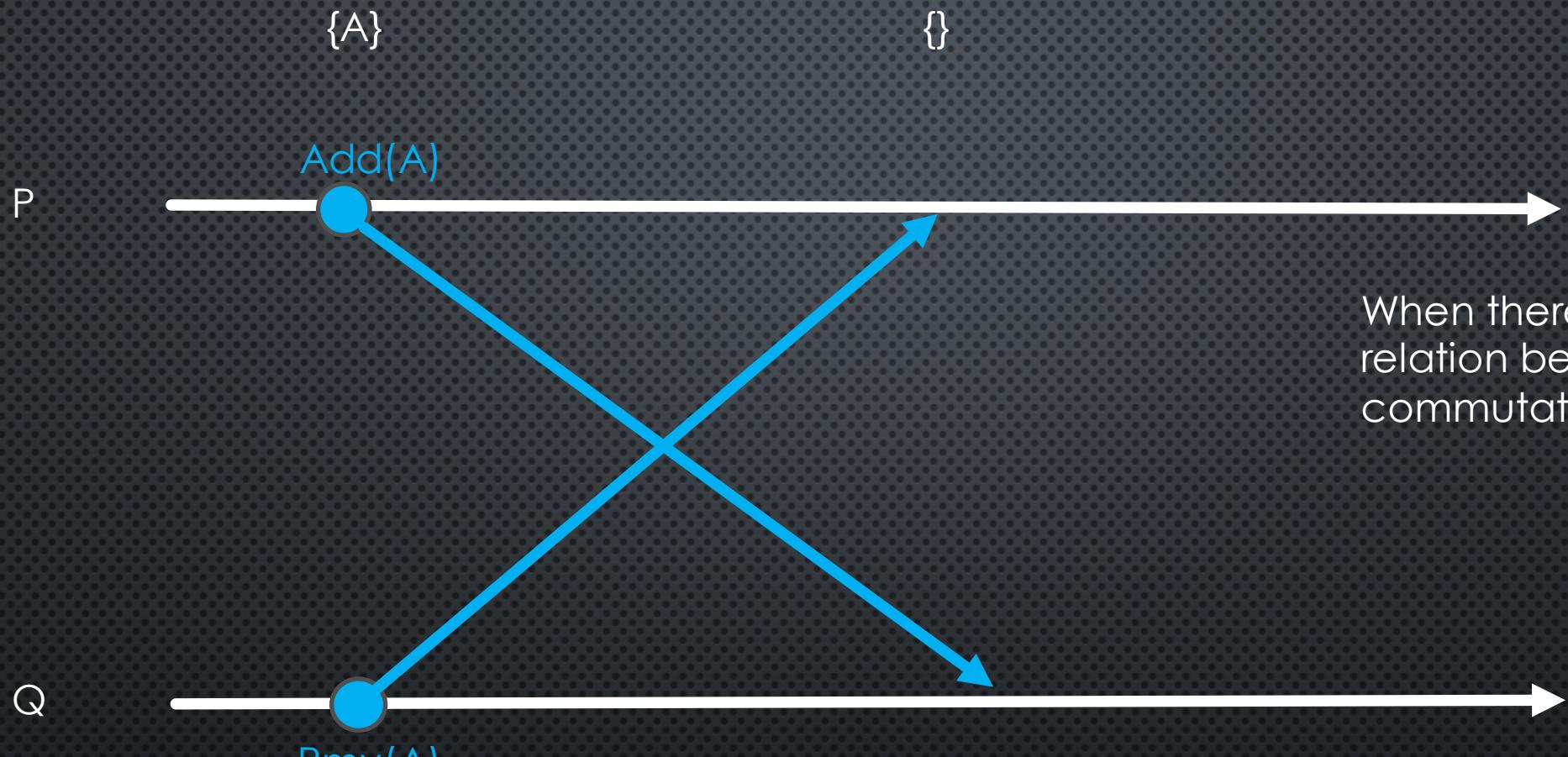
P.S. From the same sender



P.S. From different
senders i.e. covers
FIFO

Forcing delivery of operations
based on the causality
relation between their
emission

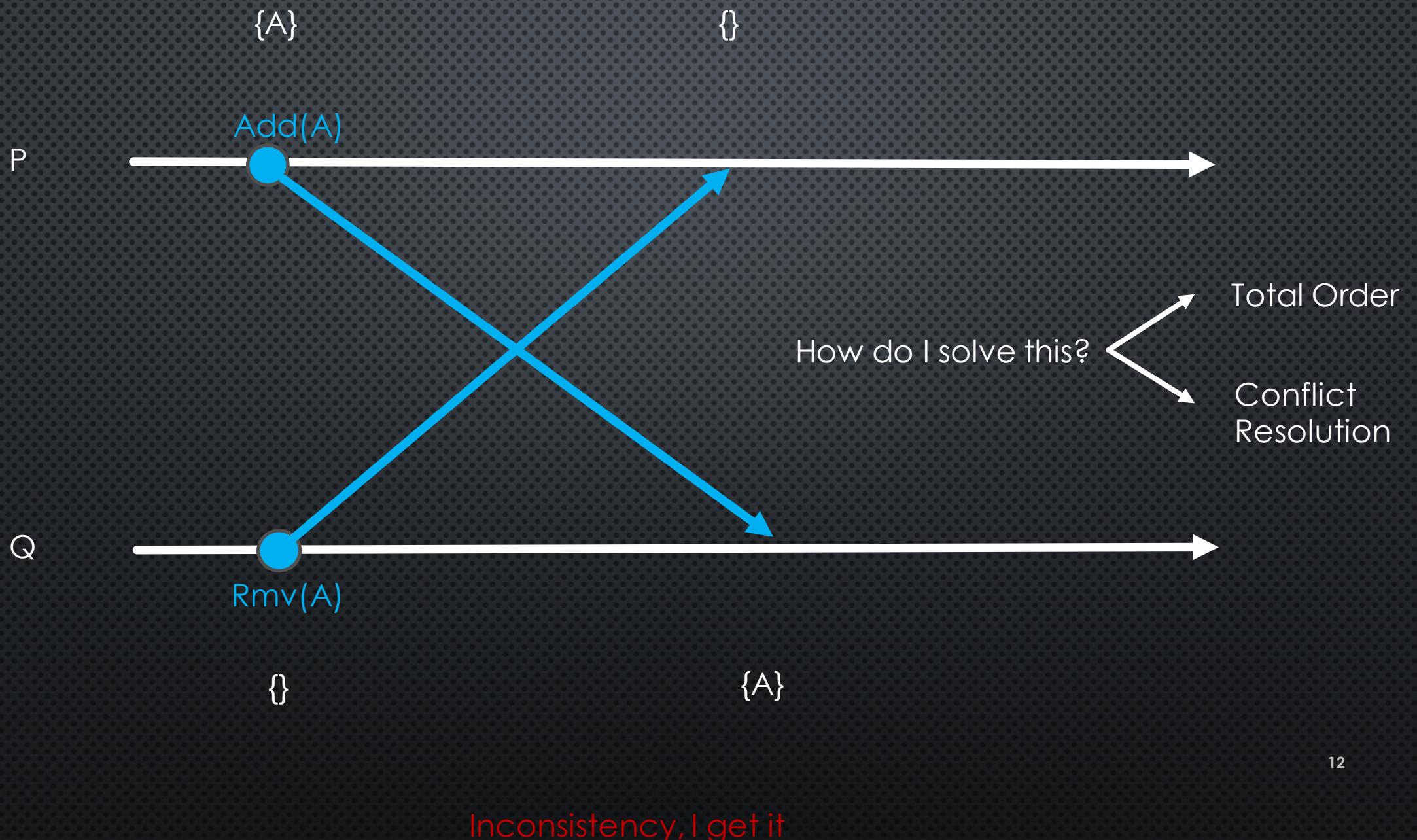




When there is no causality relation between non-commutative operations

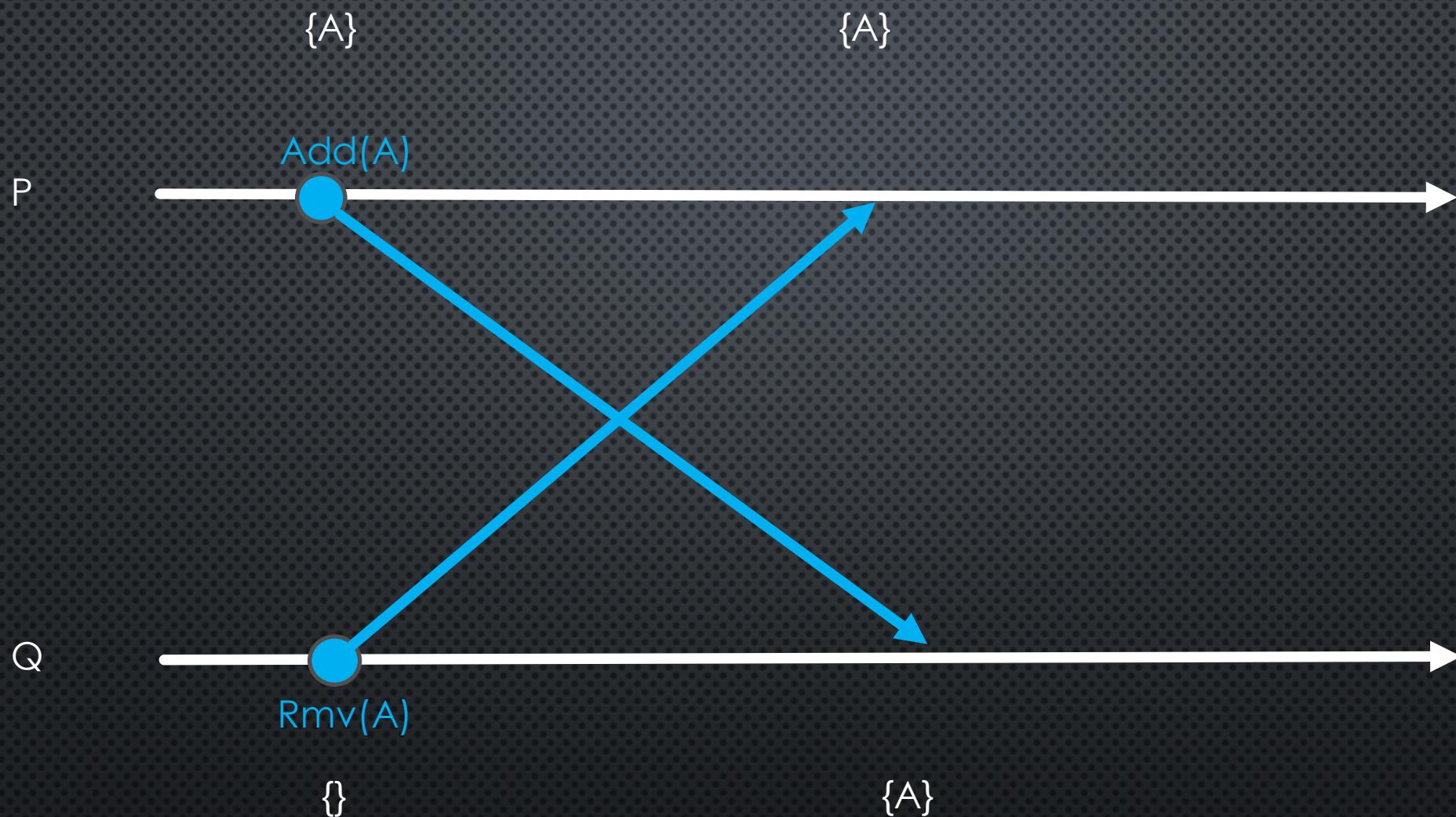
Causal order
Is NOT
enough

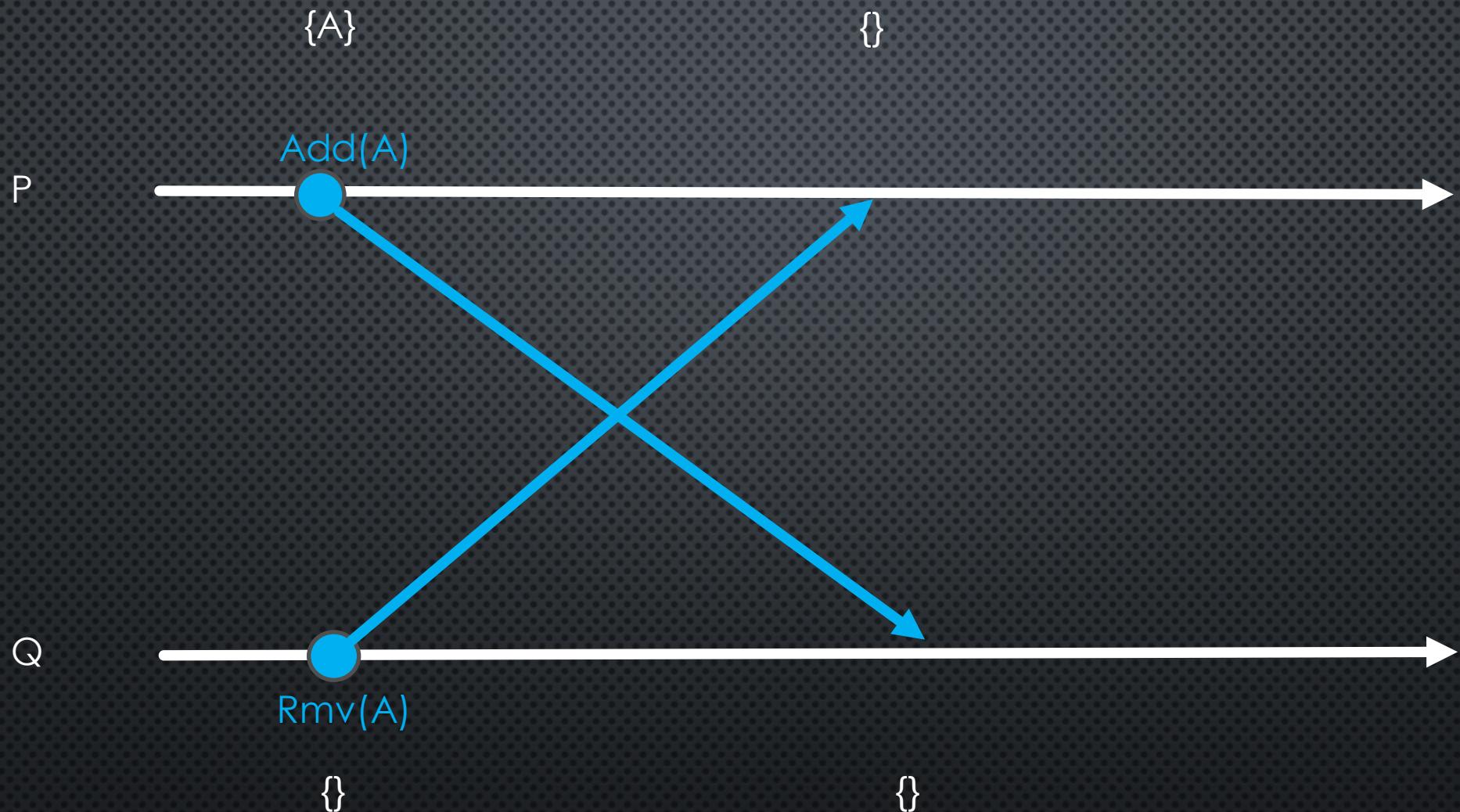
Inconsistency



WHY PICK CRDT?

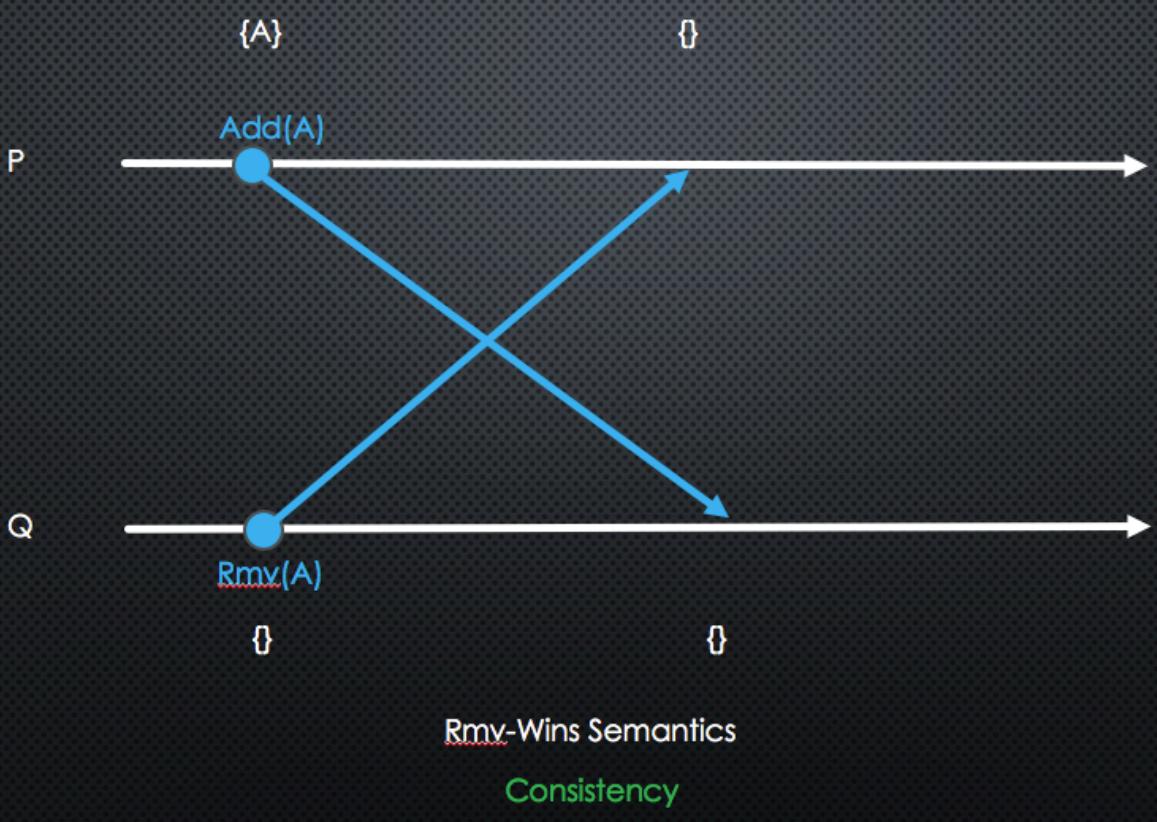
- YOU WANT AVAILABILITY AND HIGH RESPONSIVENESS
- YOU DO NOT MIND EVENTUAL CONSISTENCY
- YOU DO NOT WANT TO RESOLVE CONFLICTS MANUALLY
- YOUR SUPERVISOR IS CARLOS BAQUERO





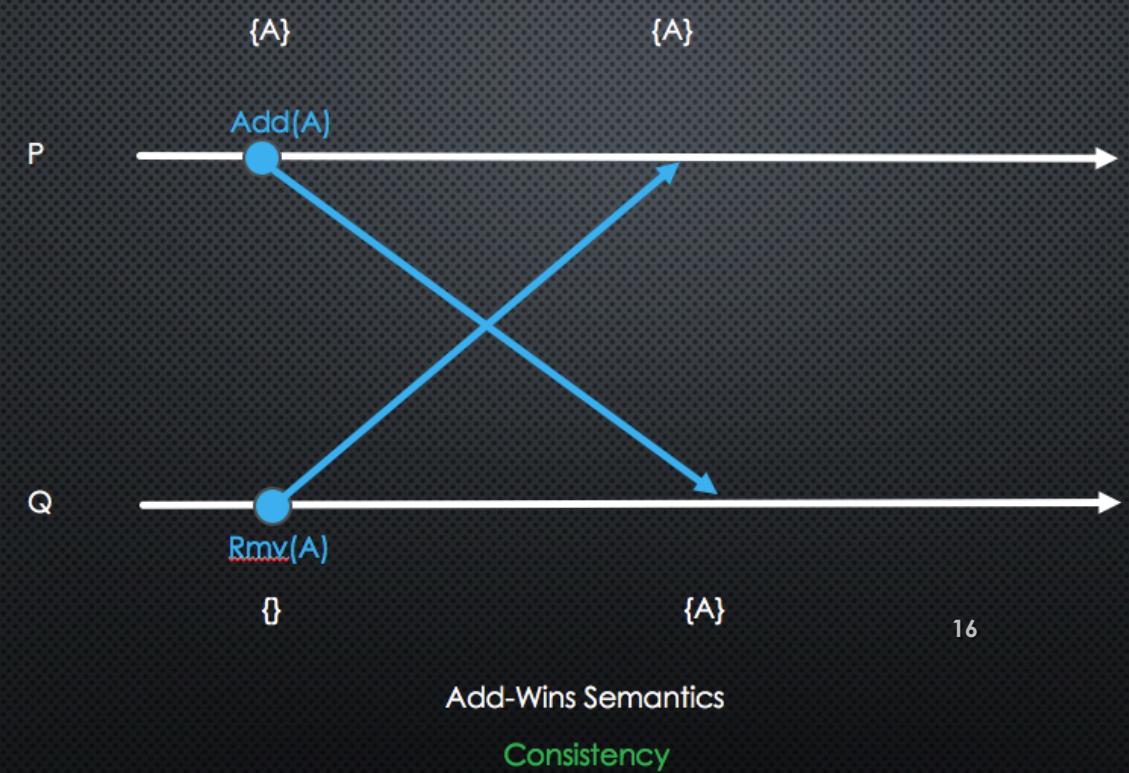
Rmv-Wins Semantics

Consistency



The “happened-before”
relation should be exposed
to the application layer

The decision is based on
The “happened-before”
Relation between concurrent
Add(A) and **Rmv(A)**



SOLUTION 1.0

- WHAT YOU NEED:

- CAUSAL ORDER



Off-the-shelf middleware

- CONFLICT RESOLUTION FOR CONCURRENT OPERATIONS



The “happened-before” relation between ops



SOLUTION 1.1

- WHAT YOU NEED:

- CAUSAL ORDER



- CONFLICT RESOLUTION FOR CONCURRENT OPERATIONS



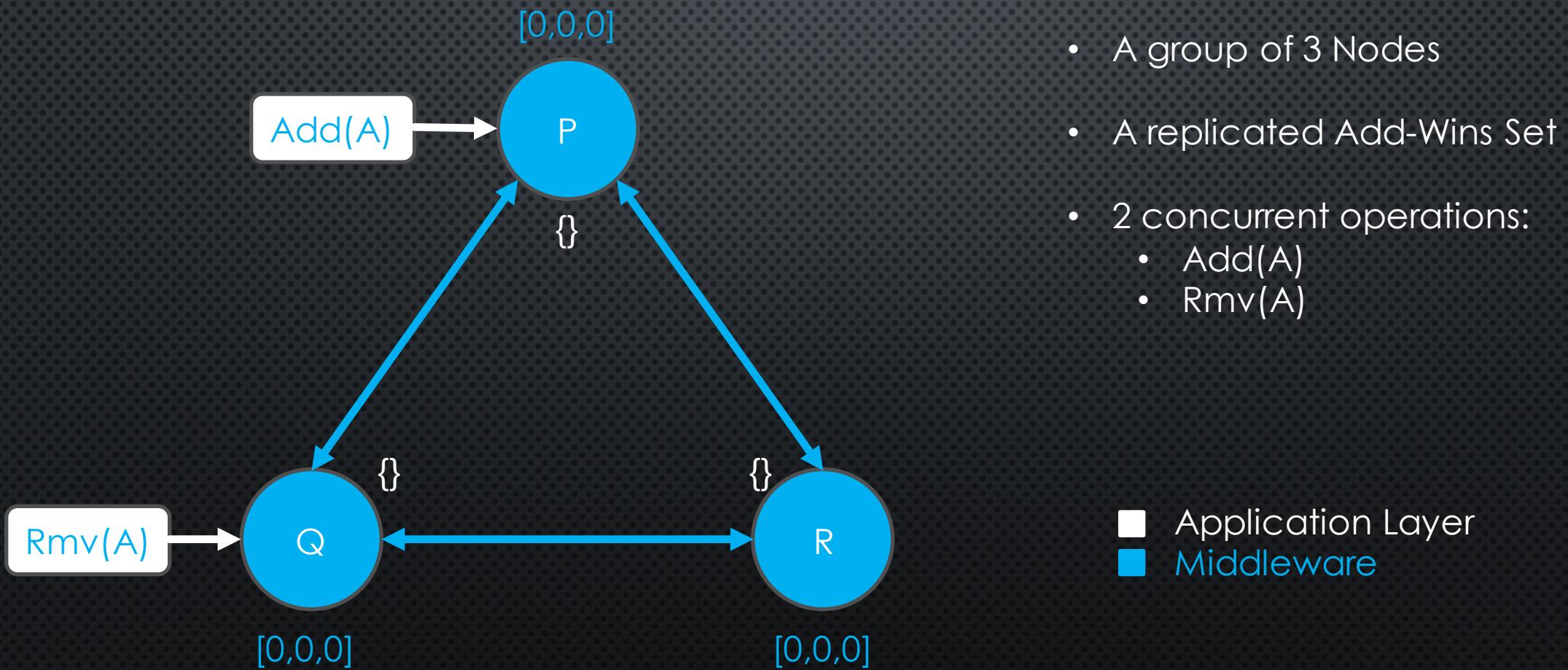
Modify Off-the-shelf middleware to expose

Off-the-shelf middleware

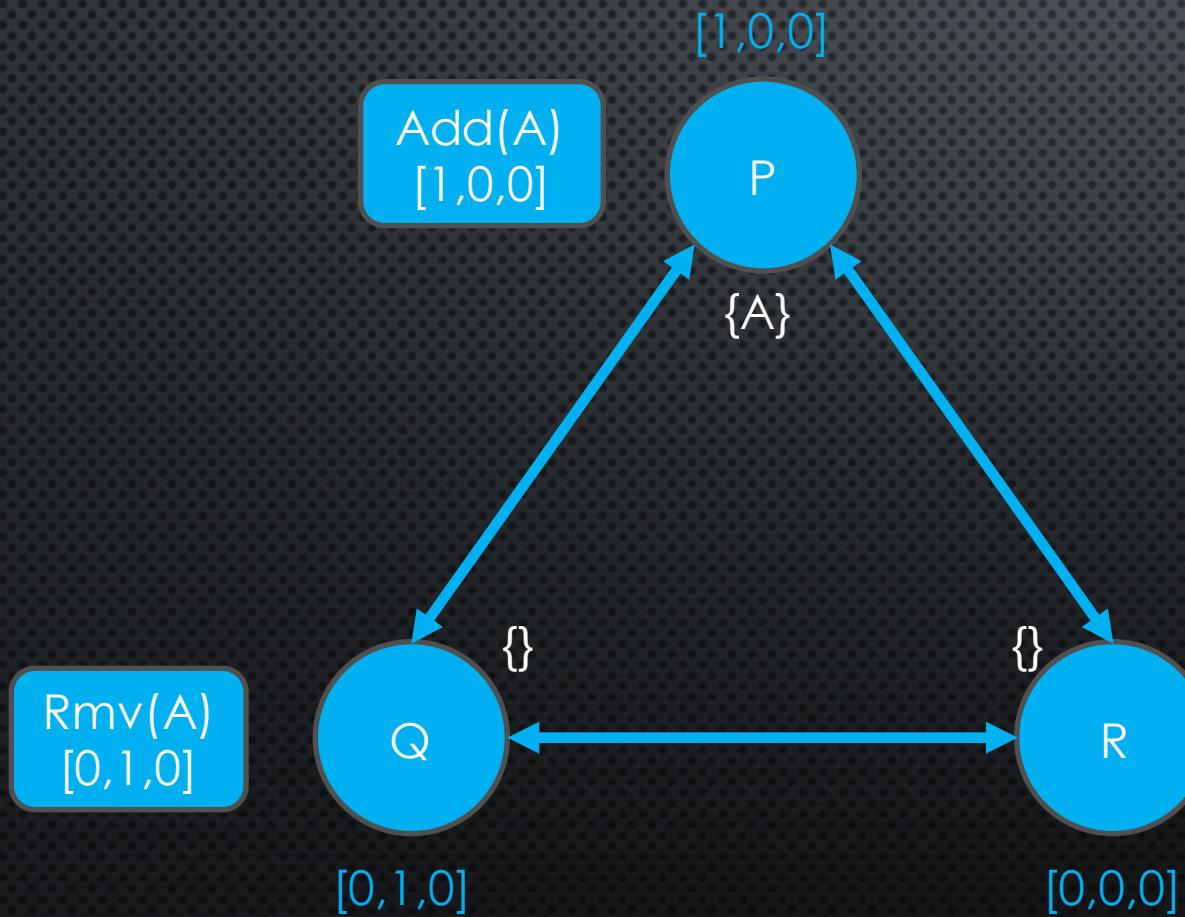
The “happened-before” relation between ops



TESTING SOLUTION 1.1



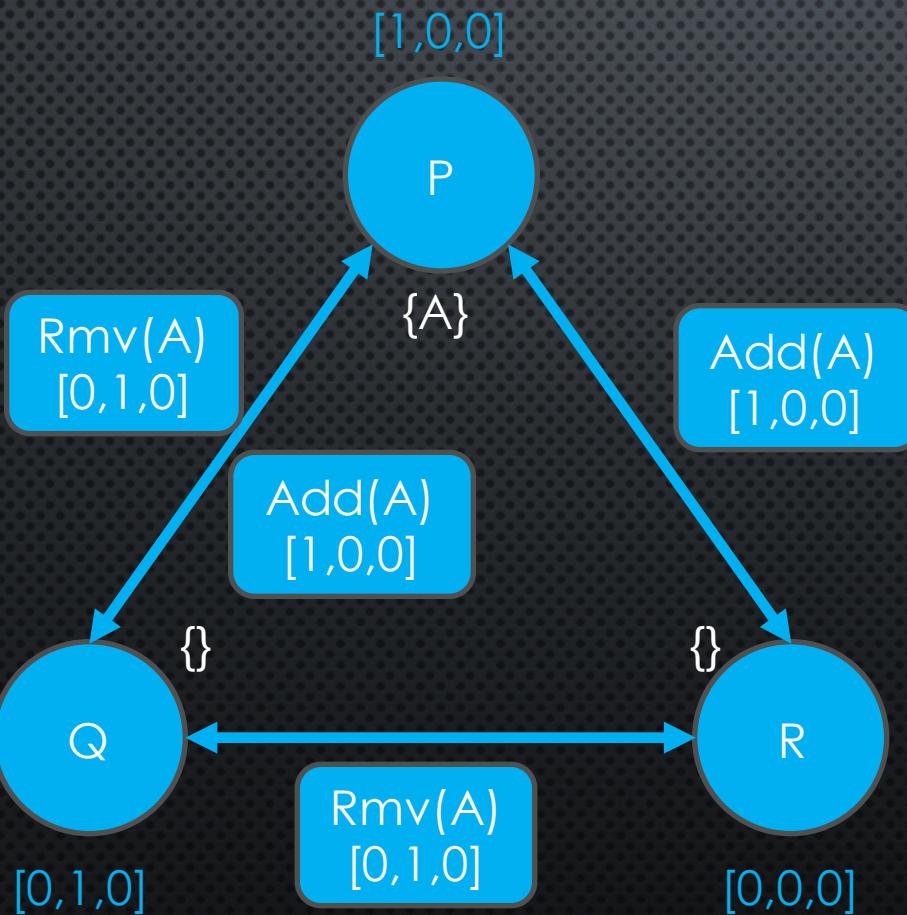
TESTING SOLUTION 1.1



Scenario 1:

- P tags the op Add(A)
- Q tags the op Rmv(A)

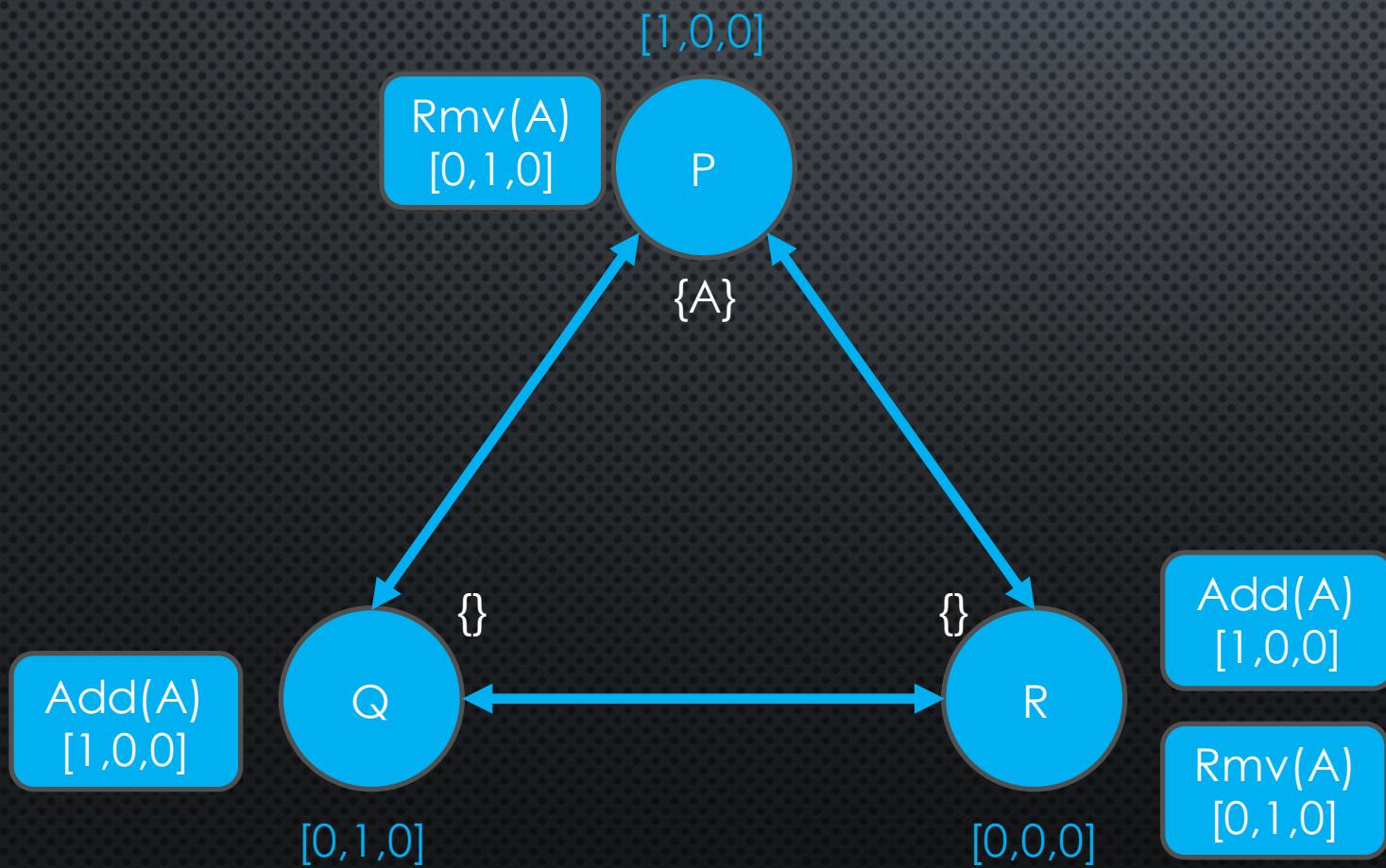
TESTING SOLUTION 1.1



Scenario 1:

- P tags the op $\text{Add}(A)$
- Q tags the op $\text{Rmv}(A)$
- P and Q bcast their tagged ops

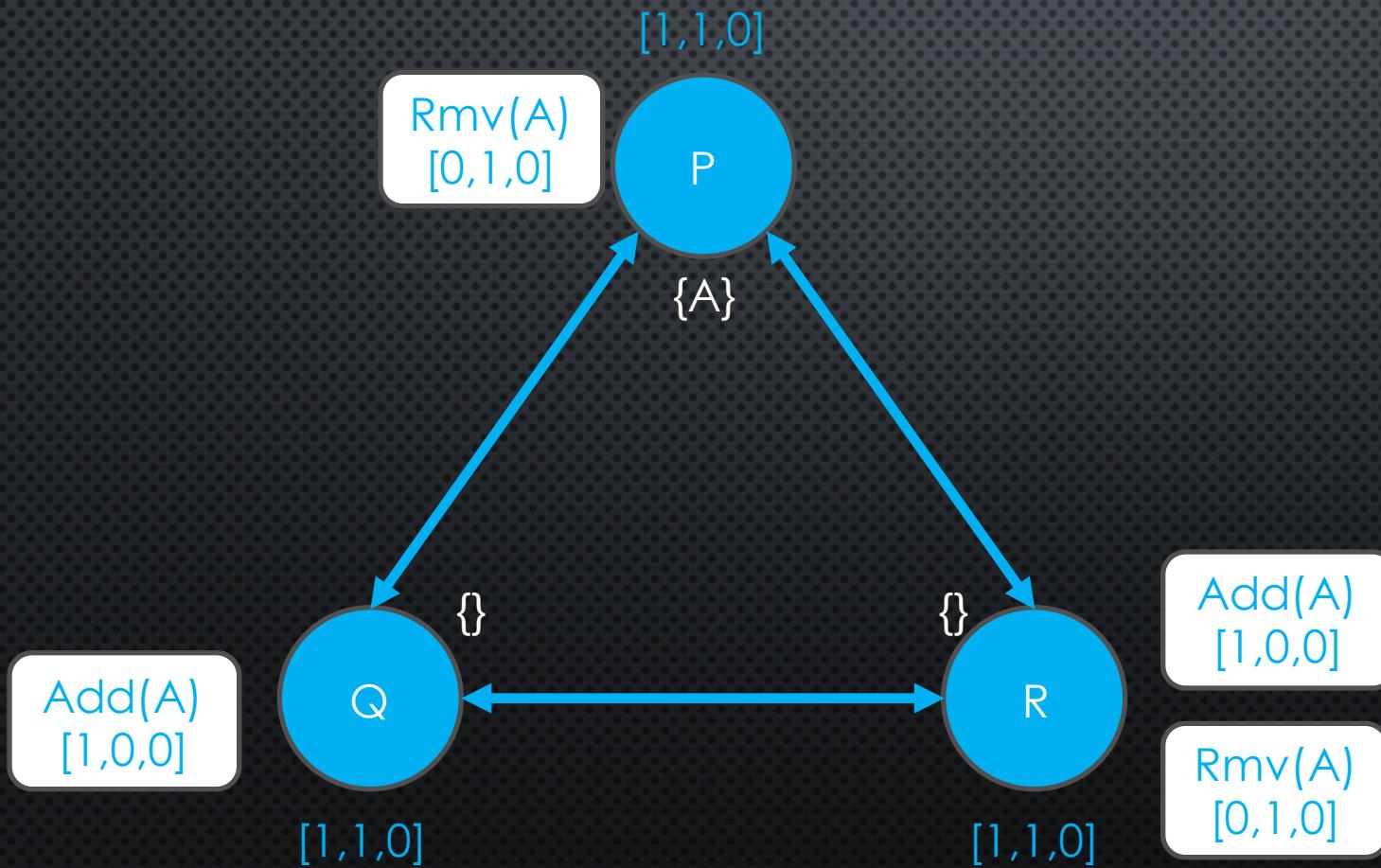
TESTING SOLUTION 1.1



Scenario 1:

- P tags the op $\text{Add}(A)$
- Q tags the op $\text{Rmv}(A)$
- P and Q bcast their tagged ops

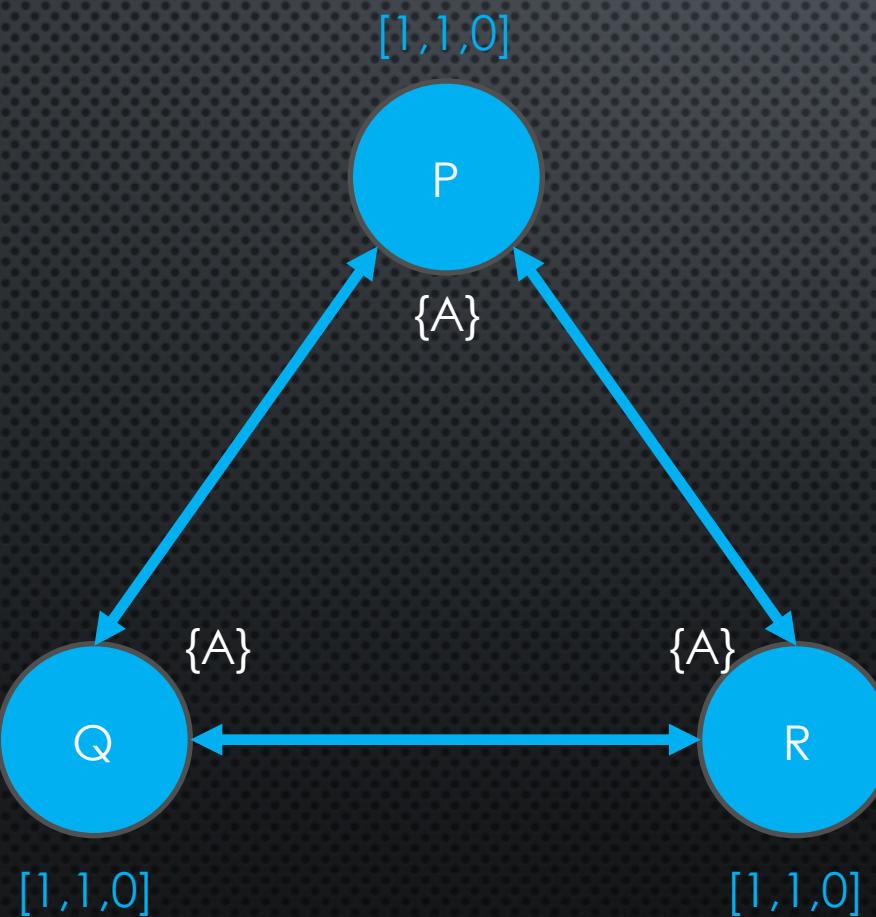
TESTING SOLUTION 1.1



Scenario 1:

- P tags the op $Add(A)$
- Q tags the op $Rmv(A)$
- P and Q bcast their tagged ops
- All nodes deliver both ops

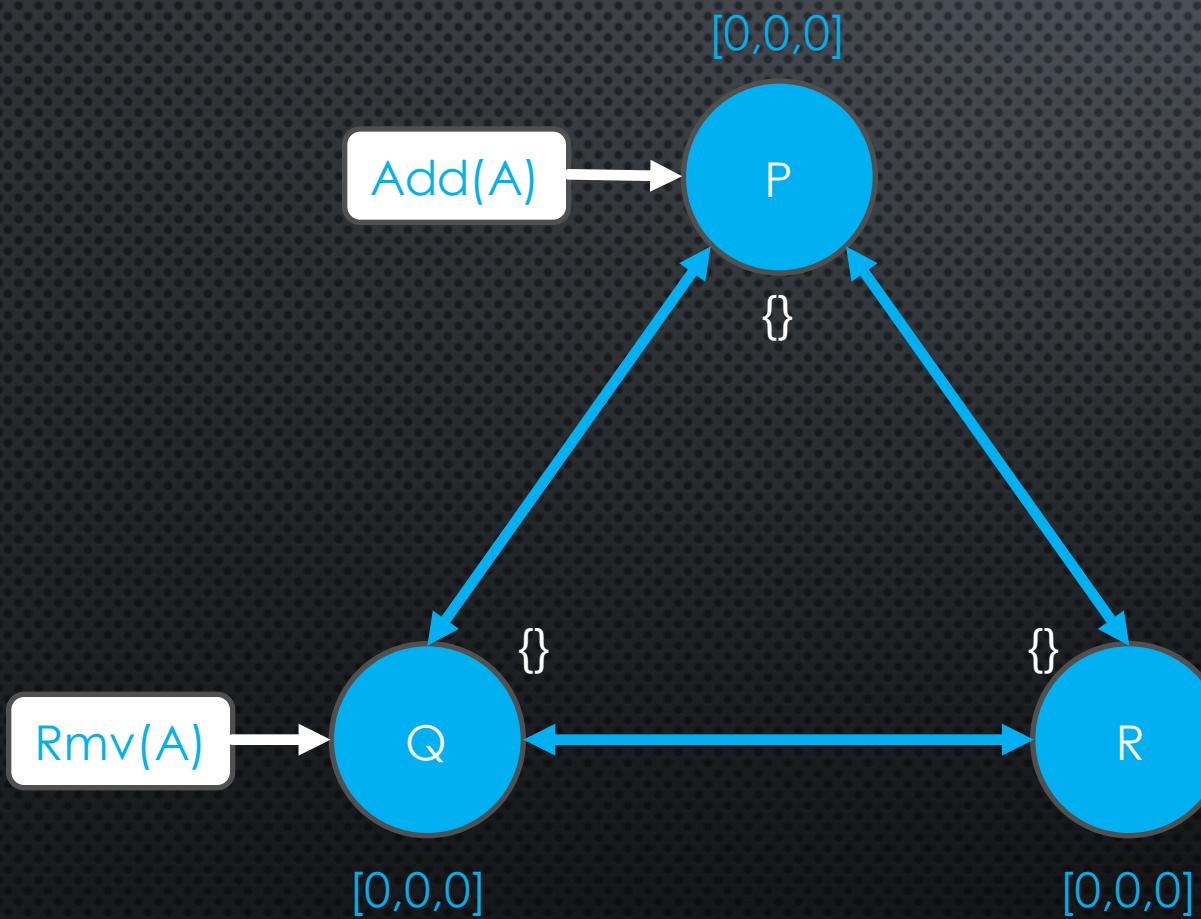
TESTING SOLUTION 1.1



Scenario 1:

- P tags the op Add(A)
- Q tags the op Rmv(A)
- P and Q bcast their tagged ops
- All nodes deliver both ops
- $[1,0,0]$ and $[0,1,0]$ are concurrent
 - Add-Wins Semantics
 - A is there

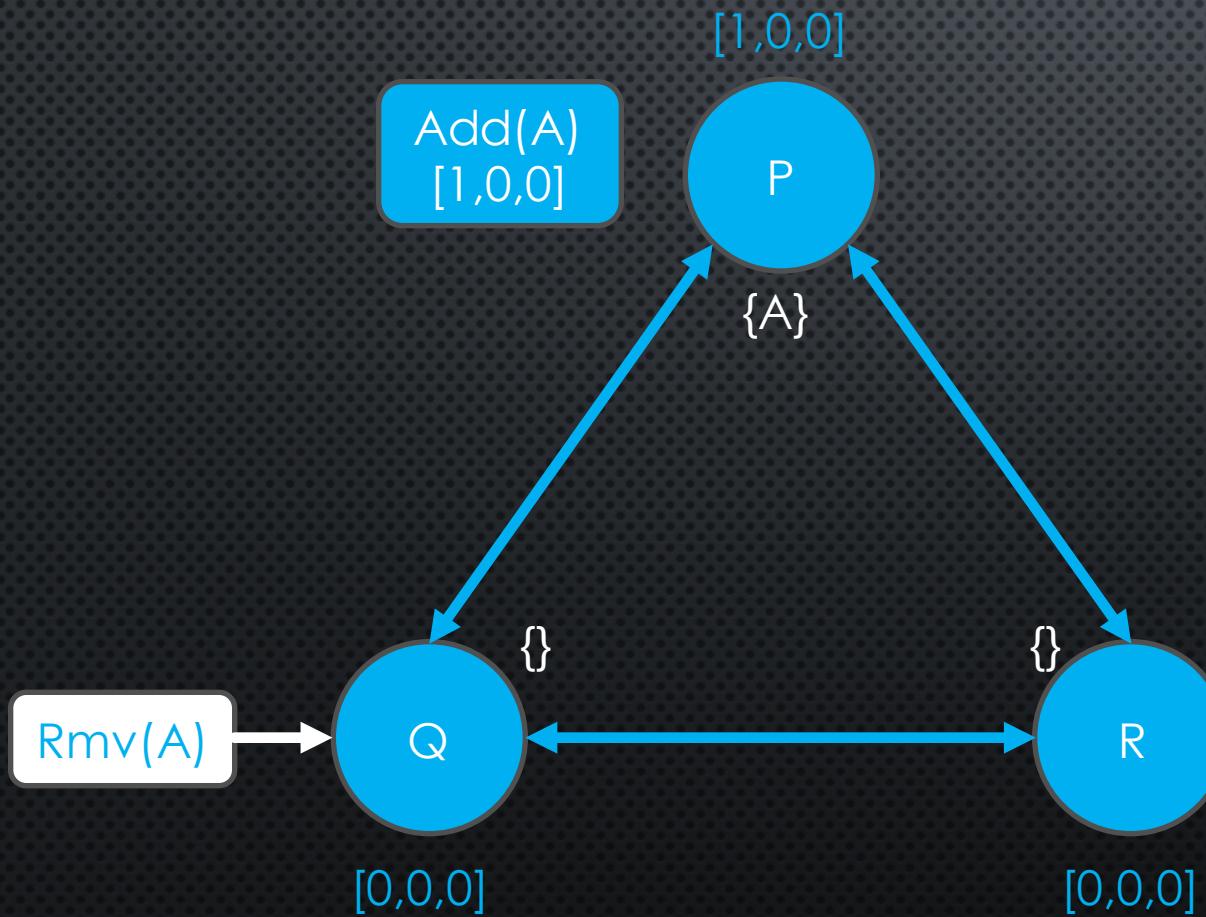
TESTING SOLUTION 1.1



- A group of 3 Nodes
- A replicated Add-Wins Set
- 2 concurrent operations:
 - Add(A)
 - Rmv(A)

■ Application Layer
■ Middleware

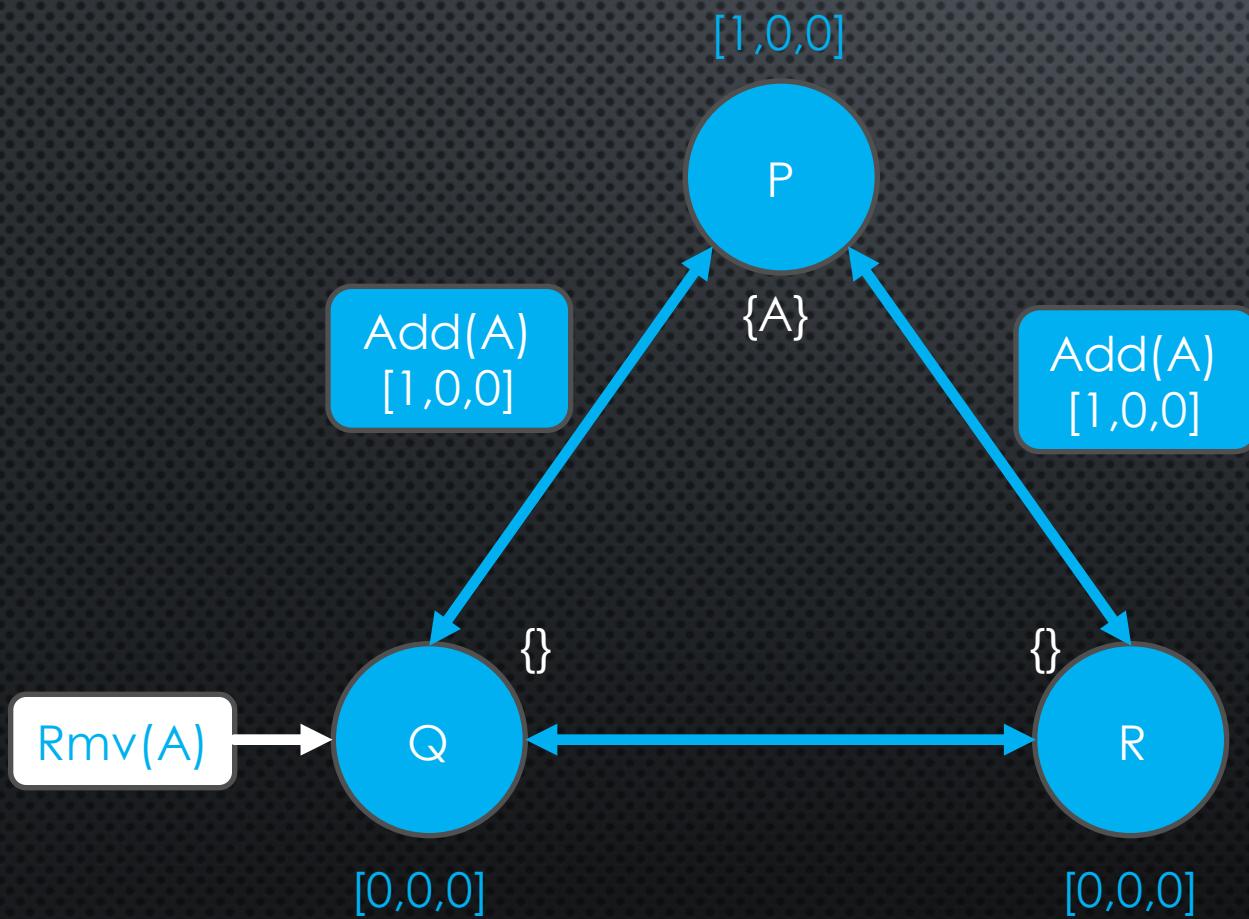
TESTING SOLUTION 1.1



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue

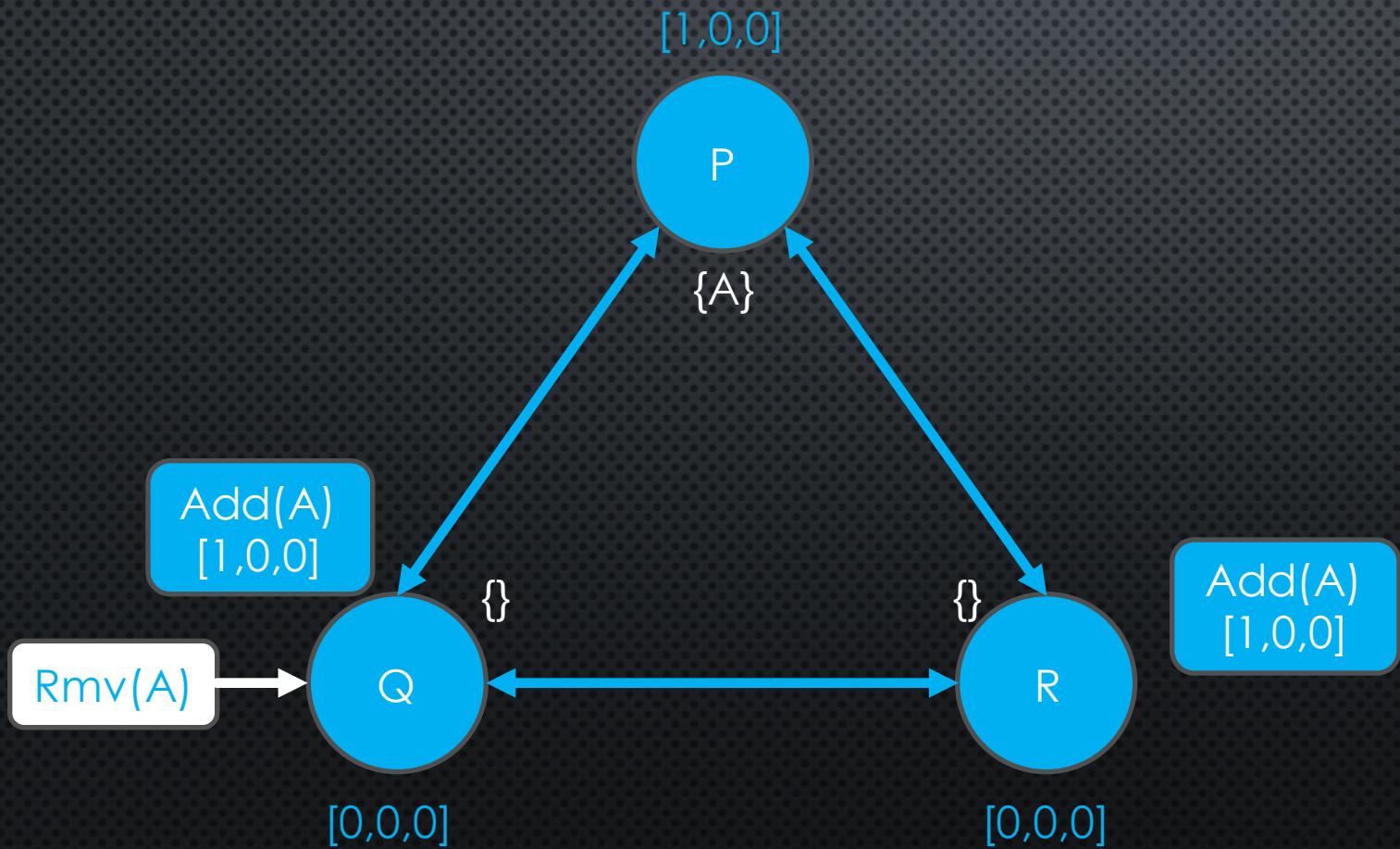
TESTING SOLUTION 1.1



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op

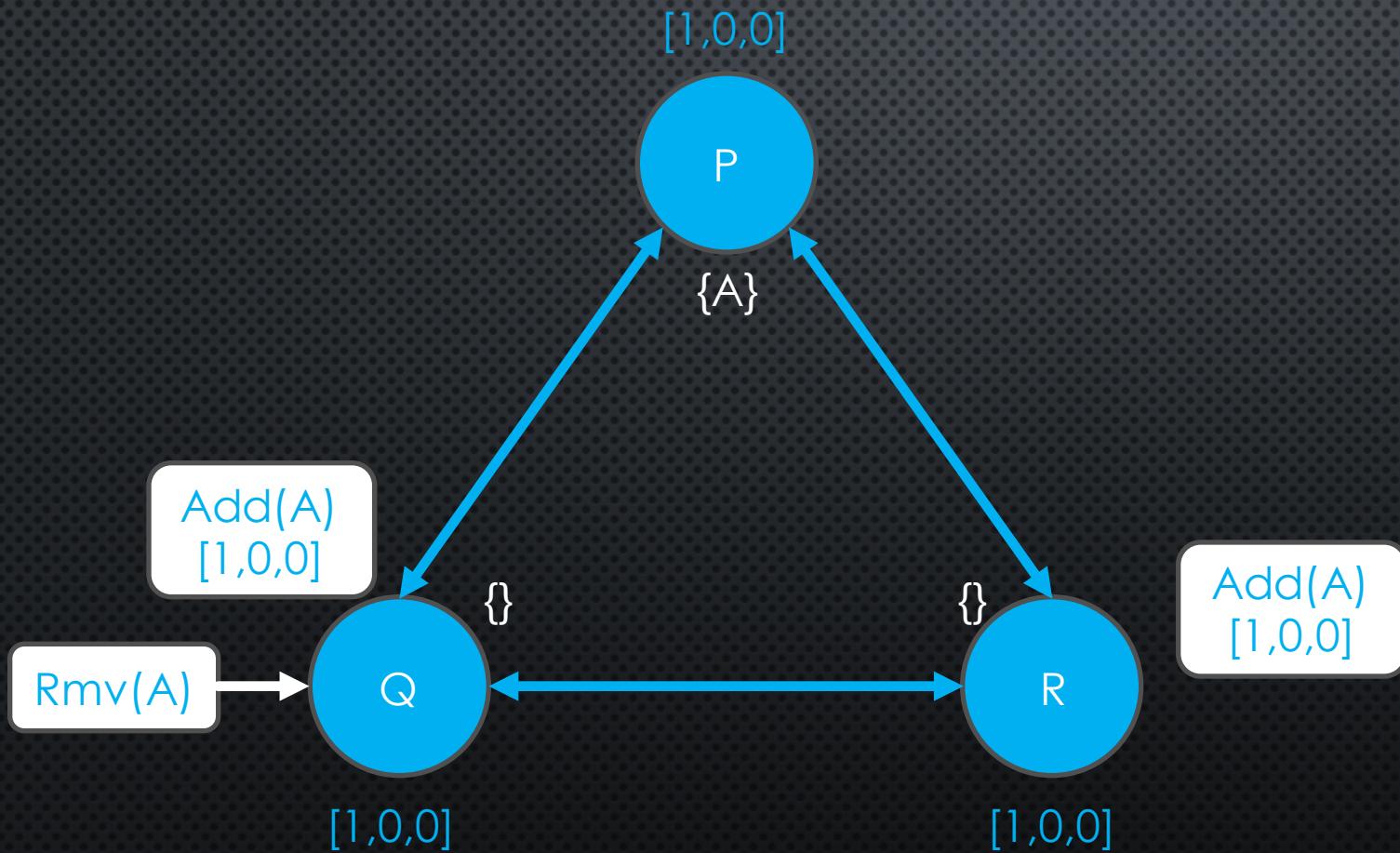
TESTING SOLUTION 1.1



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op

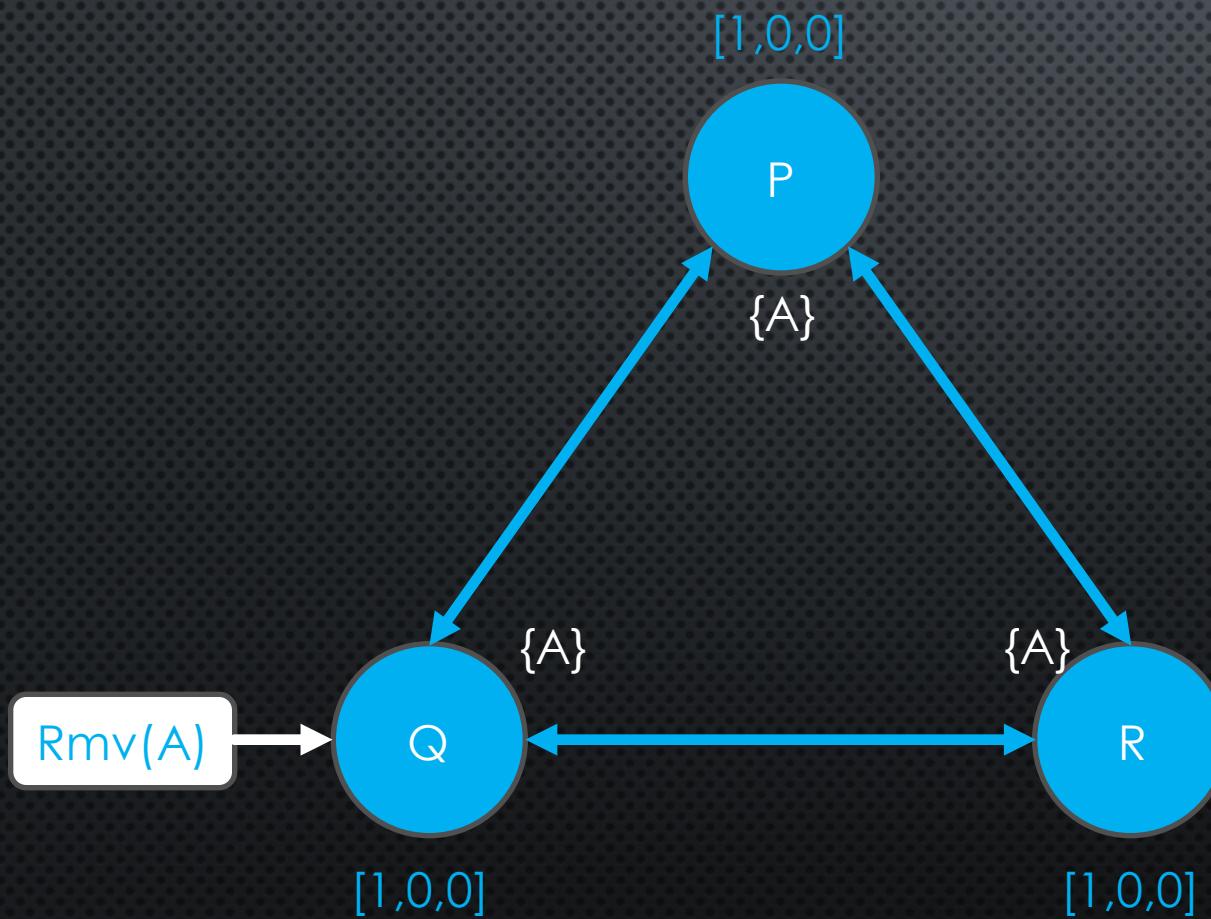
TESTING SOLUTION 1.1



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op
- All nodes deliver Add(A)

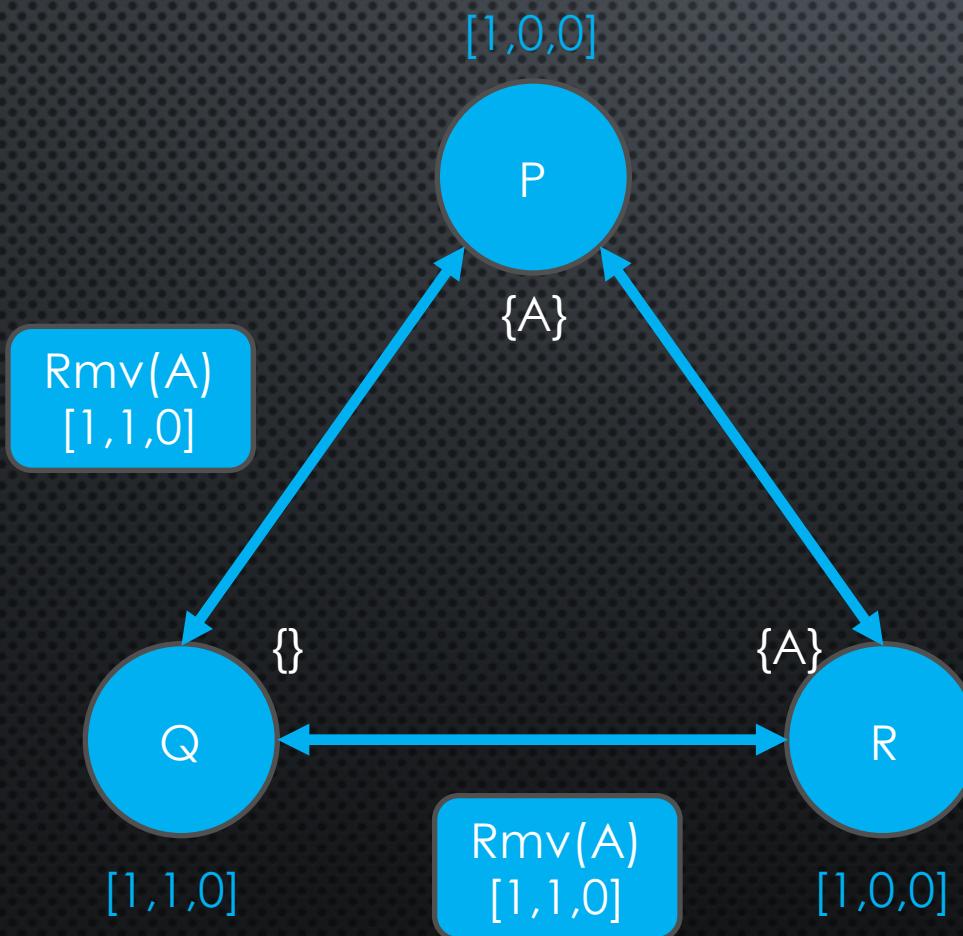
TESTING SOLUTION 1.1



Scenario 2:

- P tags the op $\text{Add}(A)$
- $\text{Rmv}(A)$ is still in the queue waiting to be processed by Q
- P broadcasts their tagged op
- All nodes deliver $\text{Add}(A)$

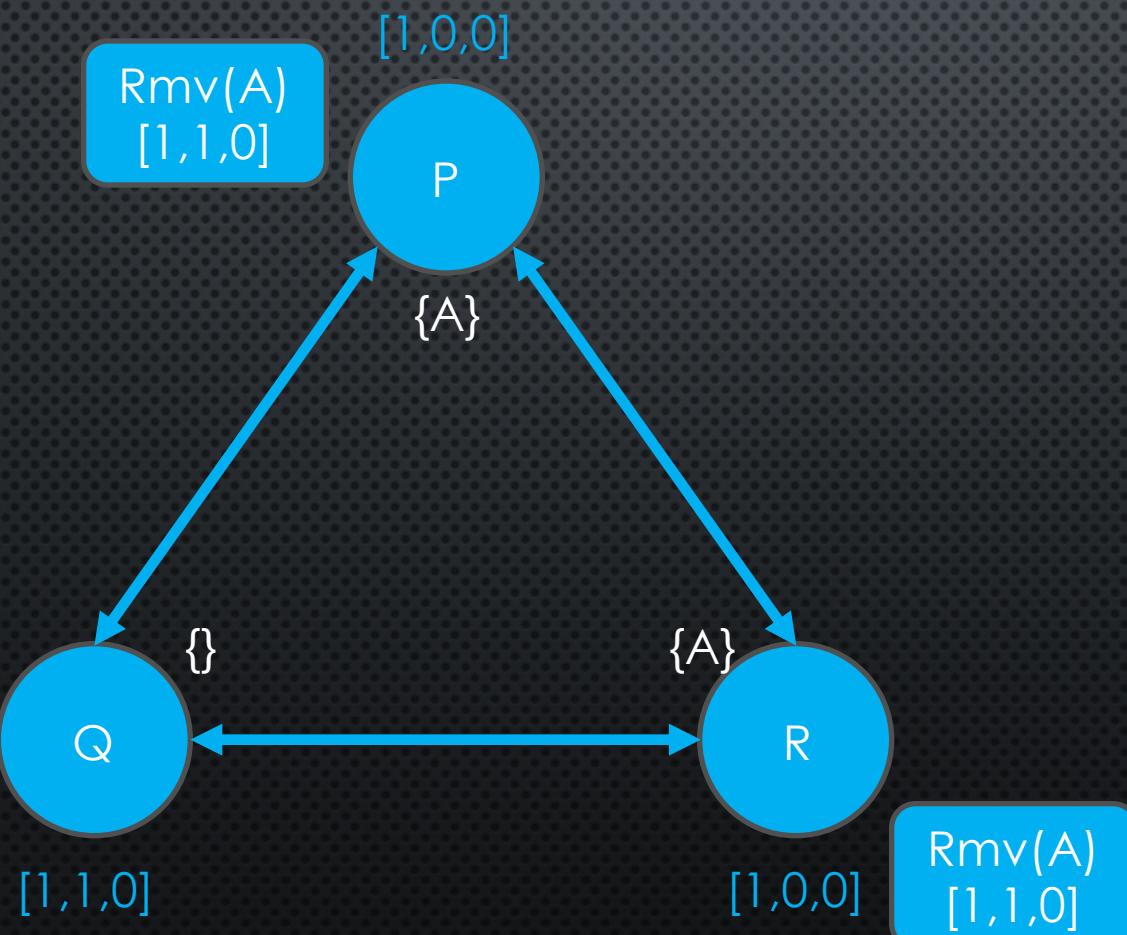
TESTING SOLUTION 1.1



Scenario 2:

- P tags the op Add(A)
- $Rmv(A)$ is still in the queue waiting to be processed by Q
- P bcasts their tagged op
- All nodes deliver Add(A)
- Q dequeues $Rmv(A)$, tags and bcasts it

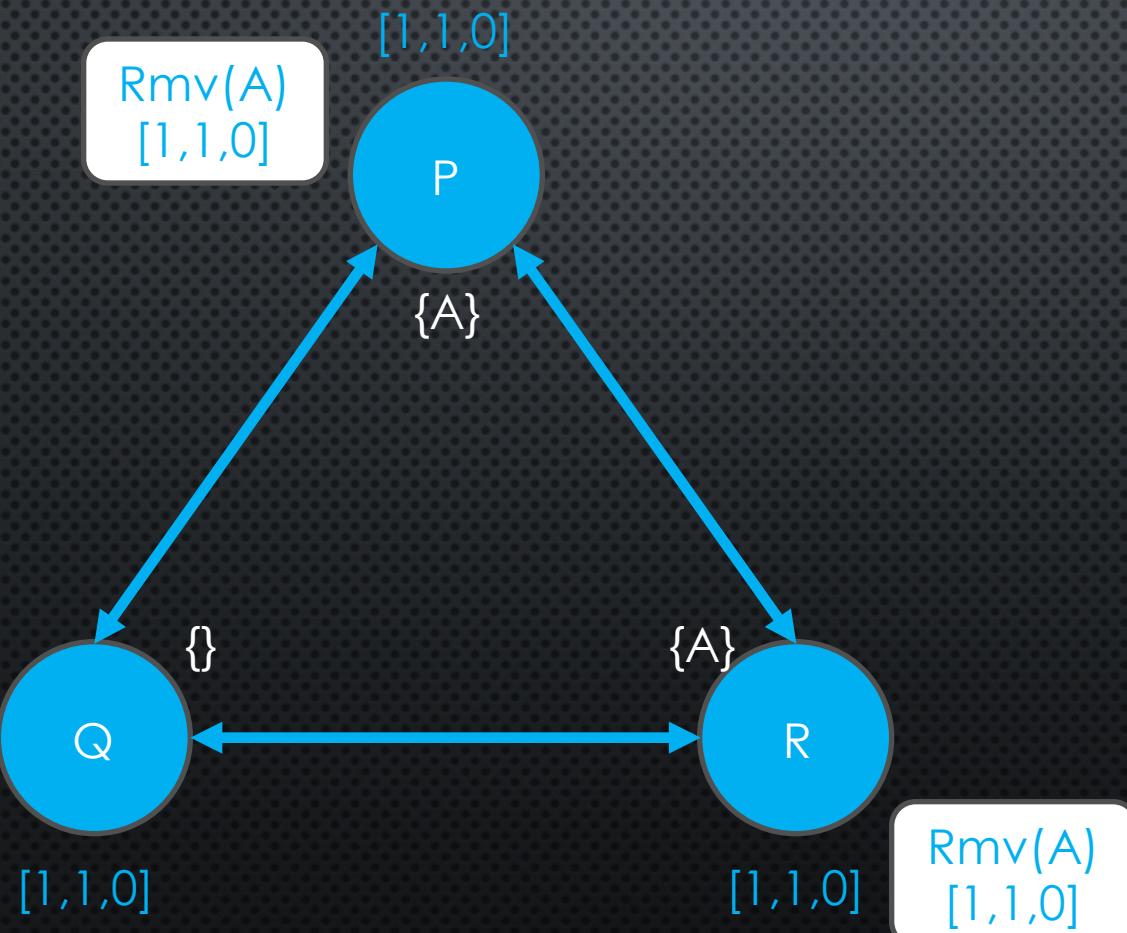
TESTING SOLUTION 1.1



Scenario 2:

- P tags the op $\text{Add}(A)$
- $\text{Rmv}(A)$ is still in the queue waiting to be processed by Q
- P bcasts their tagged op
- All nodes deliver $\text{Add}(A)$
- Q dequeues $\text{Rmv}(A)$, tags and bcasts it

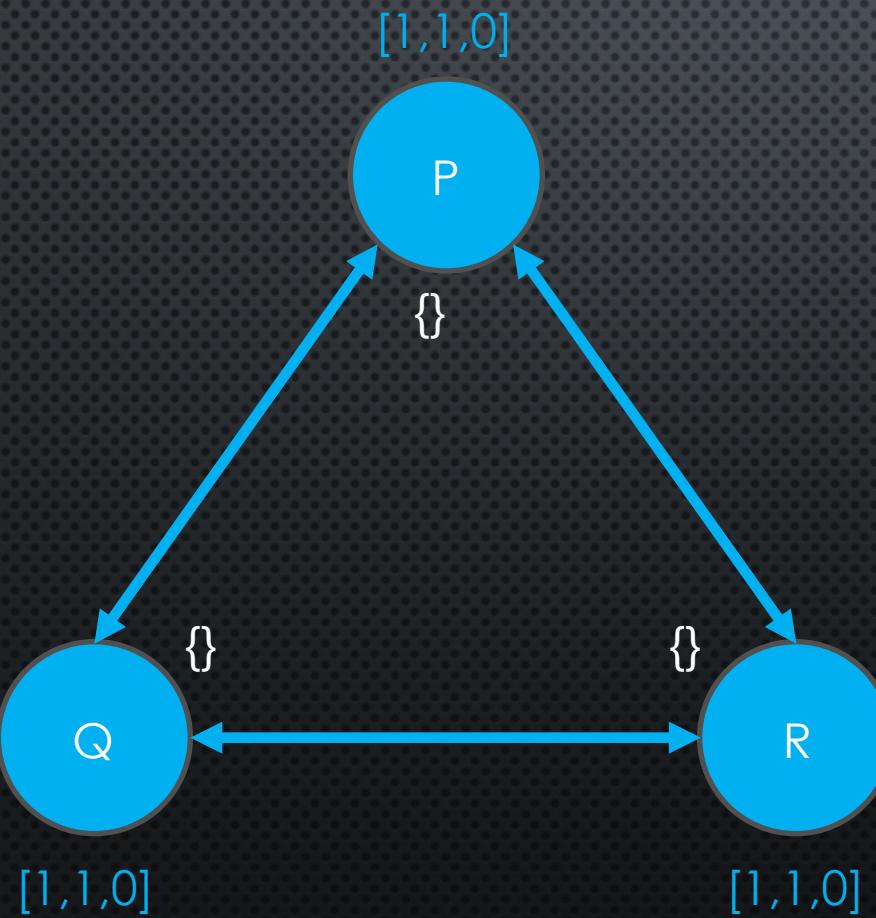
TESTING SOLUTION 1.1



Scenario 2:

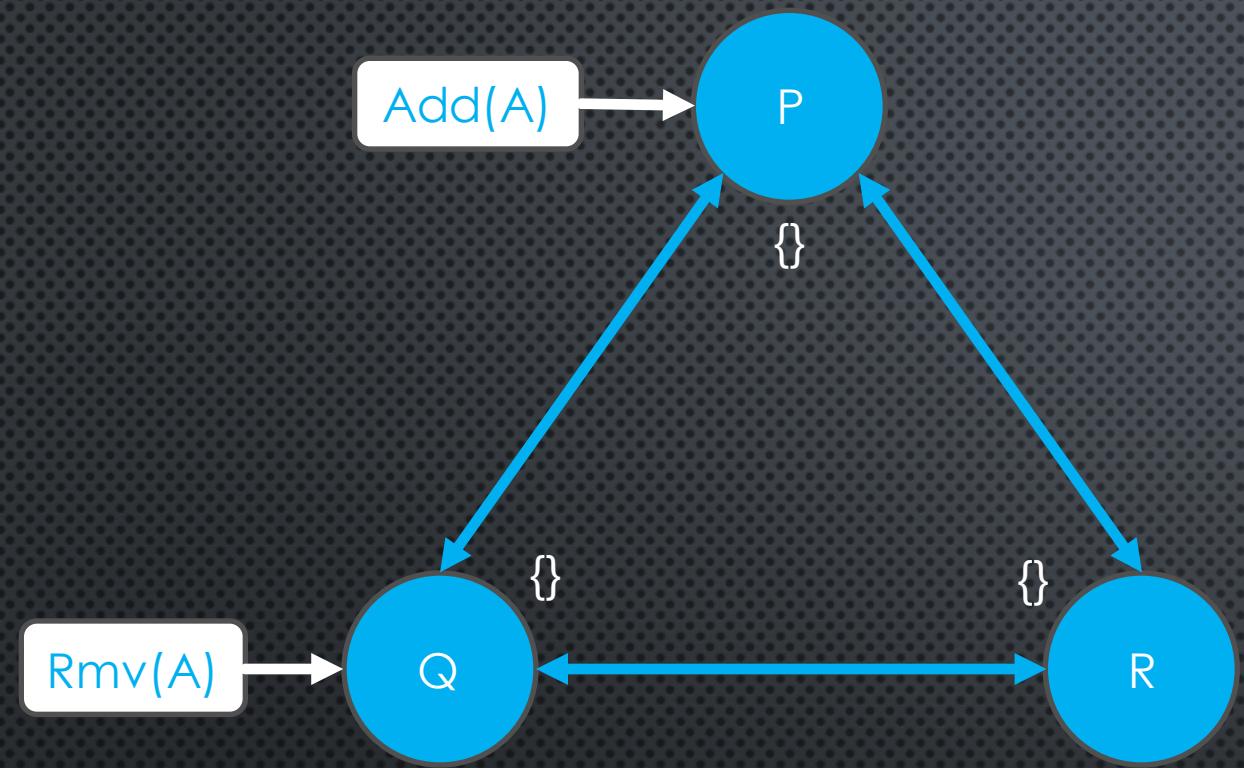
- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op
- All nodes deliver Add(A)
- Q dequeues Rmv(A) , tags and broadcasts it
- Rmv(A) delivered at P and R

TESTING SOLUTION 1.1



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P bcasts their tagged op
- All nodes deliver Add(A)
- Q dequeues Rmv(A), tags and bcasts it
- Rmv(A) delivered at P and R
- [1,0,0] happened-before [1,1,0]
 - A is not there



Scenario 1:

- Add(A) and Rmv(A) tagged as concurrent
- Final State: {A}

Classical Causal Delivery

- Does not care about concurrent operations
- Orders concurrent ops based on their delivery order
 - Could order concurrent ops as one happening before the other (Scenario 2)

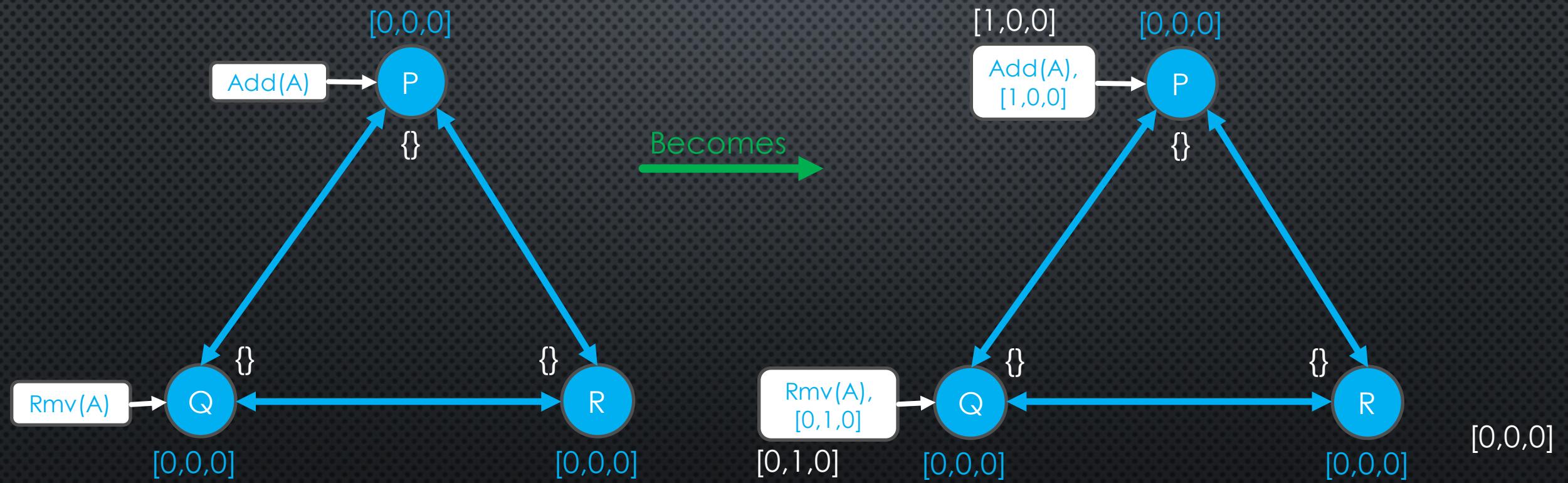
Scenario 2:

- Rmv(A) tagged as in the future of Add(A)
- Final State: {}

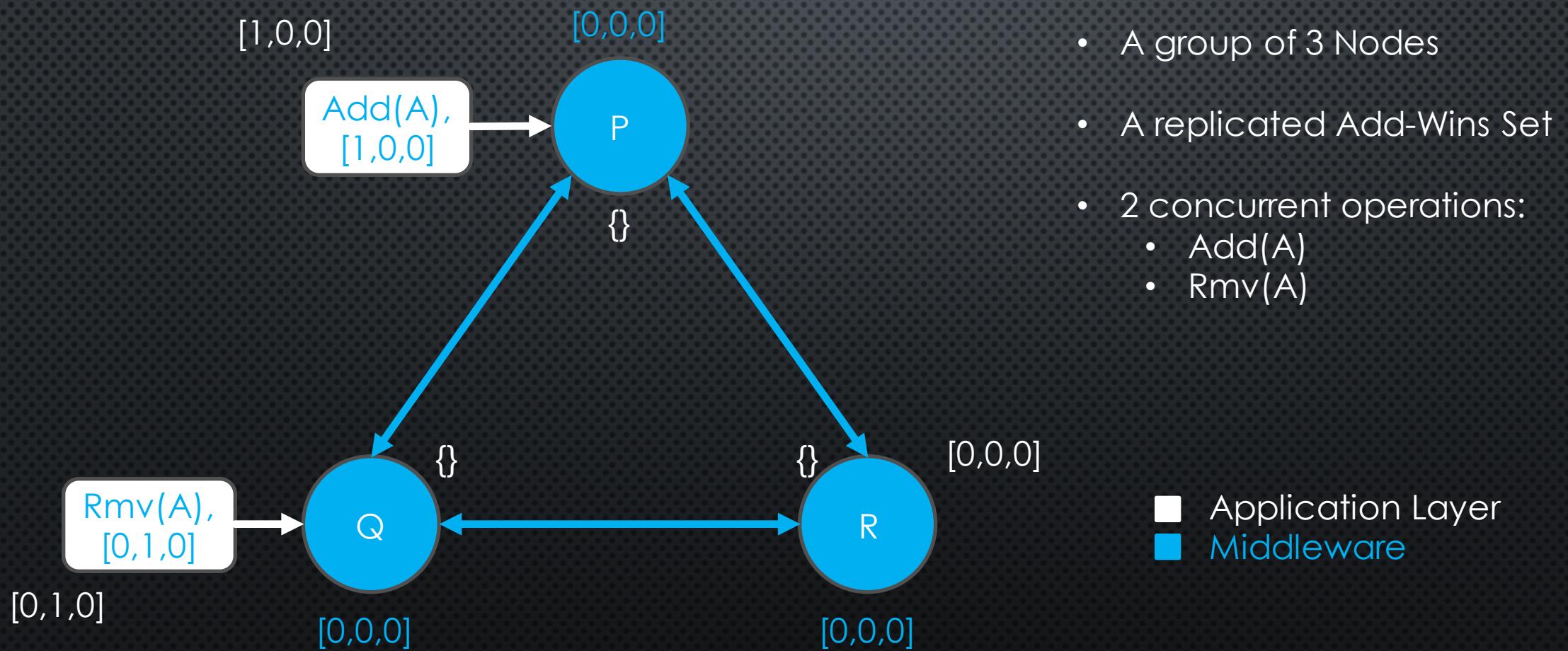
SOLUTION 2.0

- BETTER CHARACTERIZATION OF THE HAPPENED-BEFORE RELATION BETWEEN OPS
 - TAG BASED ON WHAT WAS DELIVERED AT THE APPLICATION LEVEL
- AS DELIVERY HAPPENS AT THE APPLICATION LEVEL
 - TAG AT THE APPLICATION LEVEL

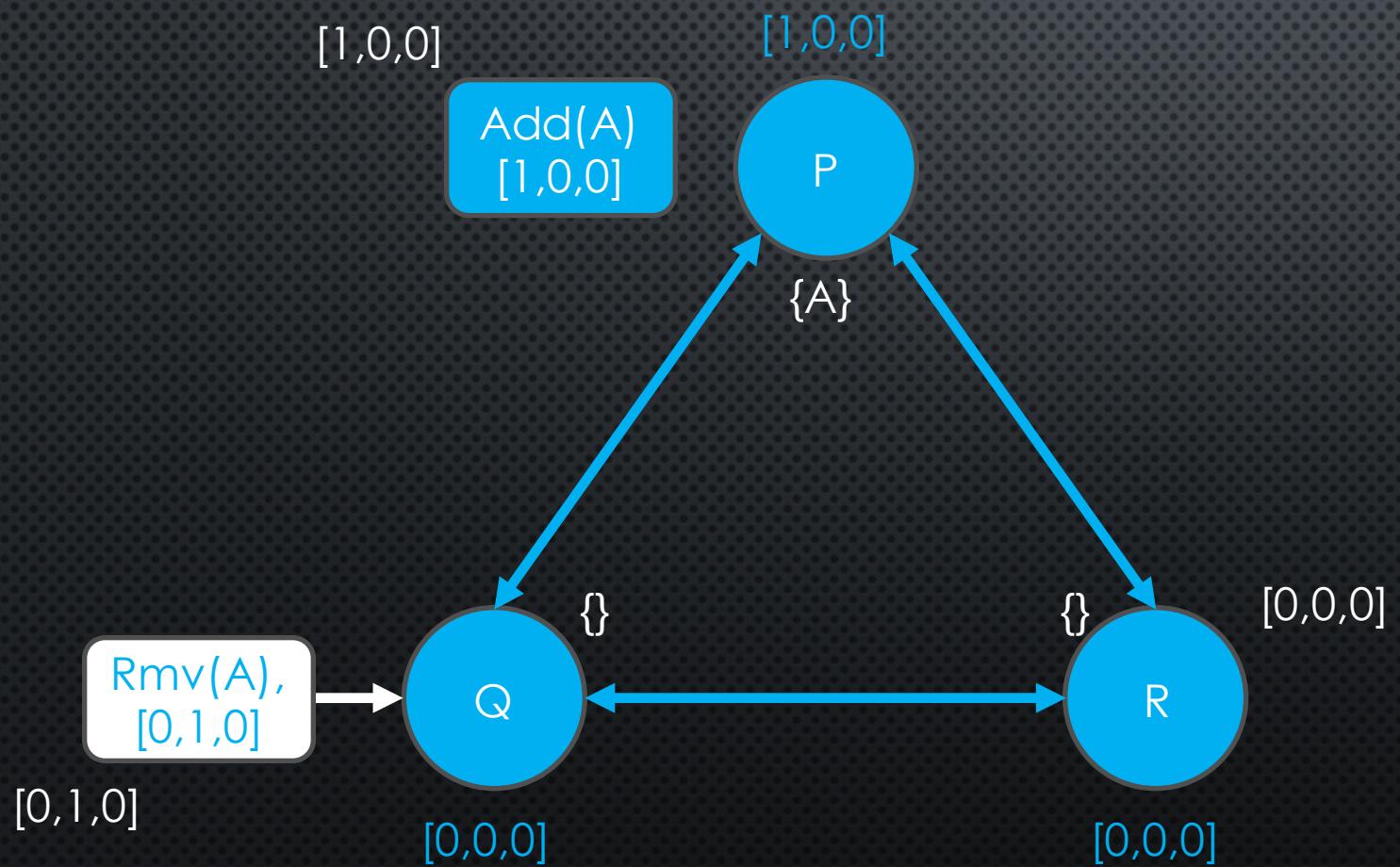
SOLUTION 2.0



TESTING SOLUTION 2.0



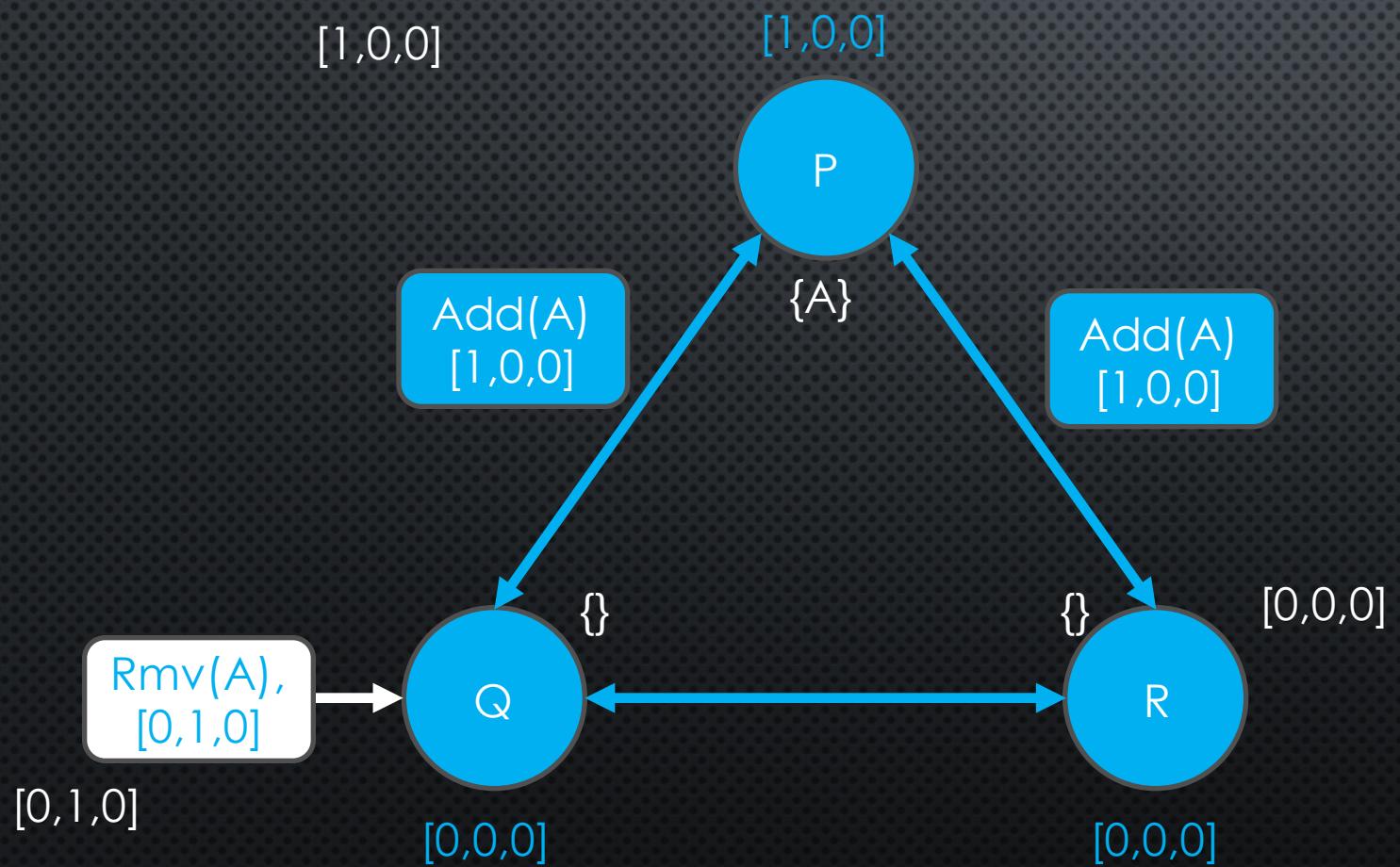
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue

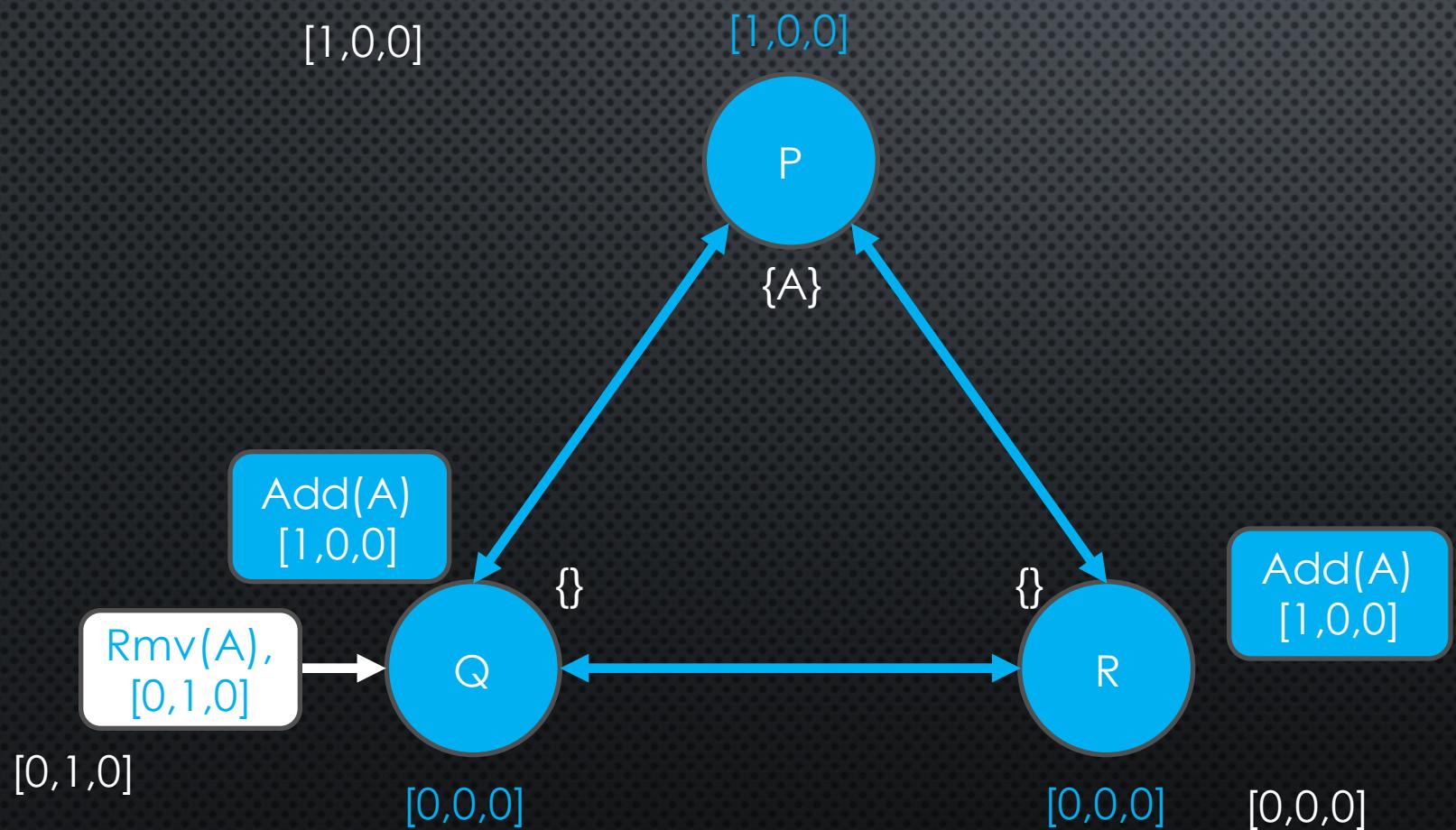
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op

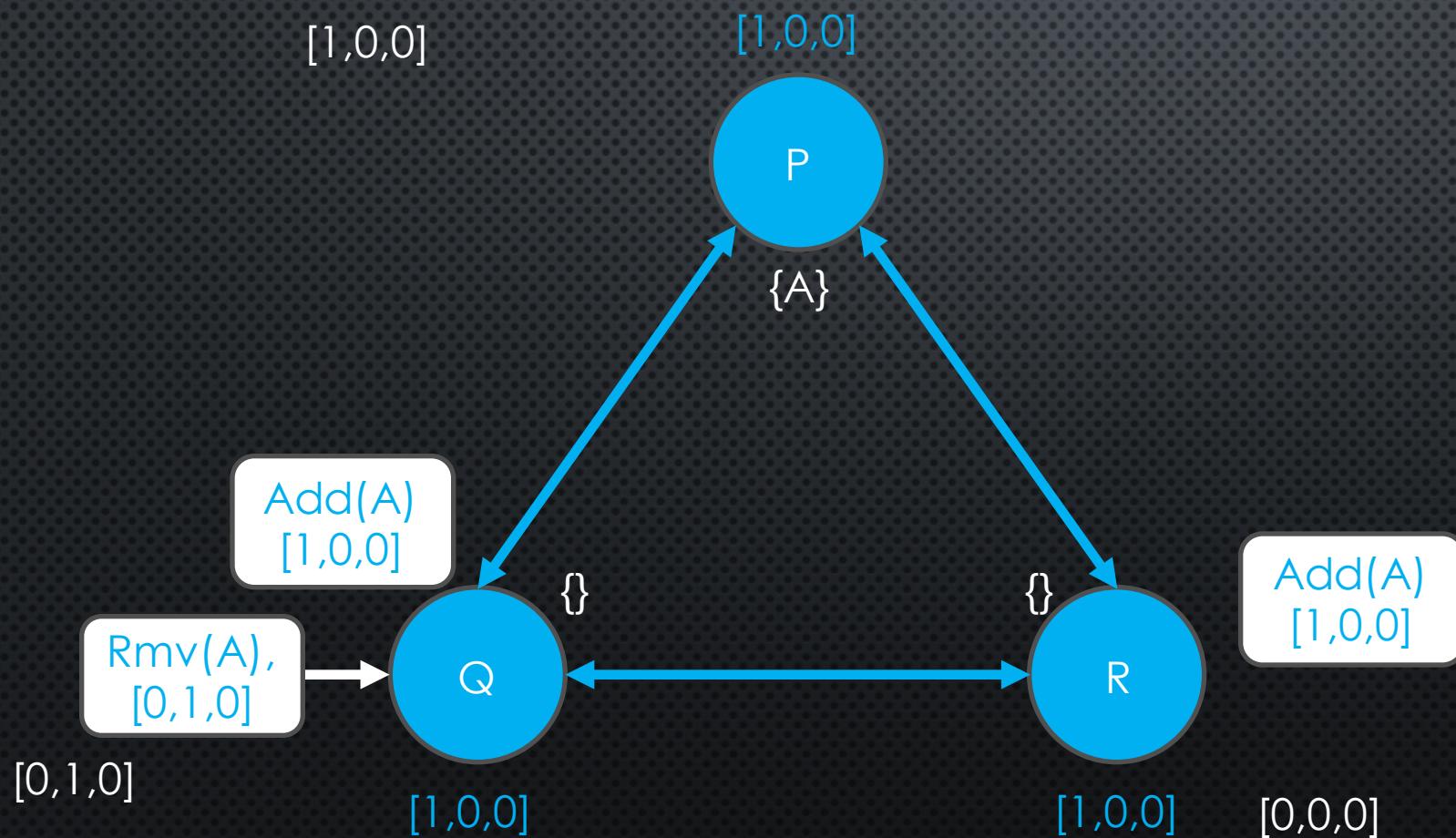
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op

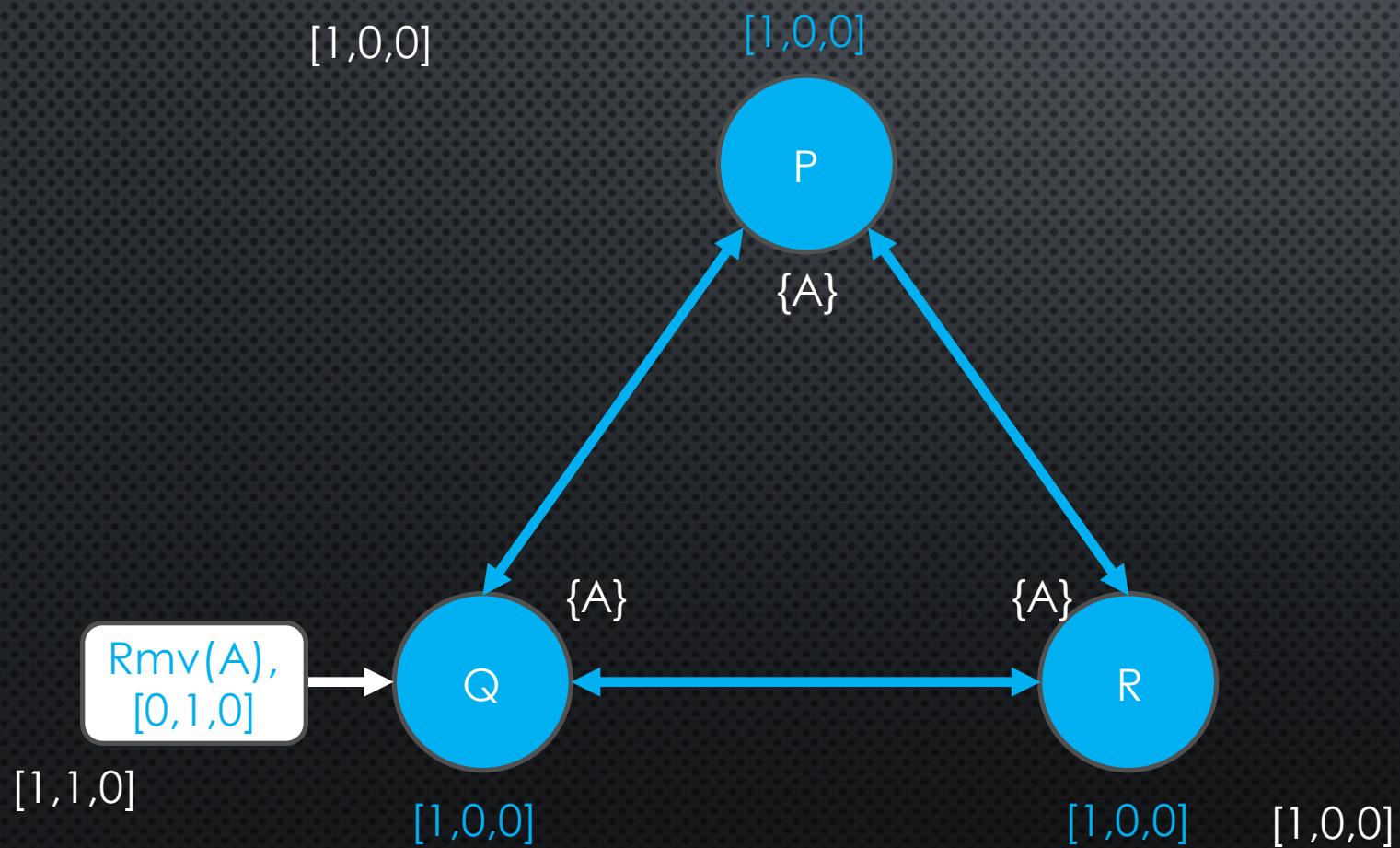
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op $\text{Add}(A)$
- $\text{Rmv}(A)$ is still in the queue waiting to be processed by Q
- P broadcasts their tagged op
- All nodes deliver $\text{Add}(A)$

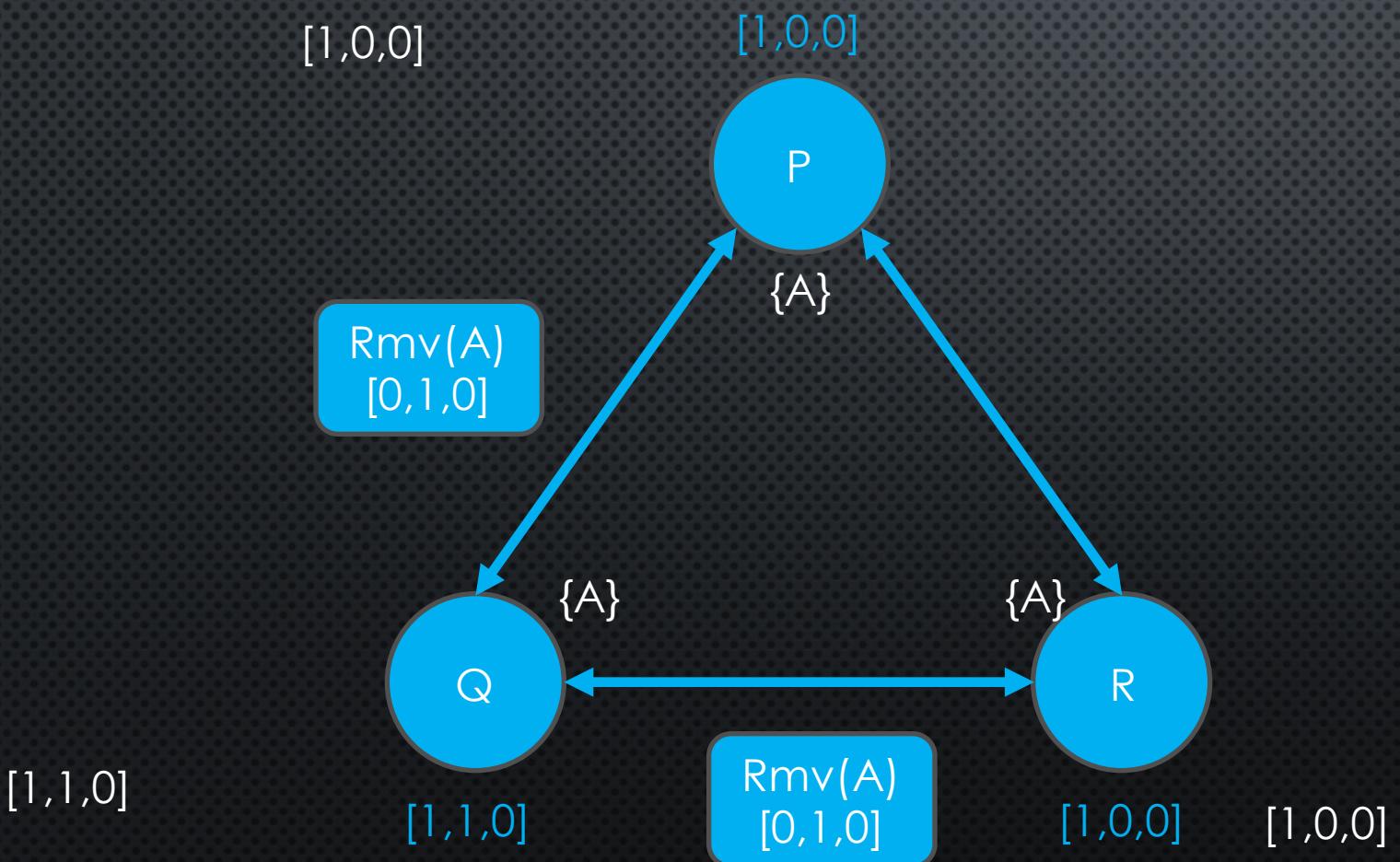
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op
- All nodes deliver Add(A)

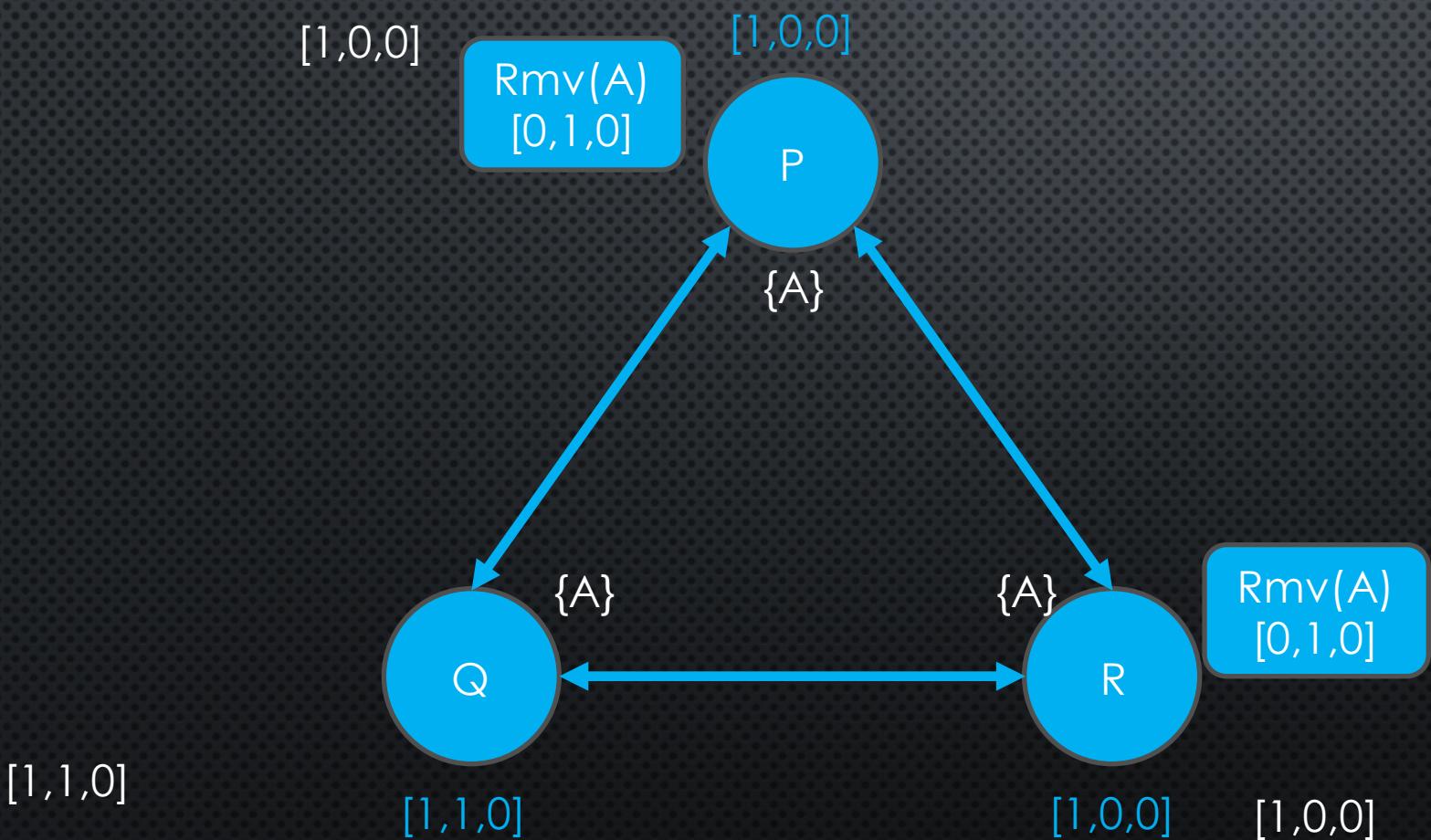
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op Add(A)
- $Rmv(A)$ is still in the queue waiting to be processed by Q
- P bcasts their tagged op
- All nodes deliver Add(A)
- Q dequeues $Rmv(A)$, tags and bcasts it

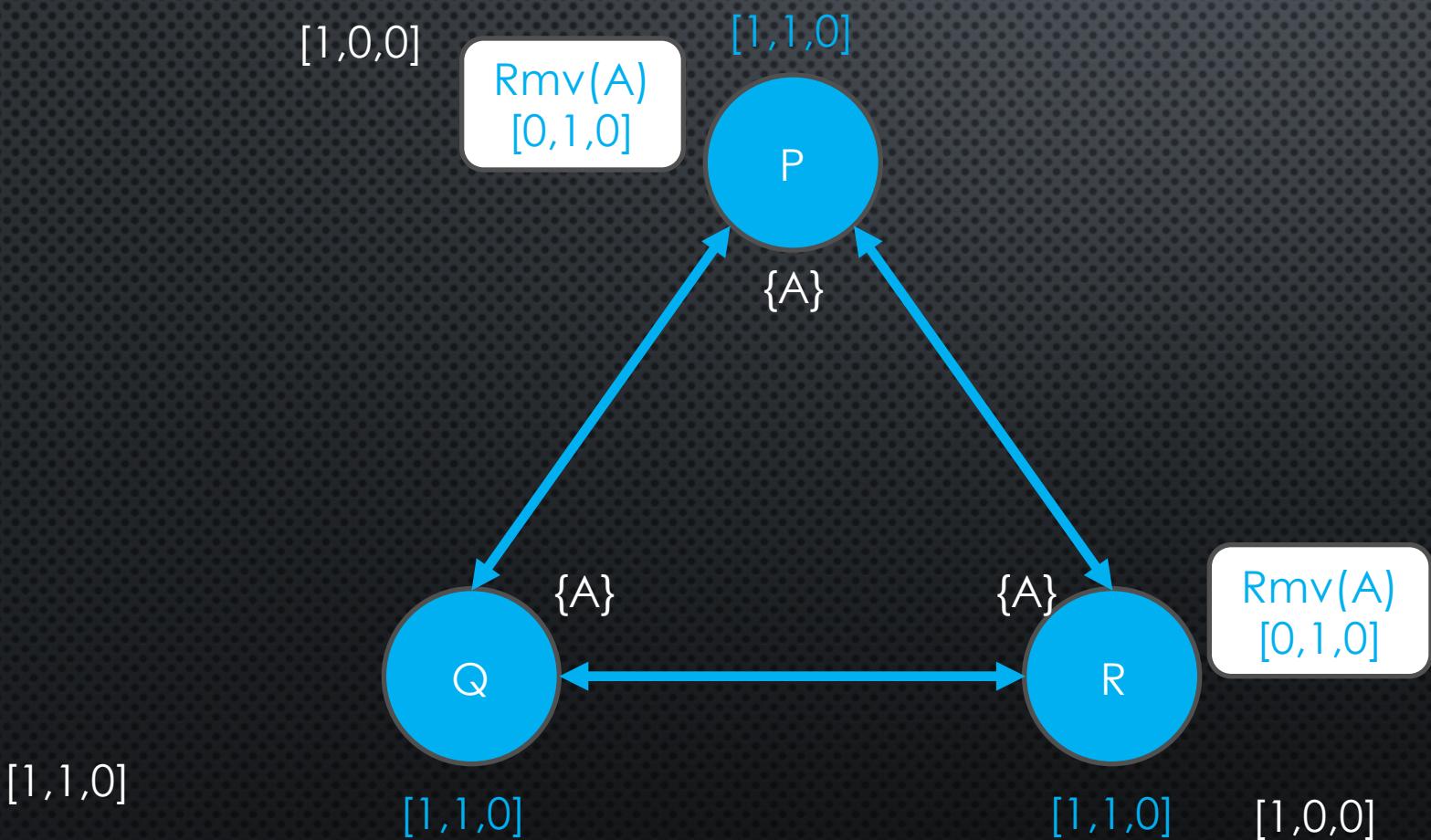
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op $\text{Add}(A)$
- $\text{Rmv}(A)$ is still in the queue waiting to be processed by Q
- P bcasts their tagged op
- All nodes deliver $\text{Add}(A)$
- Q dequeues $\text{Rmv}(A)$, tags and bcasts it

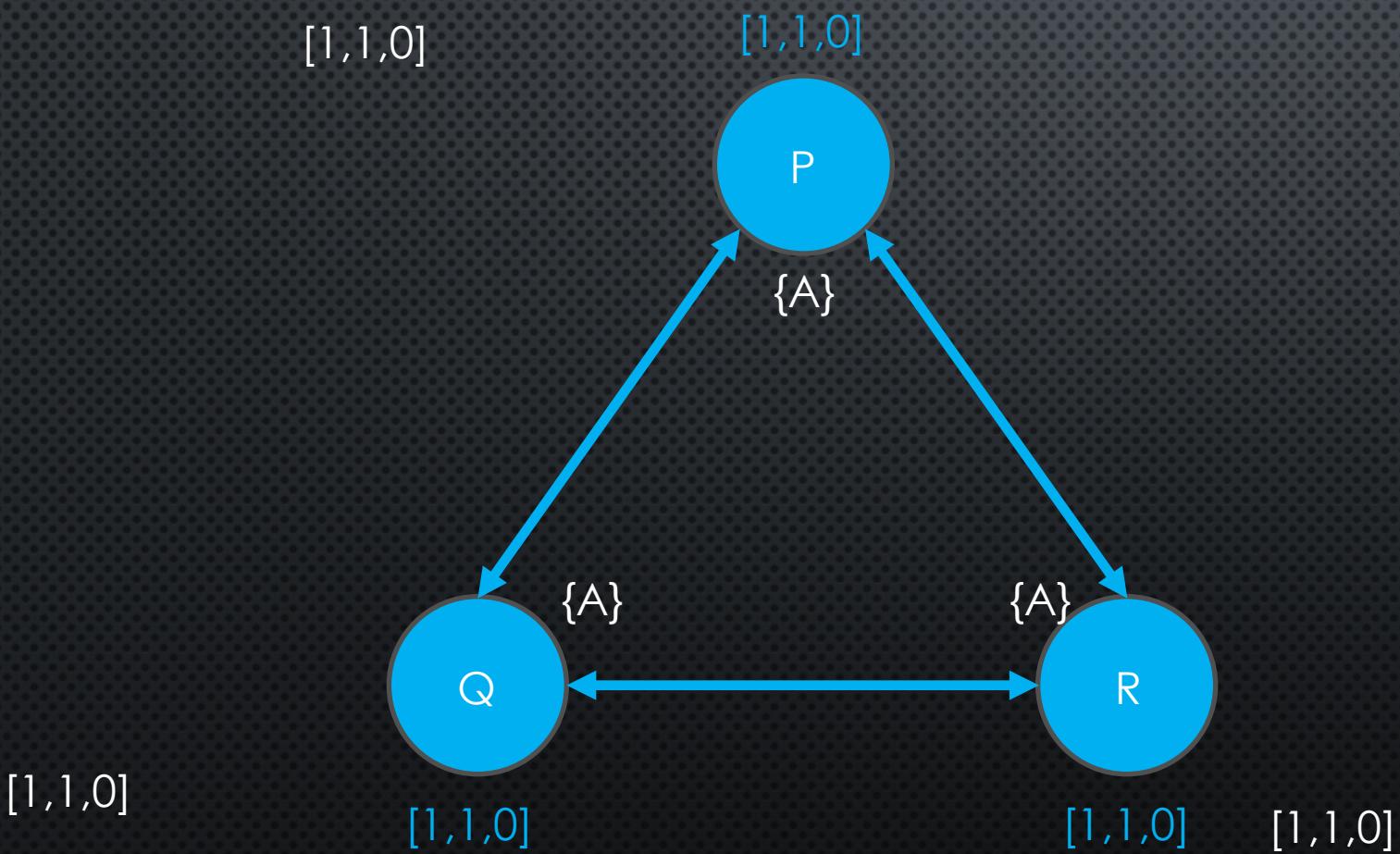
TESTING SOLUTION 2.0



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P broadcasts their tagged op
- All nodes deliver Add(A)
- Q dequeues Rmv(A), tags and broadcasts it
- Rmv(A) delivered at P and R

TESTING SOLUTION 2.0



Scenario 2:

- P tags the op Add(A)
- Rmv(A) is still in the queue waiting to be processed by Q
- P bcasts their tagged op
- All nodes deliver Add(A)
- Q dequeues Rmv(A) , tags and bcasts it
- Rmv(A) delivered at P and R
- [1,0,0] concurrent [0,1,0]
 - Add Wins: A is there

SOLUTION 2.0

- A.K.A. TAGGED CAUSAL DELIVERY
- BETTER CHARACTERIZATION OF THE HAPPENED-BEFORE RELATION BETWEEN OPS
- EXPECTED RESULT (NOT DIFFERENT FROM SCENARIO 1)
- PROBLEM SOLVED

A WELCOME SIDE EFFECT

- AS YOU NOTICED IN SOLUTION 1.1 SCENARIO 2:
 - CAUSAL DELIVERY SOMETIMES TAGS CONCURRENT OPERATIONS AS ONE HAPPENED-BEFORE THE OTHER
 - THIS LEADS TO OVER ORDERING OPERATIONS AND COULD INCLUDE EXTRA DELAY ON DELIVERY
 - THIS DOES NOT HAPPEN IN SOLUTION 2.0 (TAGGED CAUSAL DELIVERY)

CAUSAL DELIVERY VS TAGGED CAUSAL DELIVERY

- SOLUTION 1.2
 - TAGS AT THE MIDDLEWARE LAYER
 - ALL TIMESTAMP INFORMATION IN MIDDLEWARE
 - GUARANTEES CAUSAL DELIVERY
 - BAD CHARACTERIZATION OF HAPPENED-BEFORE BETWEEN OPS
 - OVER ORDERING OF CONCURRENT OPS
 - SLOWER DELIVERY OF OPS
- SOLUTION 2.0
 - TAGS AT THE APPLICATION LAYER
 - APPLICATION NEEDS TO KEEP A TIMESTAMP
 - GUARANTEES CAUSAL DELIVERY
 - GOOD CHARACTERIZATION OF HAPPENED-BEFORE BETWEEN OPS
 - NO OVER ORDERING OF CONCURRENT OPS
 - FASTER DELIVERY OF OPS

Questions ?

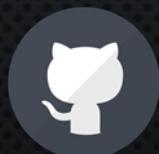
Slides: <https://bit.ly/tagged-causal>

Erlang implementation:

- Reliable Causal Broadcast: <https://github.com/gyounes/RCB>
- Tagged Reliable Causal Broadcast: https://github.com/gyounes/trcb_base



@g_unis



gyounes