

DISTRIBUTED UPDATES

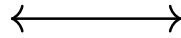
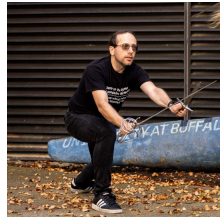
CSE 4/562: Database Systems | Lecture 23

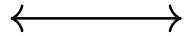
DB. Sys.: T.C.B.: Ch. 17.1-17.5

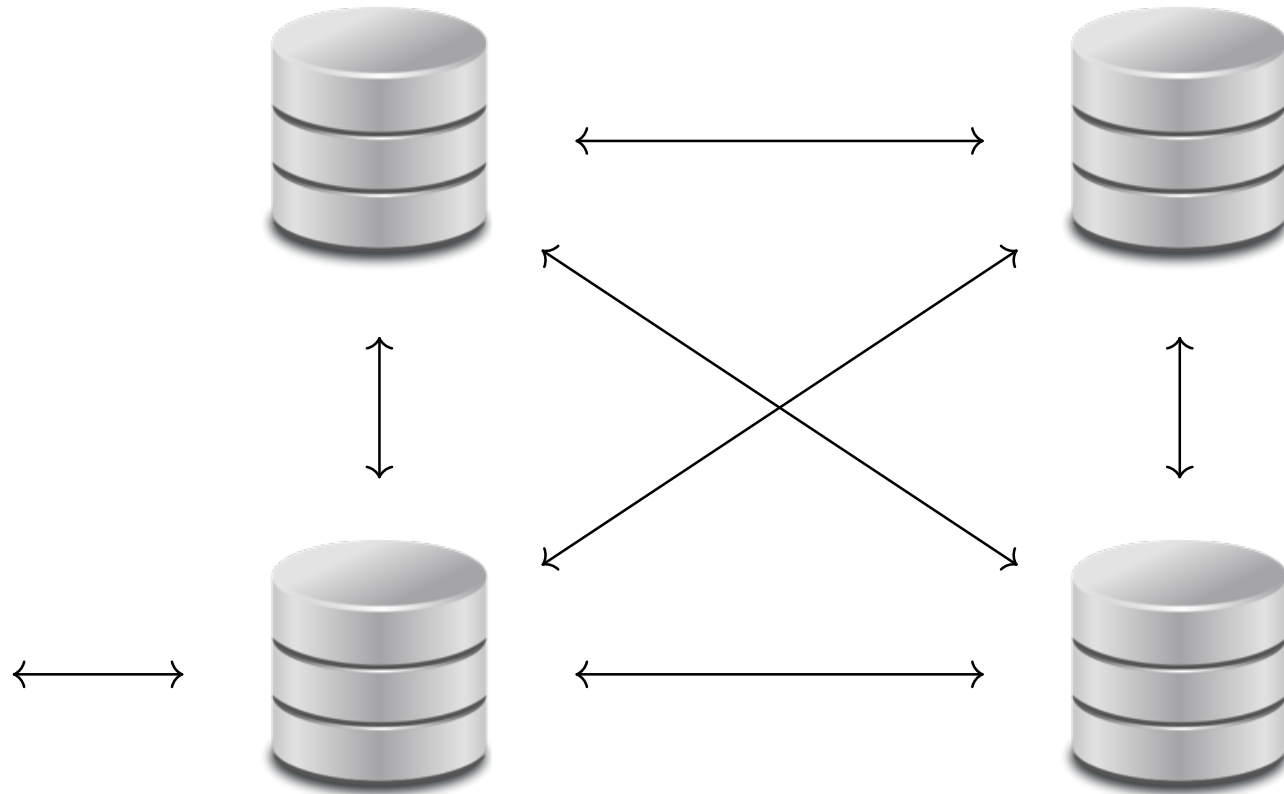
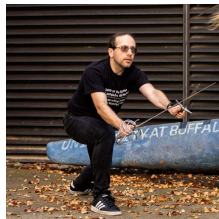
- Checkpoint 4 due Mon
- Course Reviews 0/31
 - At 28/31 I will publish a question worth no less than 10/100 pt from the final
- Final Exam
 - Will replace midterm if better

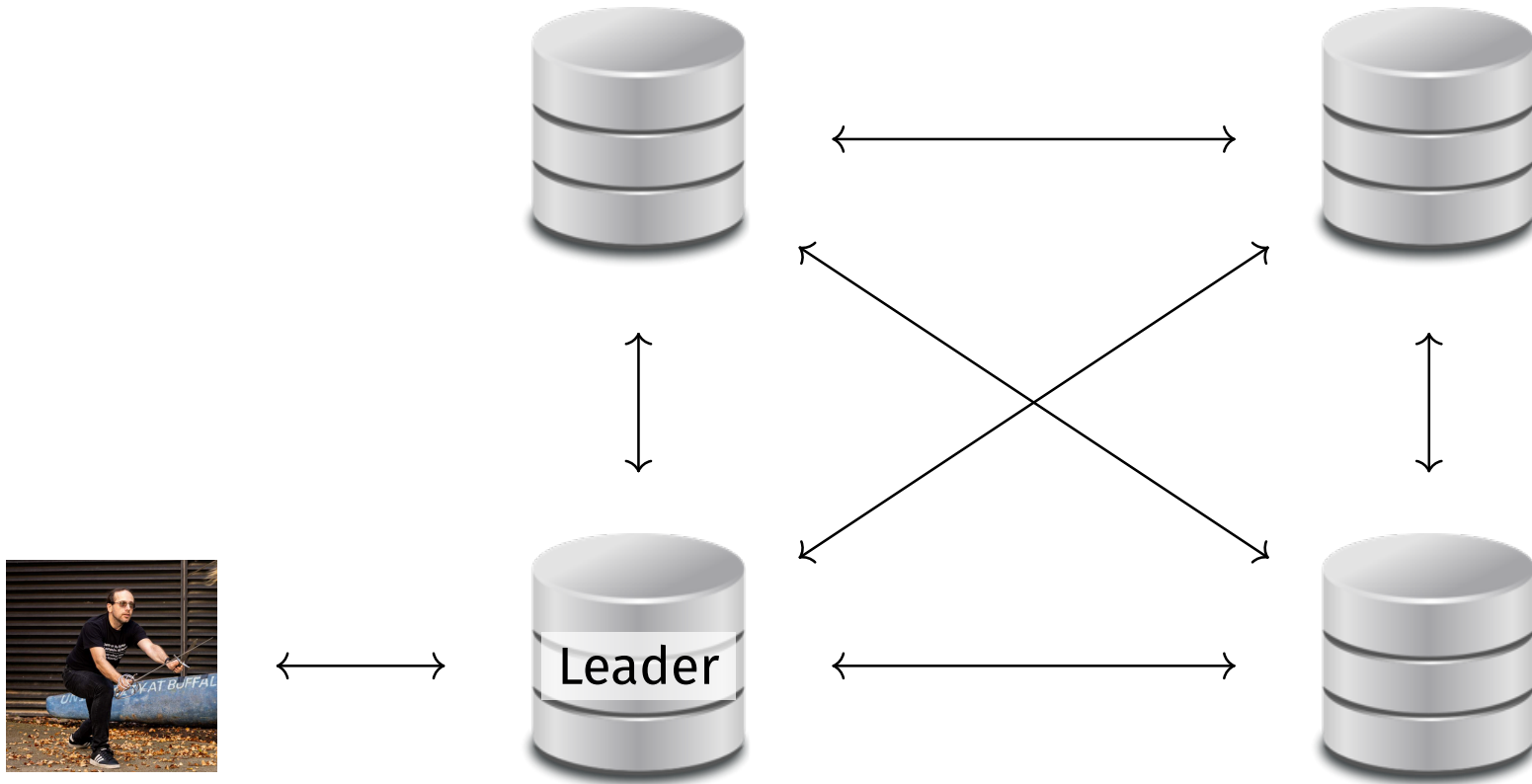
**How do we preserve ACID
in distributed databases.**

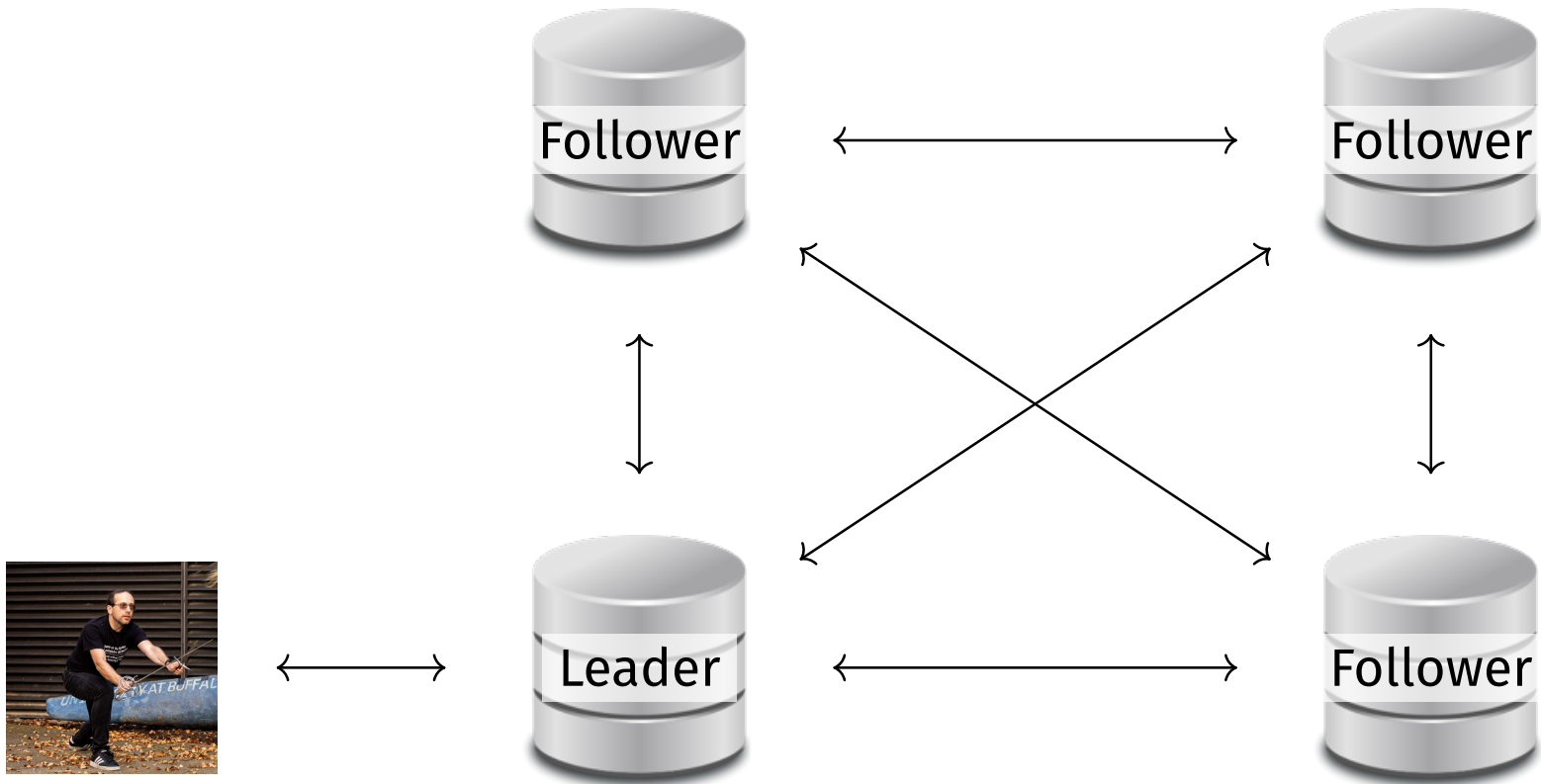


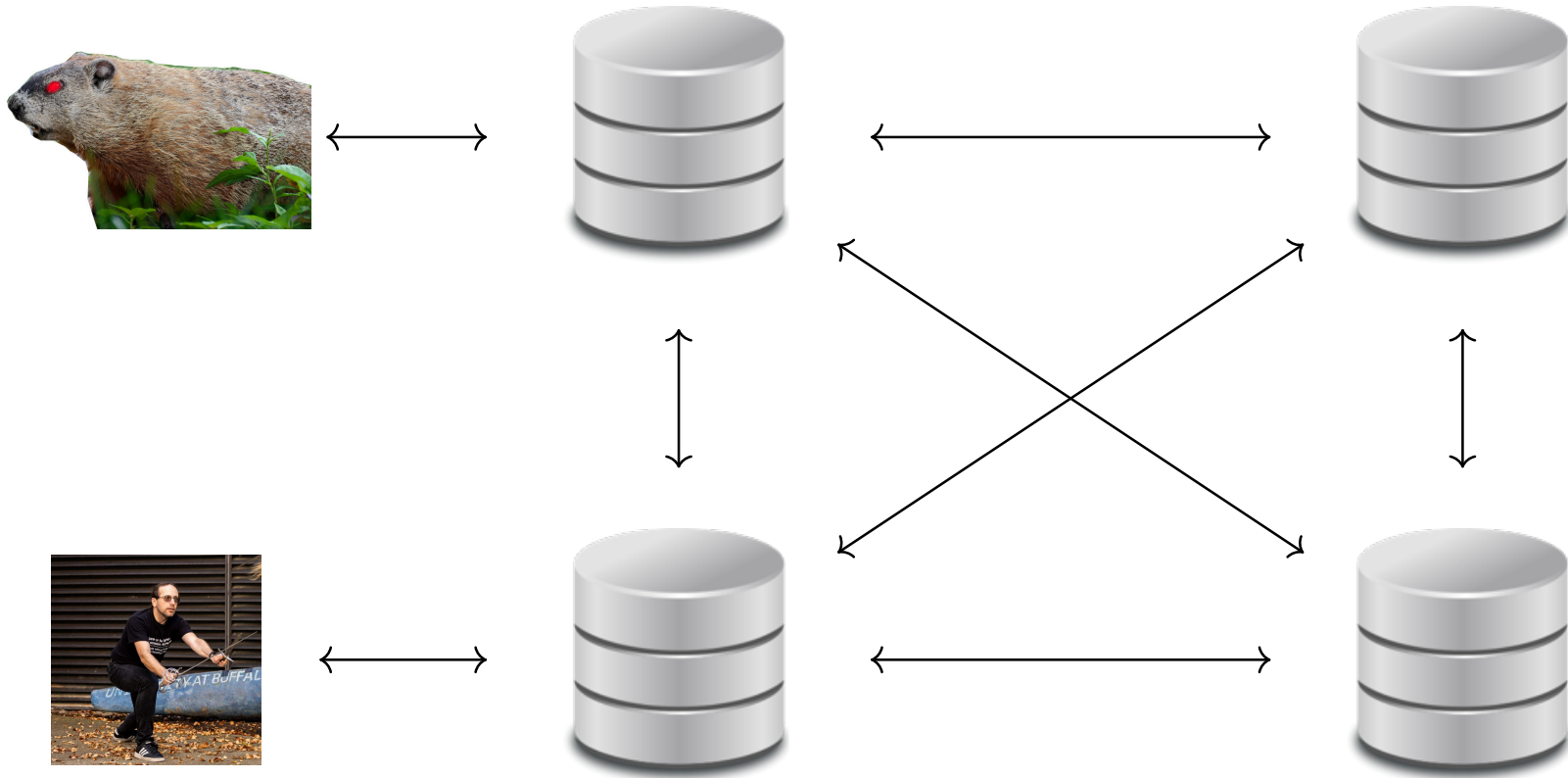












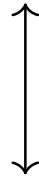
1. Each node can enforce ACID **locally**.
2. No malicious nodes.
3. Messages arrive in full or not at all.
4. No global state.

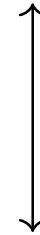
**Each node has to make
local decisions without
global state.**

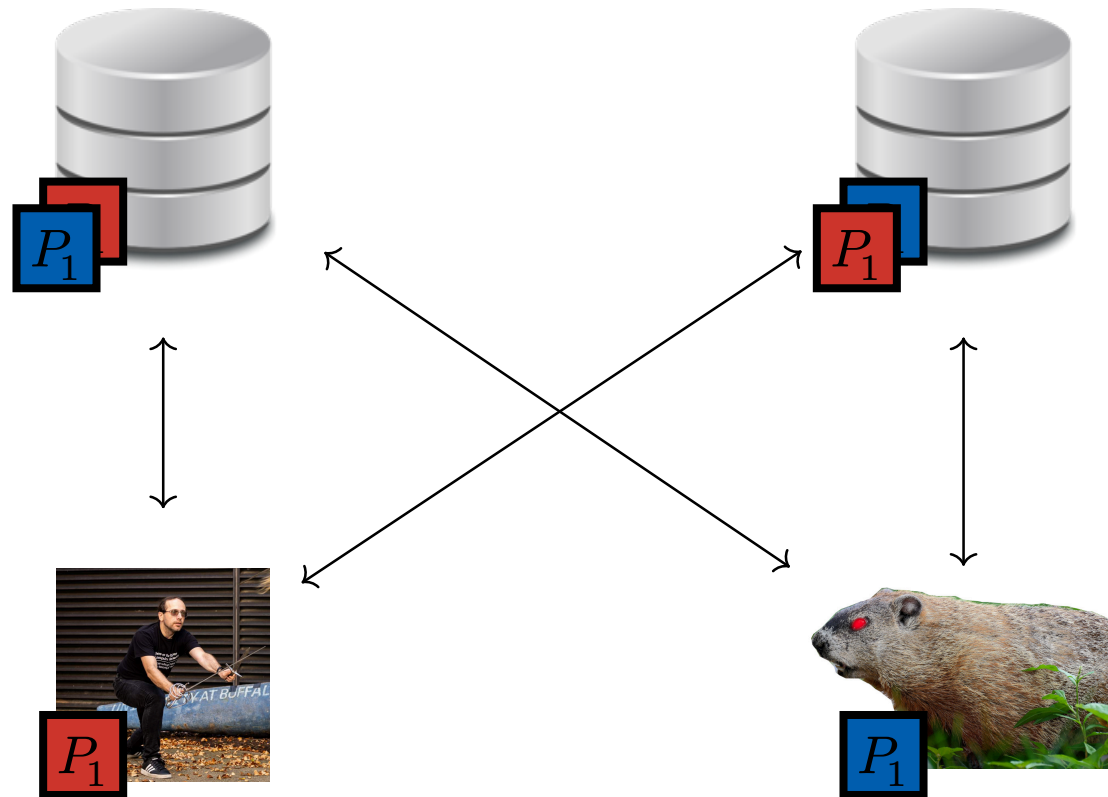


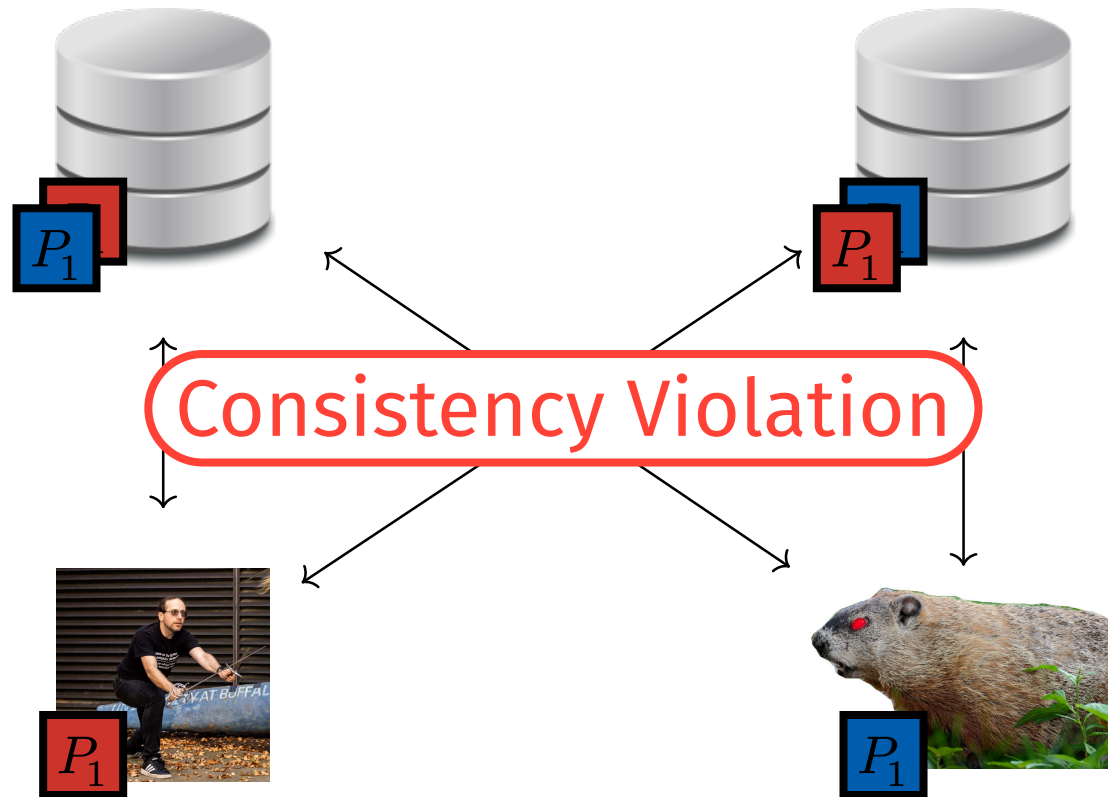
Image credit: Openclipart.org











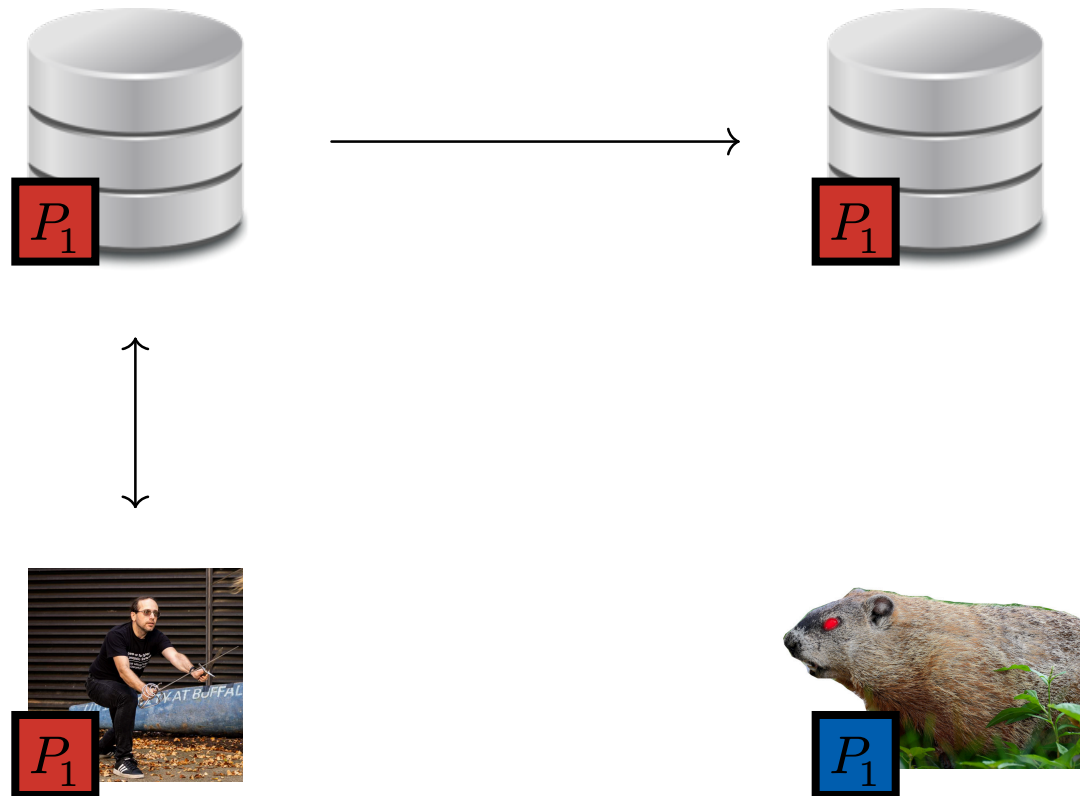
Isolation Violation

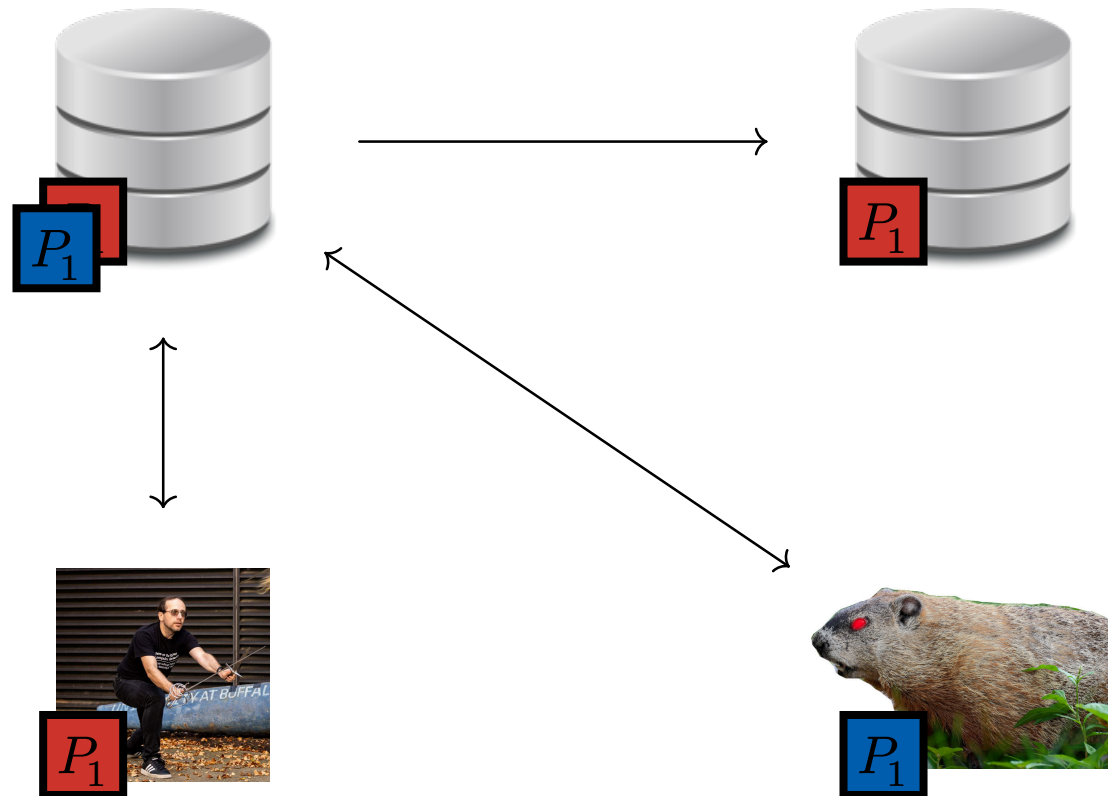
Nodes might see transaction requests in a different order.

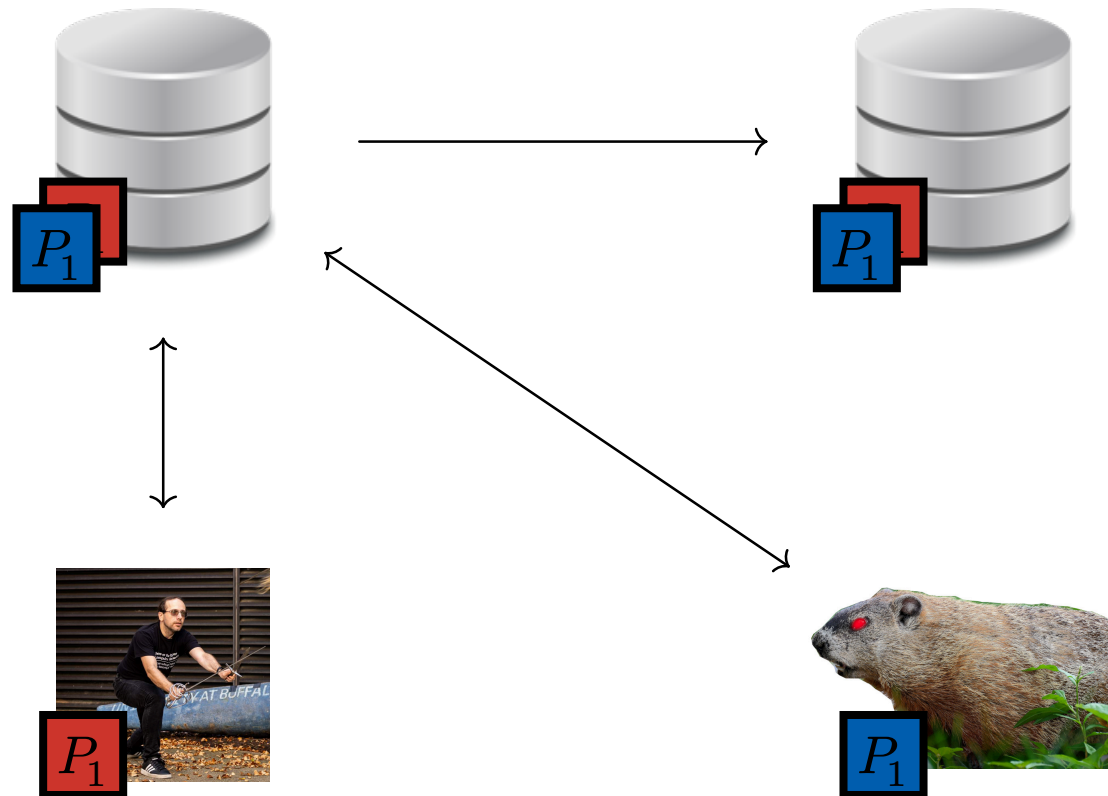
Idea 1: Route all requests through one node.











Pros

Cons

Pros

- Simple

Cons

Pros

- Simple
- Guaranteed to be correct (DB 2 sees exactly what DB 1 sees)

Cons

Pros

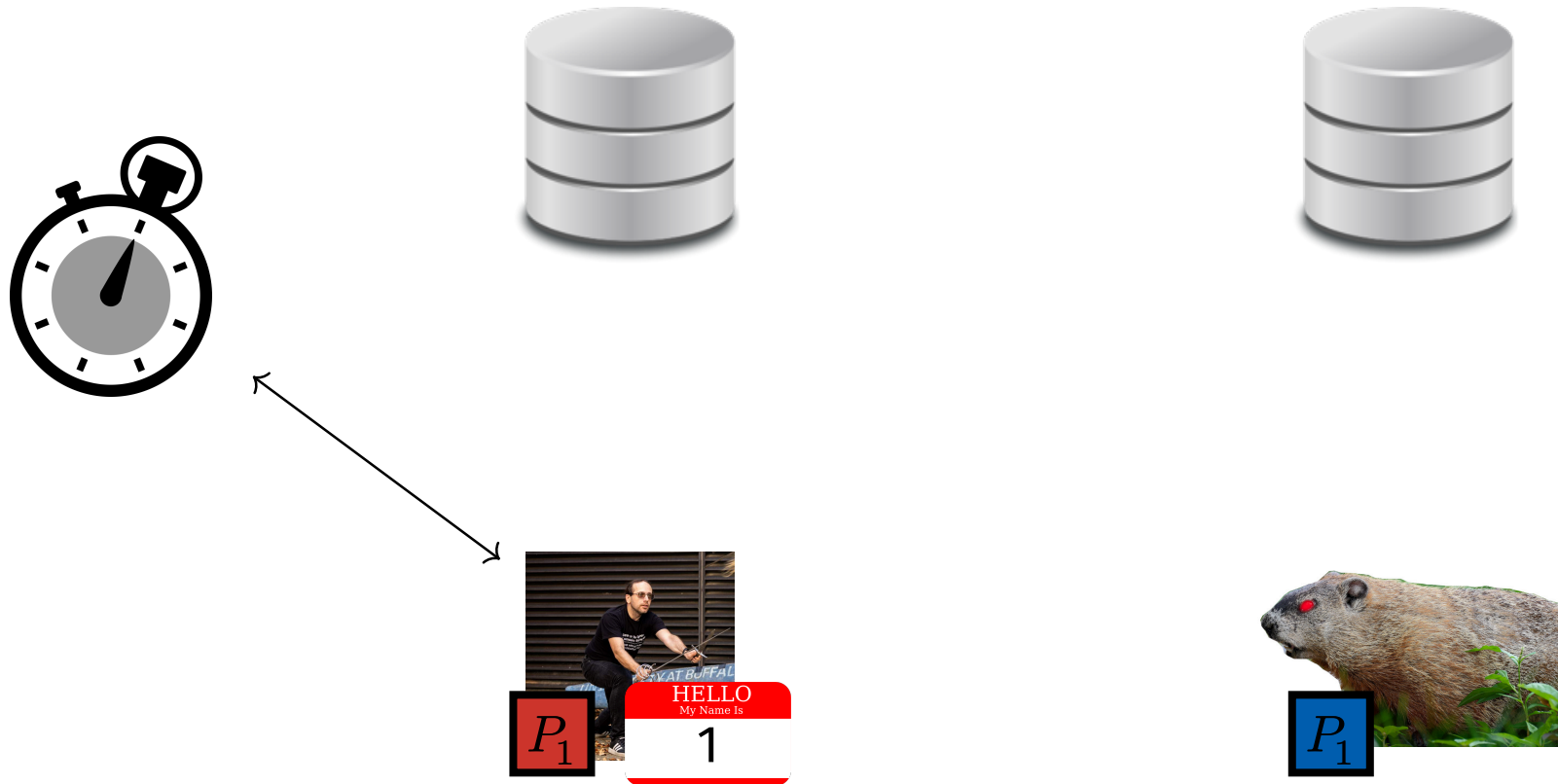
- Simple
- Guaranteed to be correct (DB 2 sees exactly what DB 1 sees)

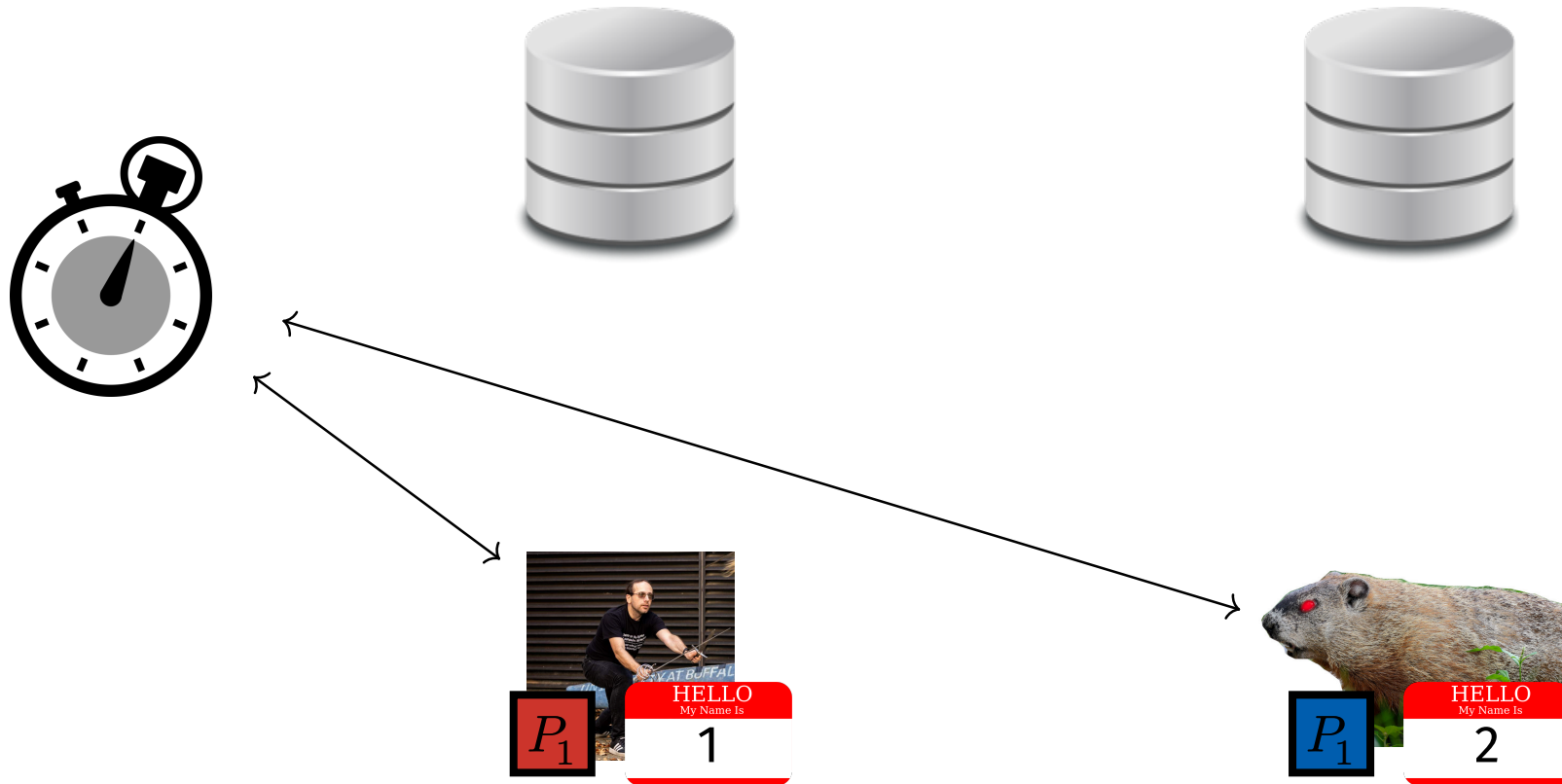
Cons

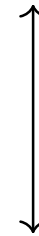
- Single Node is a Bottleneck

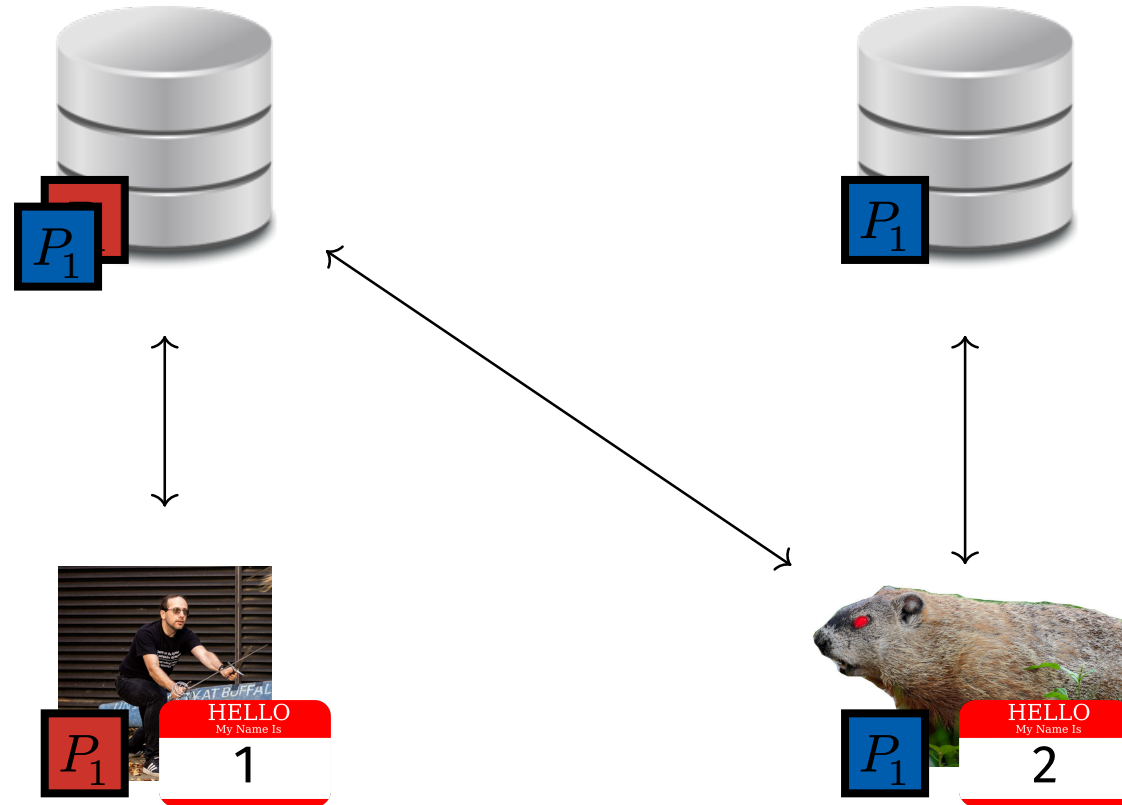
Idea 1: One node defines a serial order.

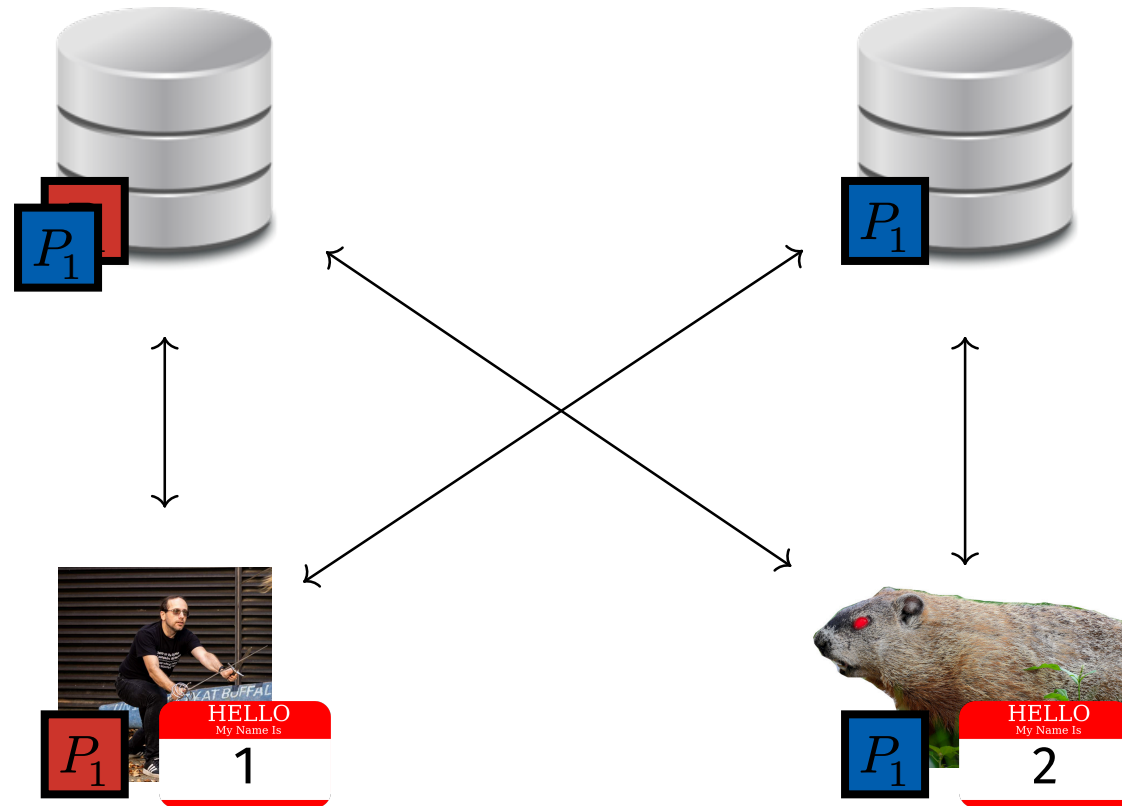


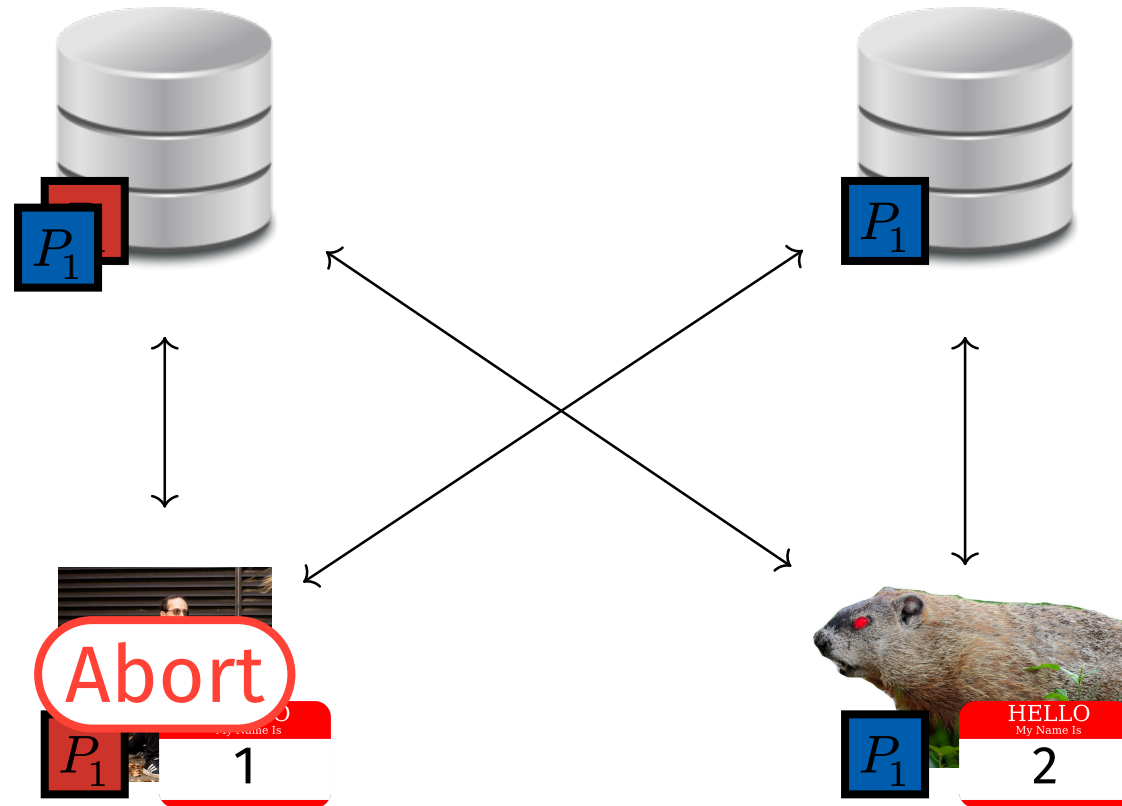












Pros

Cons

Pros

- Serving timesetamps is fast

Cons

Pros

- Serving timesetamps is fast
- No other bottlenecks

Cons

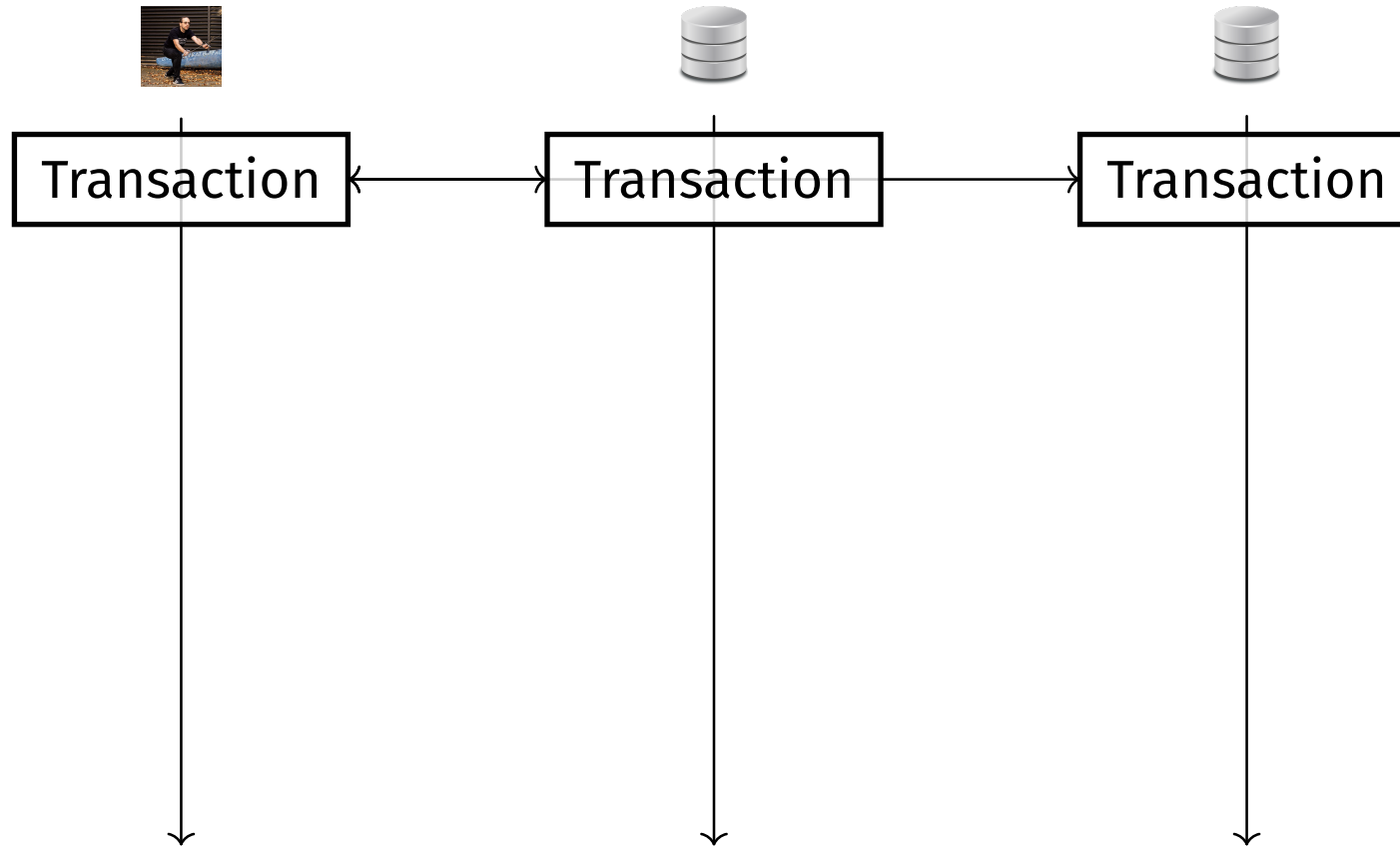
Pros

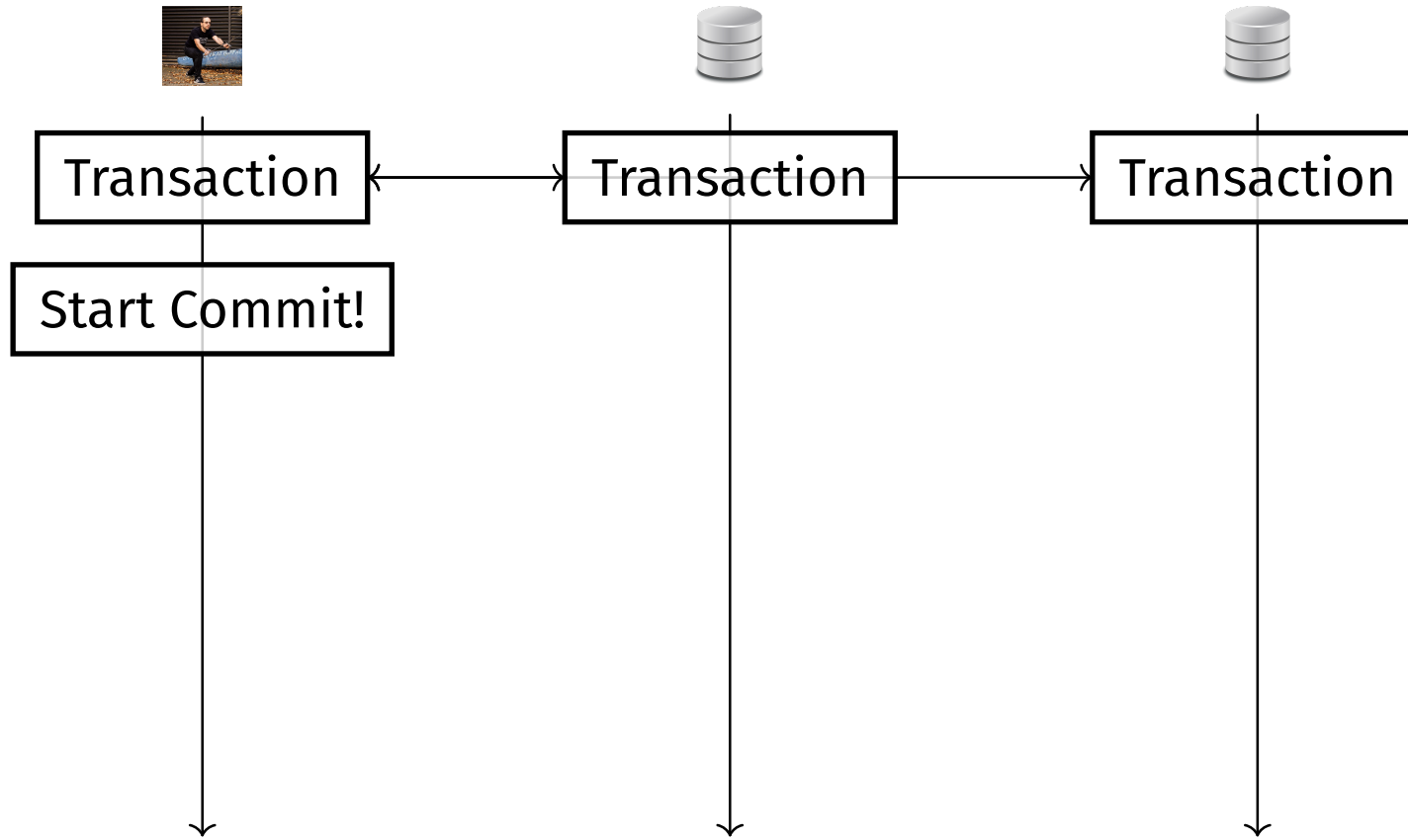
- Serving timestamps is fast
- No other bottlenecks

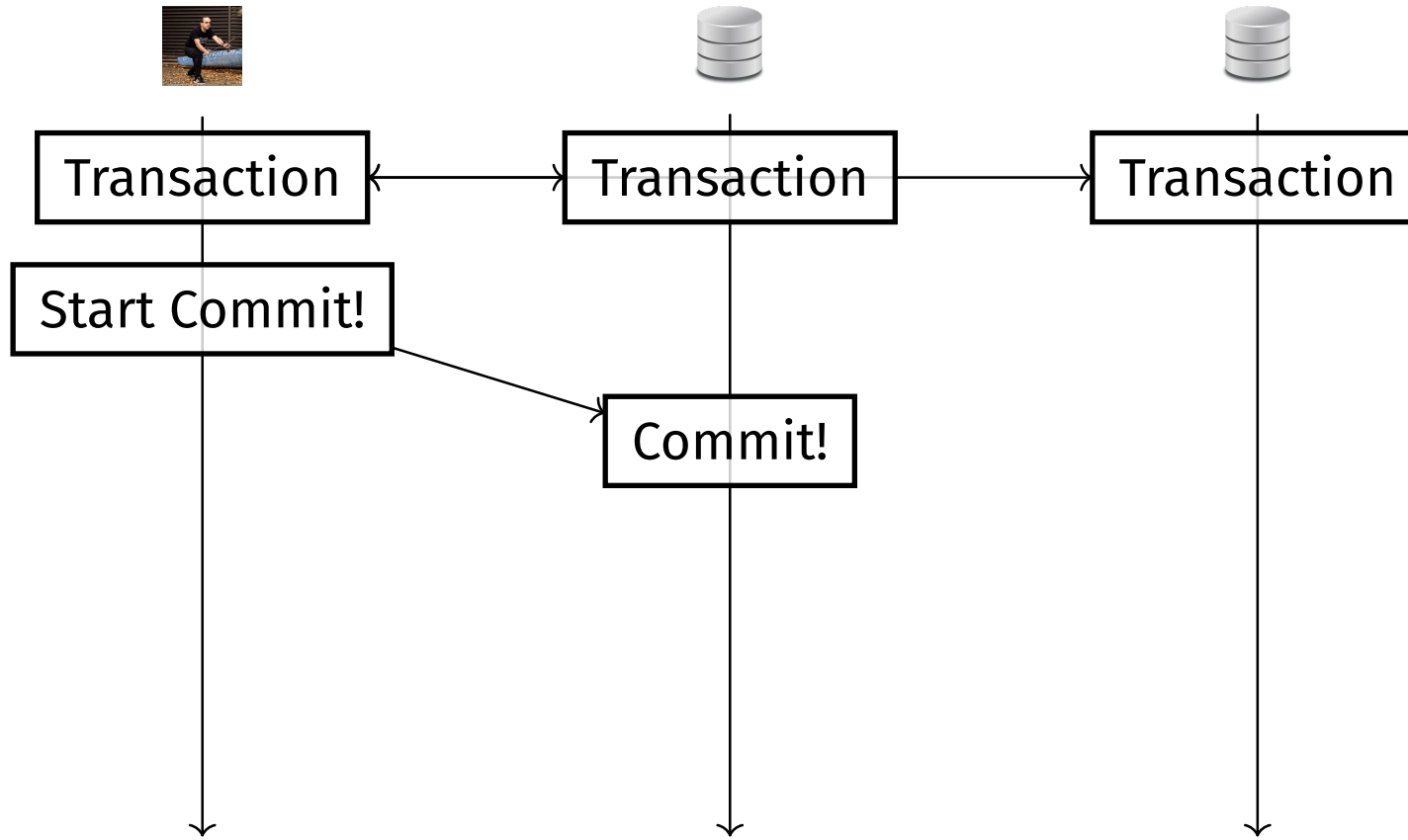
Cons

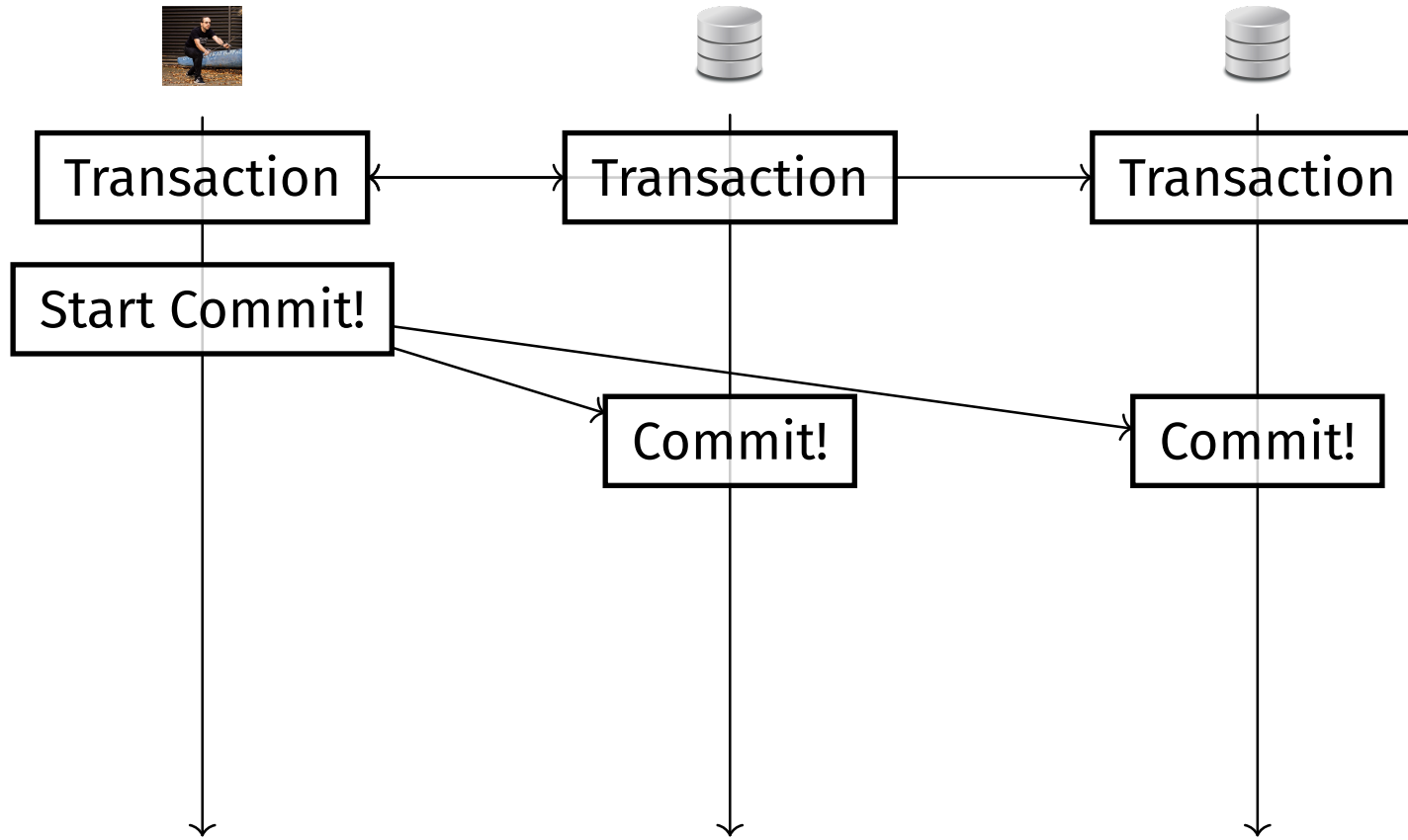
- When is a transaction safe to commit?

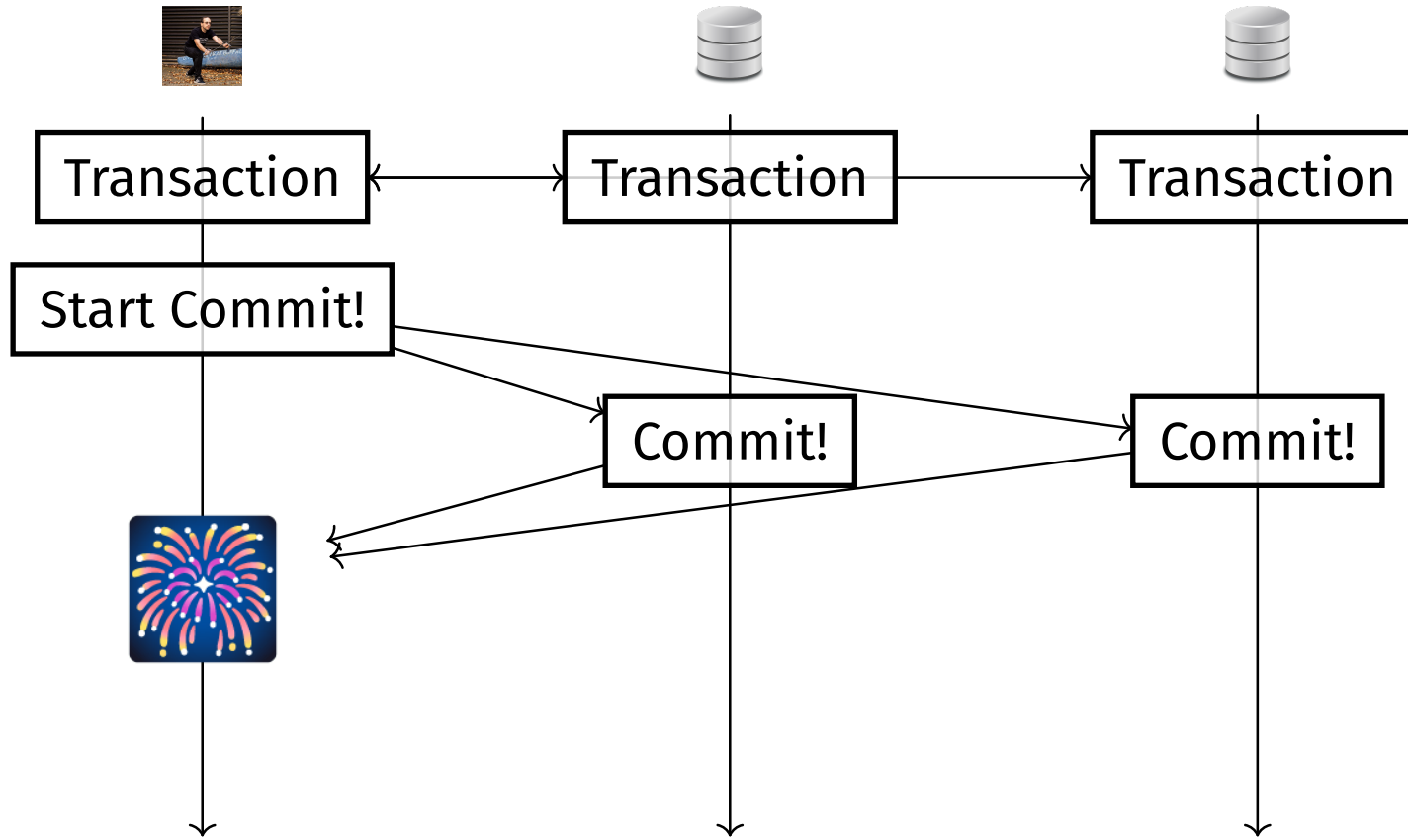




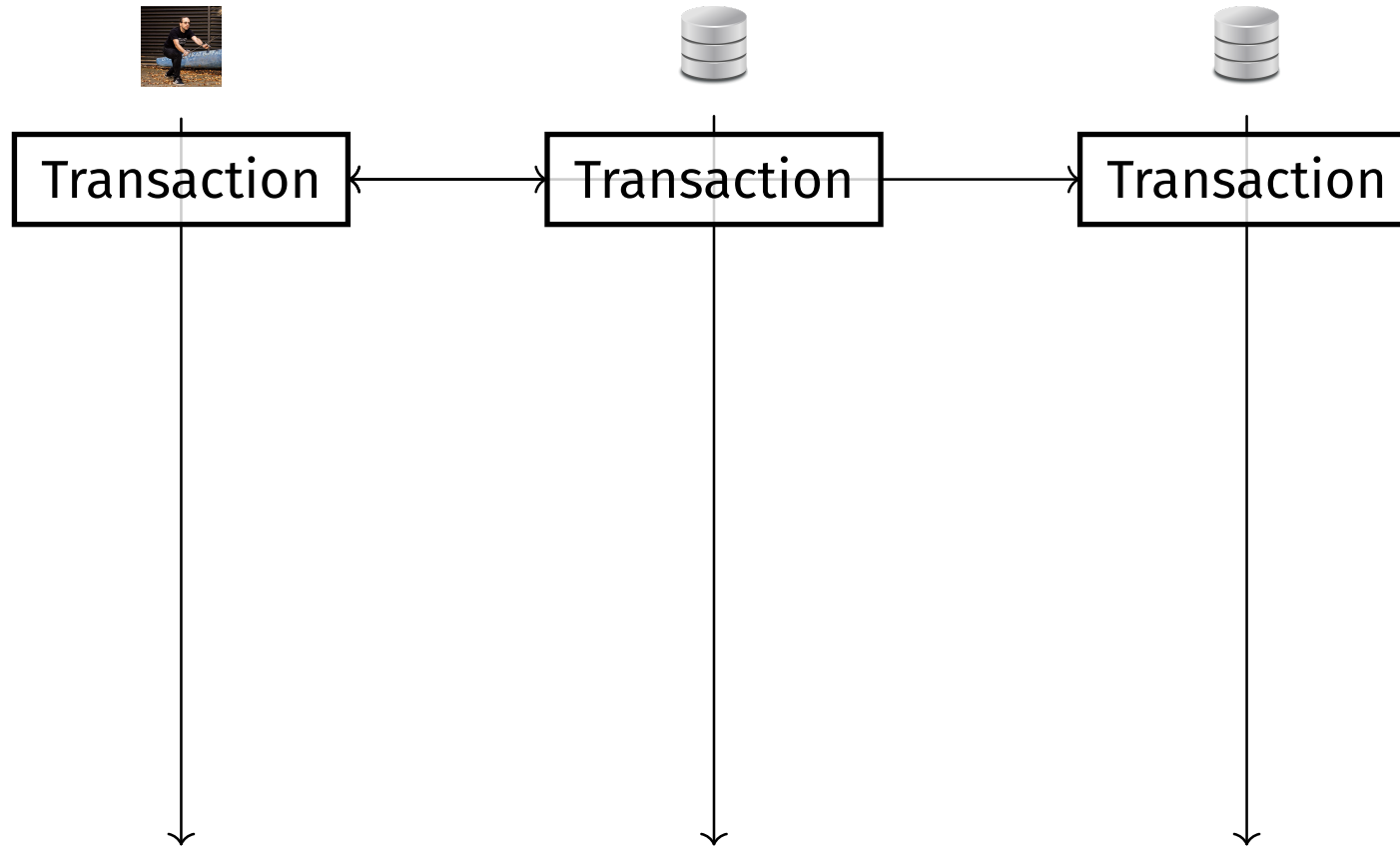


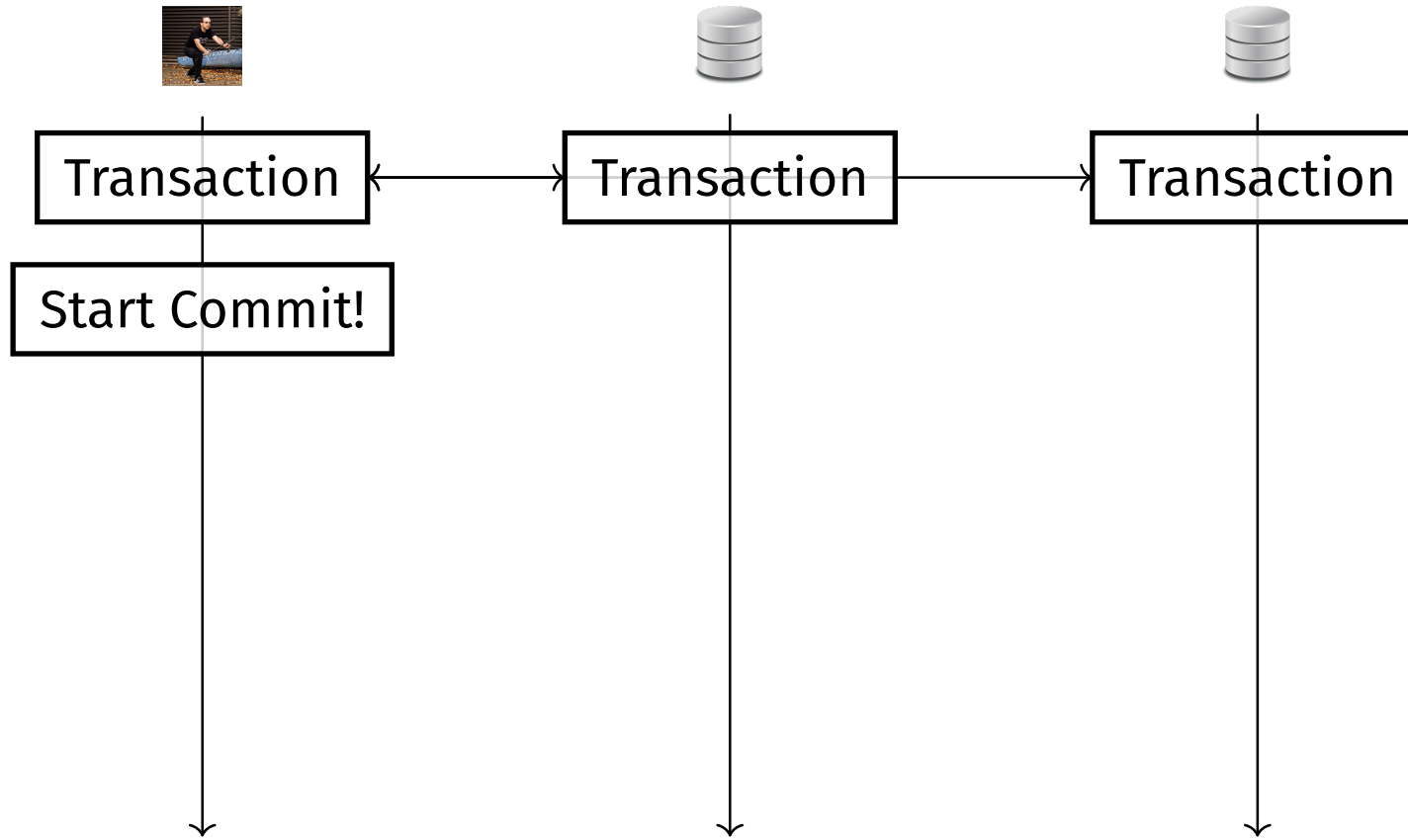


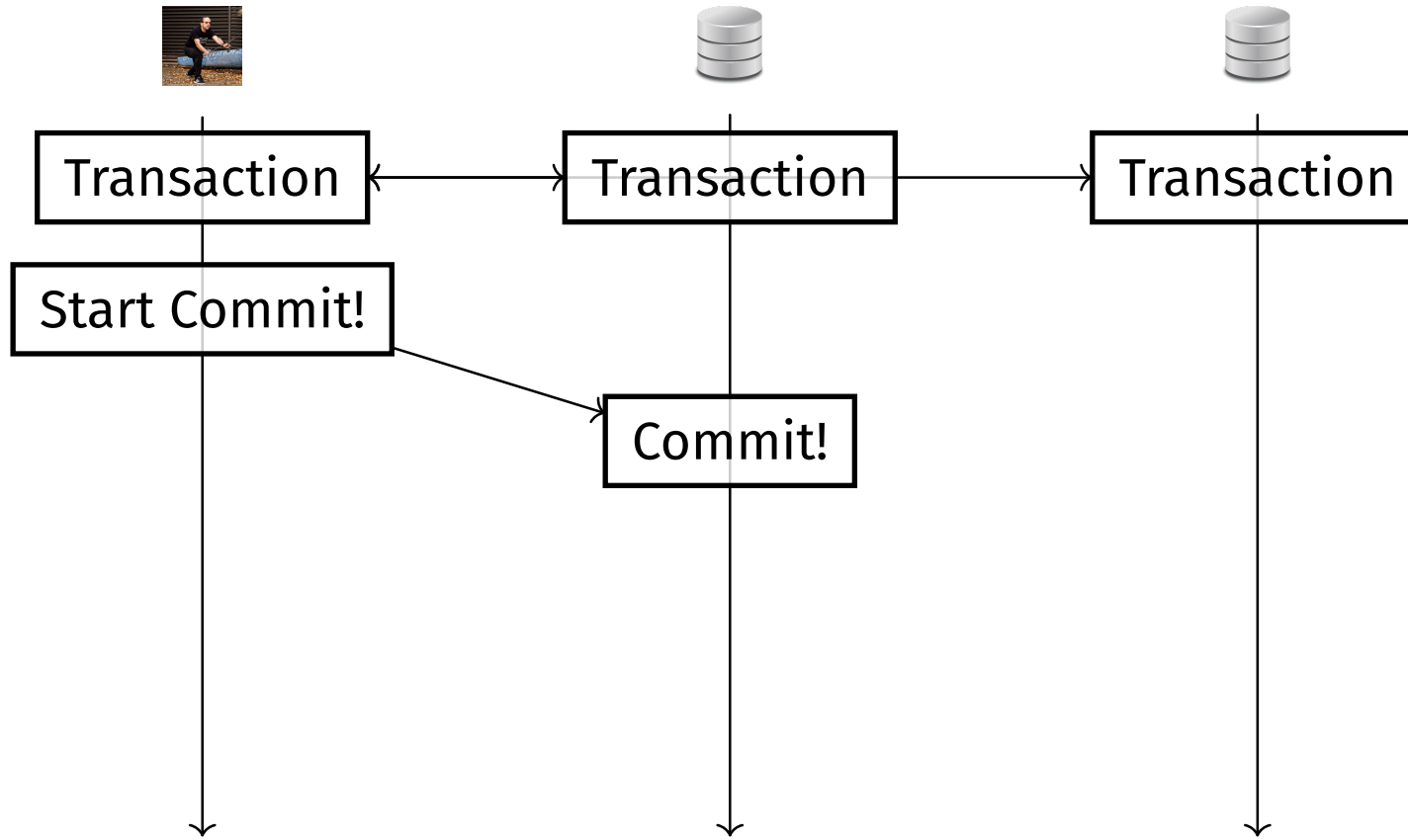


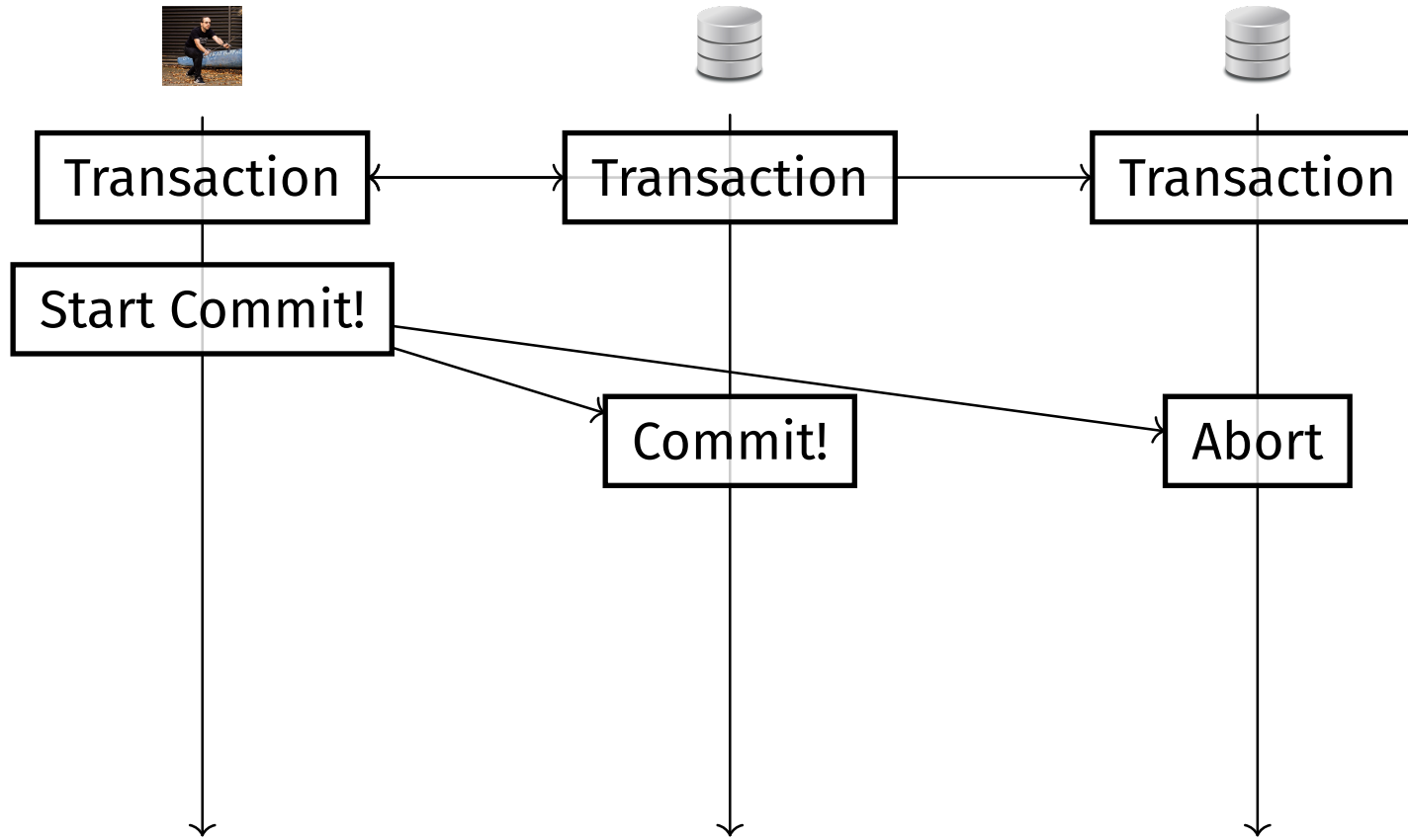


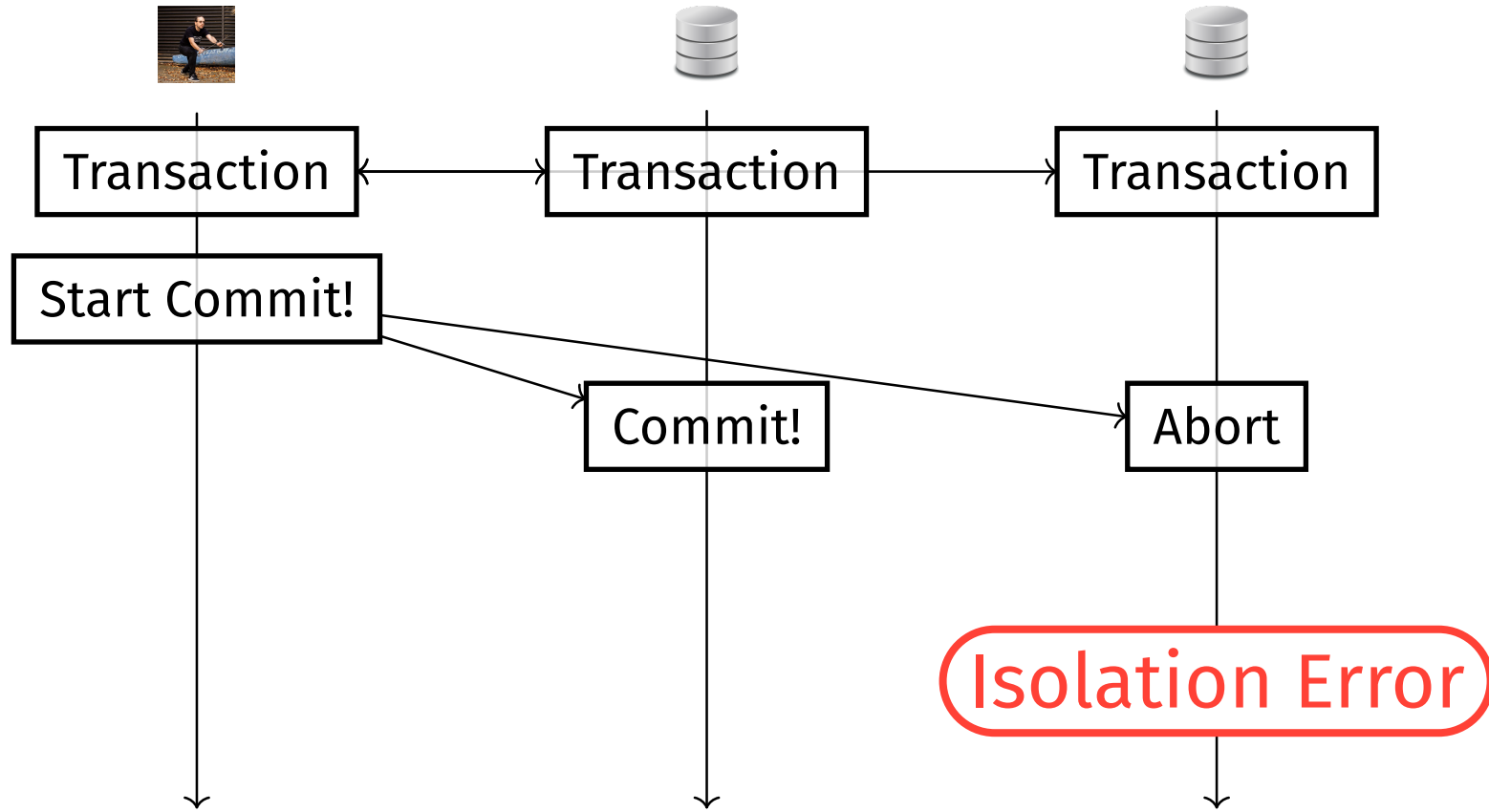








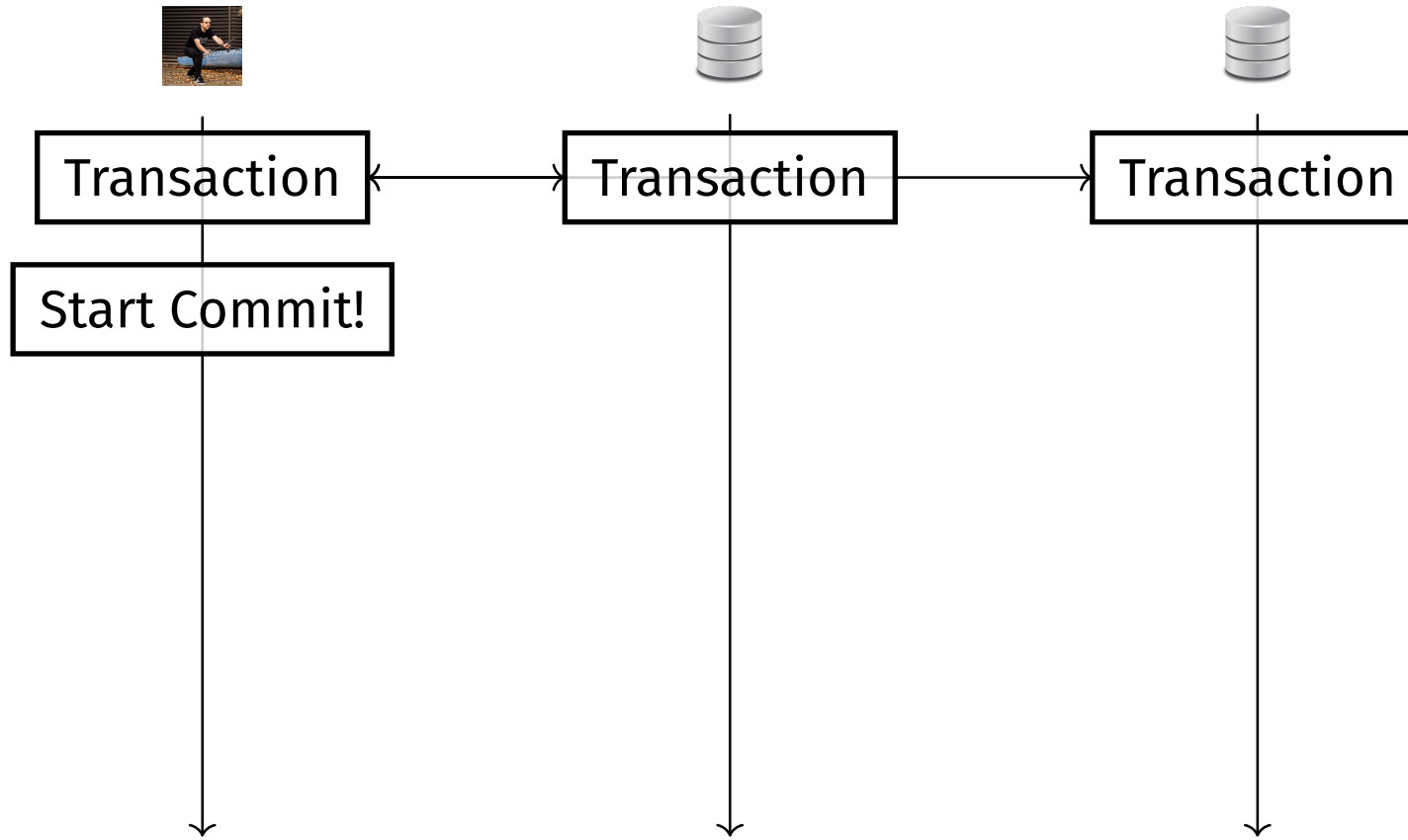


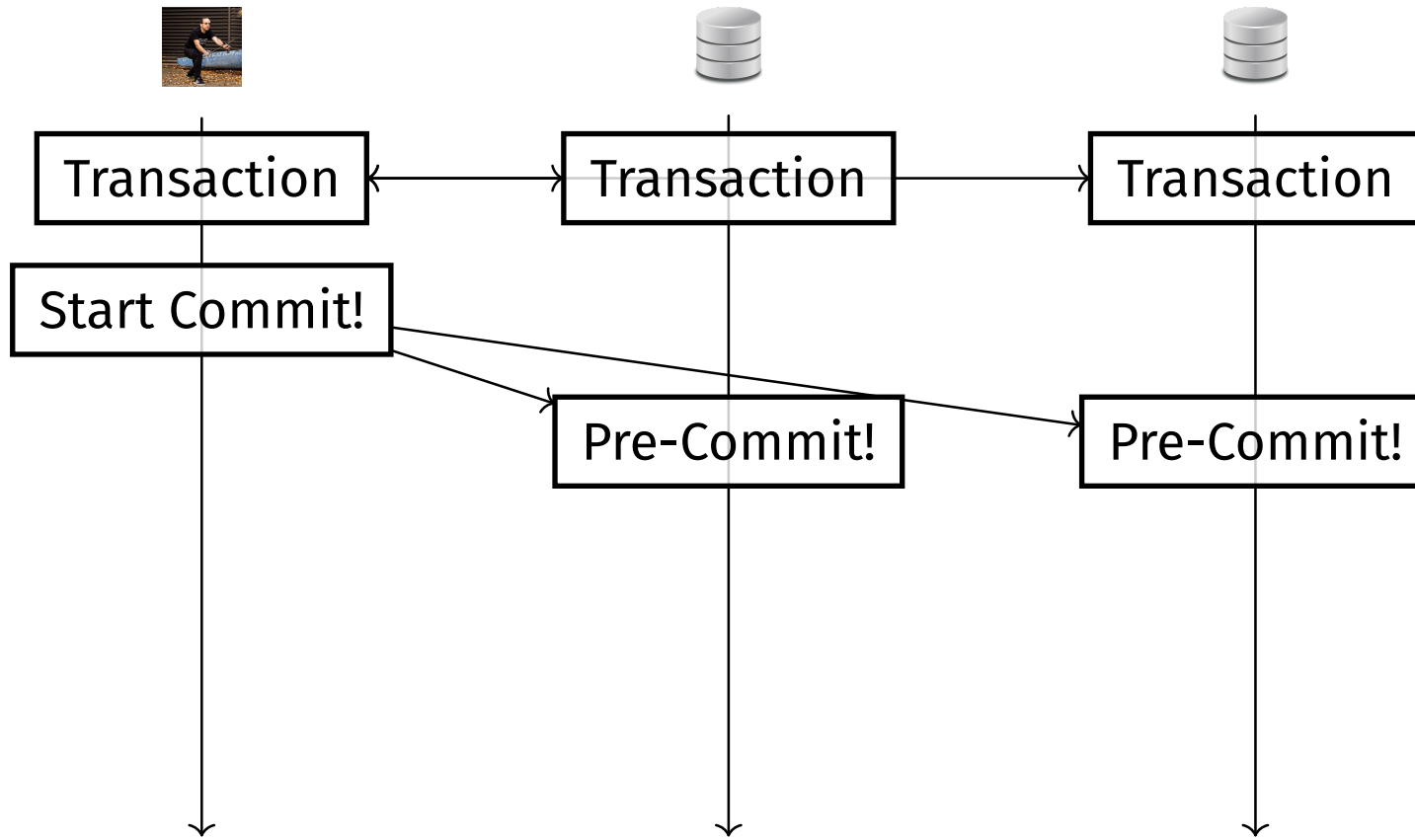


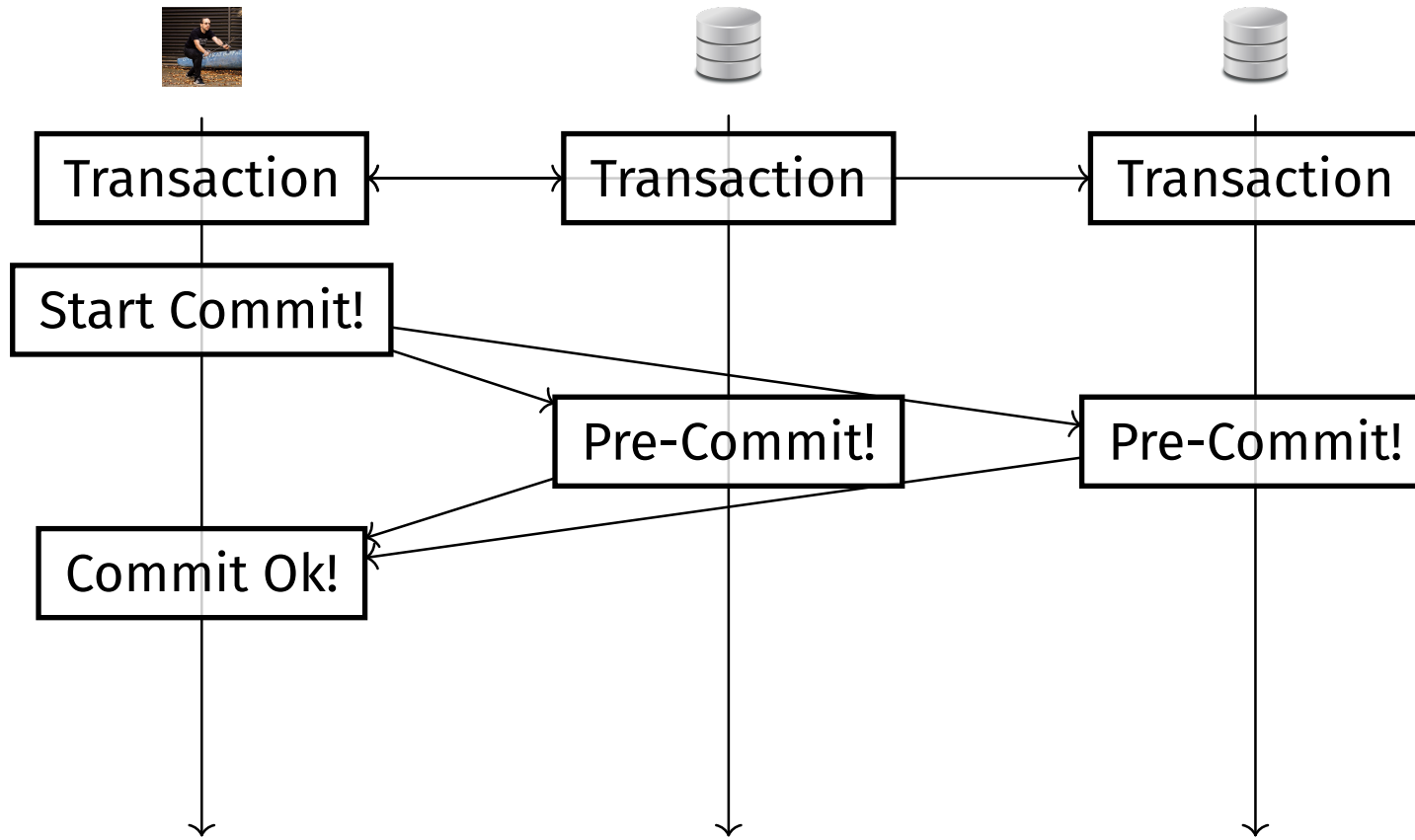
Multiple things could break during commit

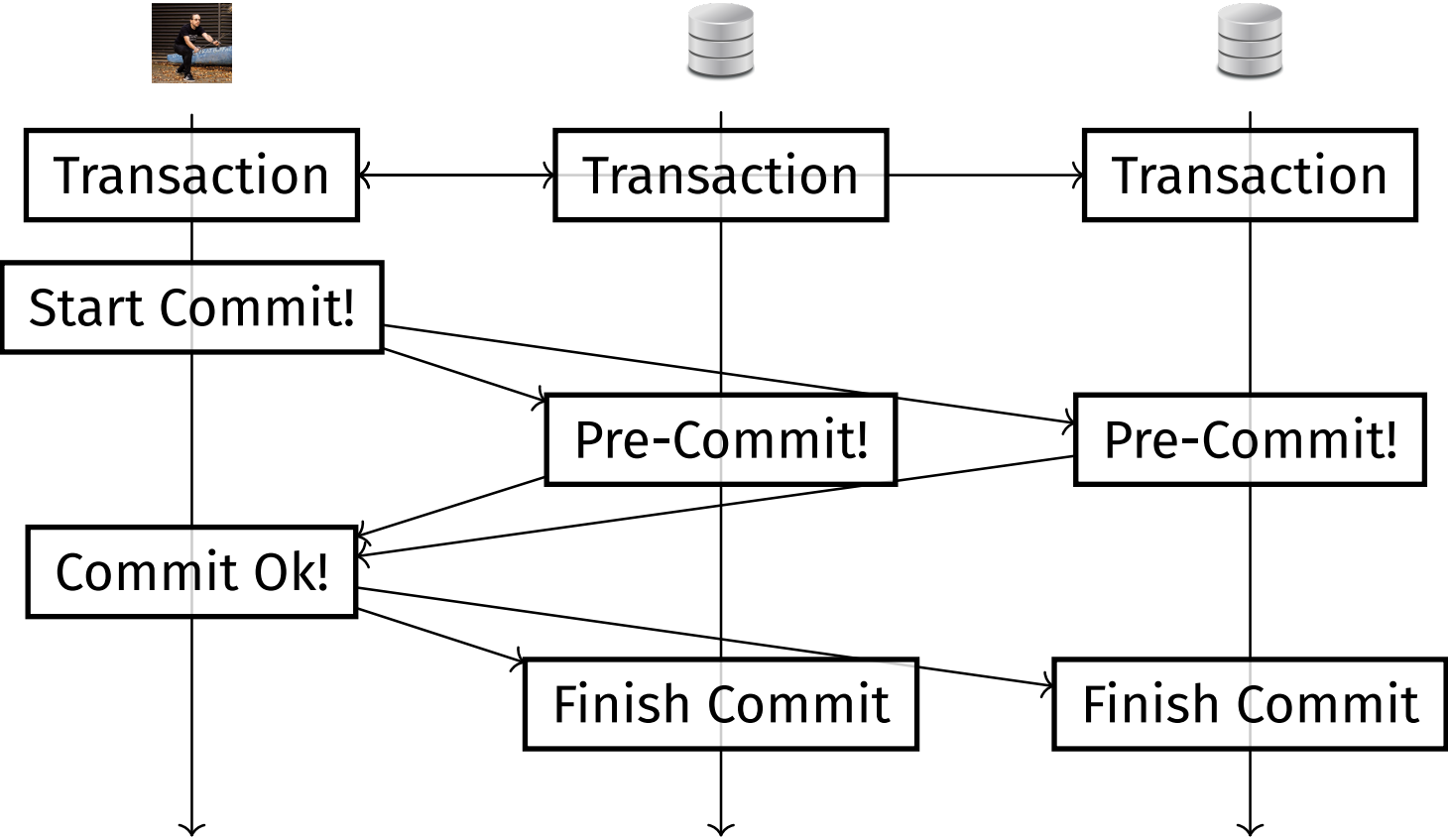
1. Optimistic concurrency control detects a concurrency violation
2. Hardware failure violates durability
3. Network failure prevents commits.

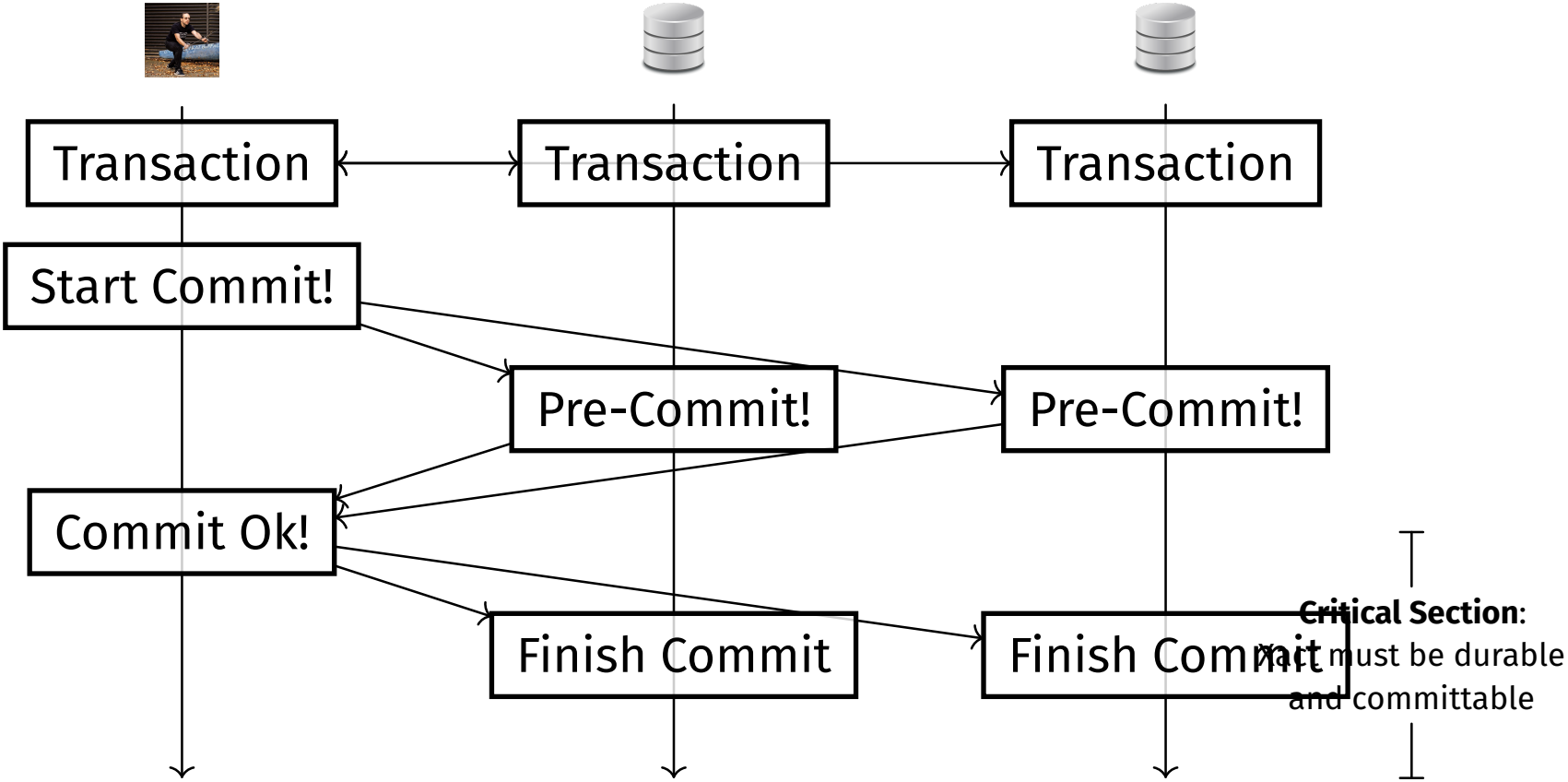
Idea: Add a validation step.

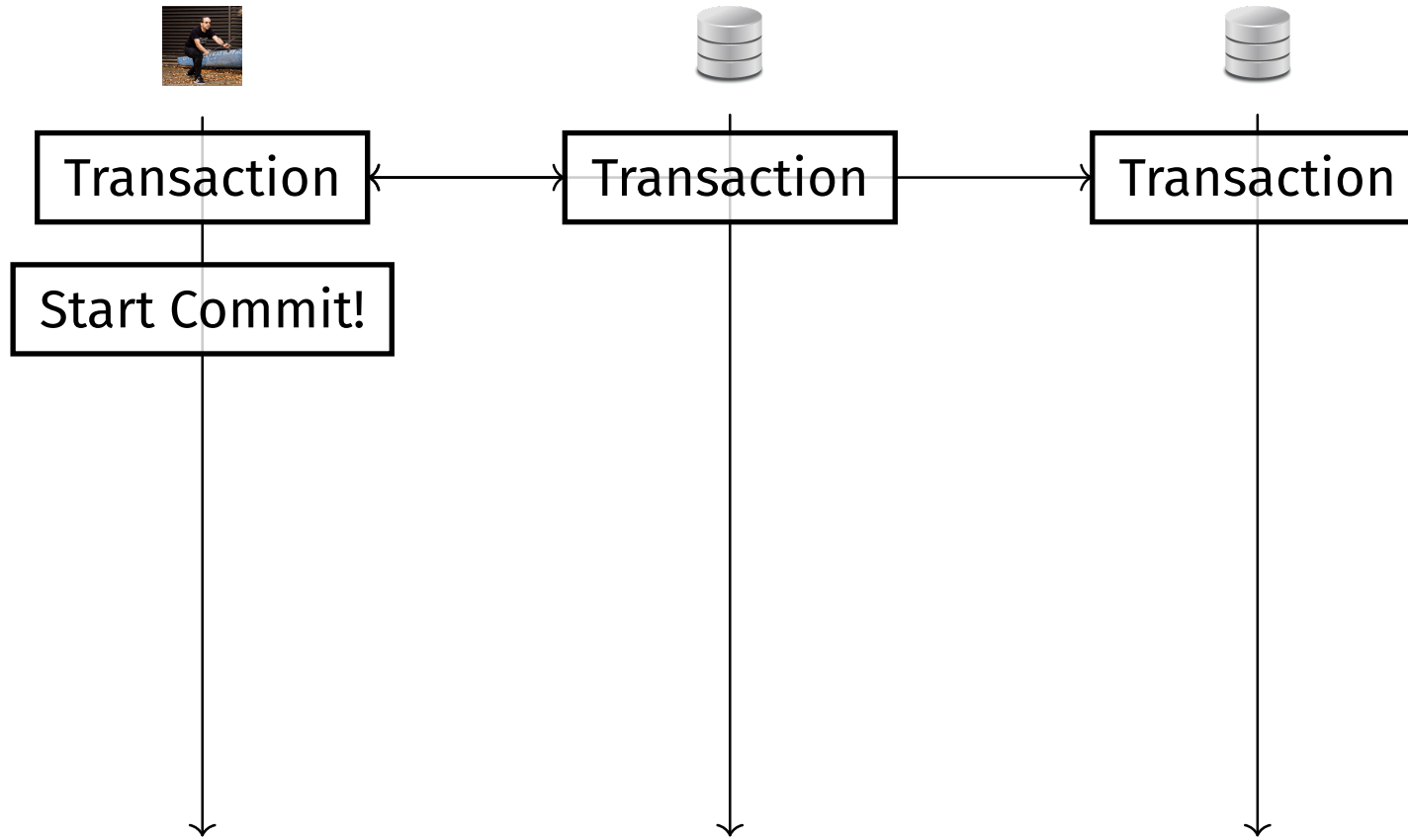


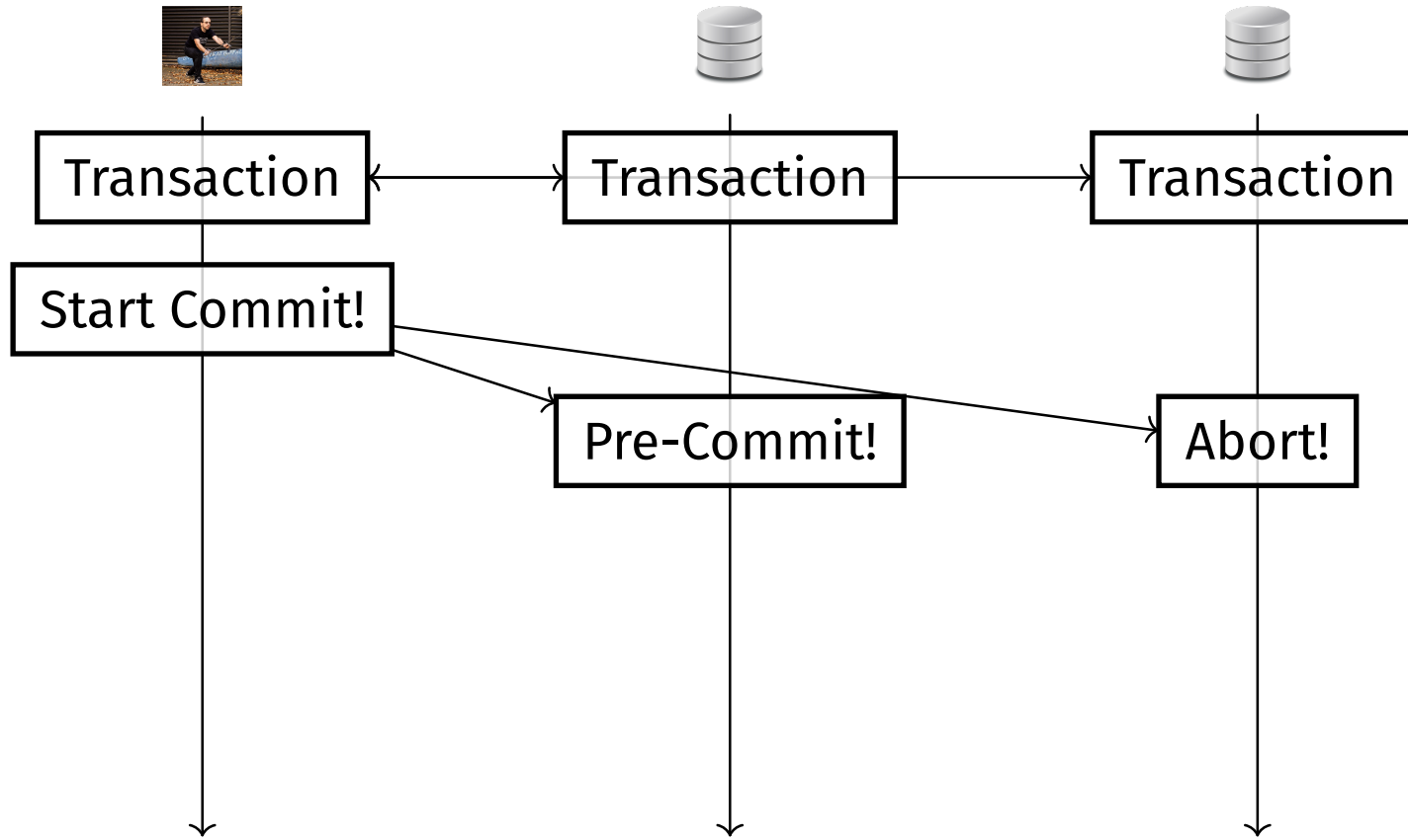


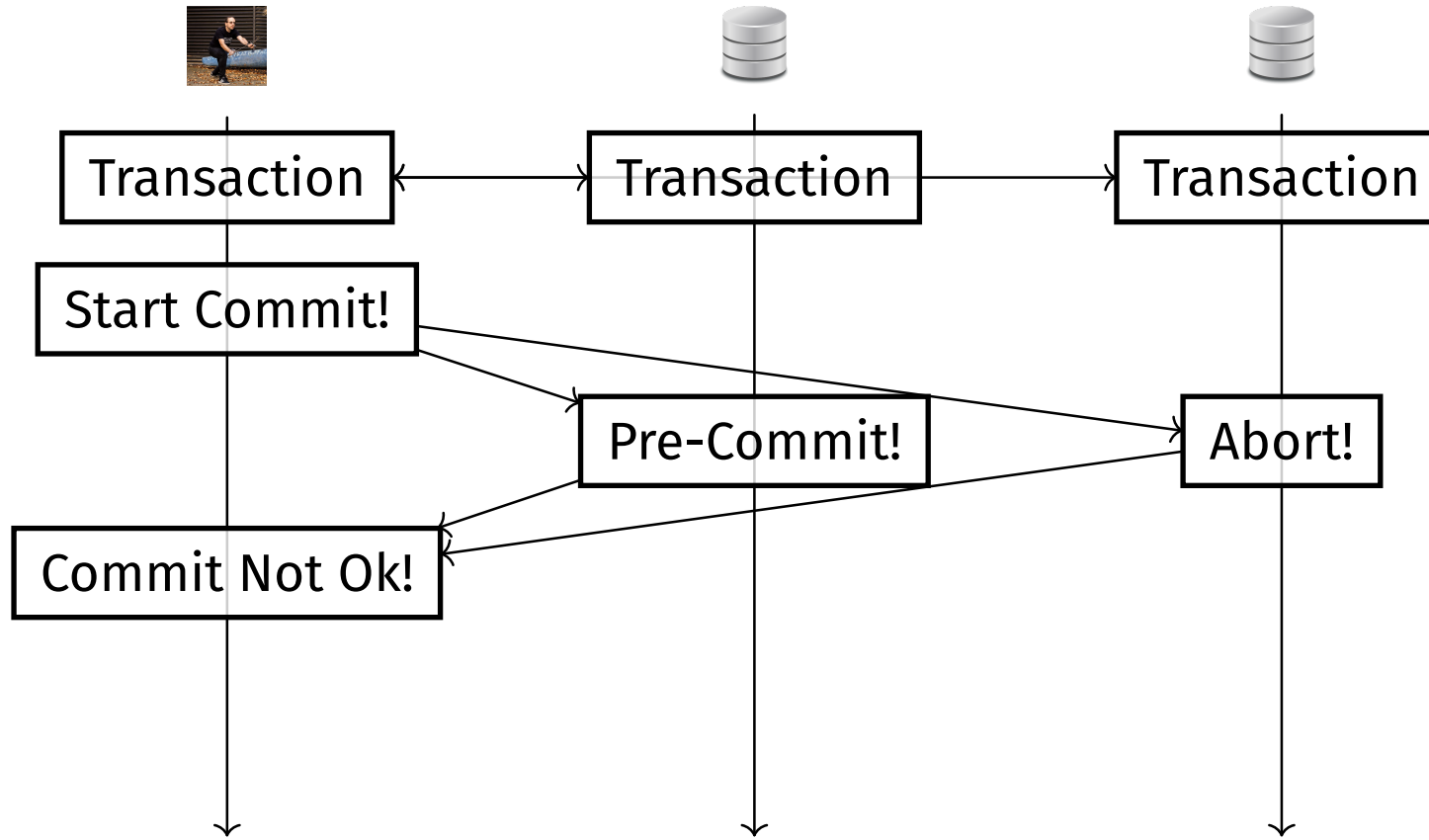


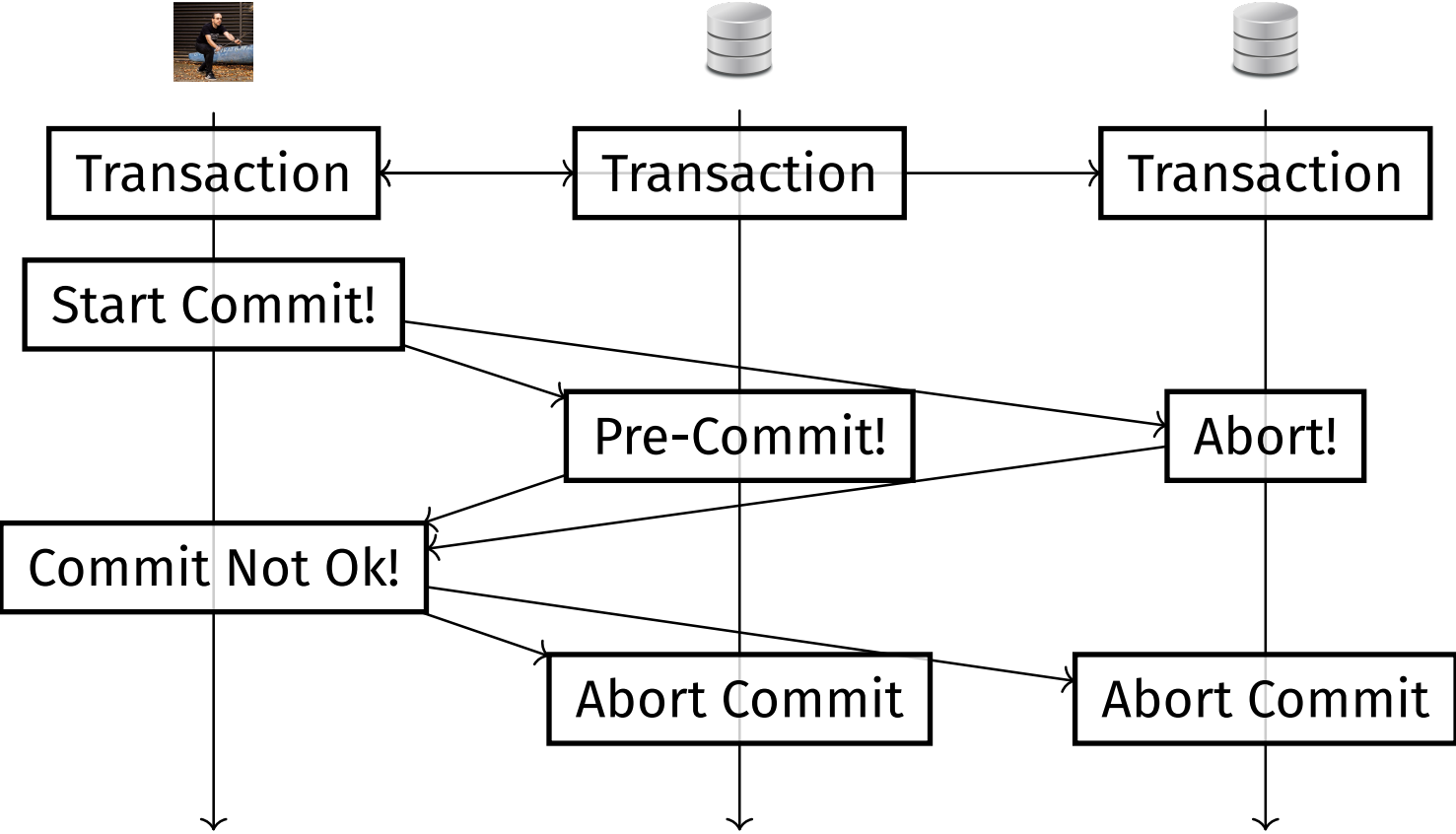












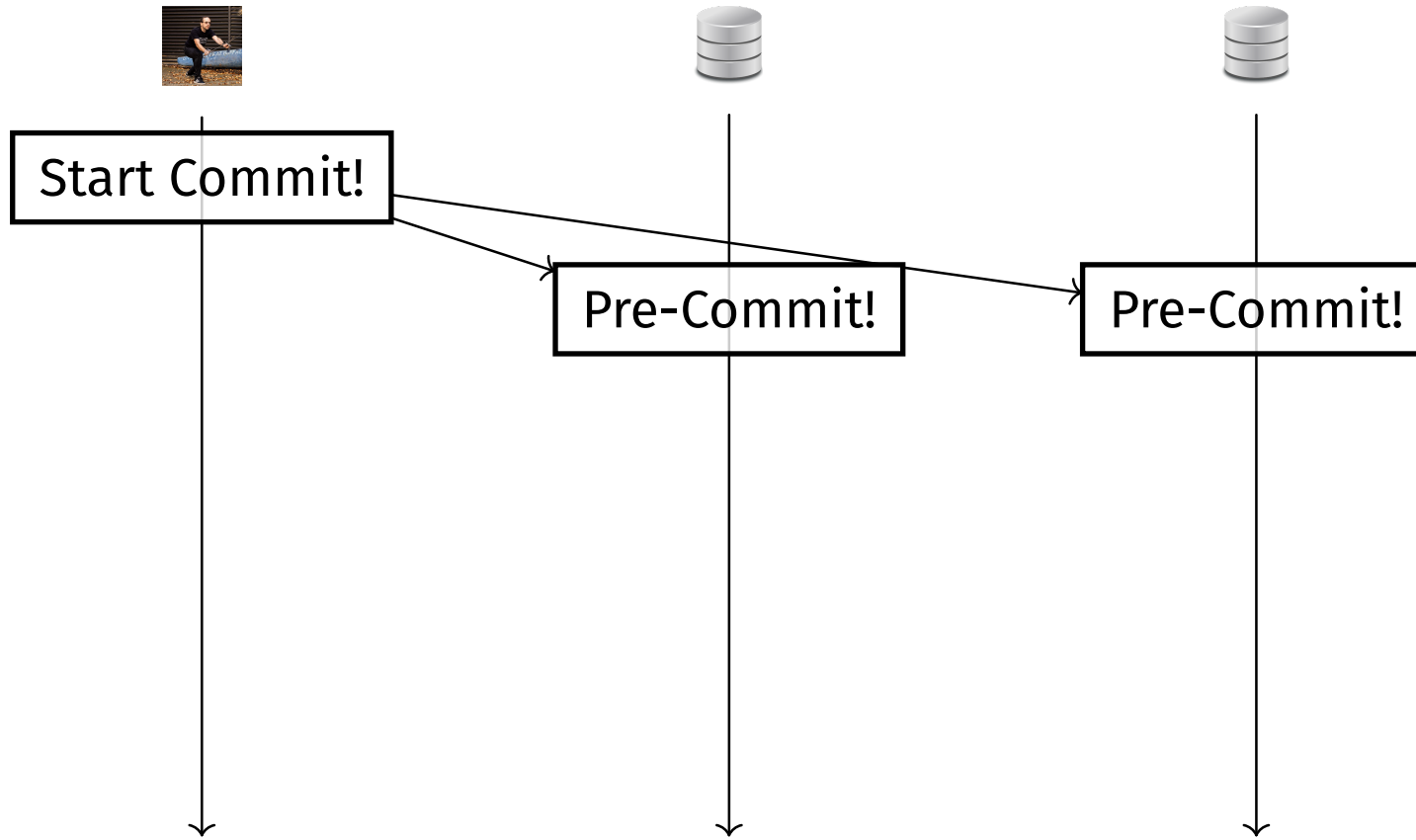
This is called 2-Phase Commit

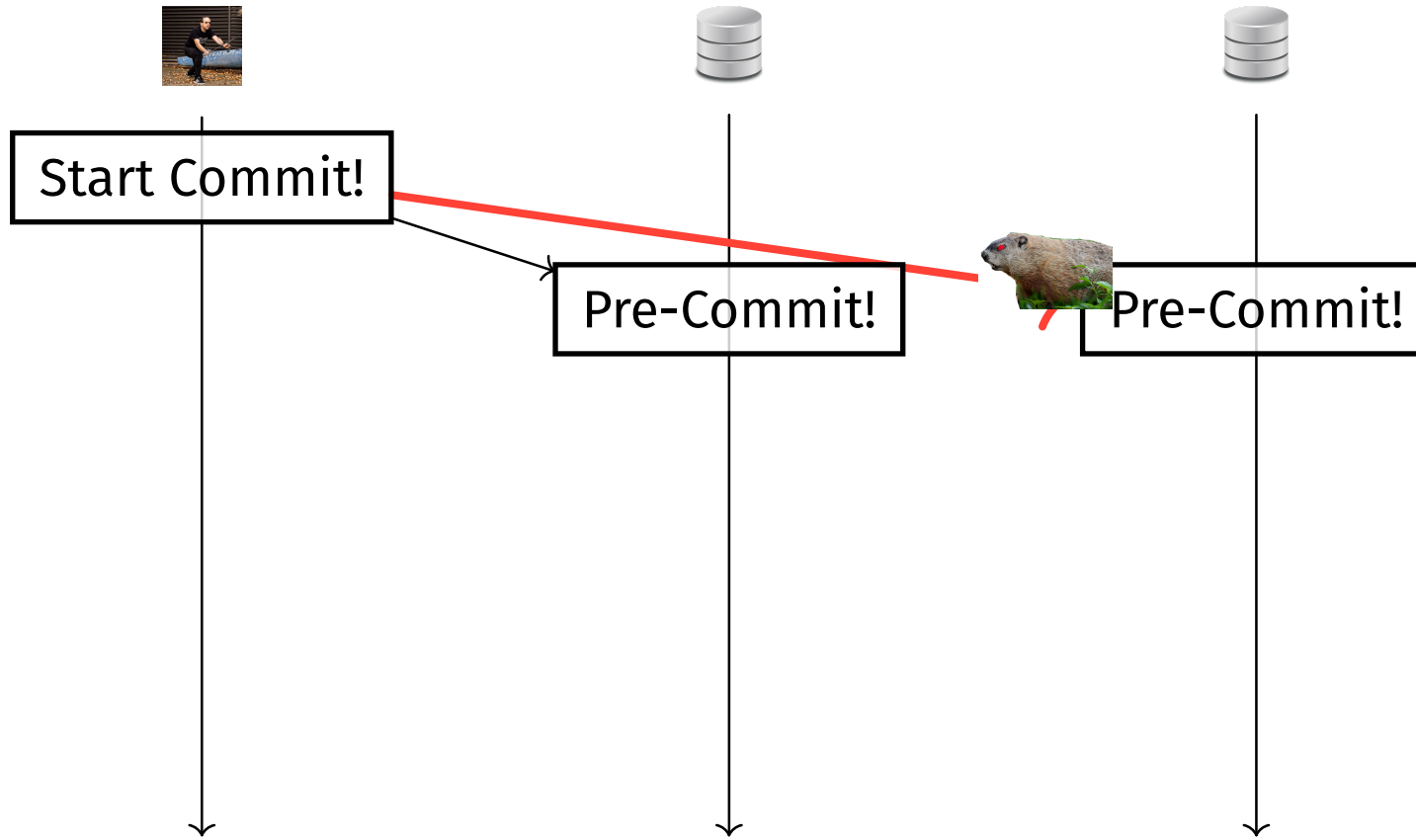
Failure Analysis

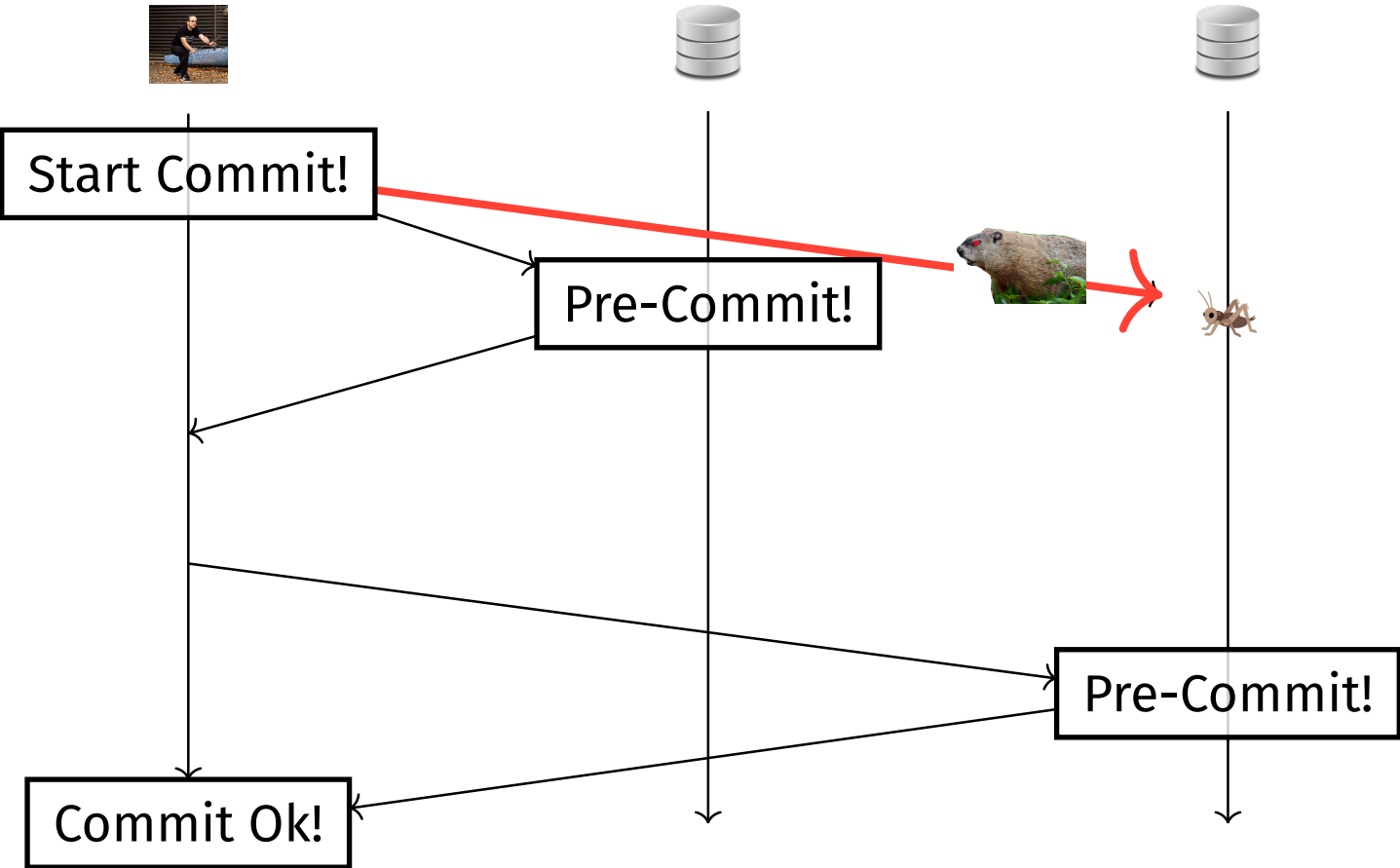
What kinds of failures can we tolerate?

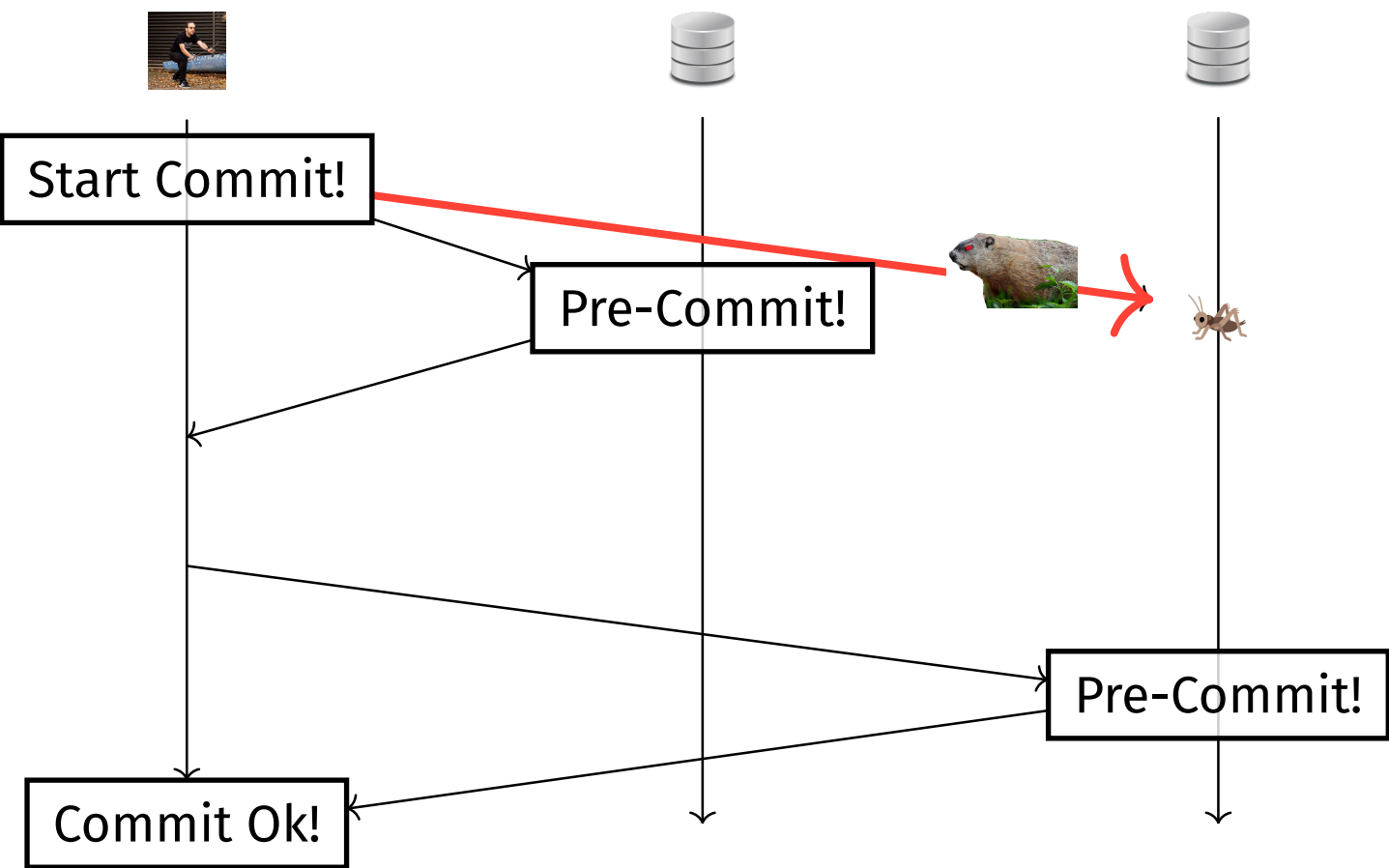
What kinds of failures can we tolerate?

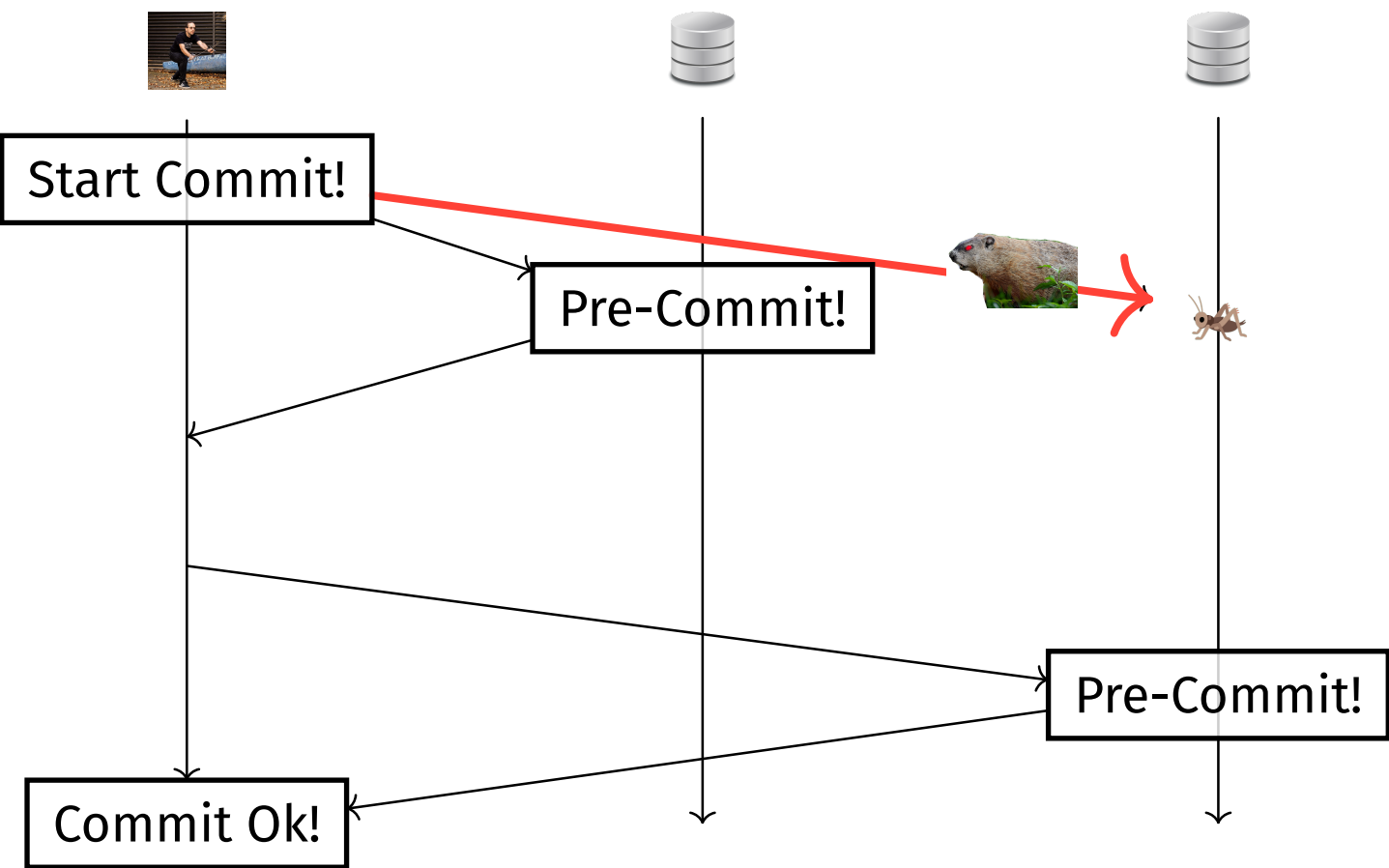
- A message is dropped
- A participant crashes recoverably (power outage)
- A participant crashes non-recoverably (marmots)

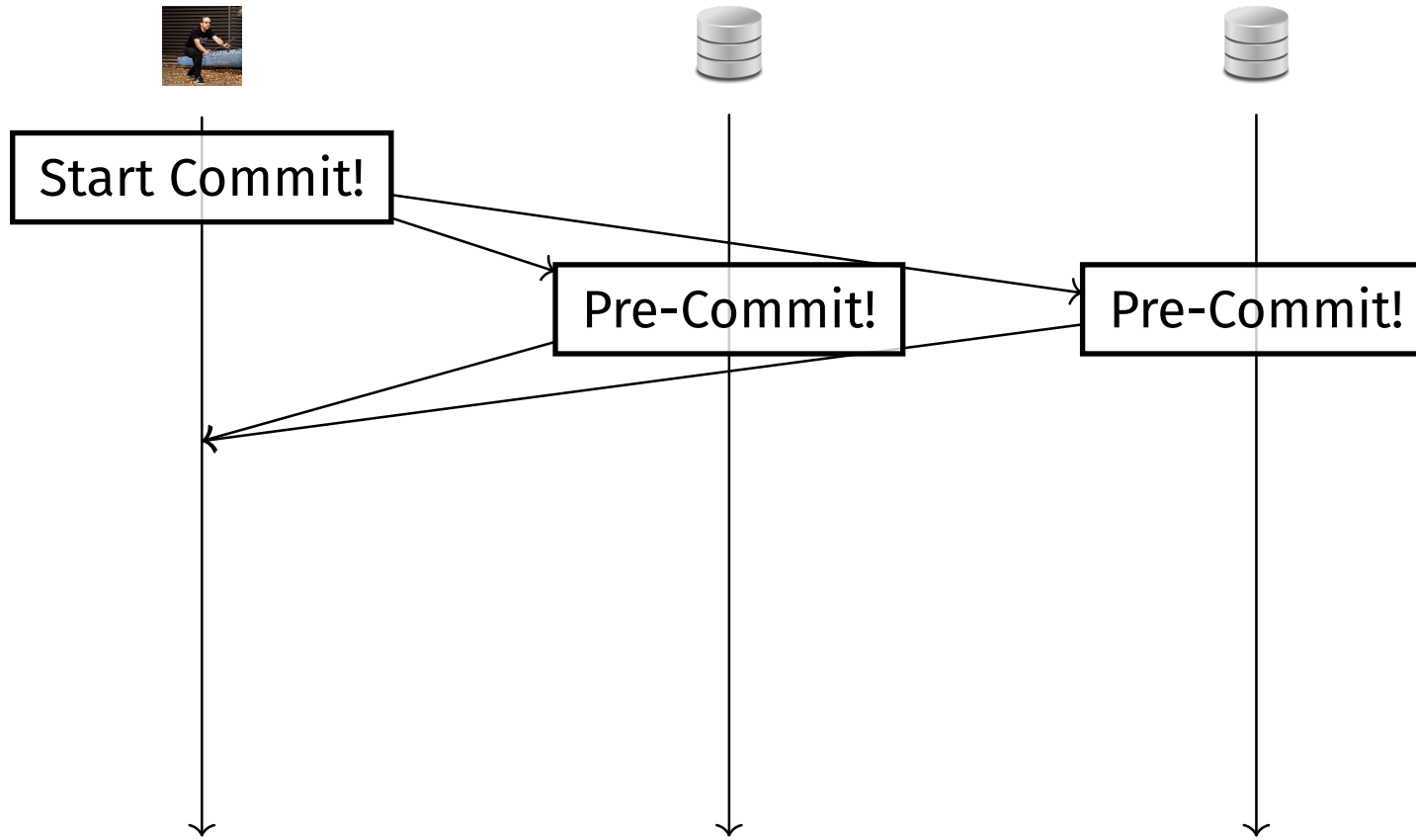


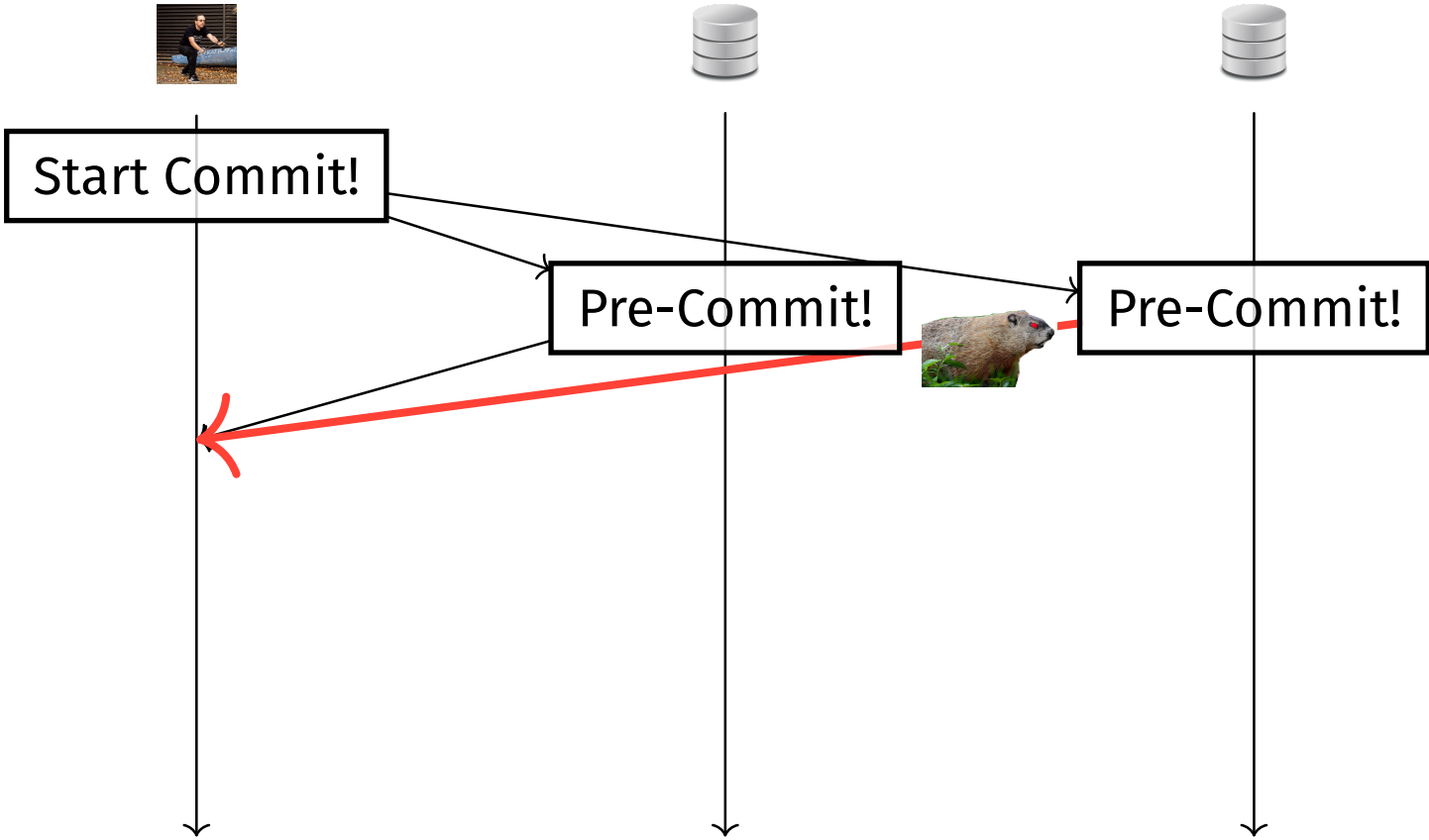


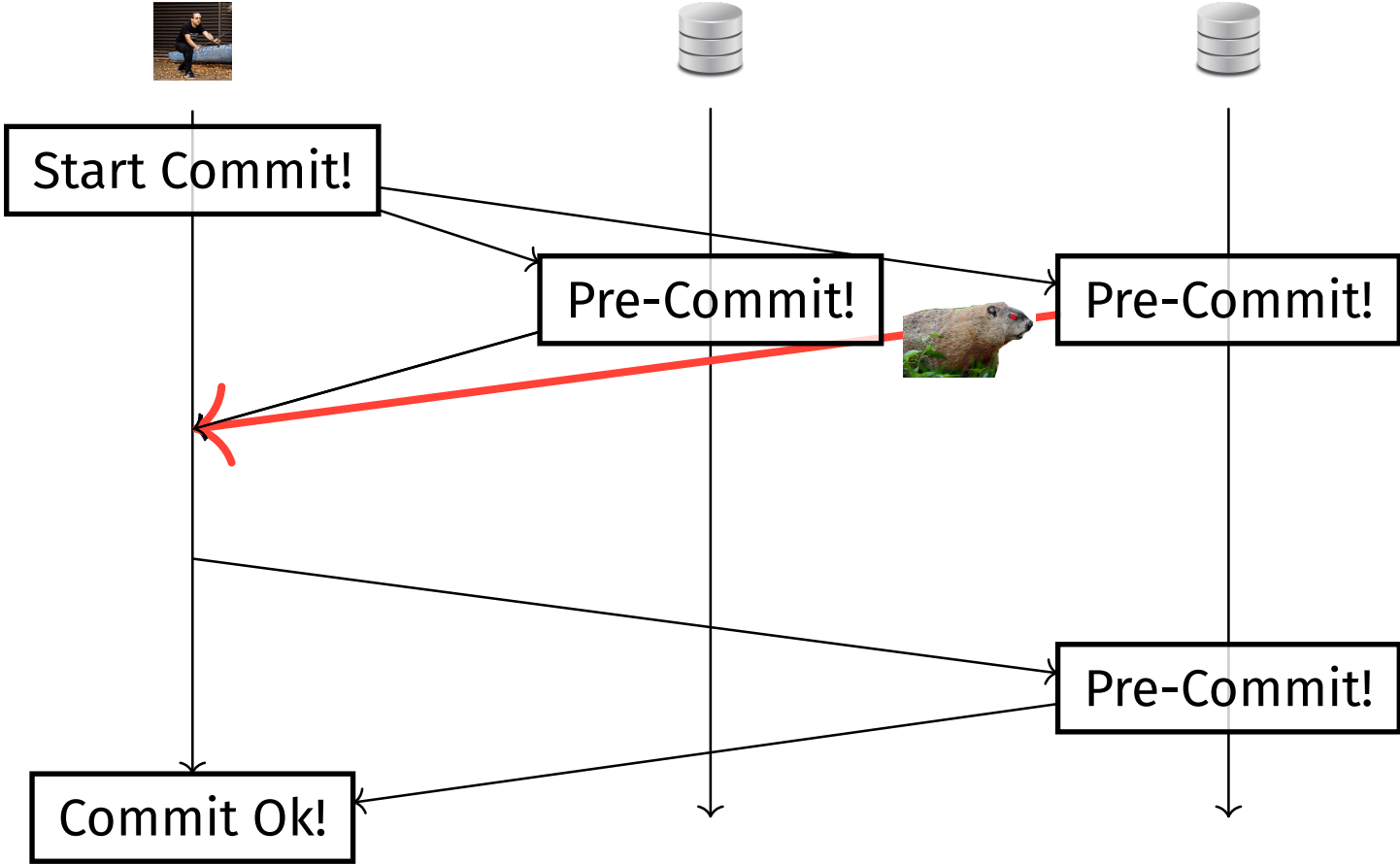


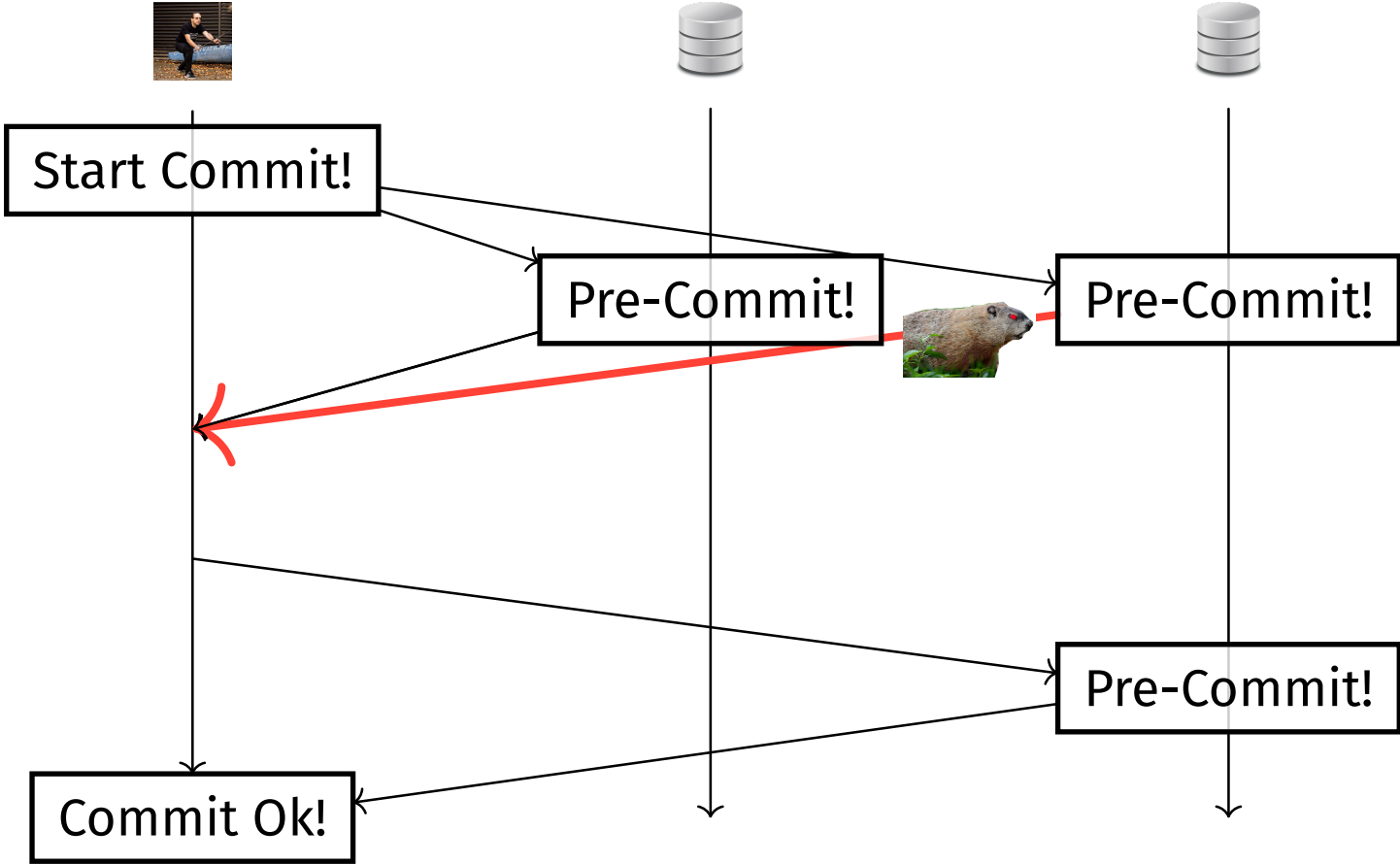


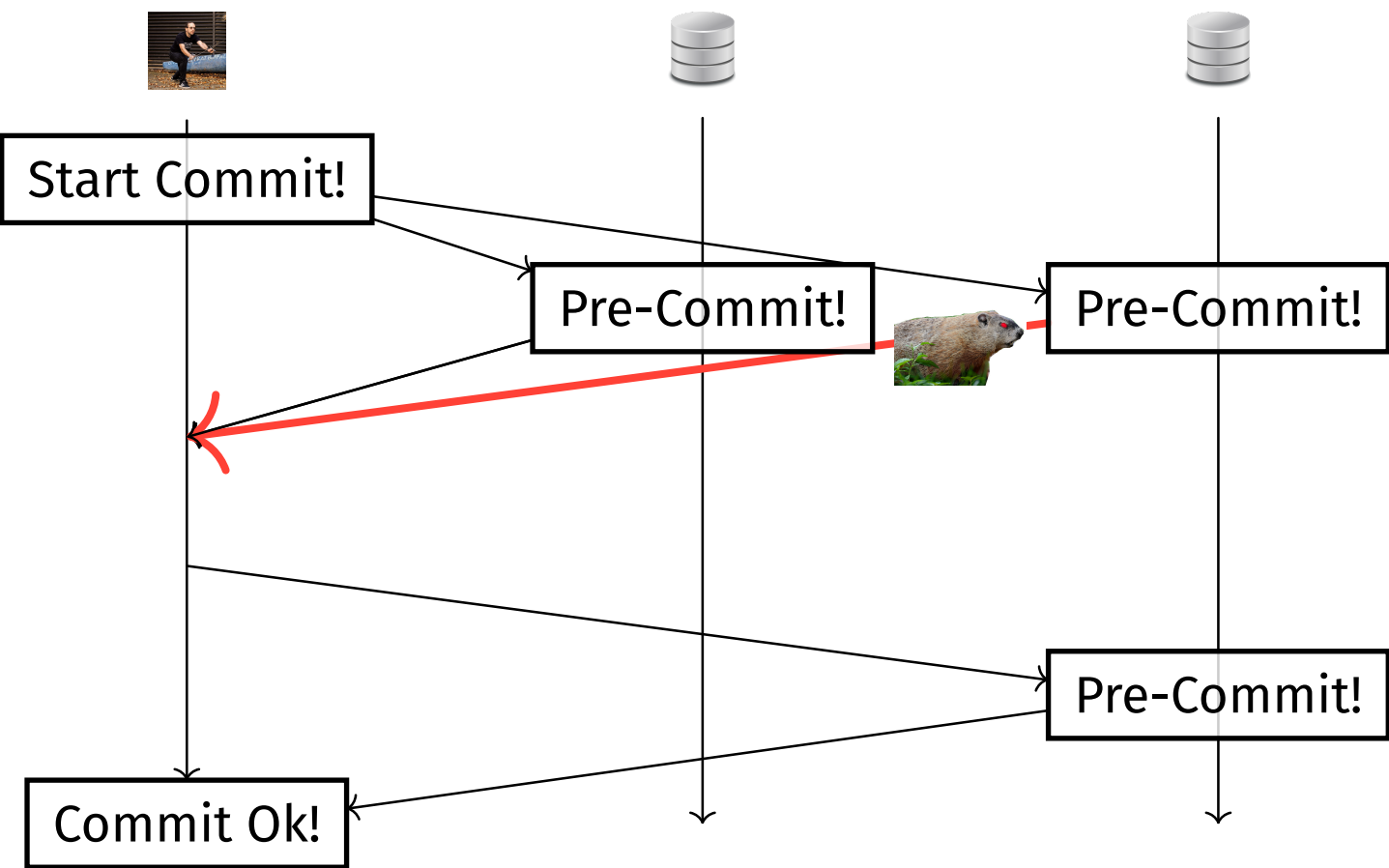












What kinds of failures can we tolerate?

- A message is dropped
- A participant crashes recoverably (power outage)
- A participant crashes non-recoverably (marmots)

1. Each node can enforce ACID **locally**.
2. No malicious nodes.
3. Messages arrive in full or not at all.
4. No global state.

If the node guarantees **Durability**,
transient failures are no different from message loss.

What kinds of failures can we tolerate?

- A message is dropped
- A participant crashes recoverably (power outage)
- A participant crashes non-recoverably (marmots)

When is the write durable?

When is the write durable?

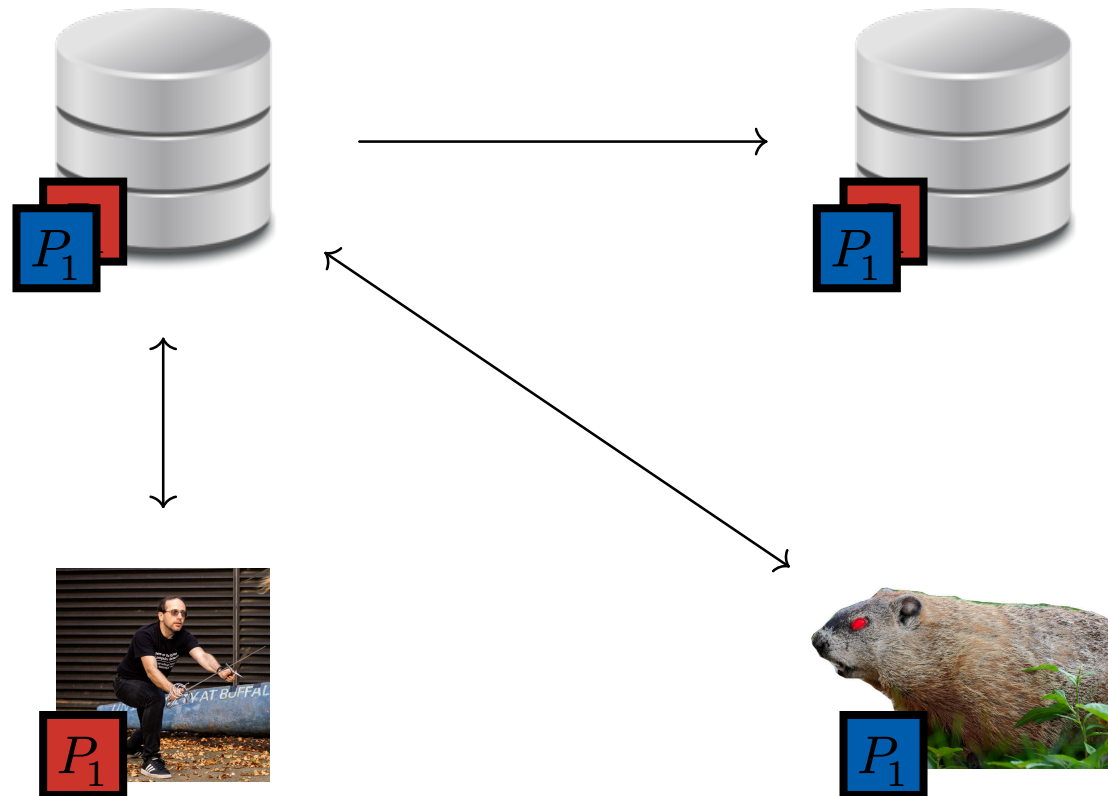
Never!

When is the write durable?

Never!

More durability = More costs

Idea: Replicate data to multiple nodes/sites/planets...



Pros

Cons

Pros

- A backup node exists to come back up when a node is eaten by marmots.

Cons

Pros

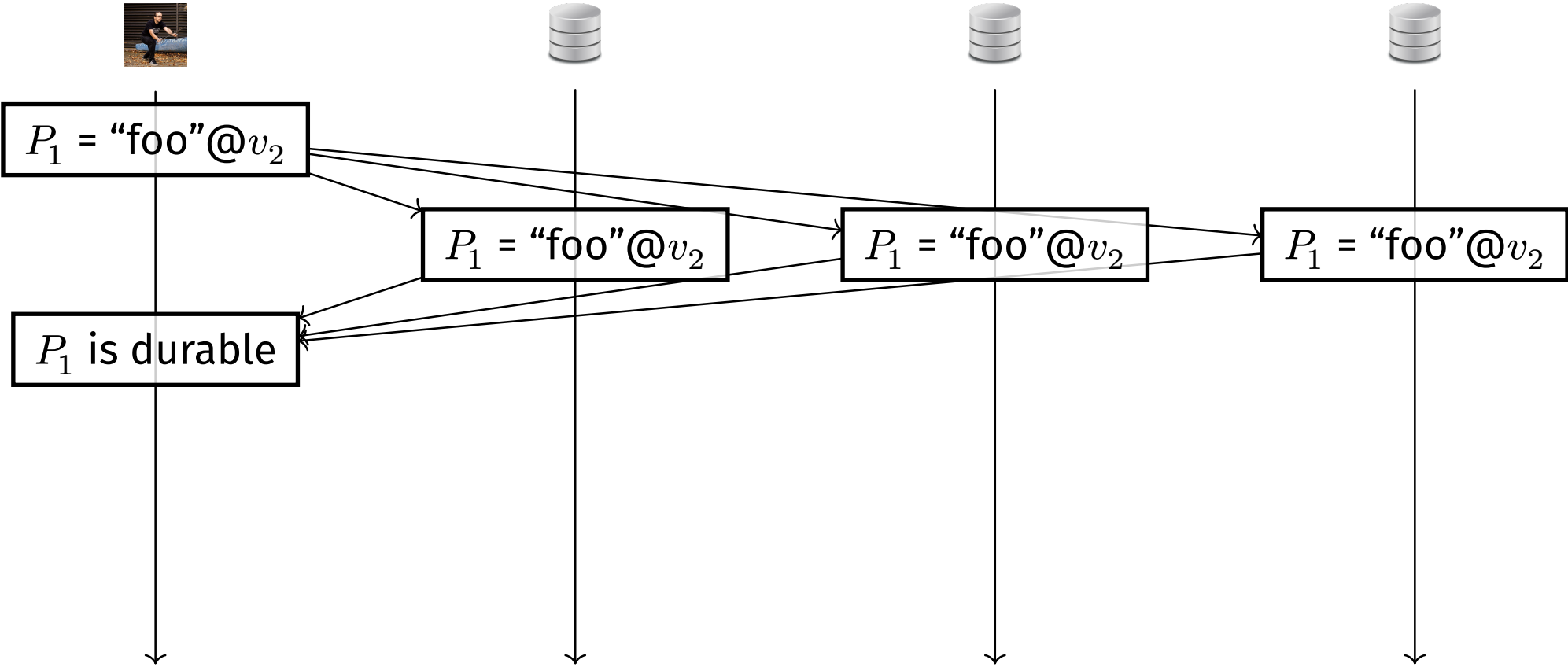
- A backup node exists to come back up when a node is eaten by marmots.

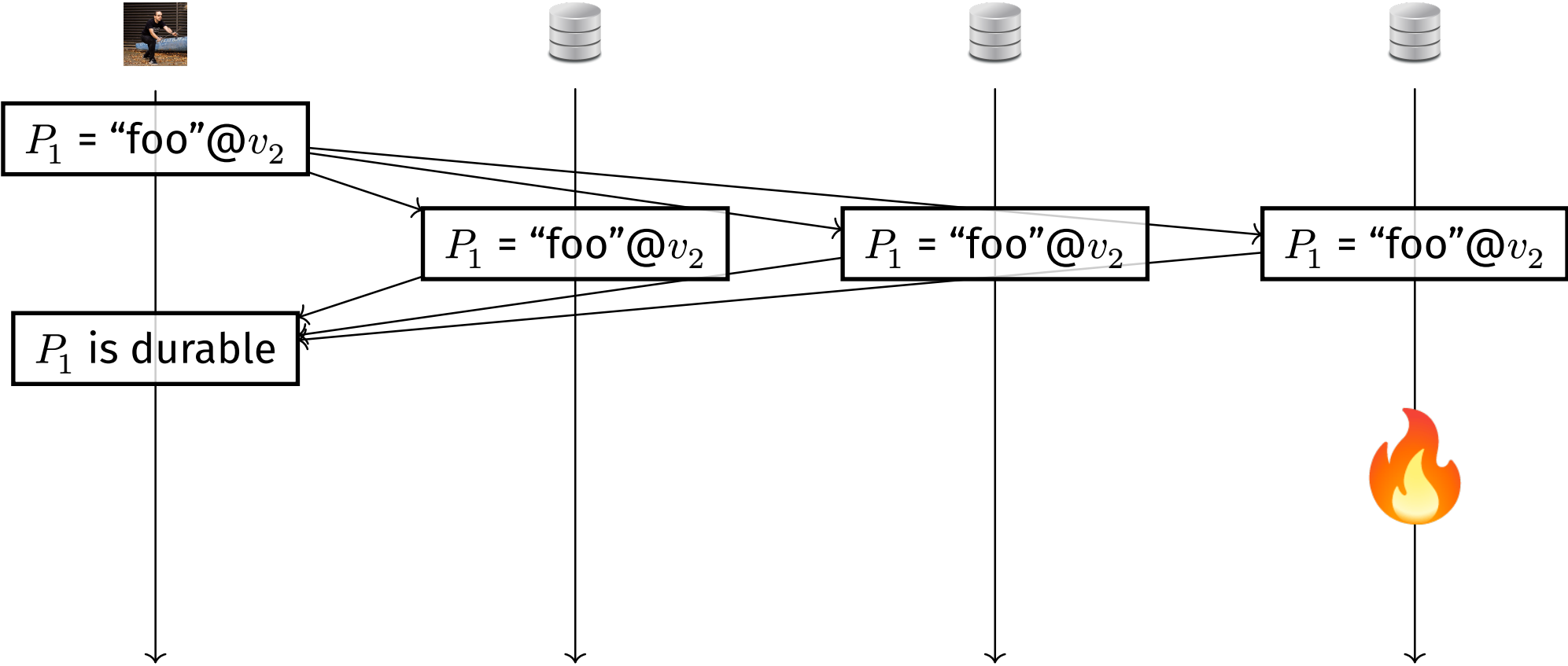
Cons

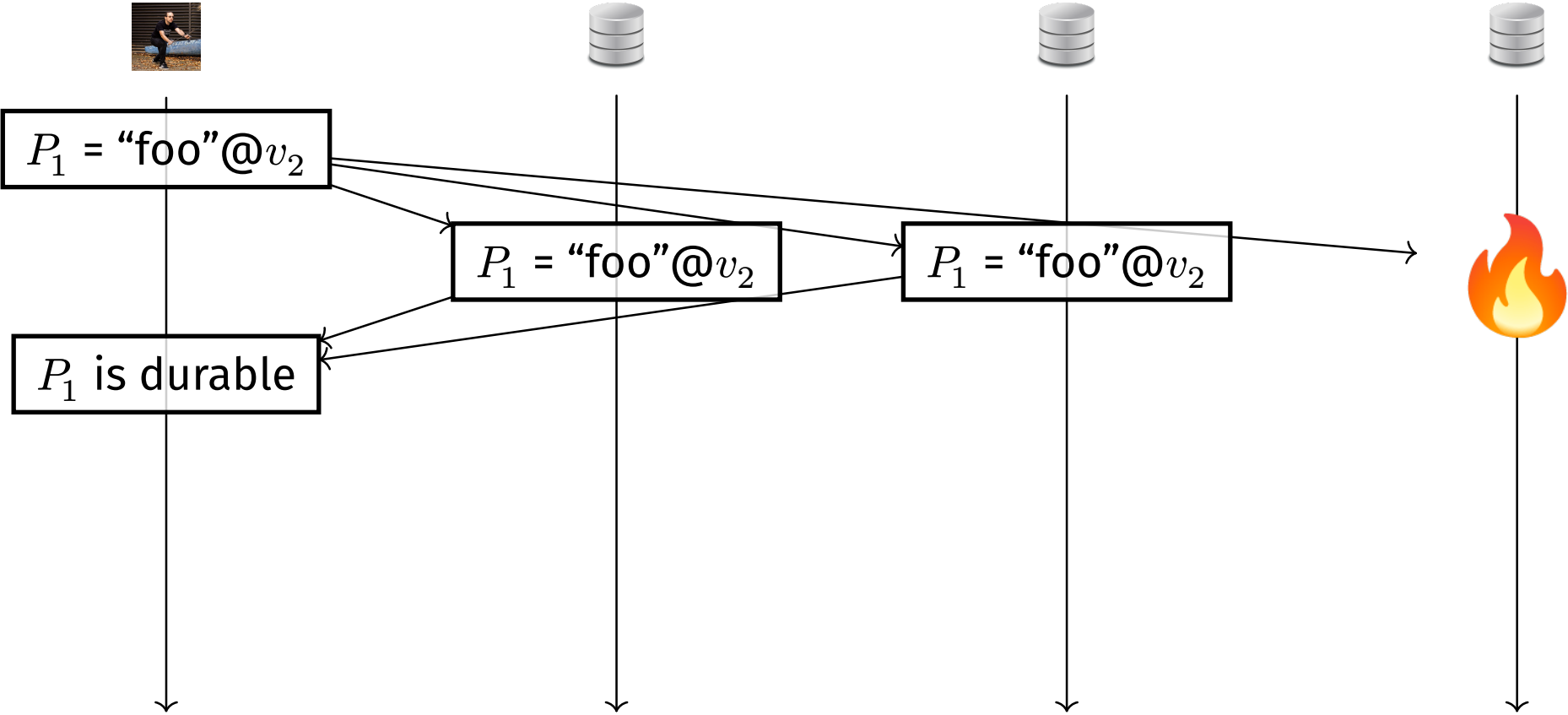
- Have to wait two message round trips (or more) to commit.

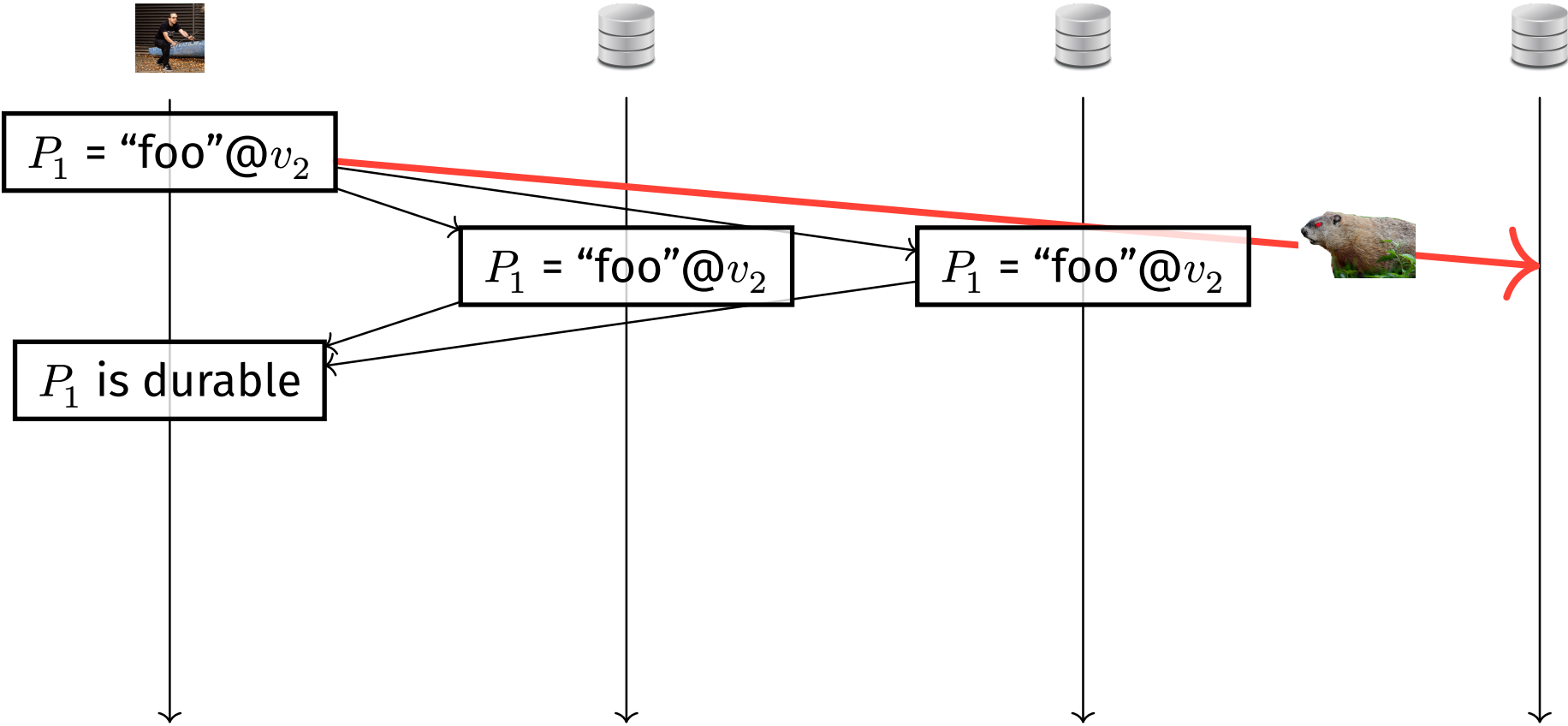
Idea 2: Perform writes in parallel

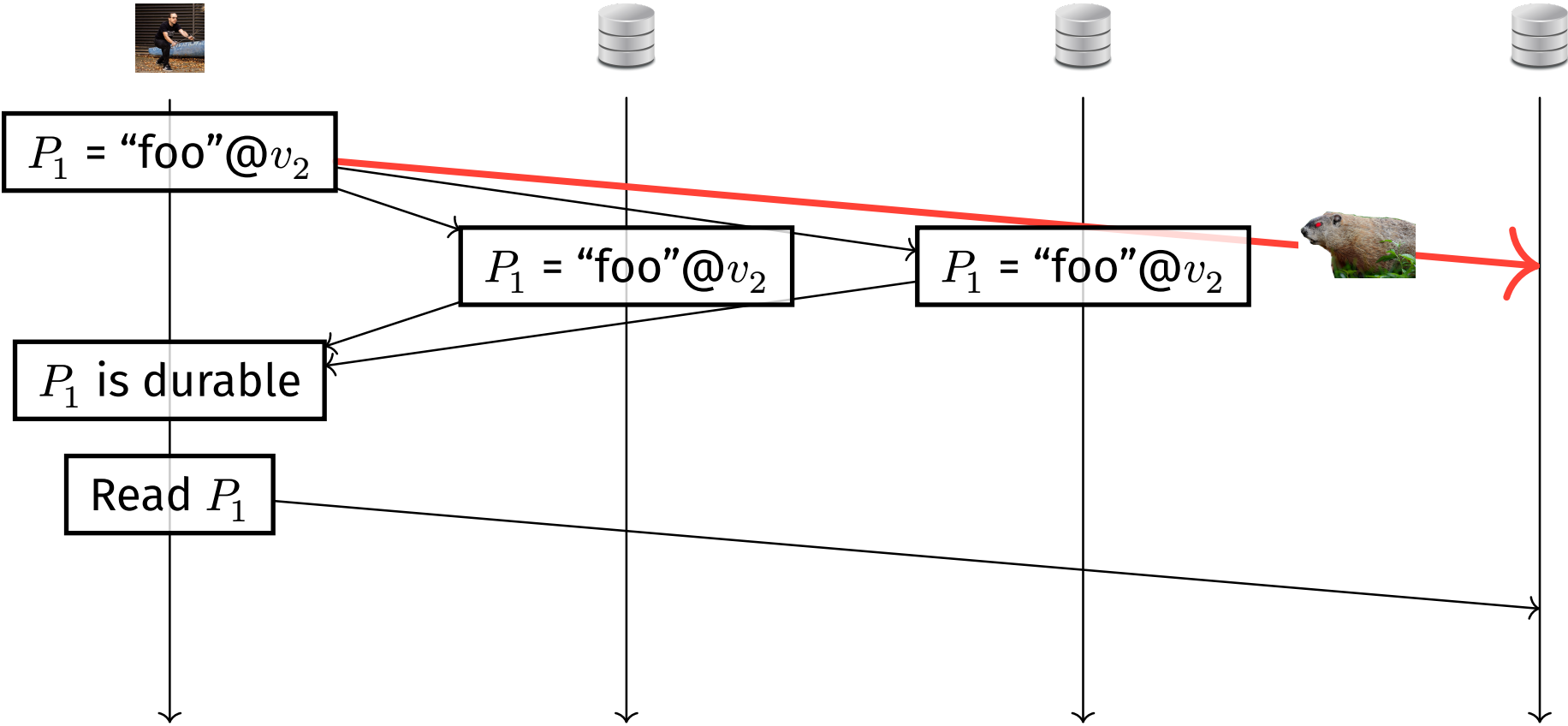
(Reminder: We're only solving for durability here)

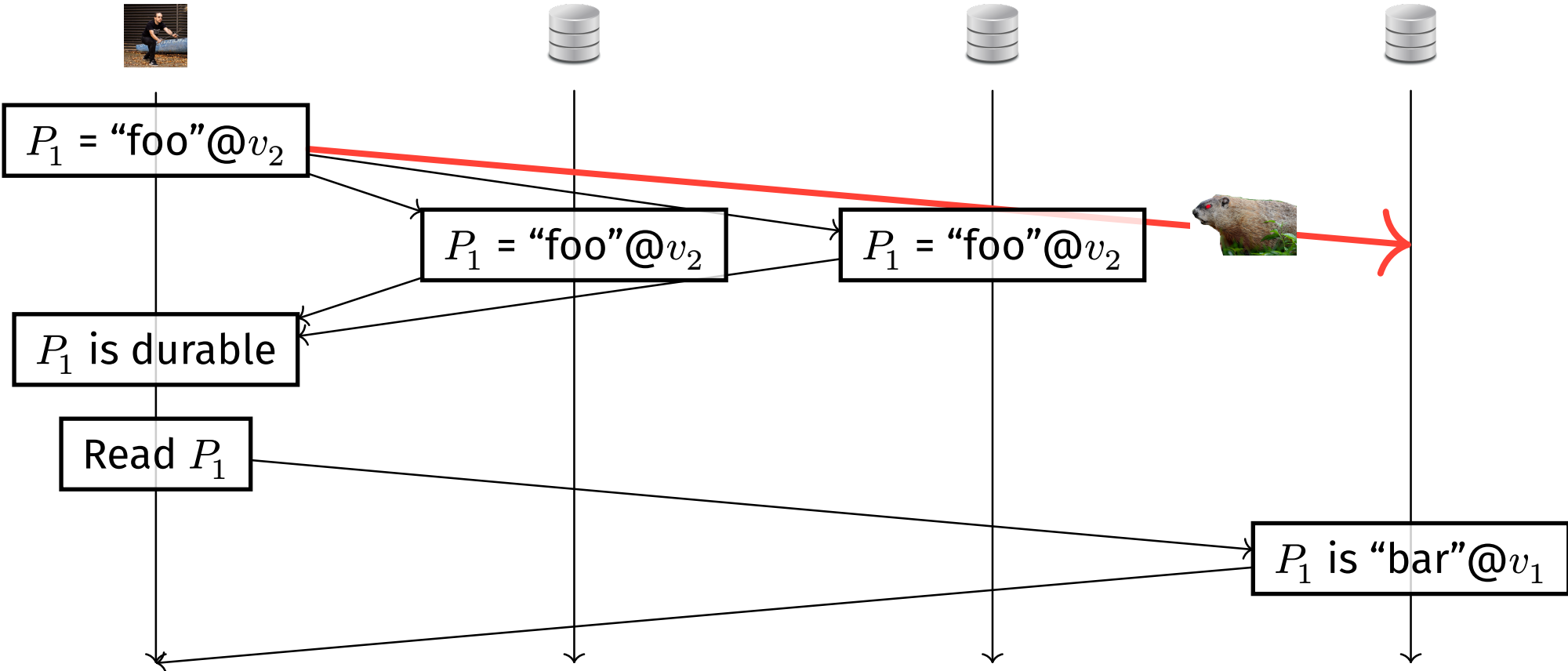


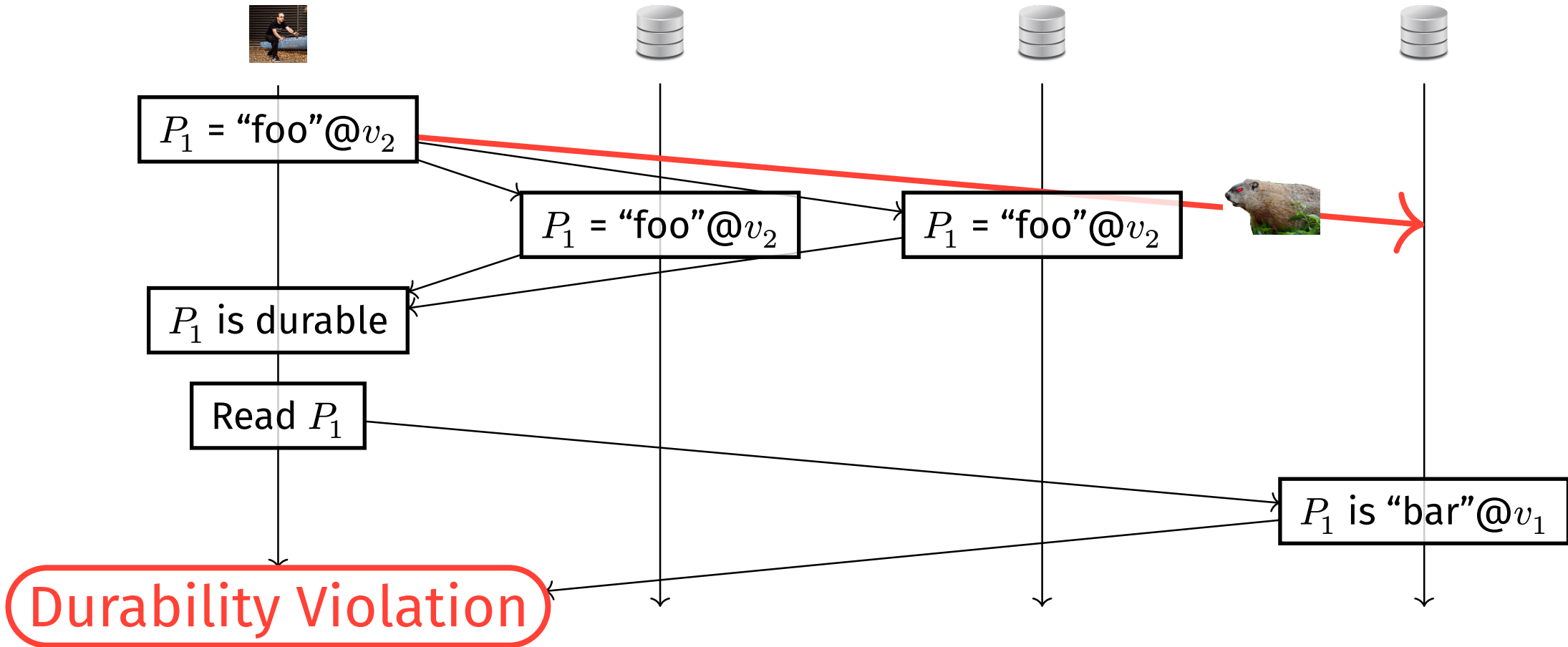


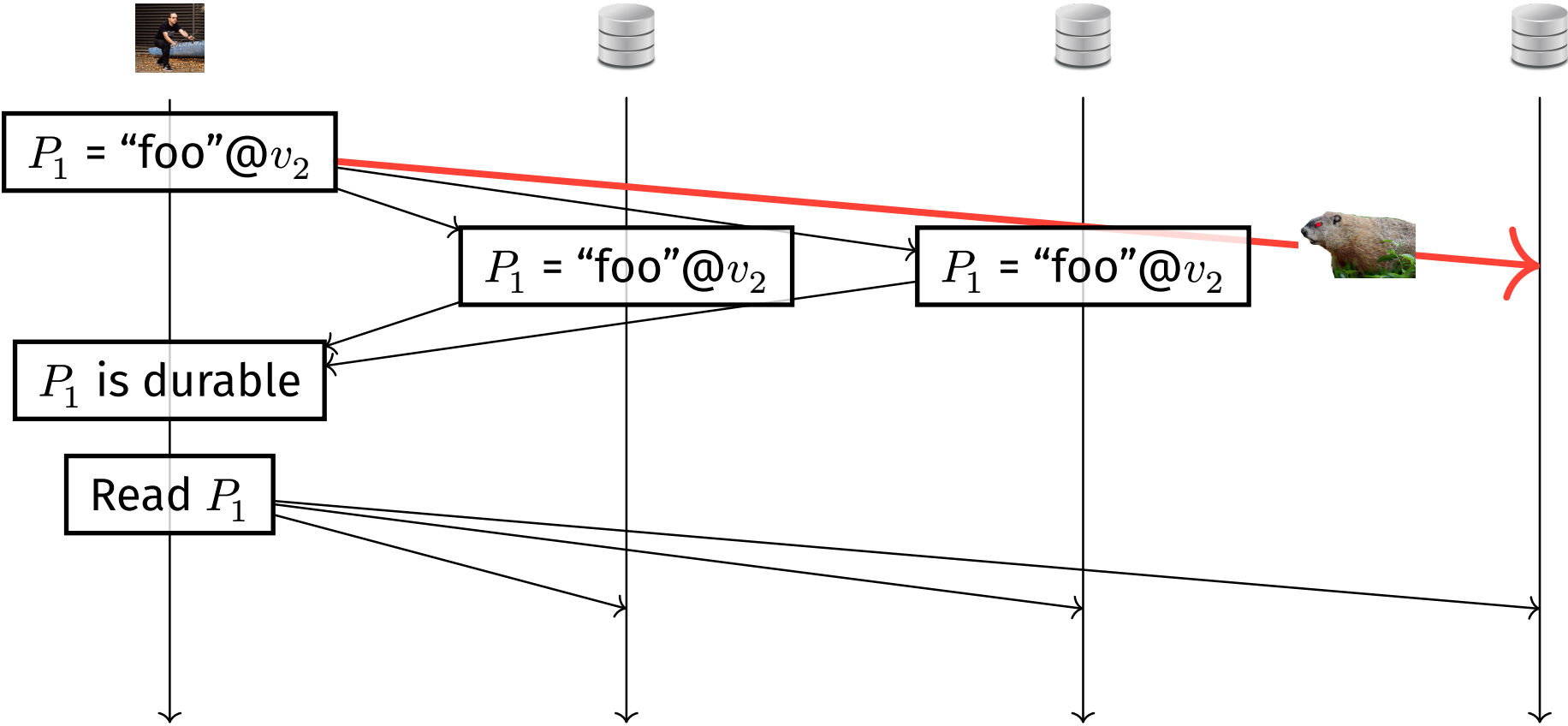


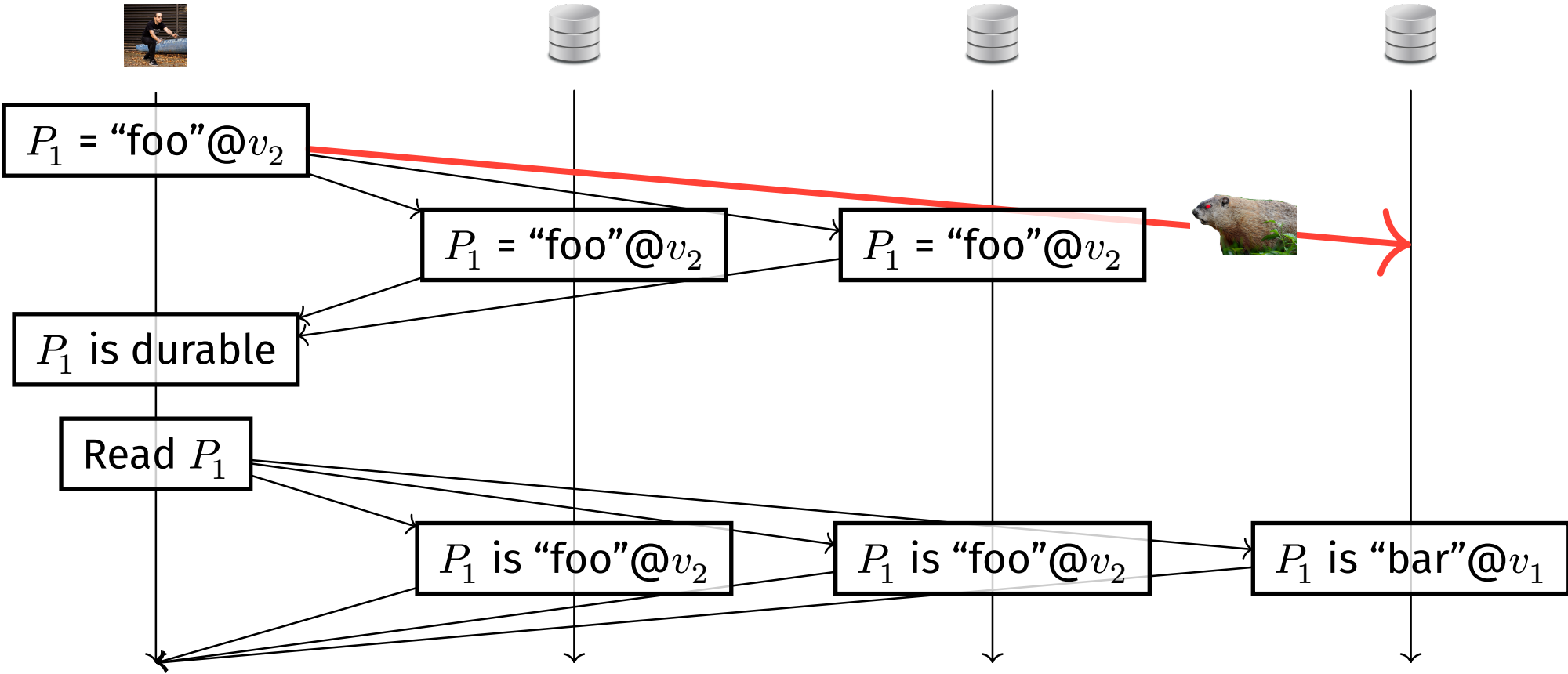


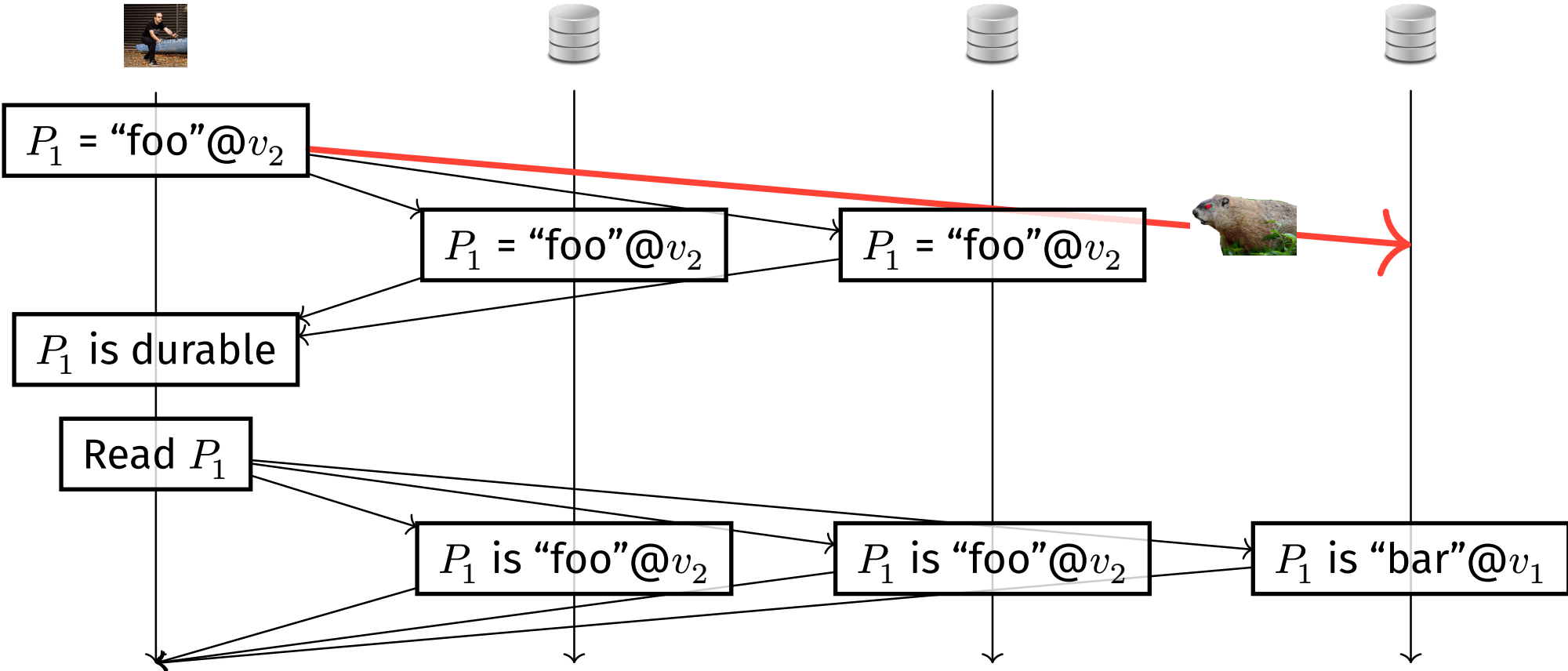












Writer

1. Write data to all replicas
2. Wait for confirmation from W replicas
3. Data is Durable?

Reader

1. Request data from all replicas
2. Receive response from R replicas

How do we guarantee the reader will see at least one copy of every write?

Pigeonhole Principle

$R + W > N$ guarantees that at least one reader and writer are on the same instance

Timeouts

What does a node assume when it fails to receive a message?

Assume transient failure and wait longer

All data remains consistent ... but unavailable while we wait, which might be indefinitely.

Assume non-recoverable failure and fail-over to backups

All data remains available ... but if it turns out to be transient, there can be inconsistencies.

C		A		P
o		v		a
n		a		r t
s		i		t o
i		l		i l
s	or	a	or	t e
t		b		i r
e		i		o a
n		l		n n
c		i		c
y		t		e
		y		

Traditionally: Pick any 2 ('P' has to be one)

C		A		P
o		v		a
n		a		r
s		i		t
i		l		i
s	or	a	during	t
t		b		i
e		i		o
n		l		n
c		i		s
y		t		
		y		

I prefer this phrasing