

# 6.033 Spring 2017

## Lecture #19

- **Distributed transactions**
  - **Availability**
  - **Replicated State Machines**

**goal:** build reliable systems from unreliable components  
the abstraction that makes that easier is

**transactions**, which provide **atomicity** and **isolation**, while not hindering **performance**

**atomicity** → **shadow copies** (simple, poor performance) or **logs** (better performance, a bit more complex)

**isolation** → **two-phase locking**

we also want transaction-based systems to be **distributed** — to run across multiple machines — and to remain **available** even through failures

**C<sub>1</sub>** **write<sub>1</sub>(X)**

**S<sub>1</sub>**

**C<sub>2</sub>** **write<sub>2</sub>(X)**

**S<sub>2</sub>**

(replica of S<sub>1</sub>)

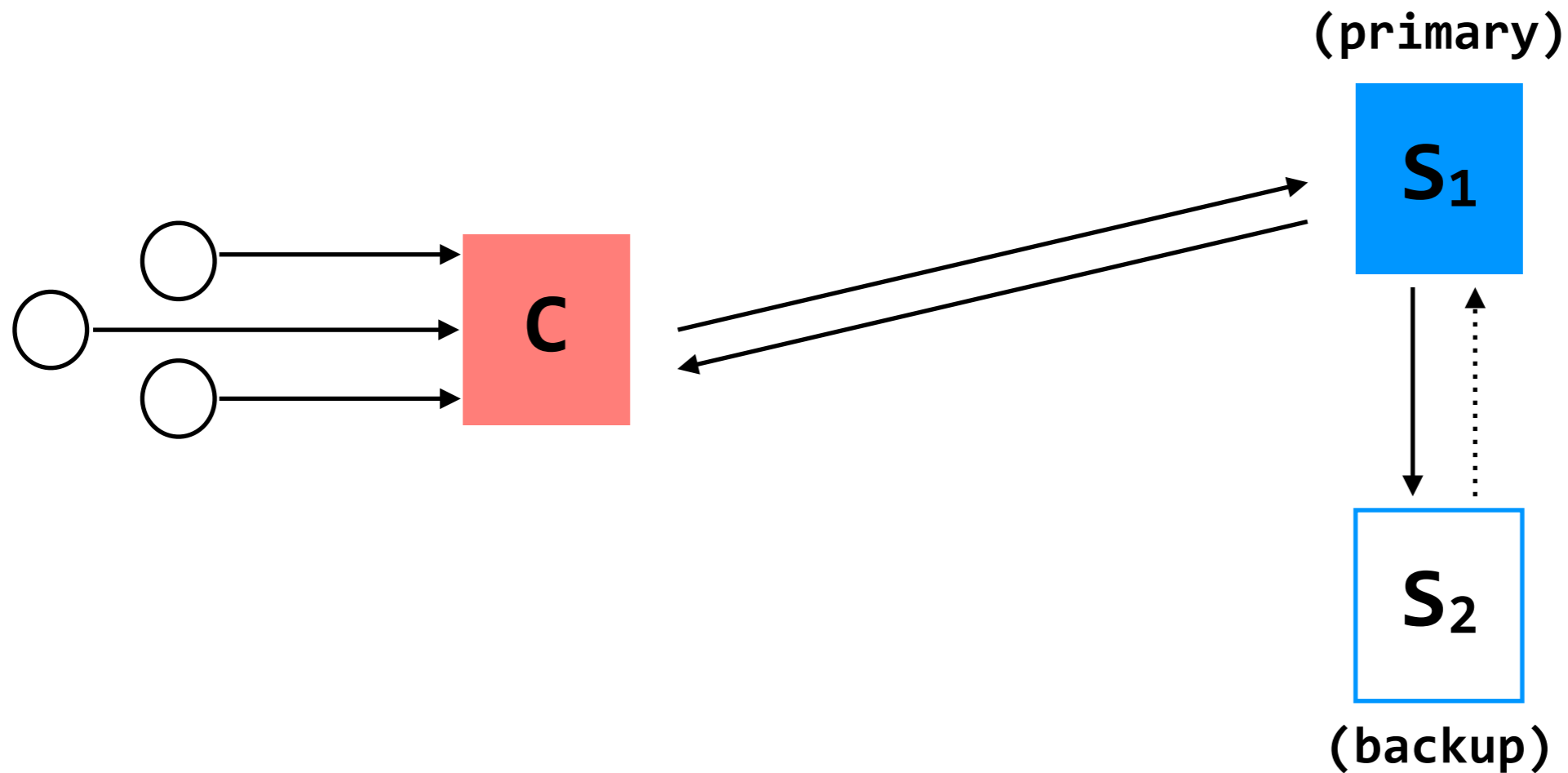
**C<sub>1</sub>**

**C<sub>2</sub>**

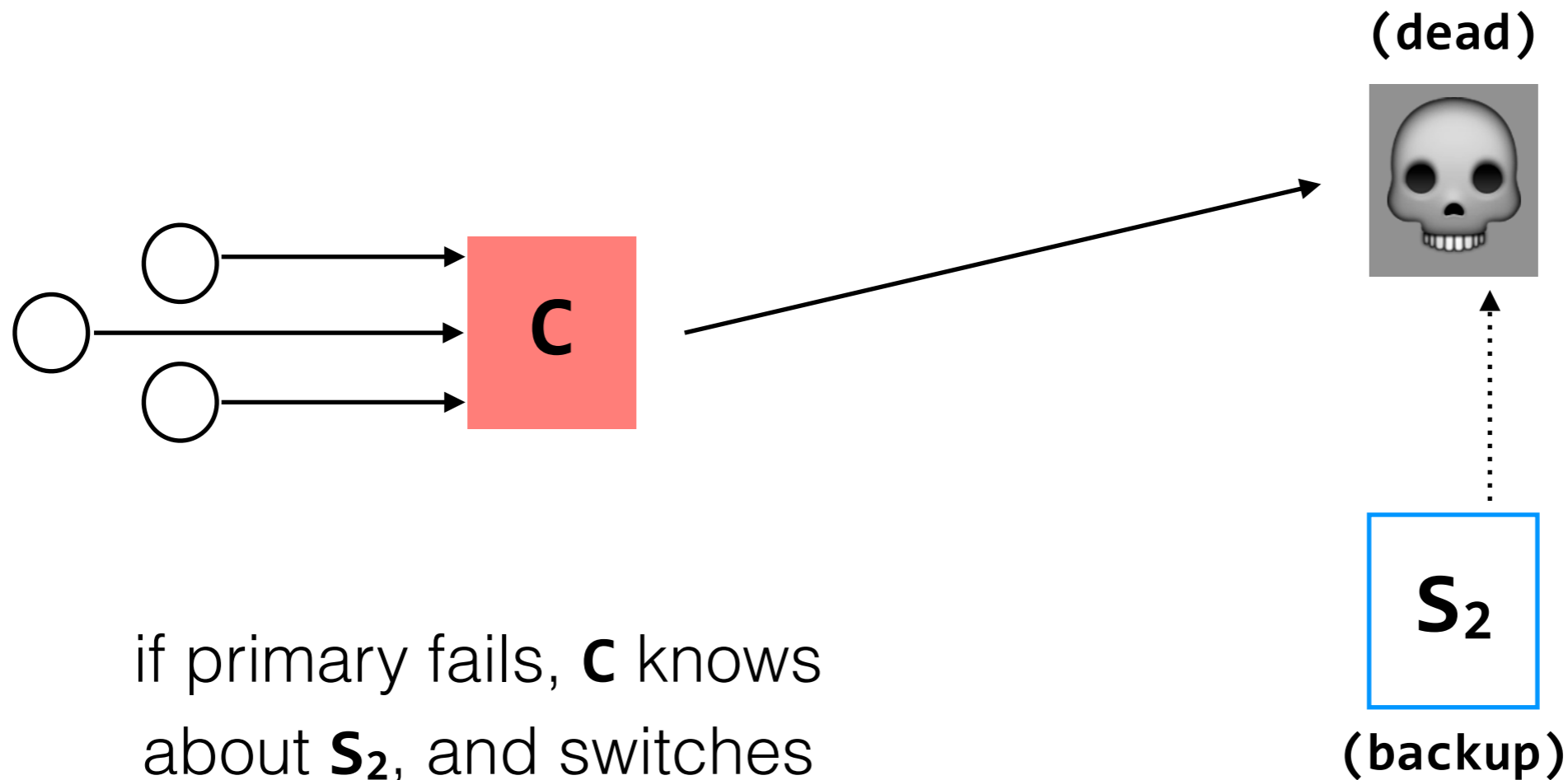
**S<sub>1</sub>** write<sub>1</sub>(X)  
write<sub>2</sub>(X)

**S<sub>2</sub>** write<sub>2</sub>(X)  
write<sub>1</sub>(X)  
(replica of S<sub>1</sub>)

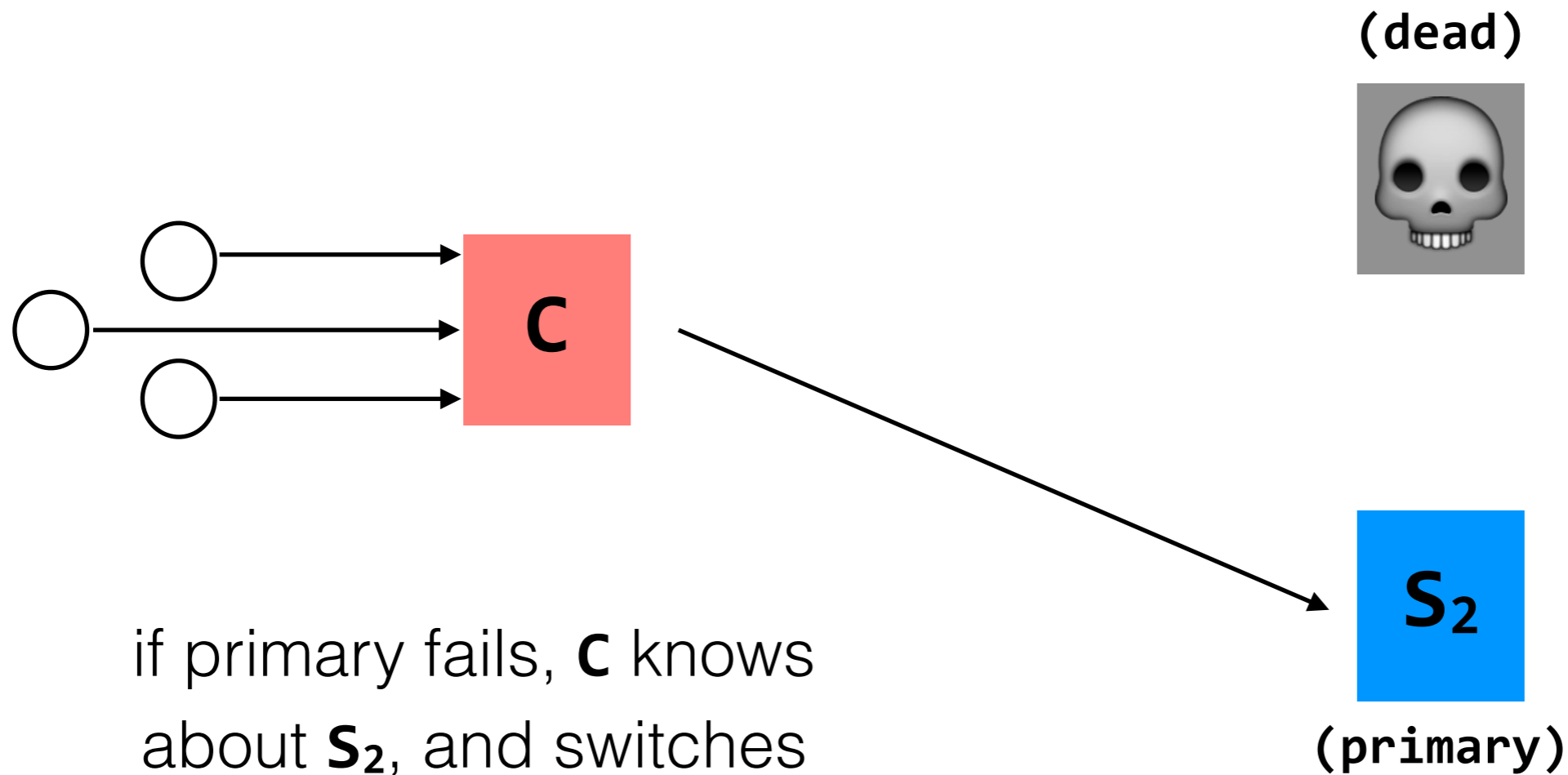
**problem:** replica servers can become inconsistent



**attempt:** coordinators communicate with primary servers, who communicate with backup servers

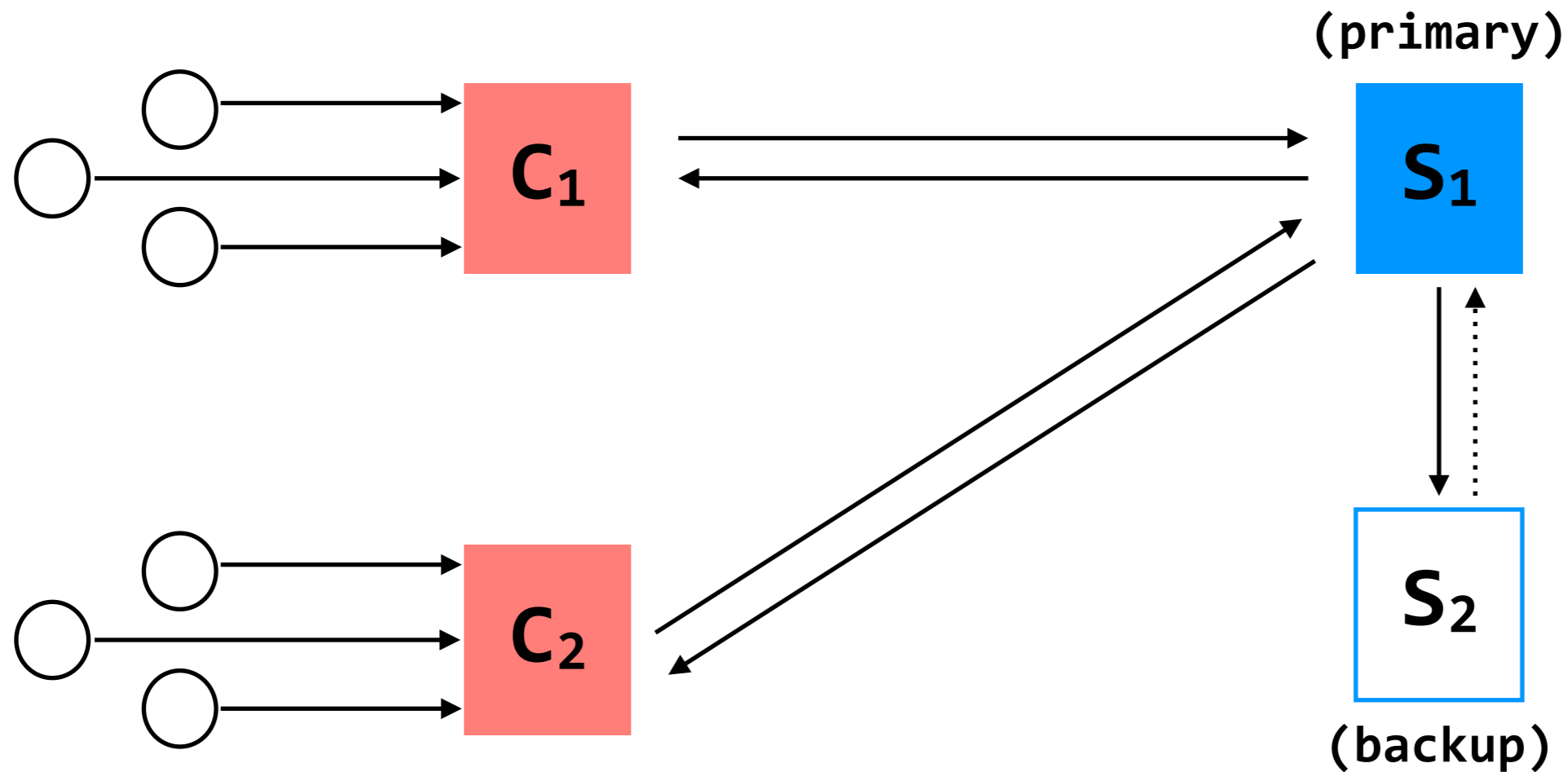


**attempt:** coordinators communicate with primary servers, who communicate with backup servers



**attempt:** coordinators communicate with primary servers, who communicate with backup servers

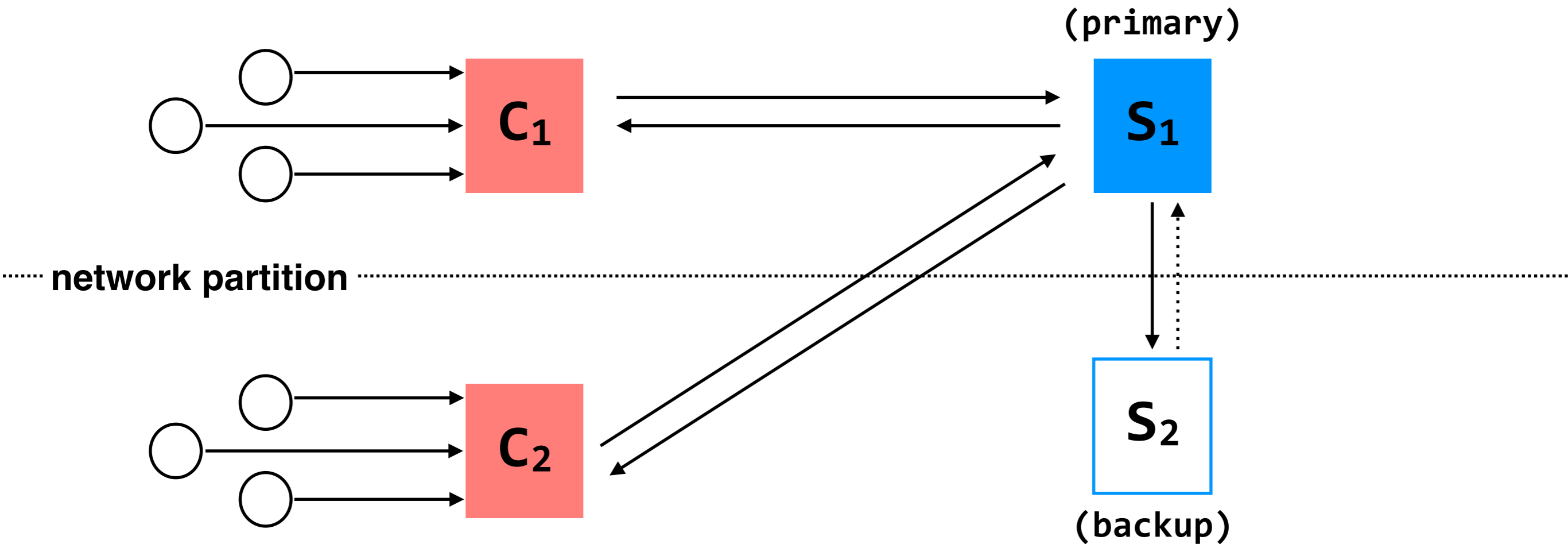
# multiple coordinators + the network = problems



**attempt:** coordinators communicate with primary servers, who communicate with backup servers

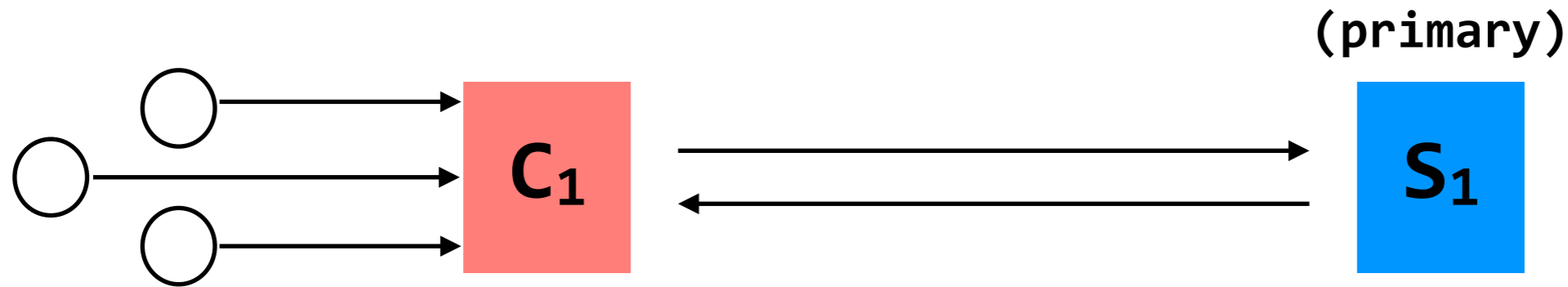


# multiple coordinators + the network = problems

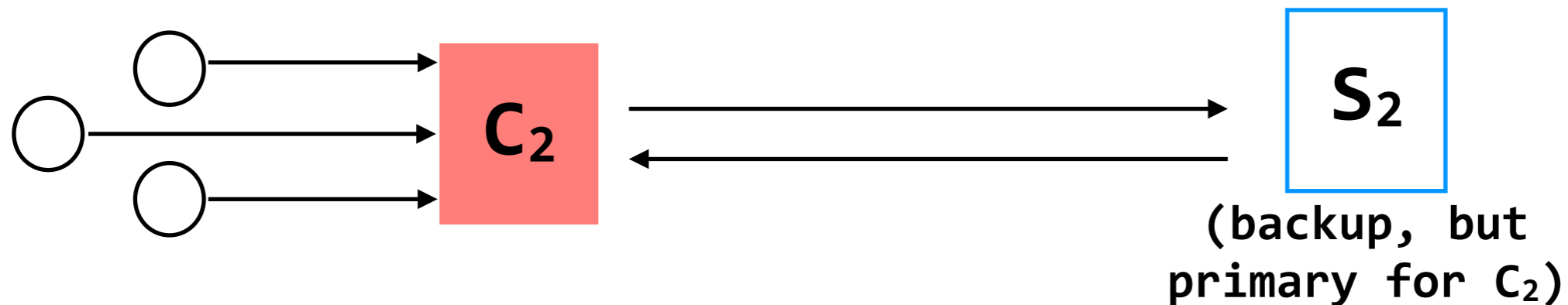


**attempt:** coordinators communicate with primary servers, who communicate with backup servers

# multiple coordinators + the network = problems



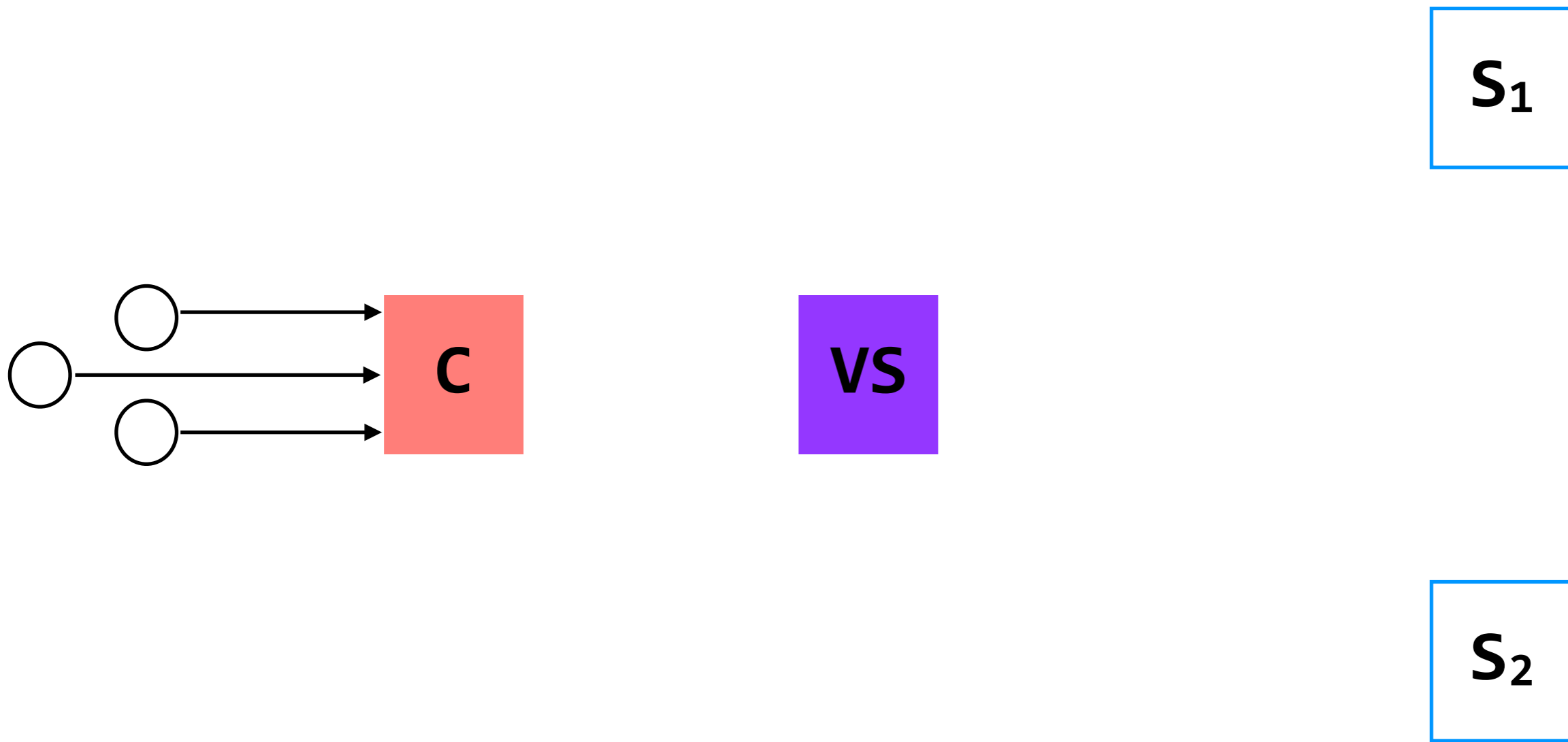
..... network partition .....



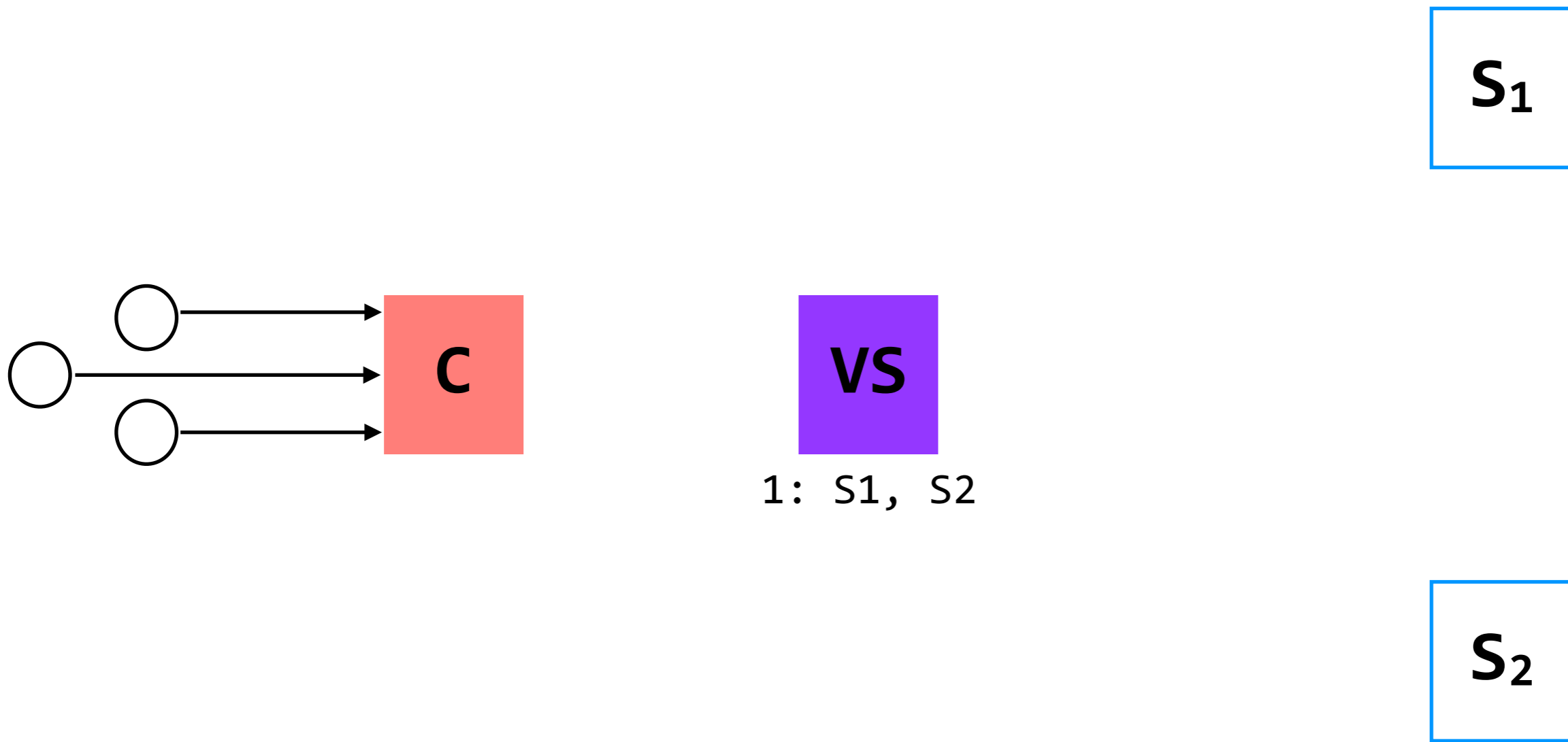
$C_1$  and  $C_2$  are using different primaries;

$S_1$  and  $S_2$  are no longer consistent

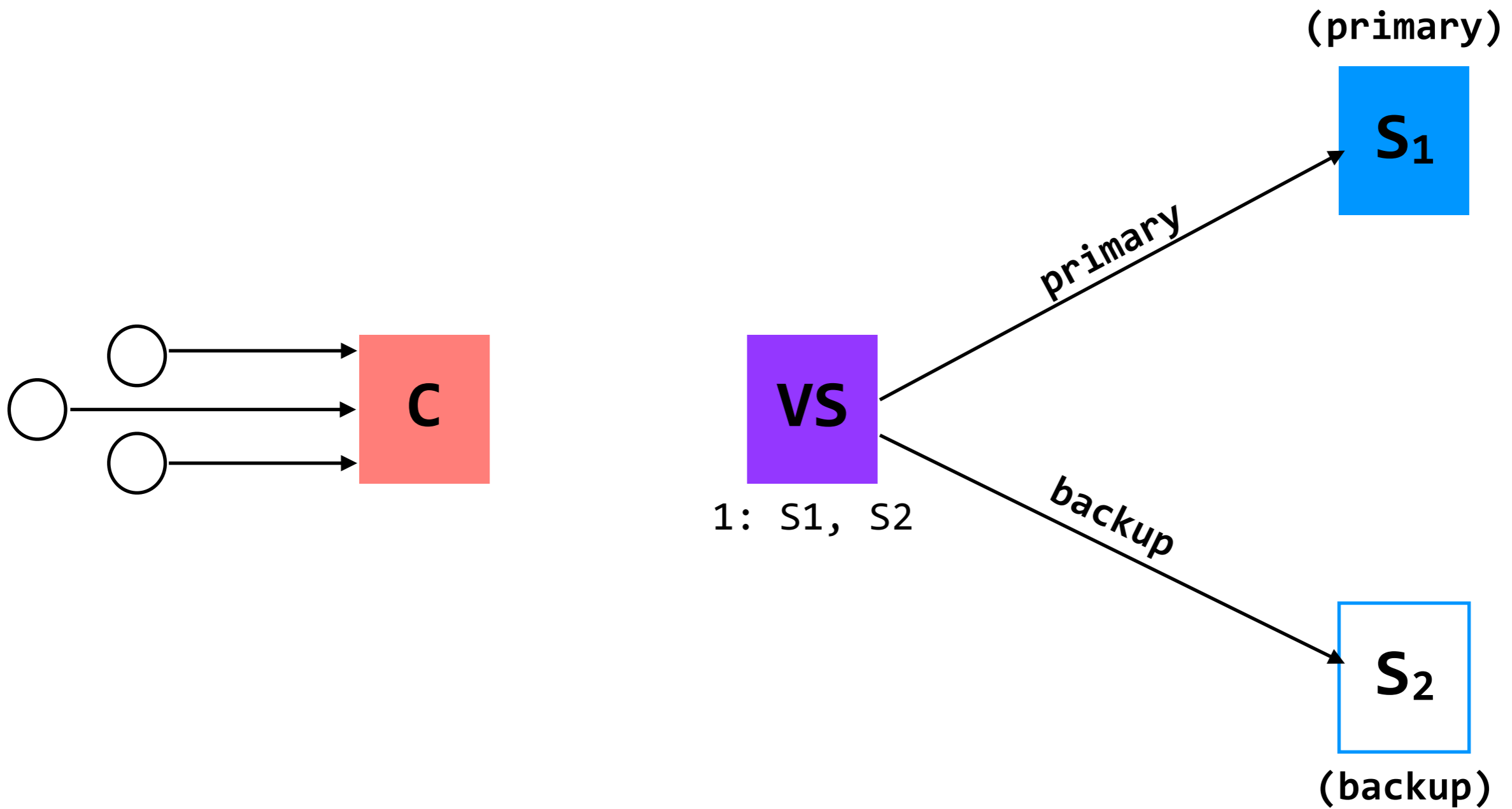
**attempt:** coordinators communicate with primary servers, who communicate with backup servers



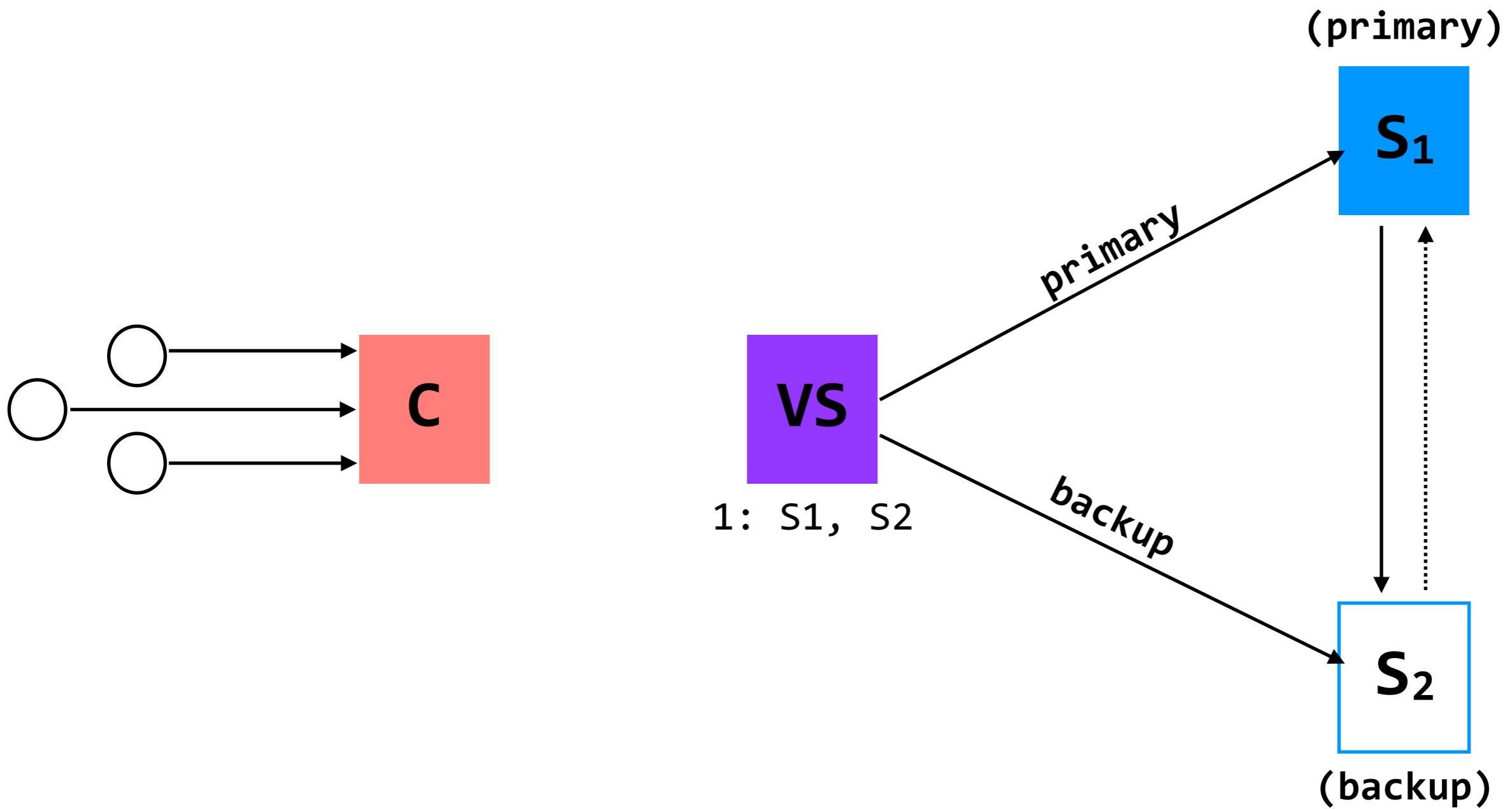
use a **view server**, which determines which replica is the primary



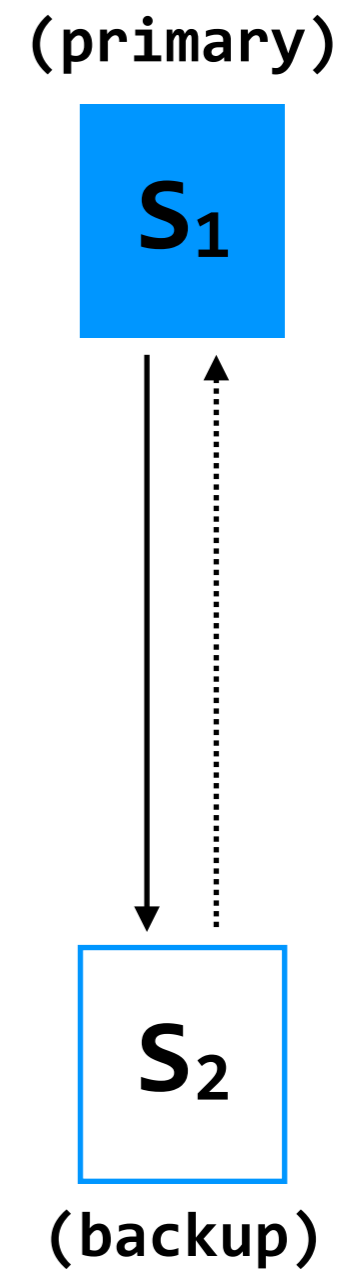
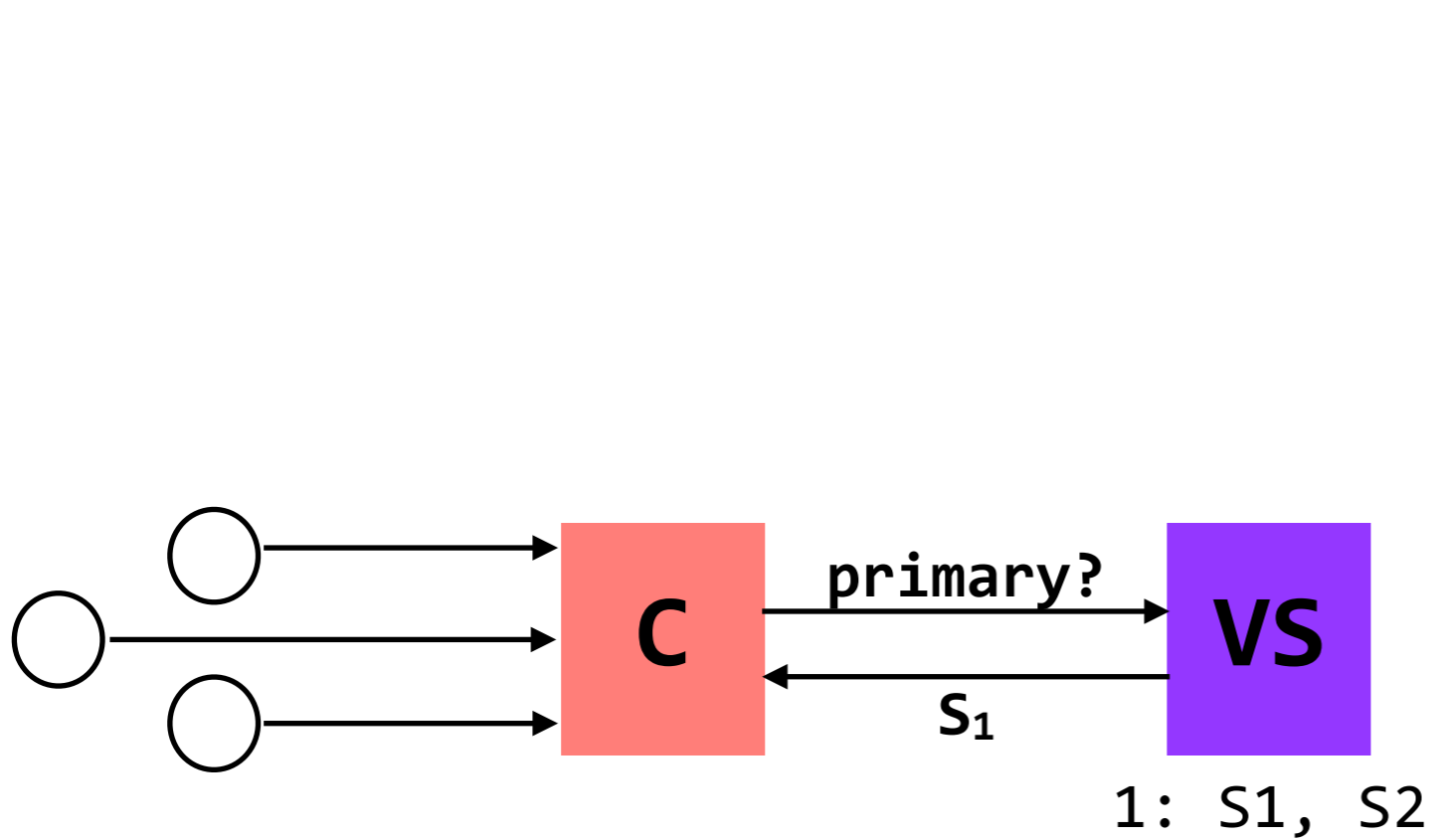
use a **view server**, which determines which replica is the primary



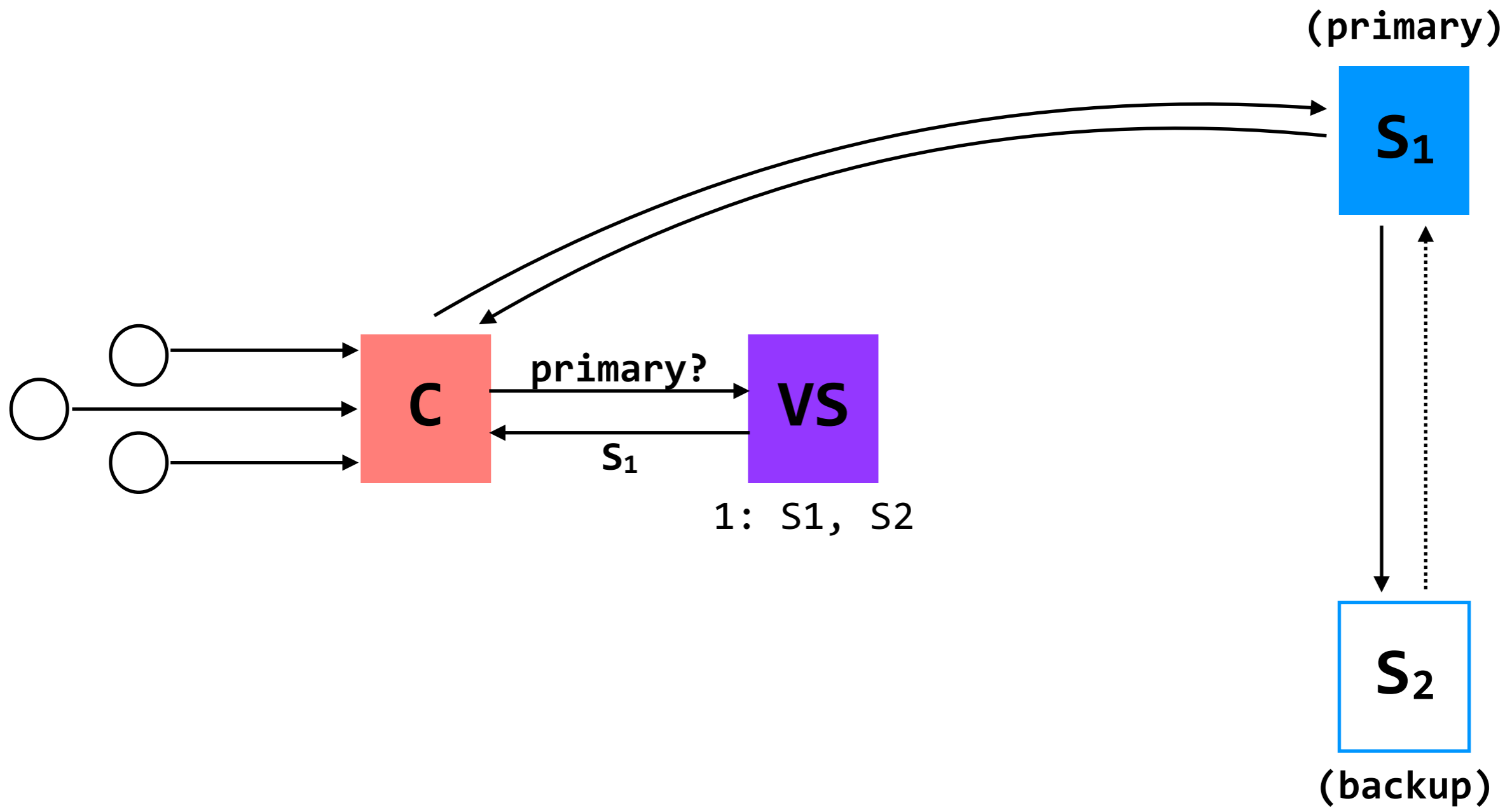
use a **view server**, which determines which replica is the primary



use a **view server**, which determines which replica is the primary

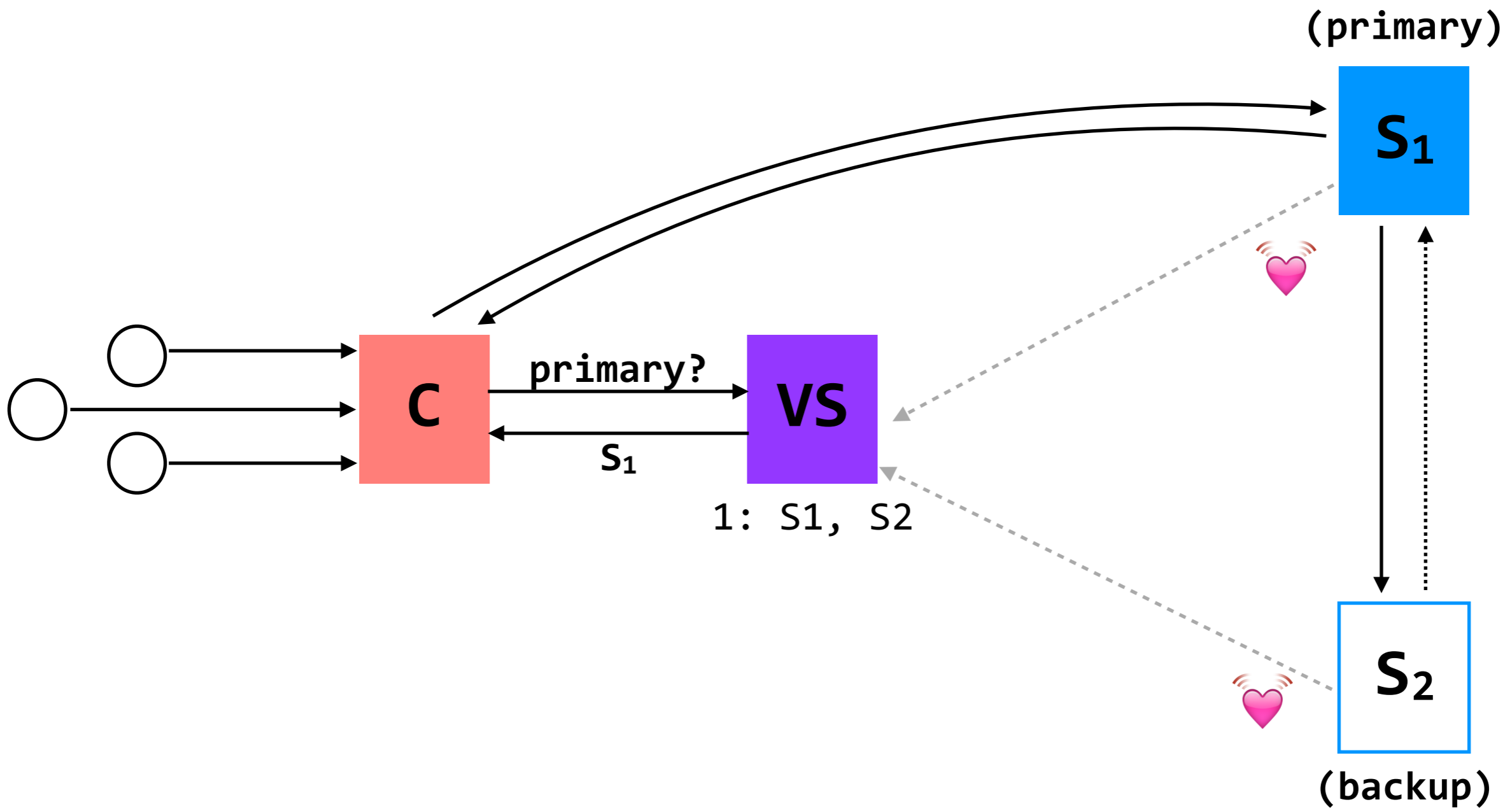


use a **view server**, which determines which replica is the primary



use a **view server**, which determines which replica is the primary

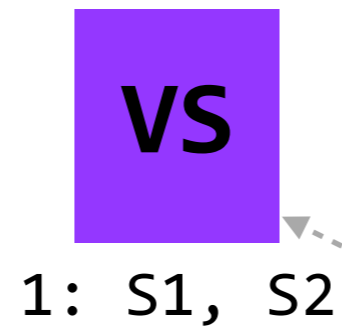
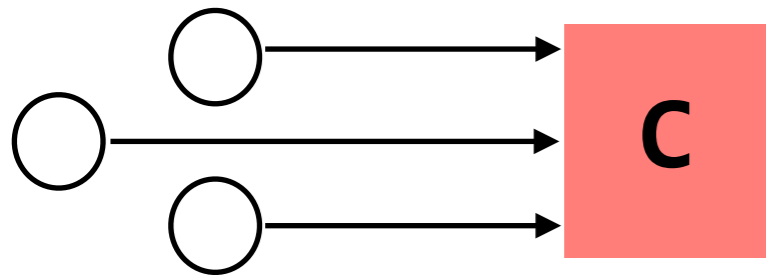




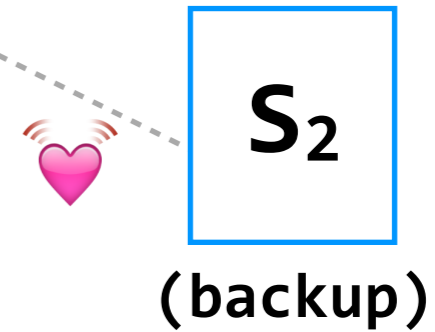
use a **view server**, which determines which replica is the primary

# handling primary failure

(dead)

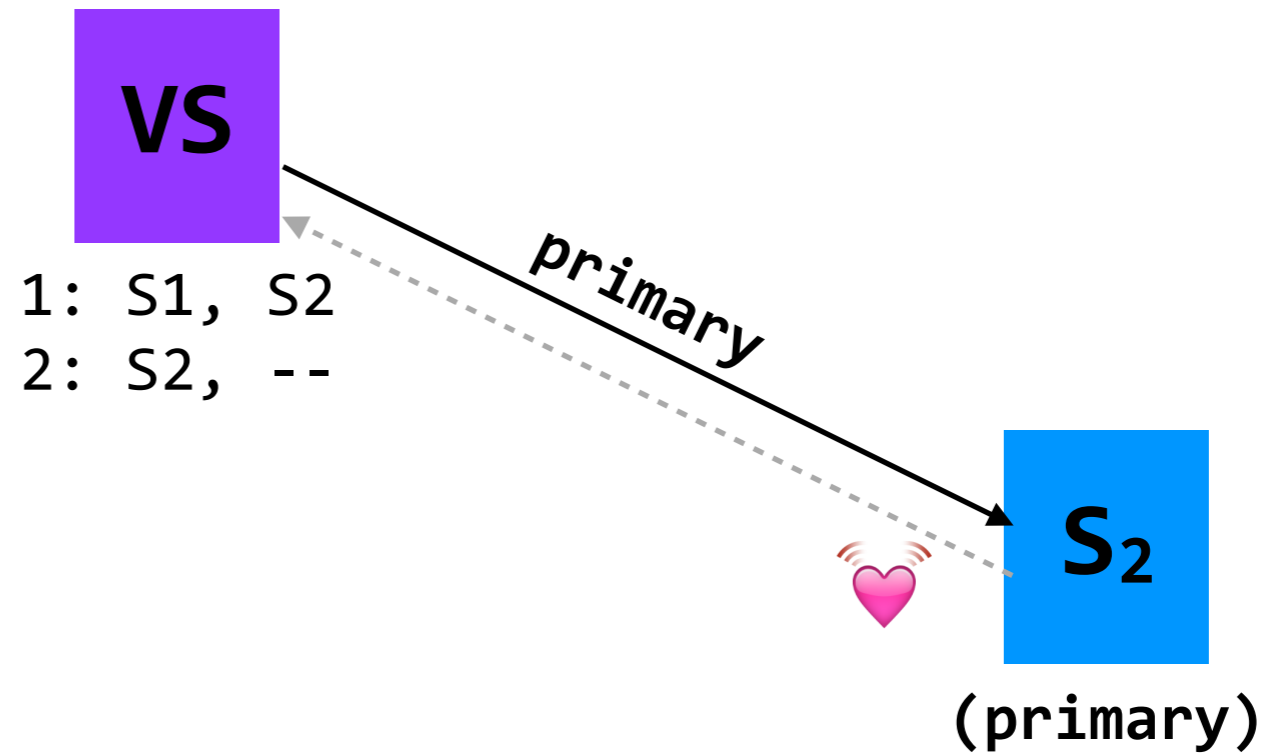
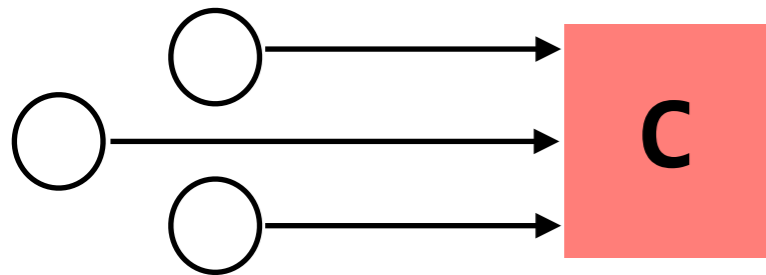


lack of pings indicates  
to **VS** that **S<sub>1</sub>** is down



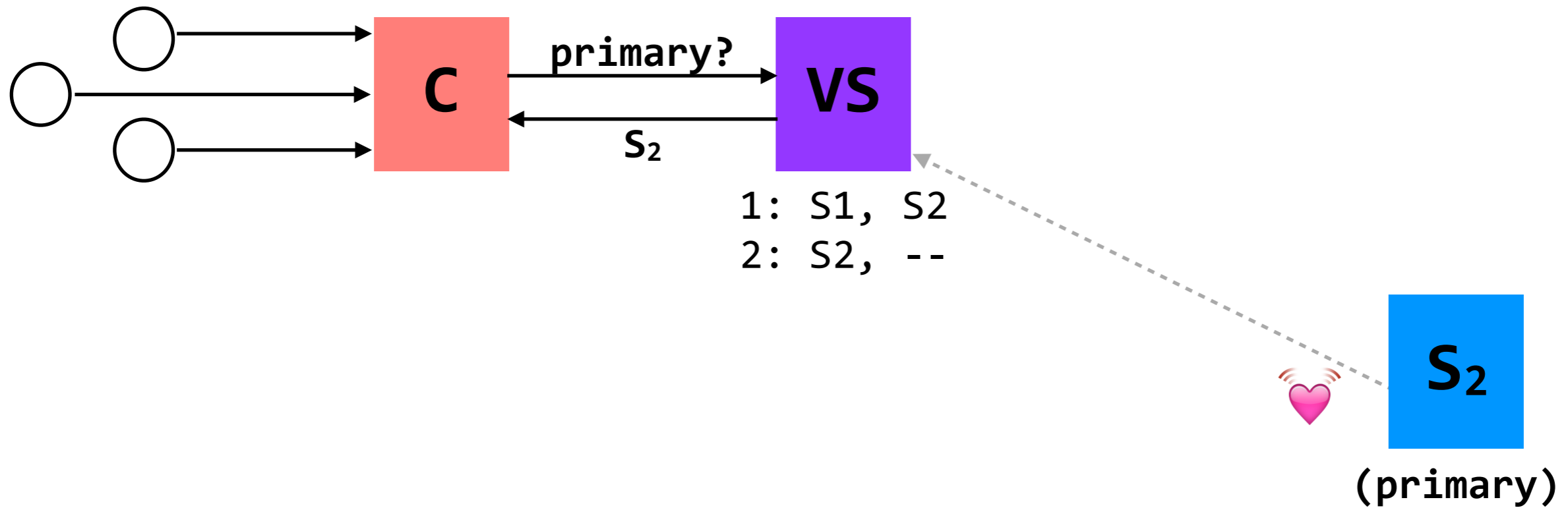
# handling primary failure

(dead)



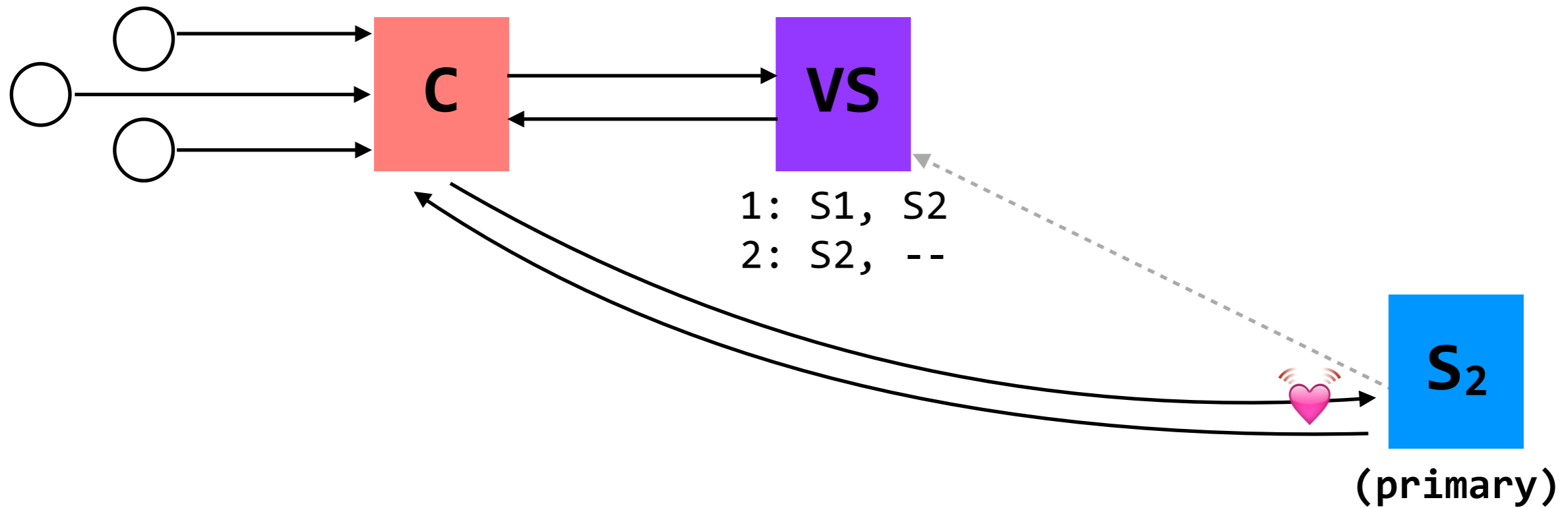
# handling primary failure

(dead)



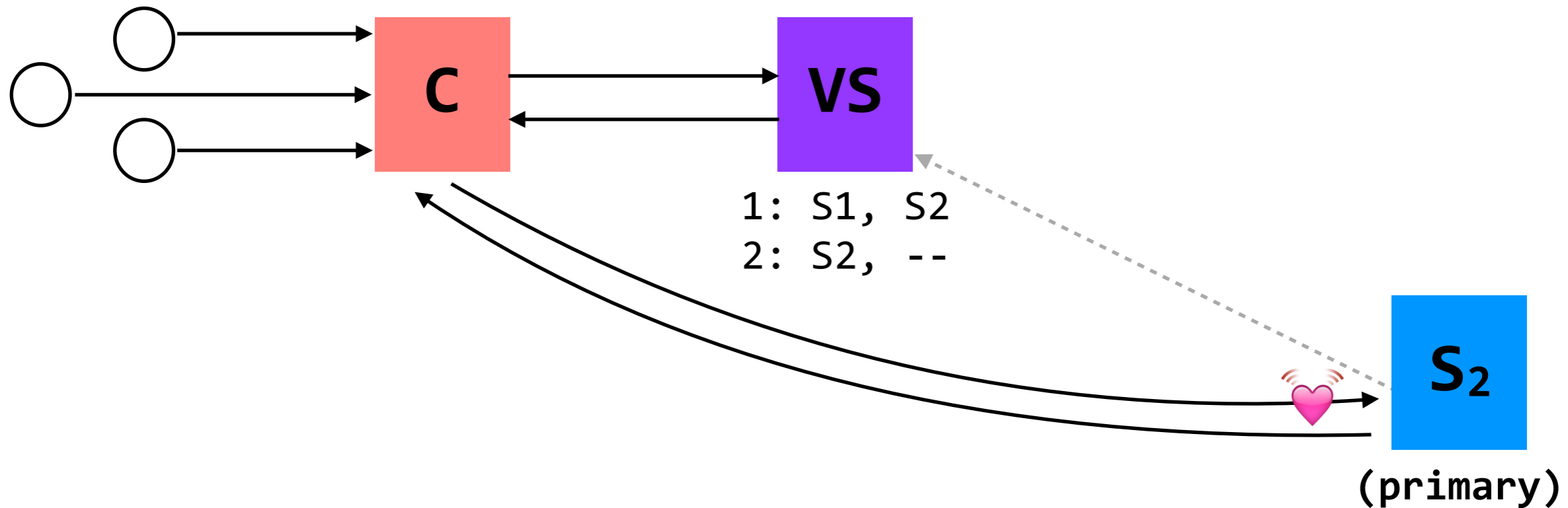
# handling primary failure

(dead)



# handling primary failure

(dead)



**before  $S_2$  knows it's primary, it will reject any requests from clients**

(and if clients had contacted  $S_1$  after it failed but before it was deemed dead, they would have received no response)

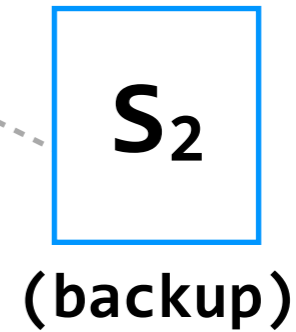
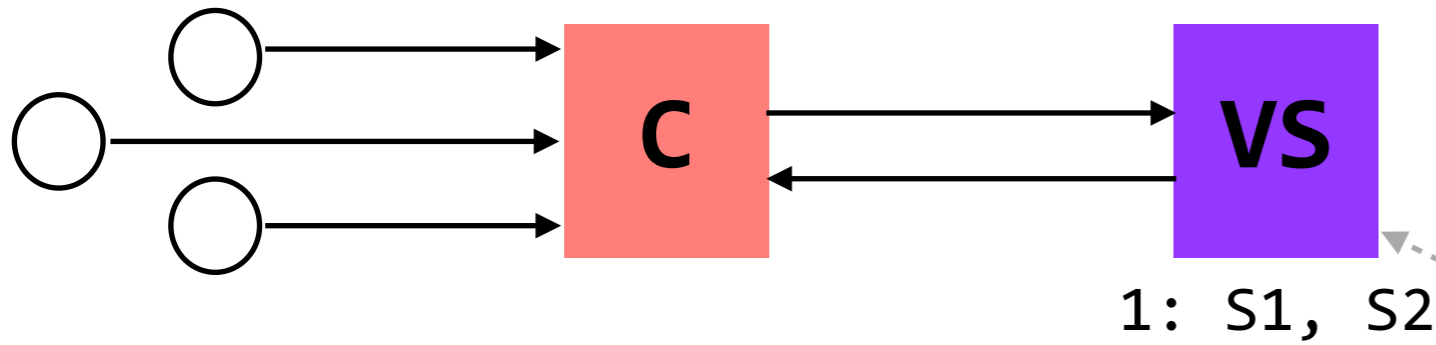
# handling primary failure due to partition

network partition

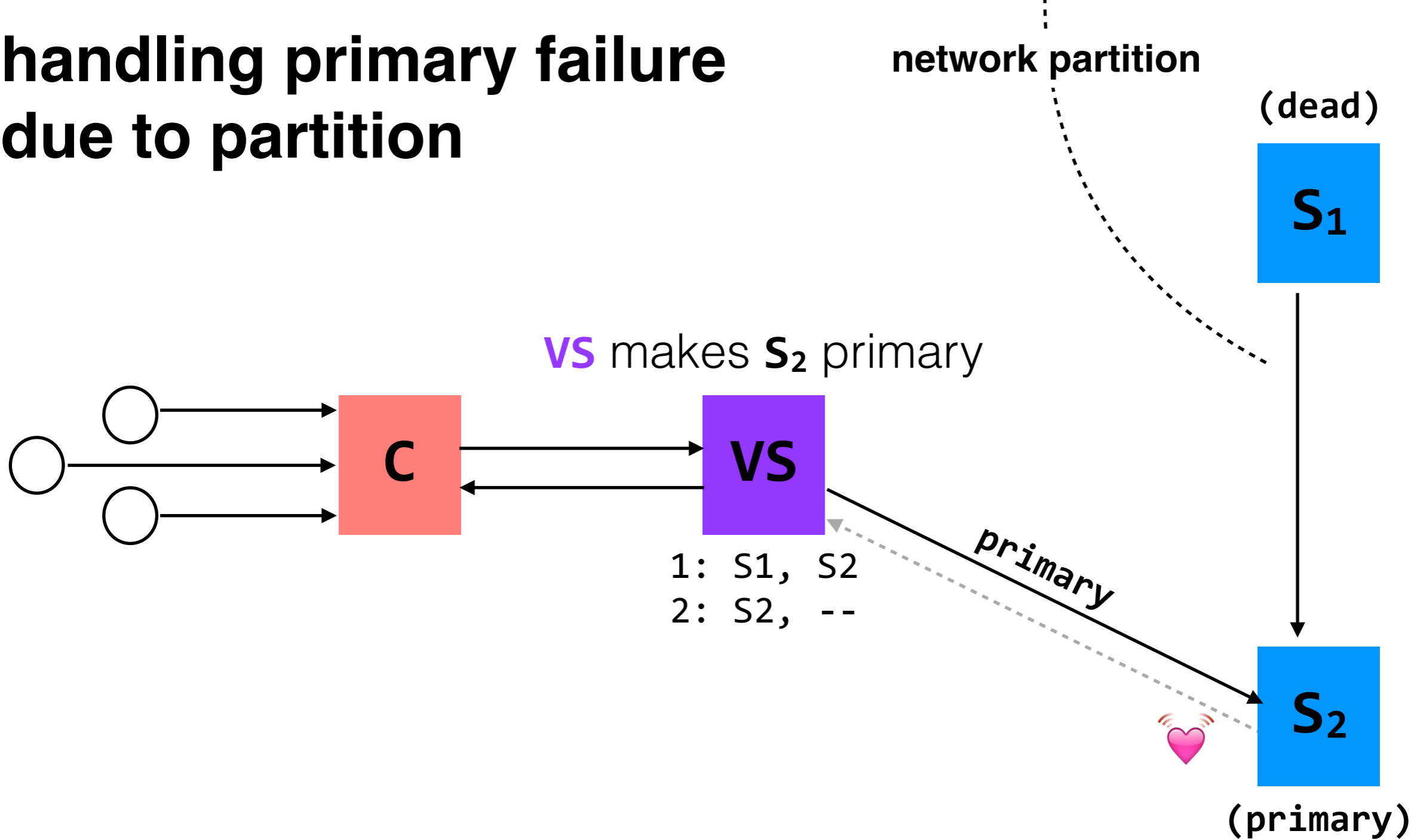
(dead)



lack of pings indicates  
to **VS** that **S<sub>1</sub>** is down

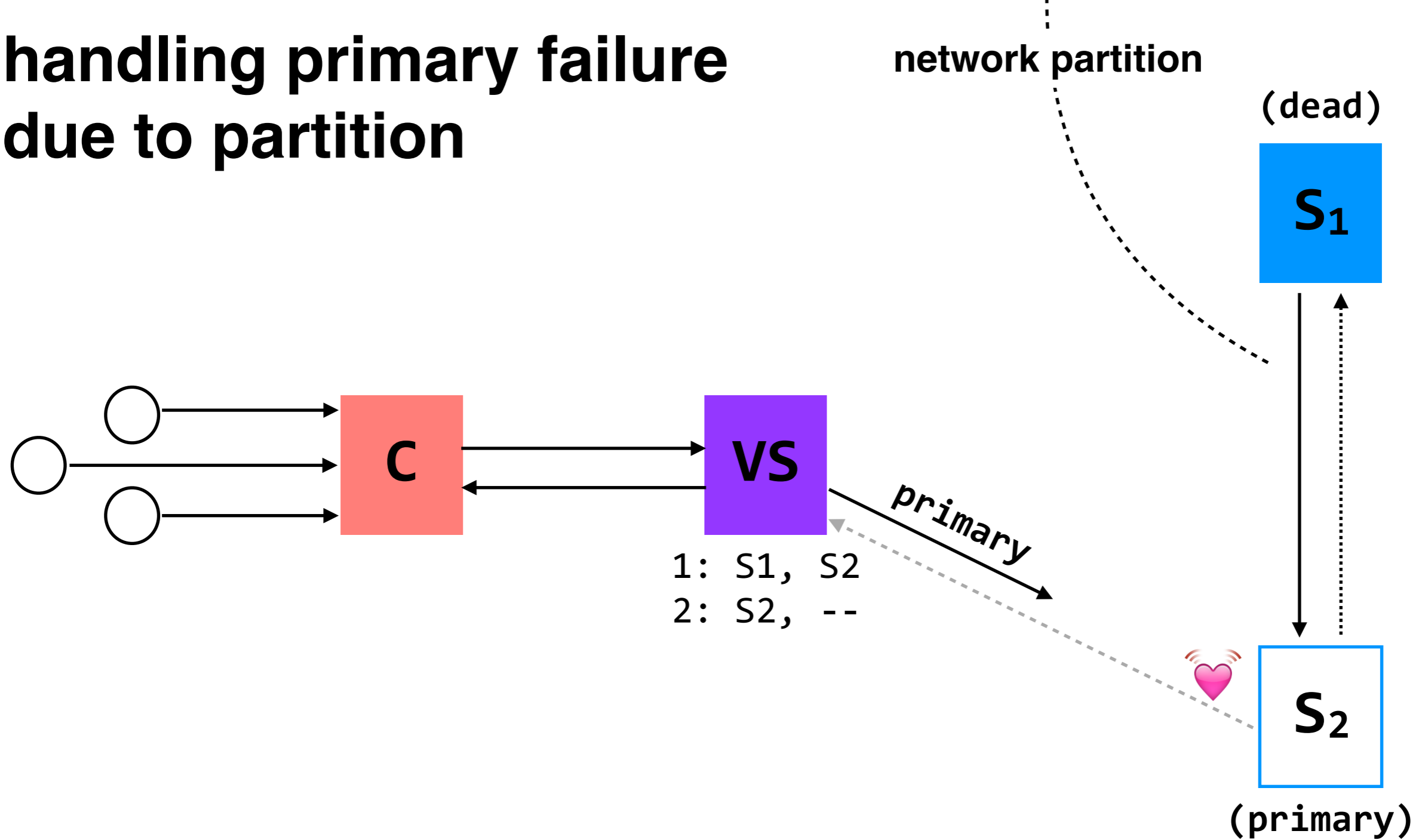


# handling primary failure due to partition



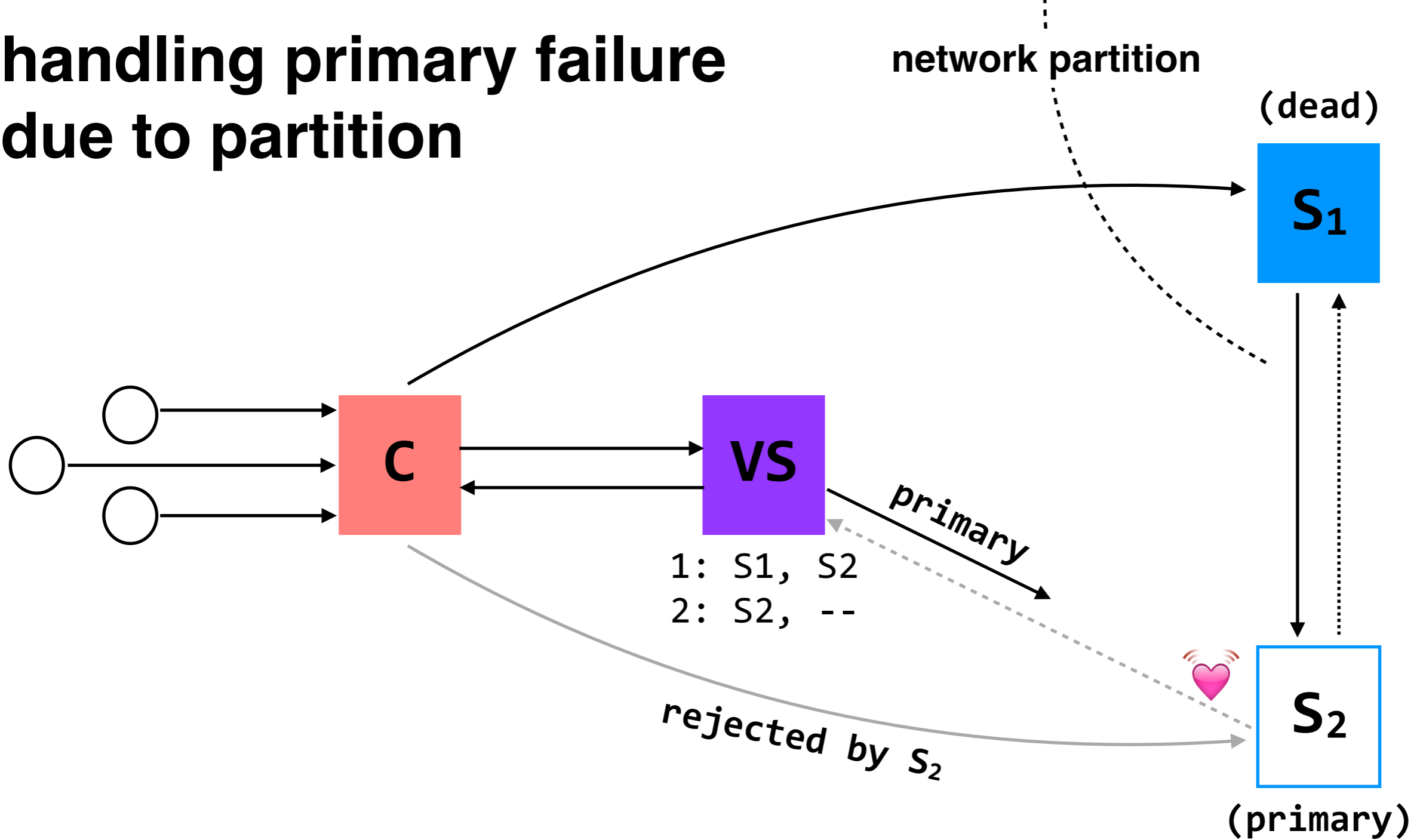


# handling primary failure due to partition



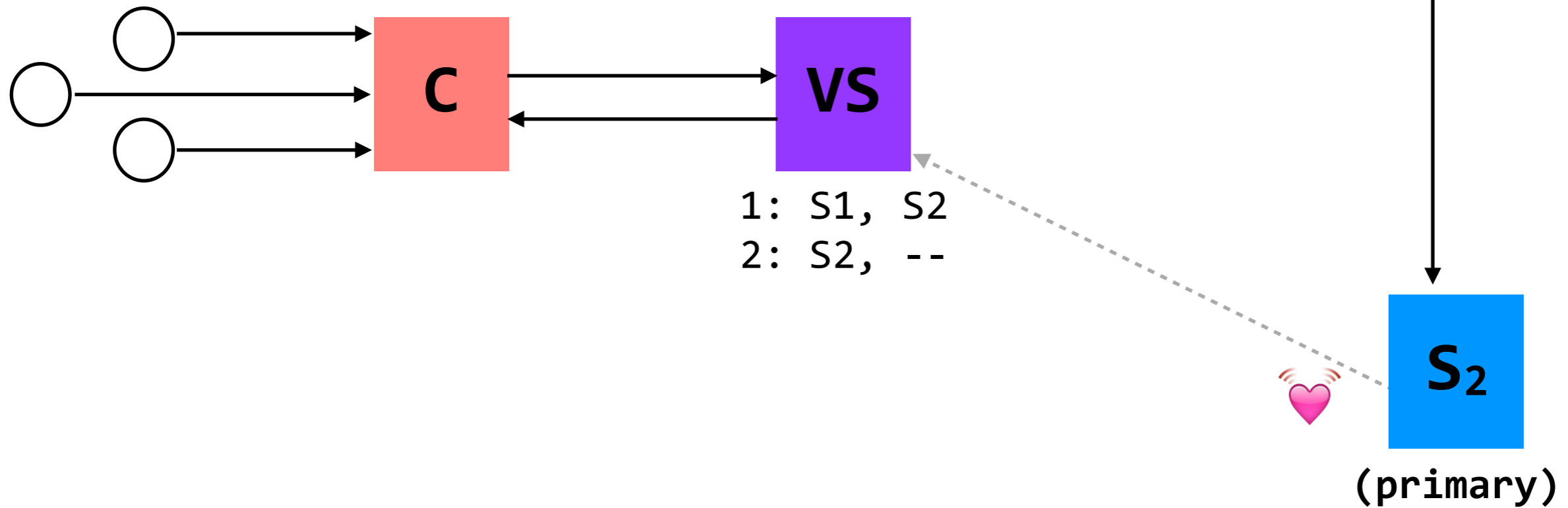
**problem:** what happens before  $S_2$  knows  
it's the primary?

# handling primary failure due to partition



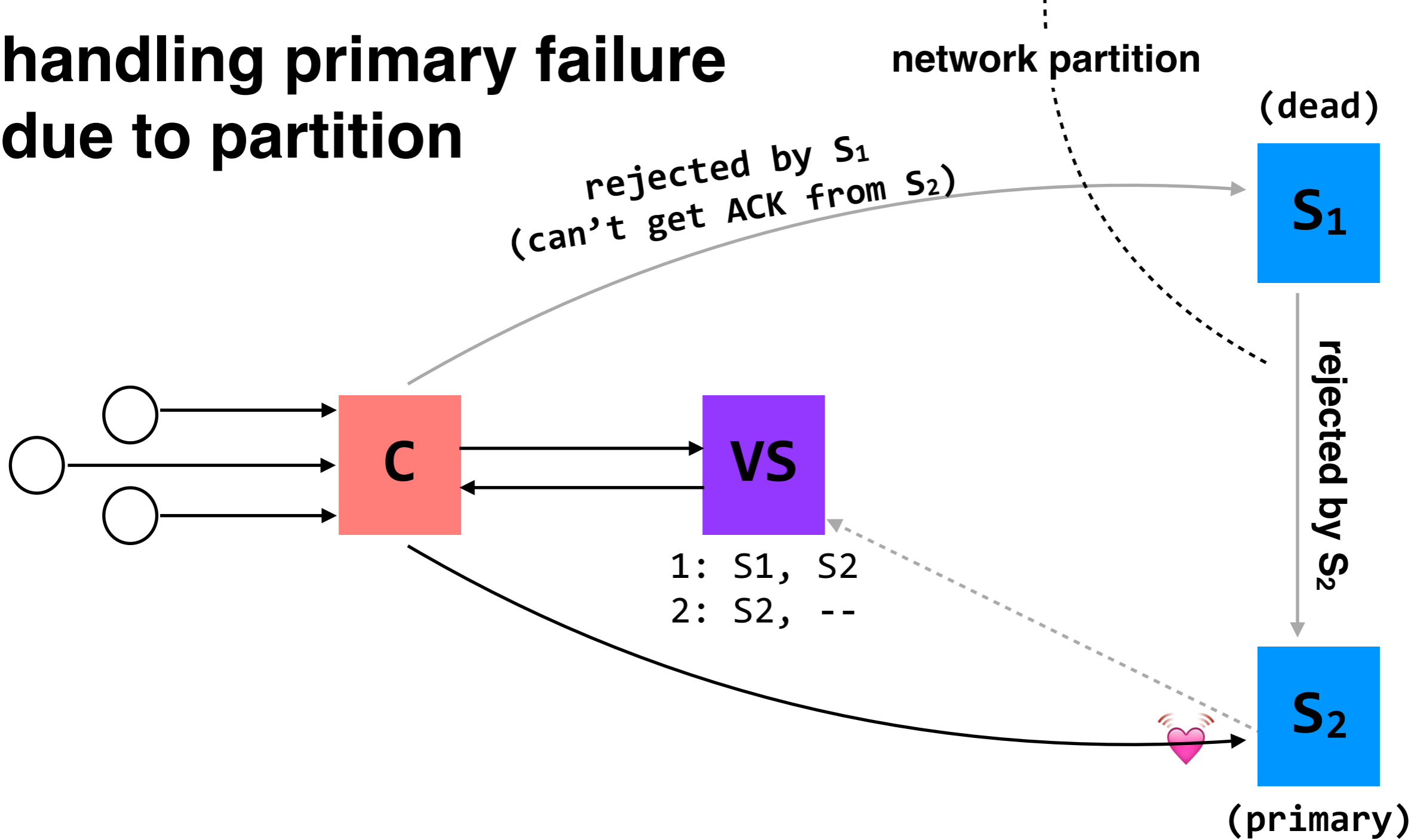
**it's okay! S<sub>2</sub> will act as backup**  
(accept updates from S<sub>1</sub>, reject coordinator requests)

# handling primary failure due to partition

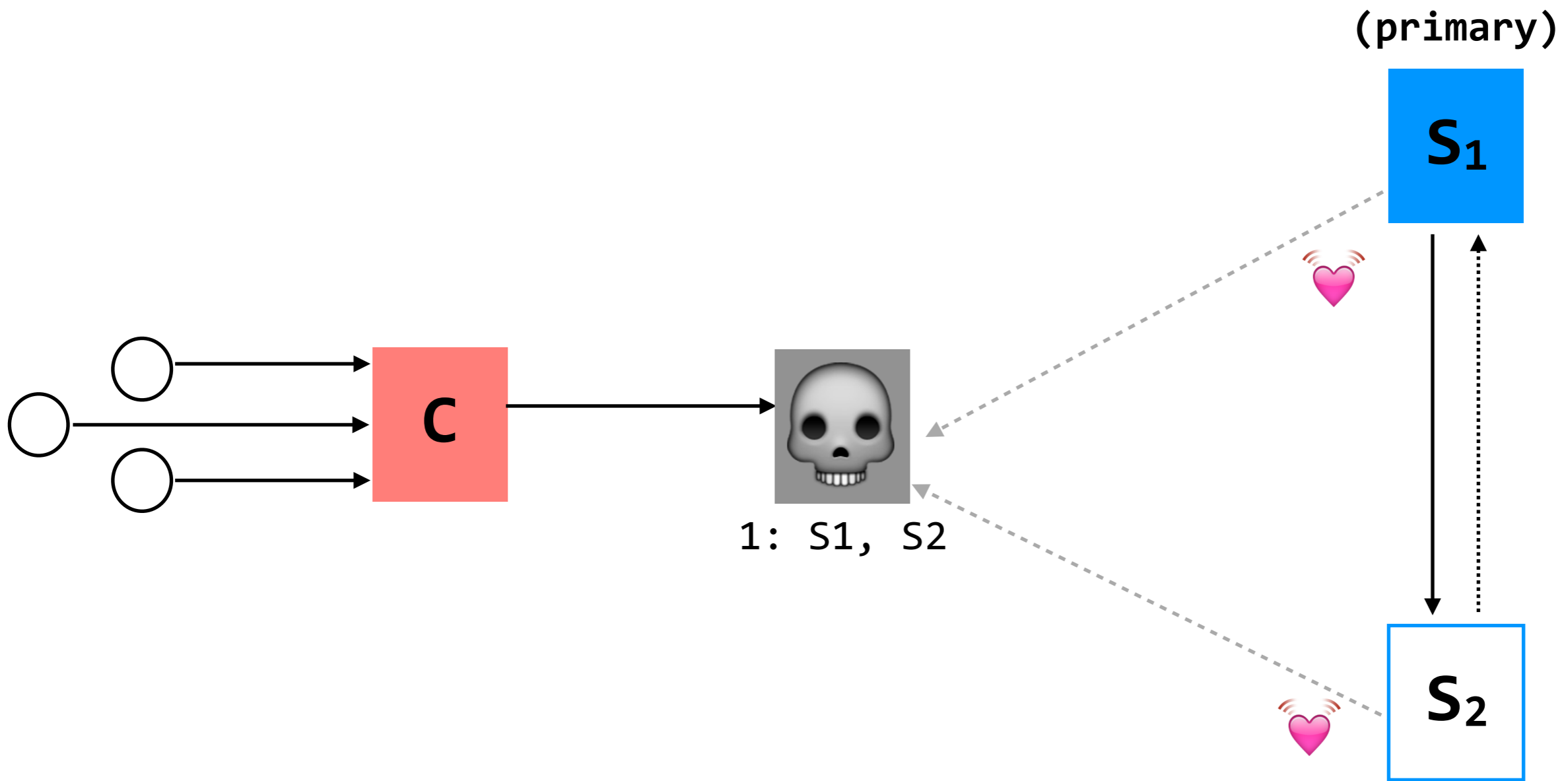


**problem:** what happens after  $S_2$  knows it's the primary, but  $S_1$  also thinks it is?

# handling primary failure due to partition



**also okay! S<sub>1</sub> won't be able to act as primary**  
(can't accept client requests because it won't get ACKs from S<sub>2</sub>)



**problem:** what if view server fails?

**go to recitation tomorrow and find out!**

- **Replicated state machines (RSMs)** provide **single-copy consistency**: operations complete as if there is a single copy of the data, though internally there are replicas.
- RSMs use a **primary-backup** mechanism for replication. The **view server** ensures that only one replica acts as the primary. It can also recruit new backups after servers fail.
- To extend this model to handle view-server failures, we need a mechanism to provide **distributed consensus**; see tomorrow's recitation (on RAFT).