

6.033 Spring 2018

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

C₁ **write₁(X)**

S₁

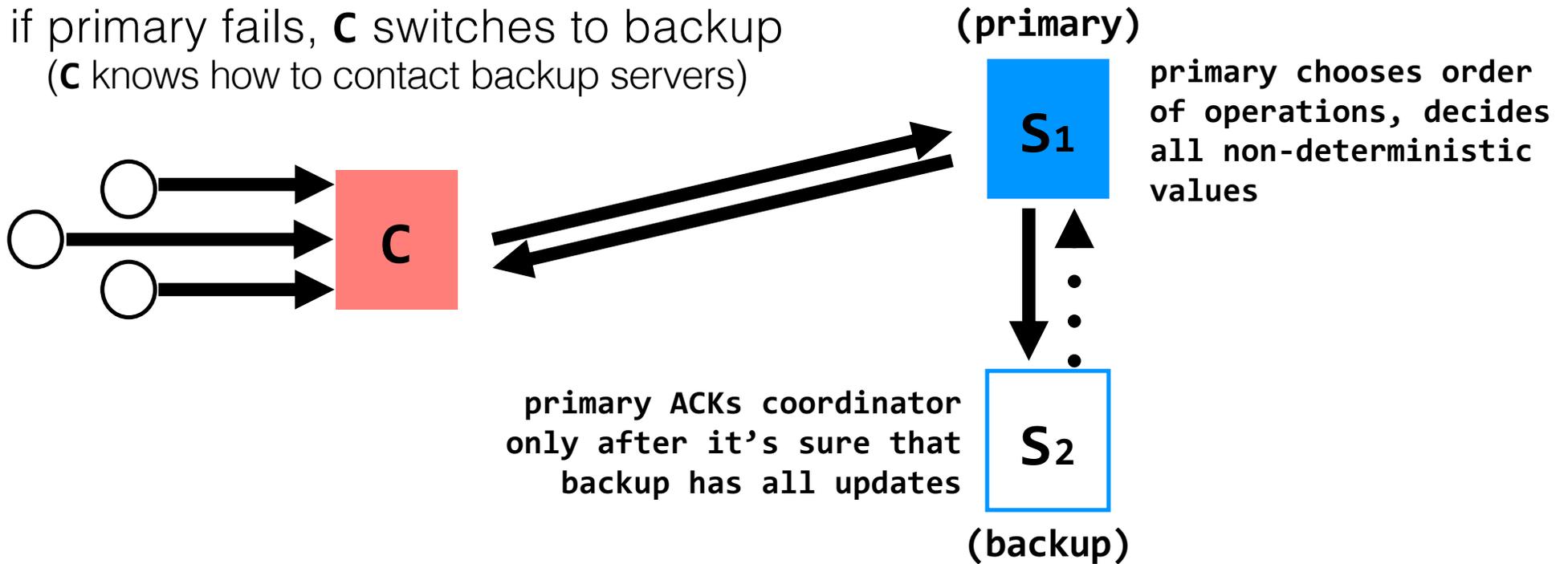
C₂ **write₂(X)**

S₂

(replica of S₁)

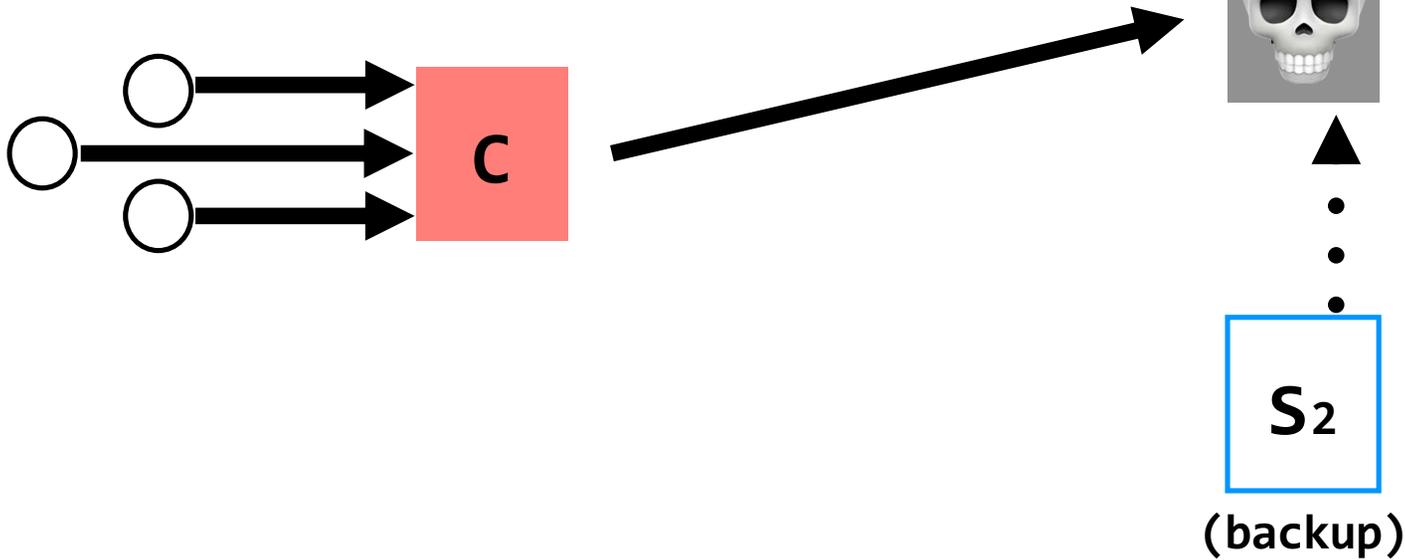
problem: replica servers can become inconsistent

if primary fails, **C** switches to backup
(**C** knows how to contact backup servers)



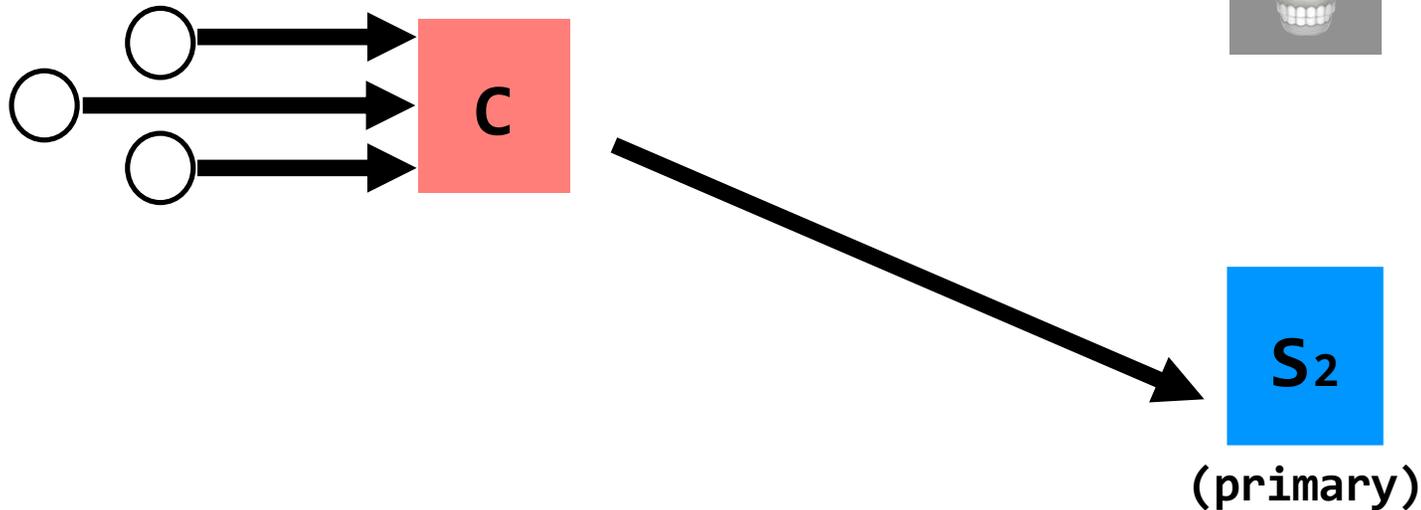
attempt: coordinators communicate with primary servers, who communicate with backup servers

if primary fails, **C** switches to backup
(**C** knows how to contact backup servers)



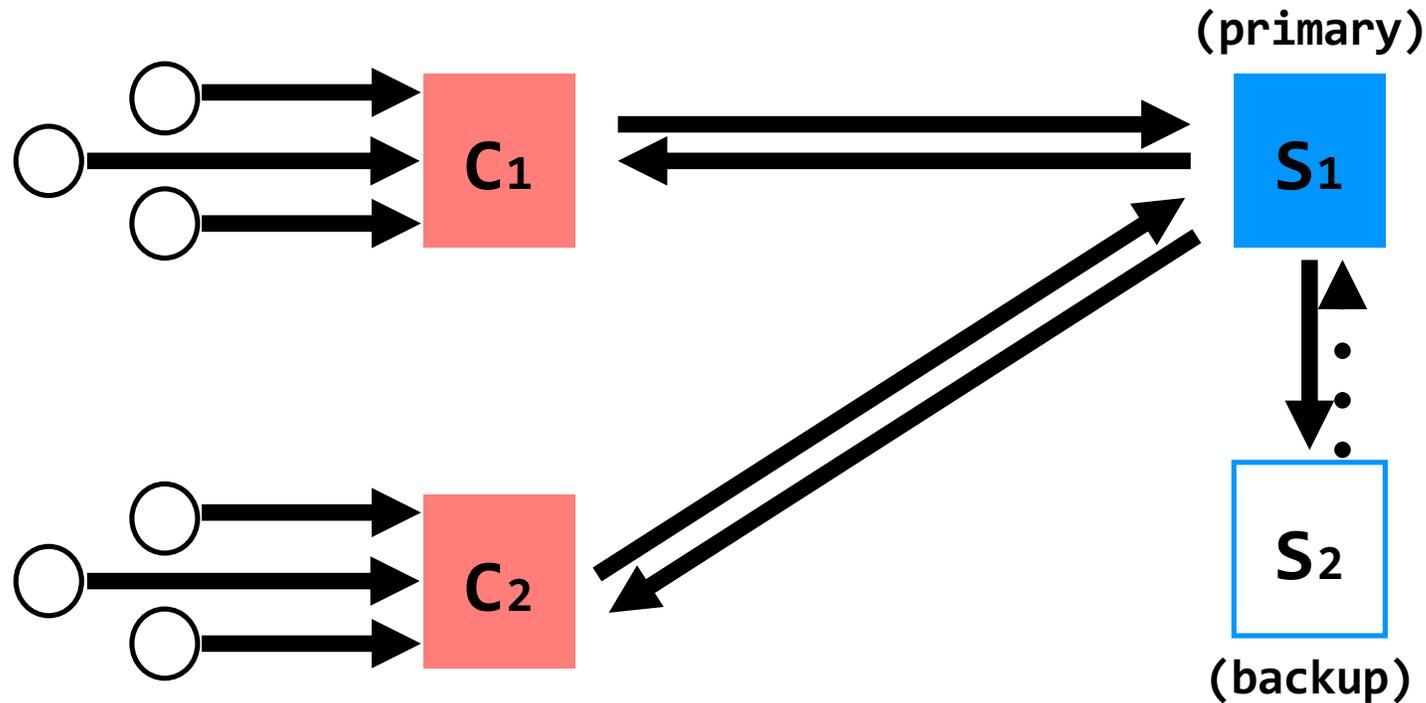
attempt: coordinators communicate with primary servers, who communicate with backup servers

if primary fails, **C** switches to backup
(**C** knows how to contact backup servers)



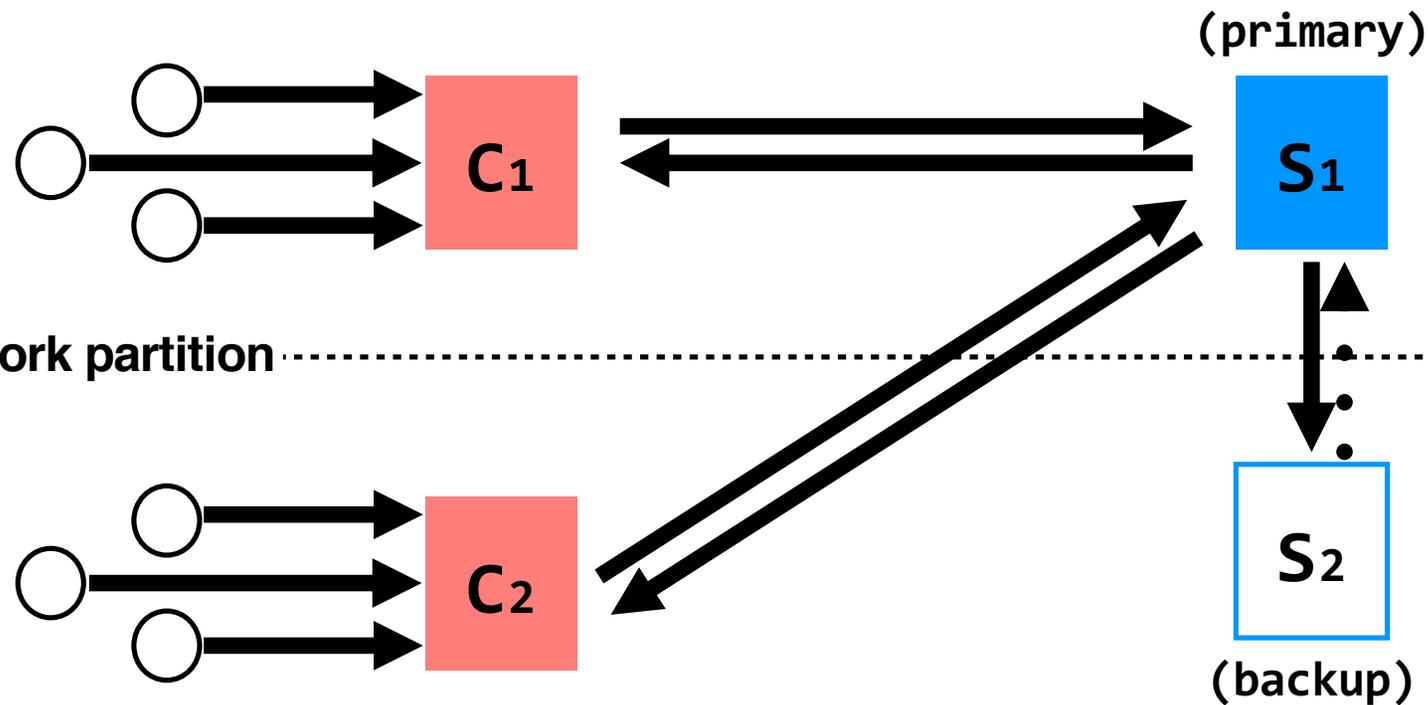
attempt: coordinators communicate with primary servers, who communicate with backup servers

multiple coordinators + the network = problems



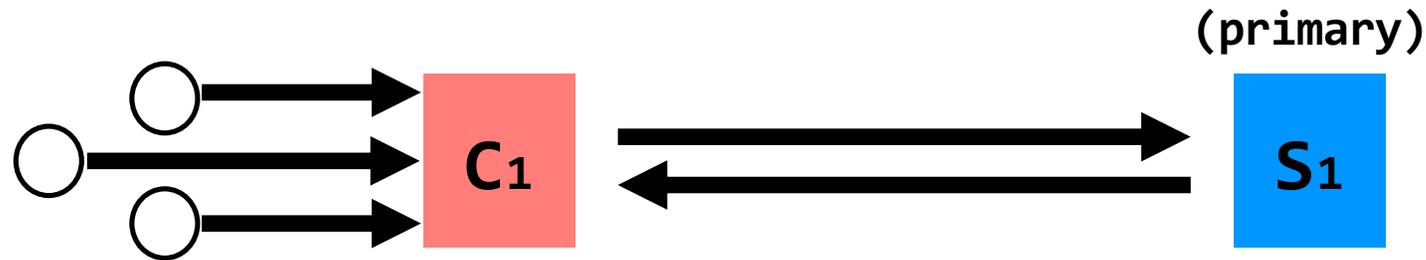
attempt: coordinators communicate with primary servers, who communicate with backup servers

multiple coordinators + the network = problems

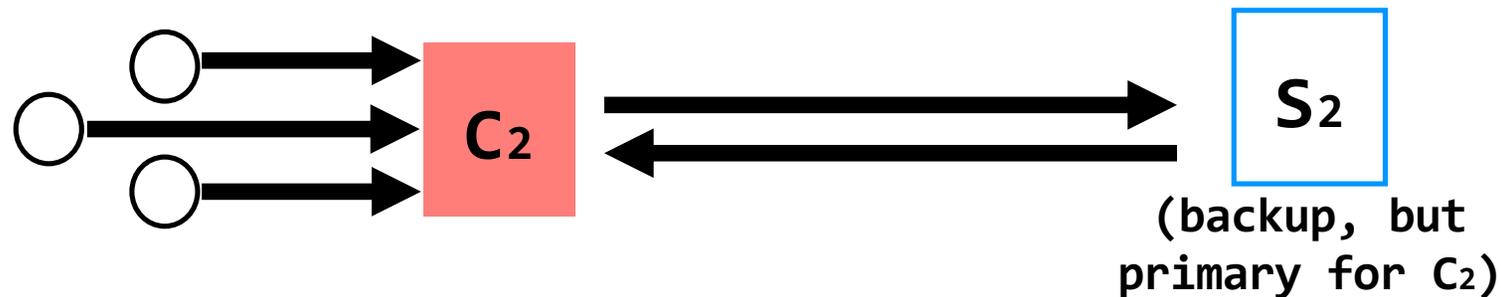


attempt: coordinators communicate with primary servers, who communicate with backup servers

multiple coordinators + the network = problems

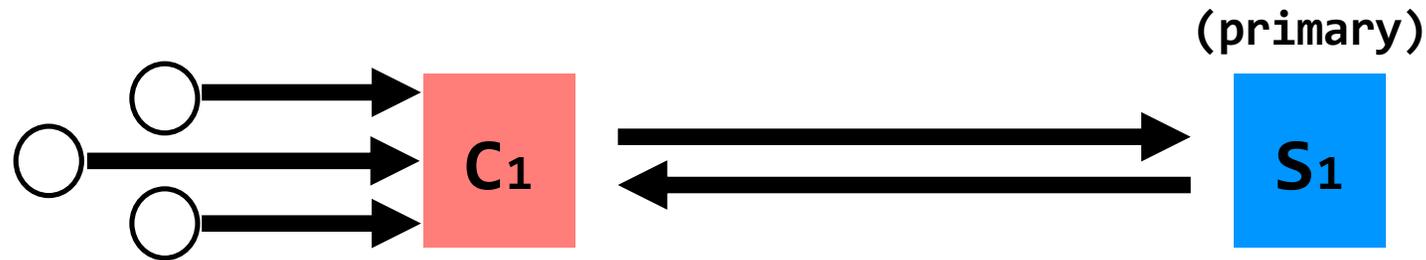


... network partition

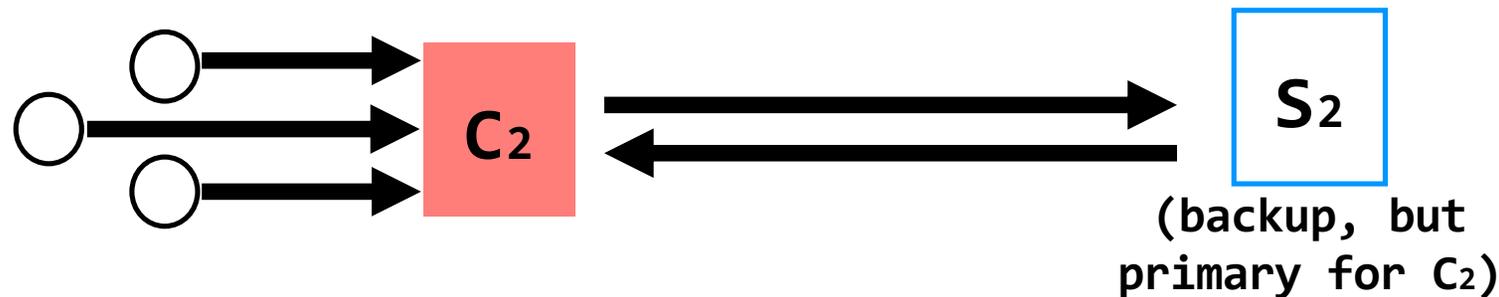


attempt: coordinators communicate with primary servers, who communicate with backup servers

multiple coordinators + the network = problems

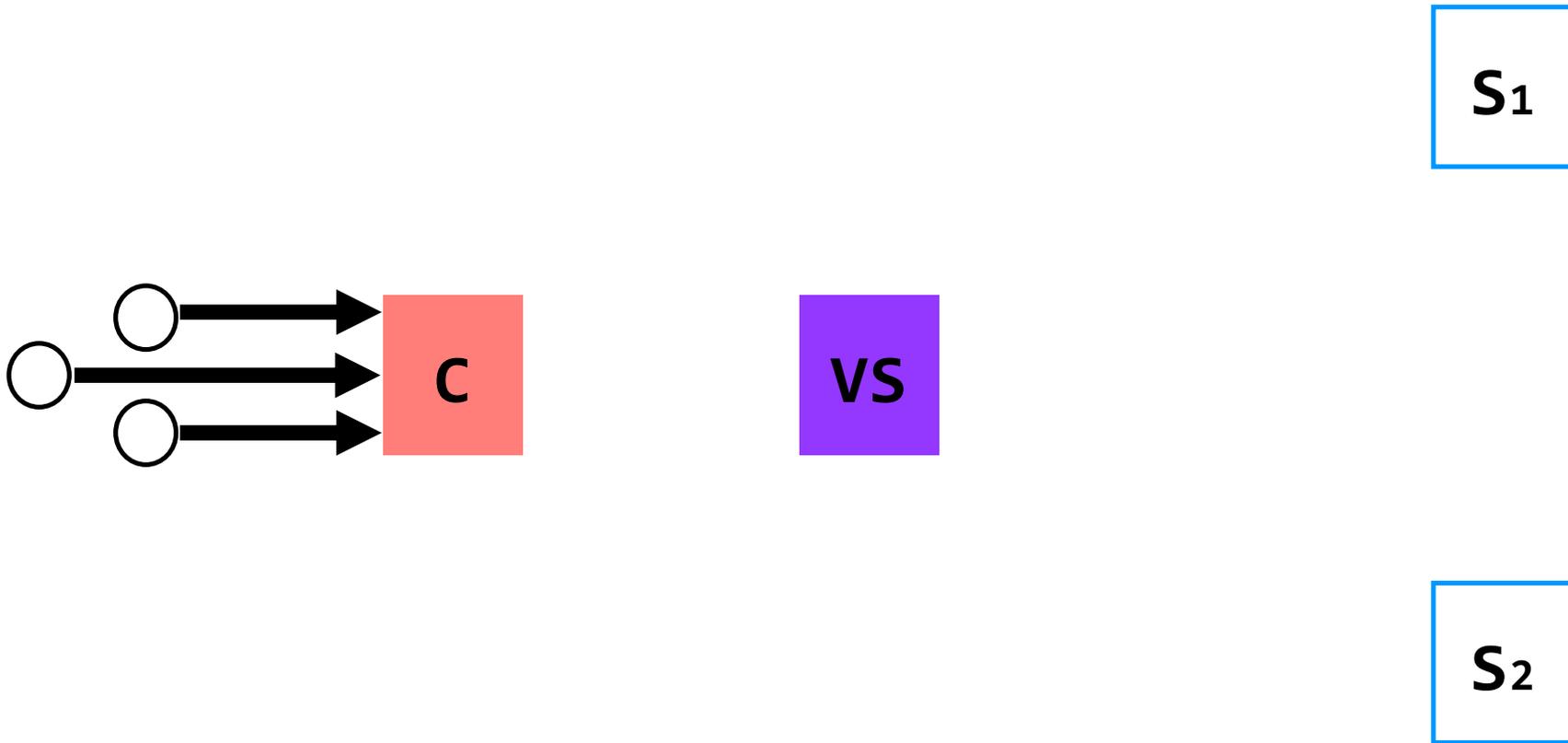


... network partition

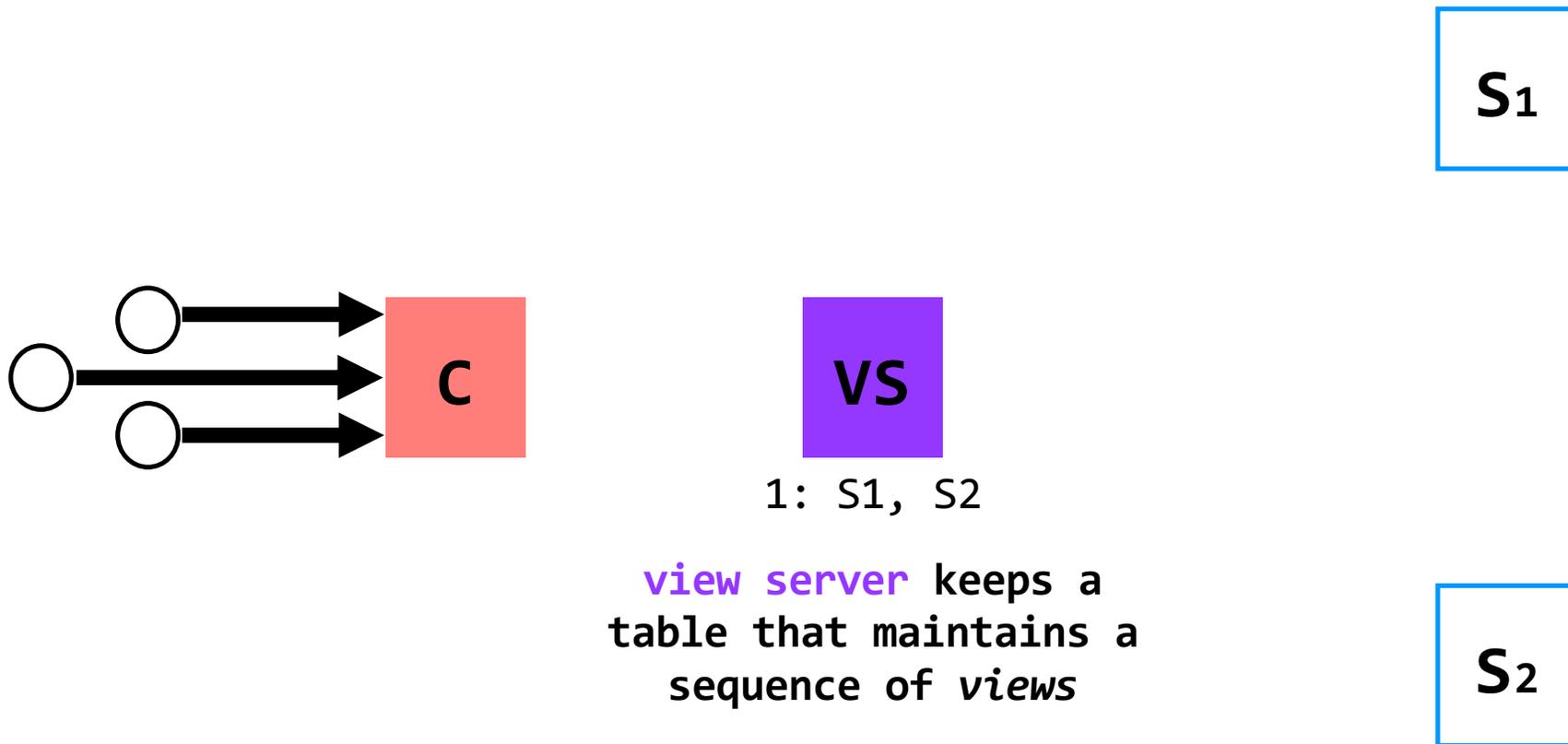


C₁ and **C₂** are using different primaries;
S₁ and **S₂** are no longer consistent

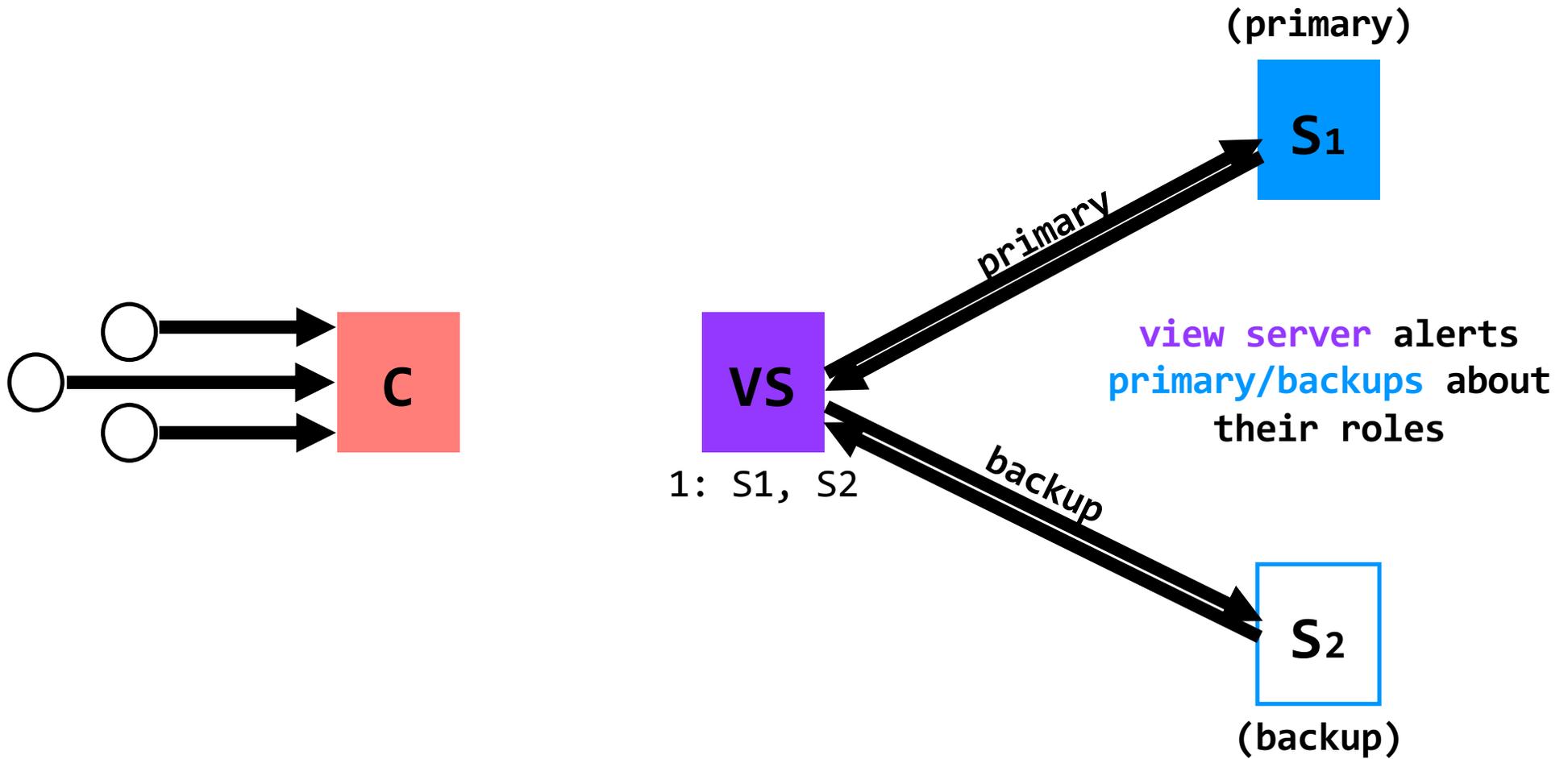
attempt: coordinators communicate with primary servers, who communicate with backup servers



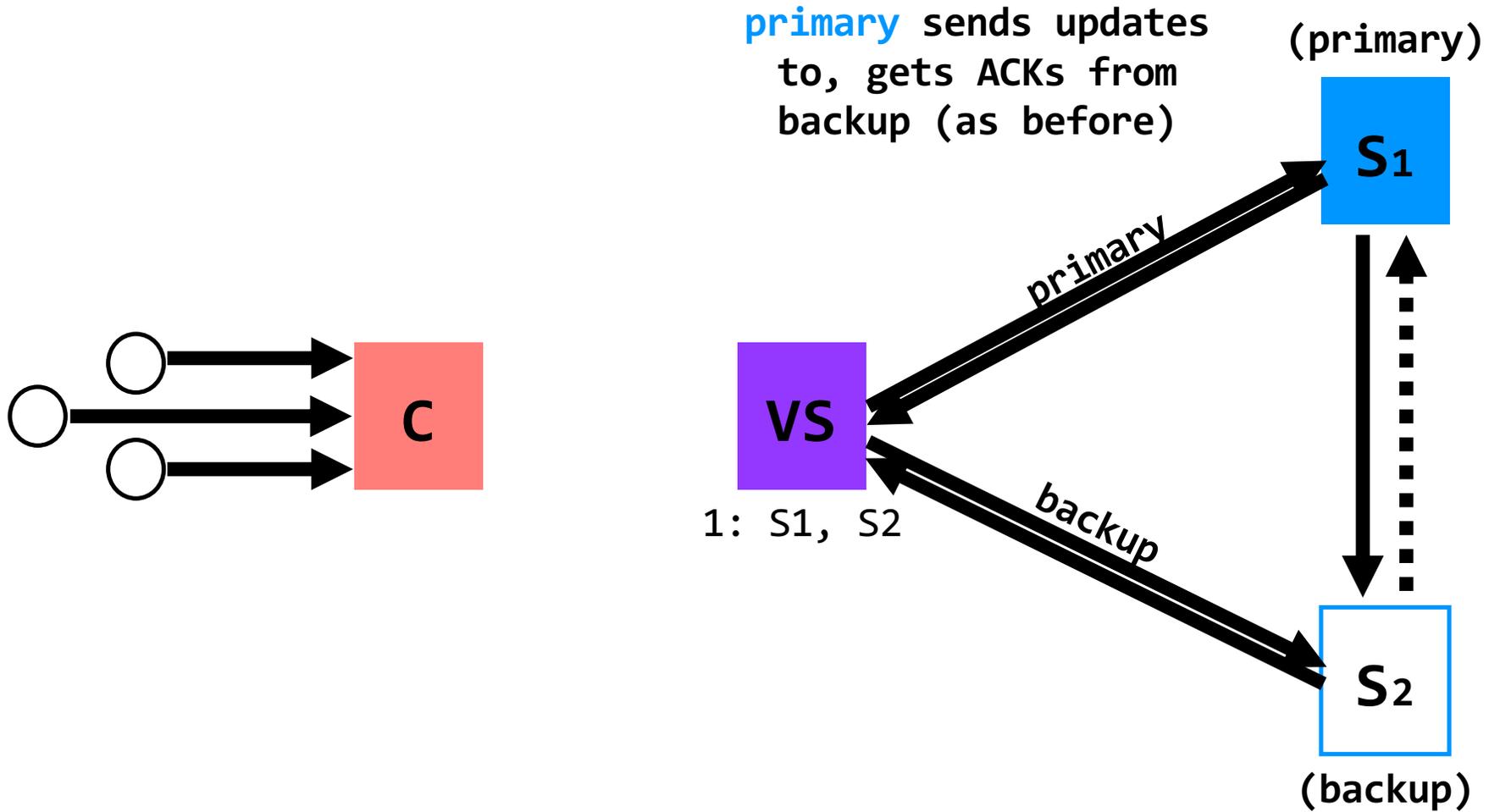
use a **view server**, which determines which replica is the primary



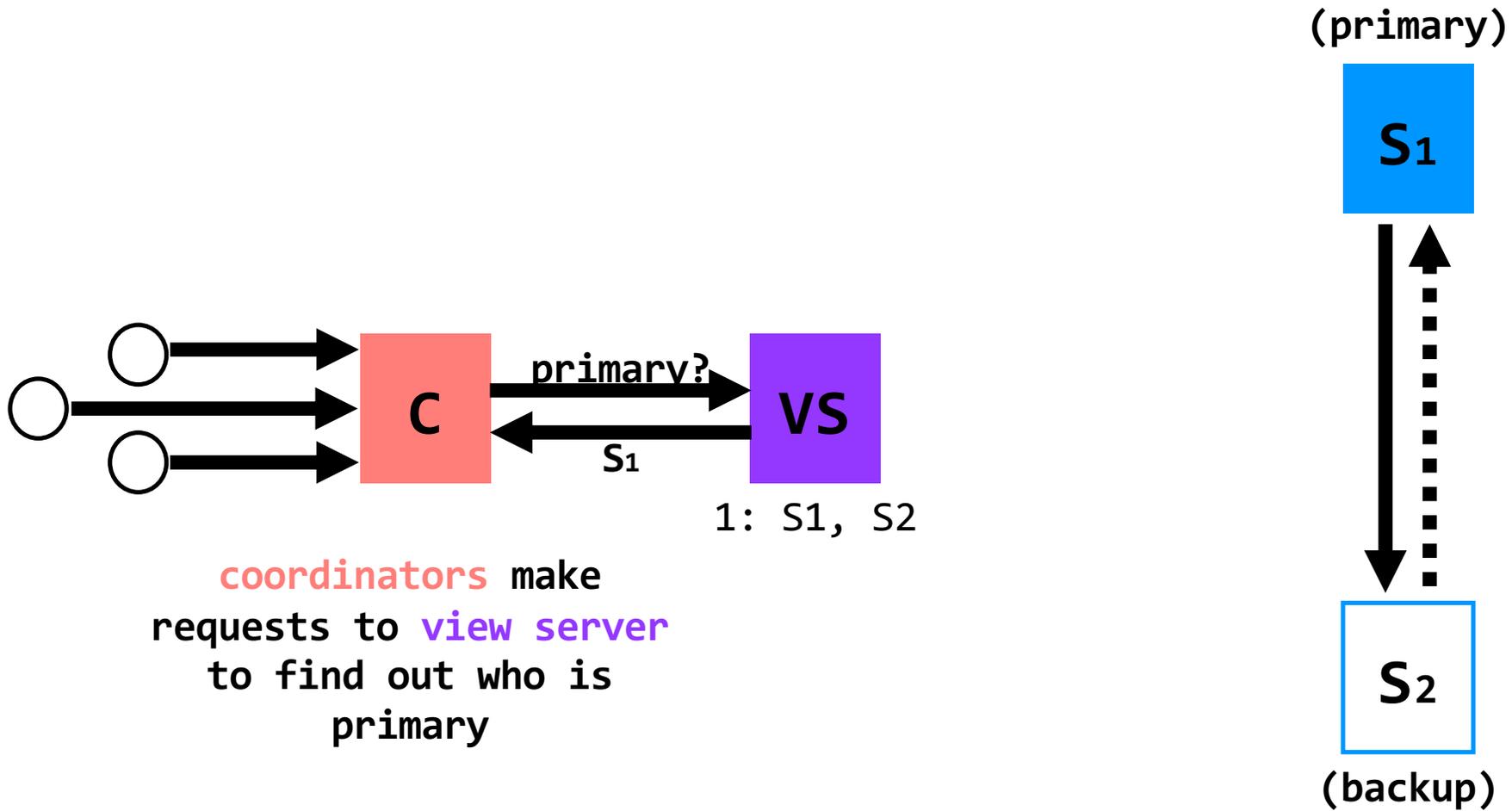
use a **view server**, which determines which replica is the primary



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

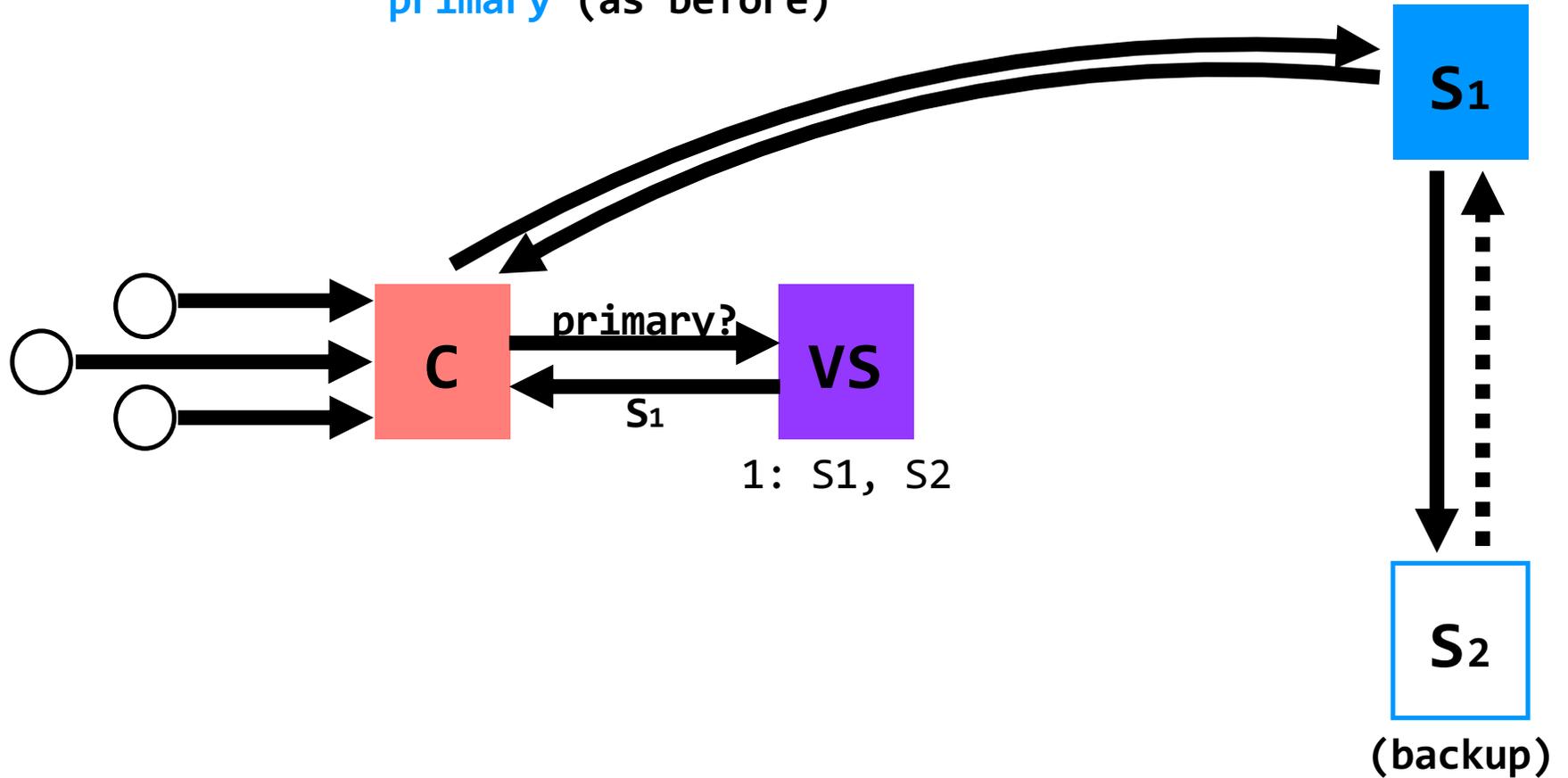


use a **view server**, which determines which replica is the primary

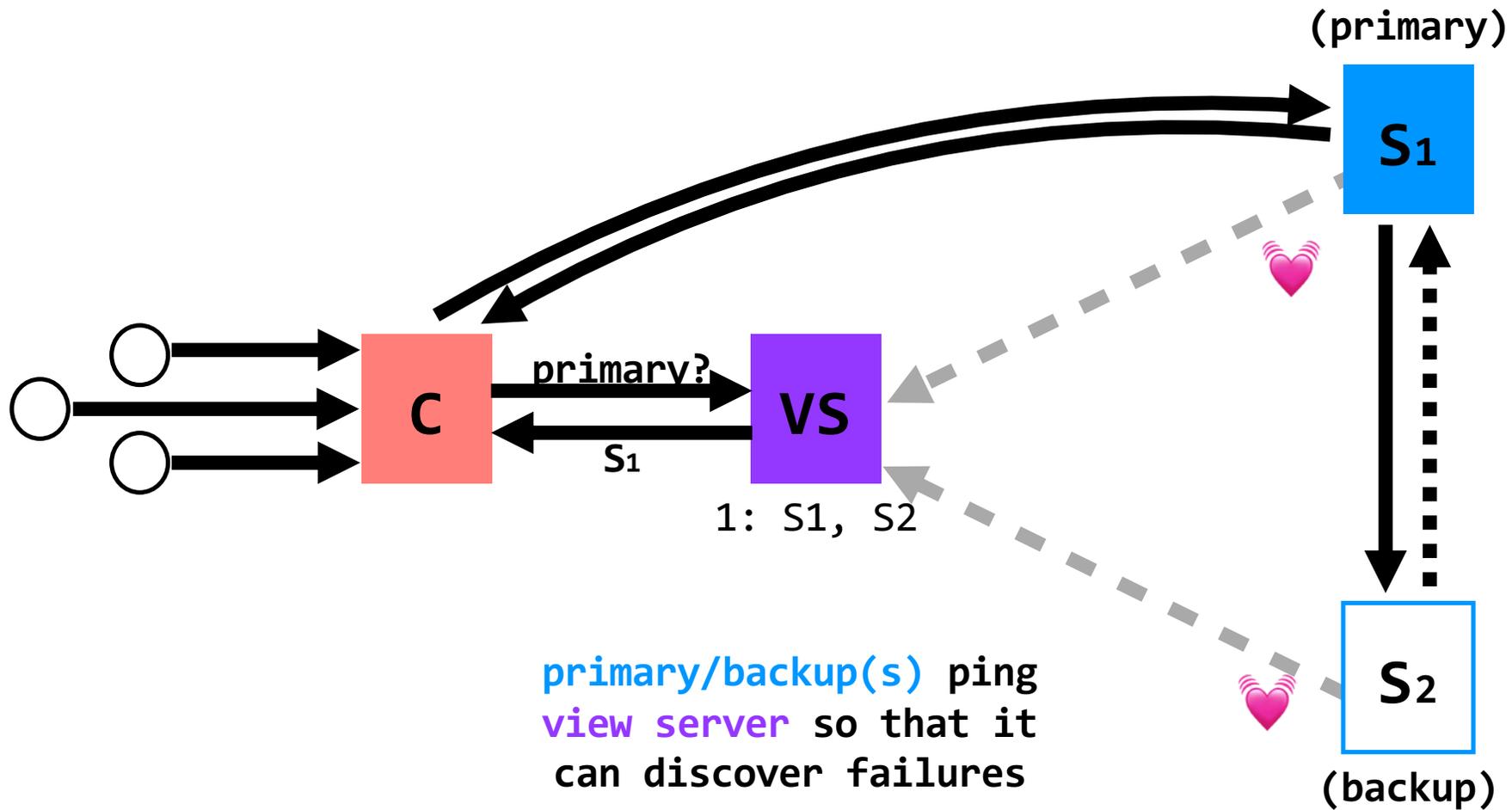


use a **view server**, which determines which replica is the primary

coordinators contact
primary (as before)

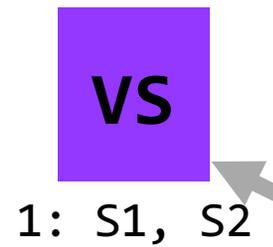
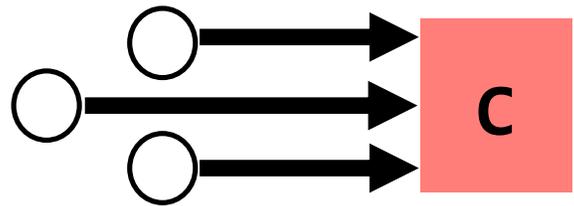


use a **view server**, which determines which replica is the primary

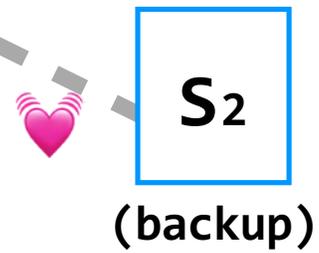


use a **view server**, which determines which replica is the primary

handling primary failure



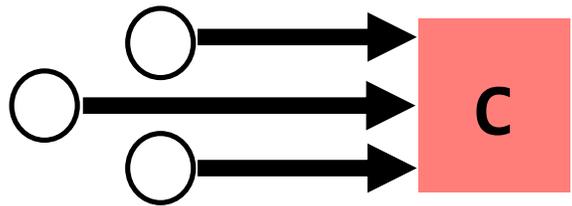
lack of pings indicates to **VS** that **S1** is down



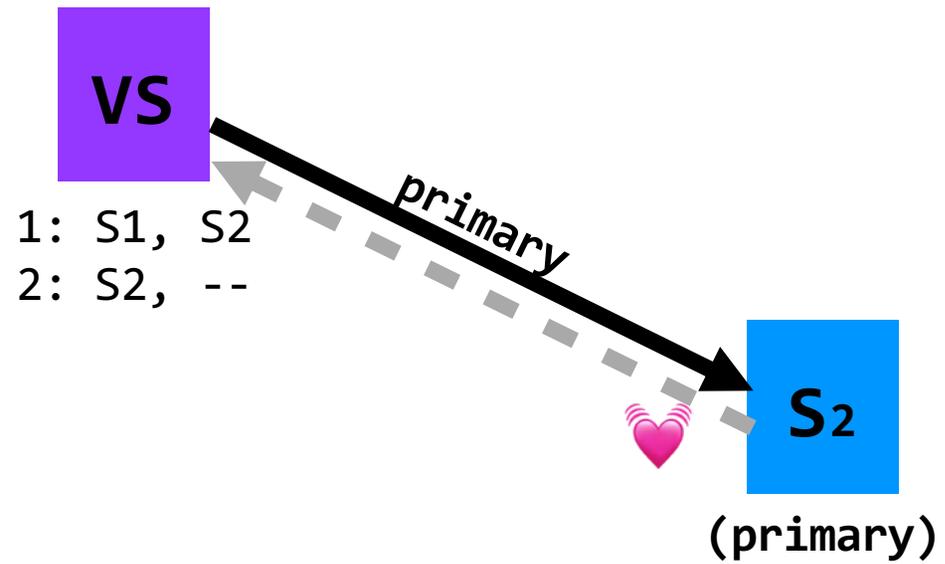
(dead)



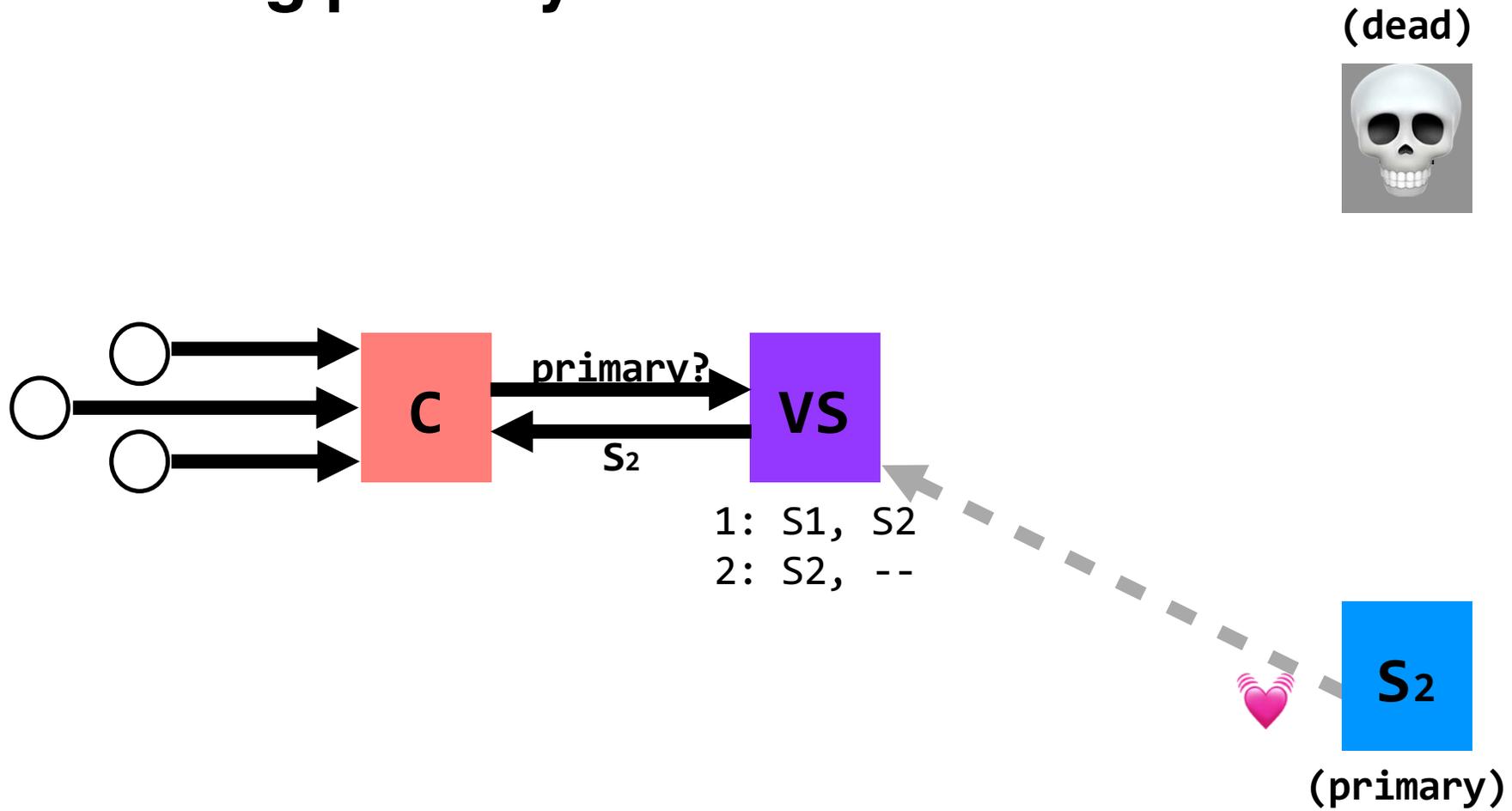
handling primary failure



(dead)

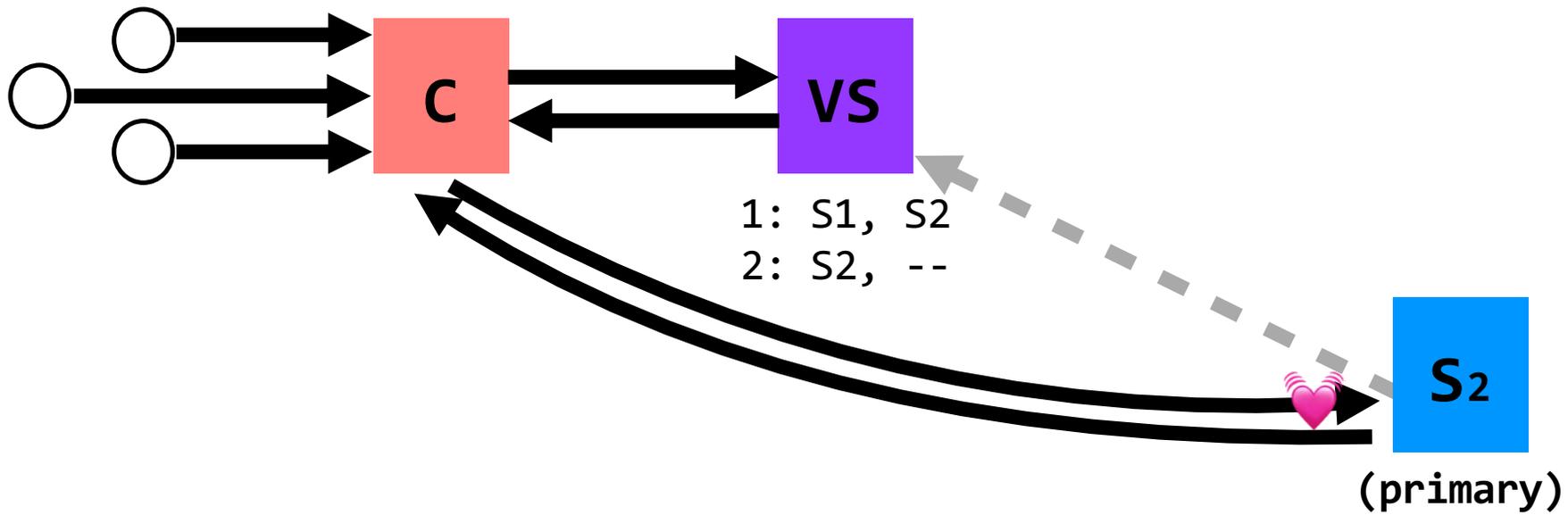


handling primary failure

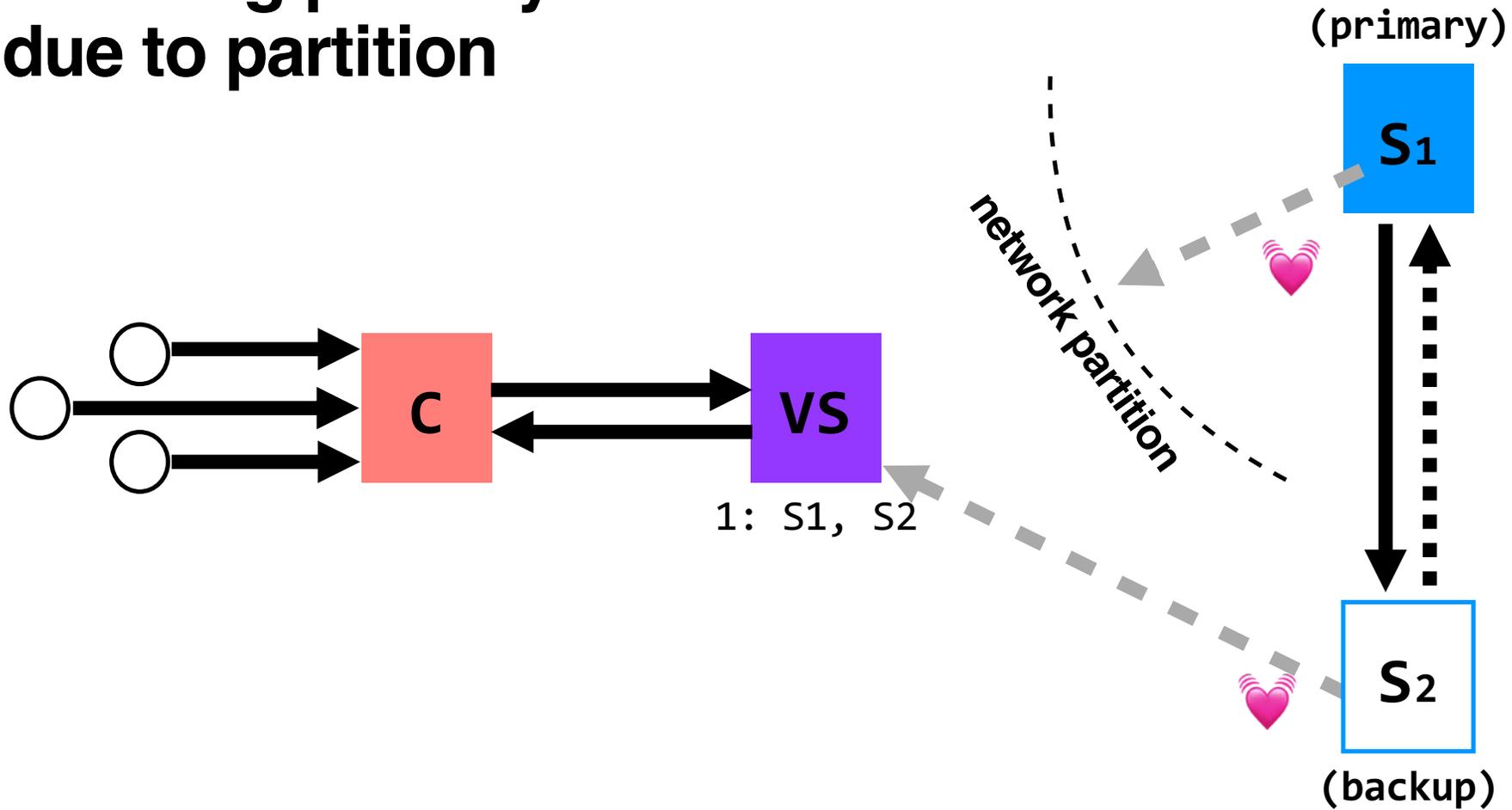


handling primary failure

(dead)

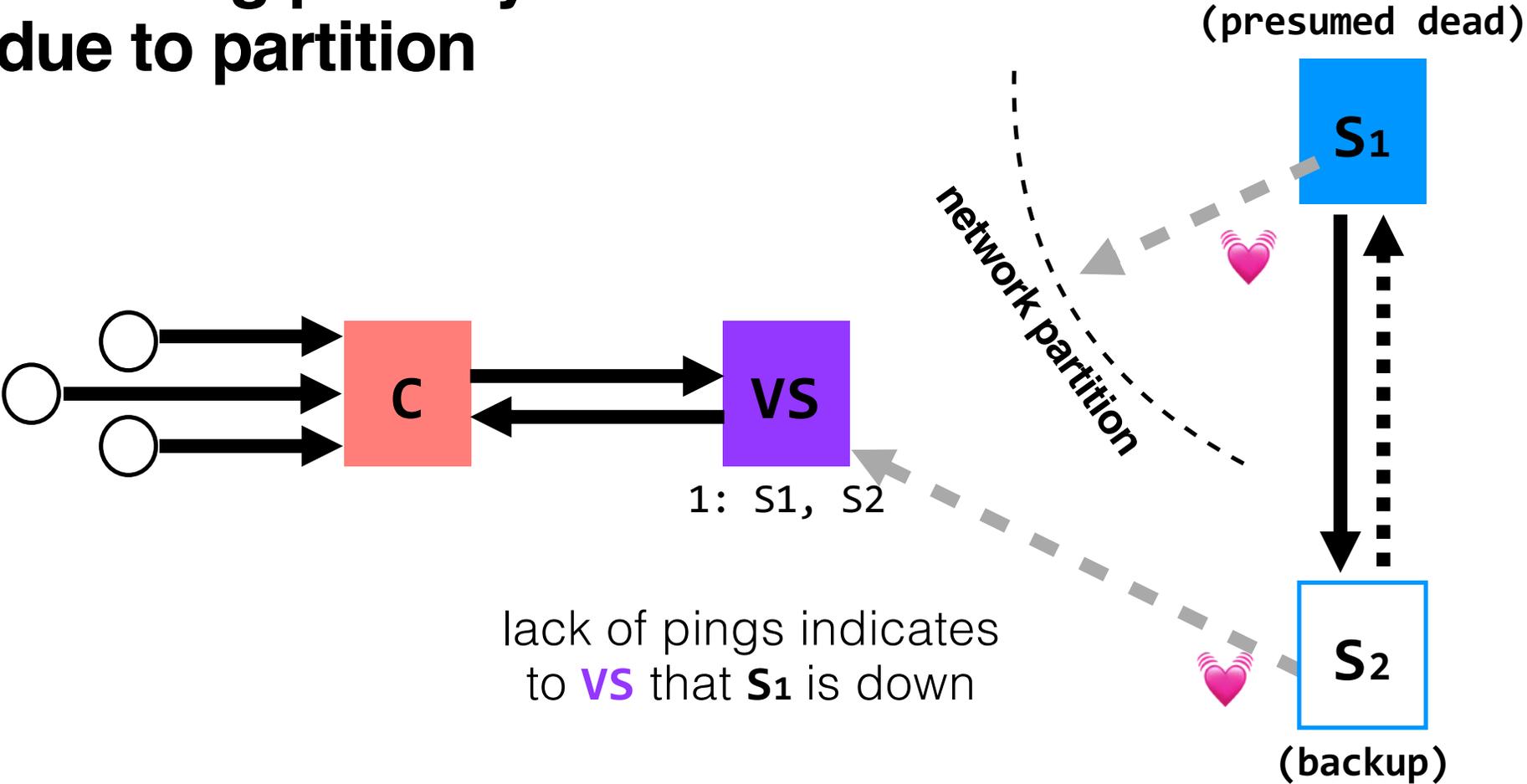


handling primary failure due to partition

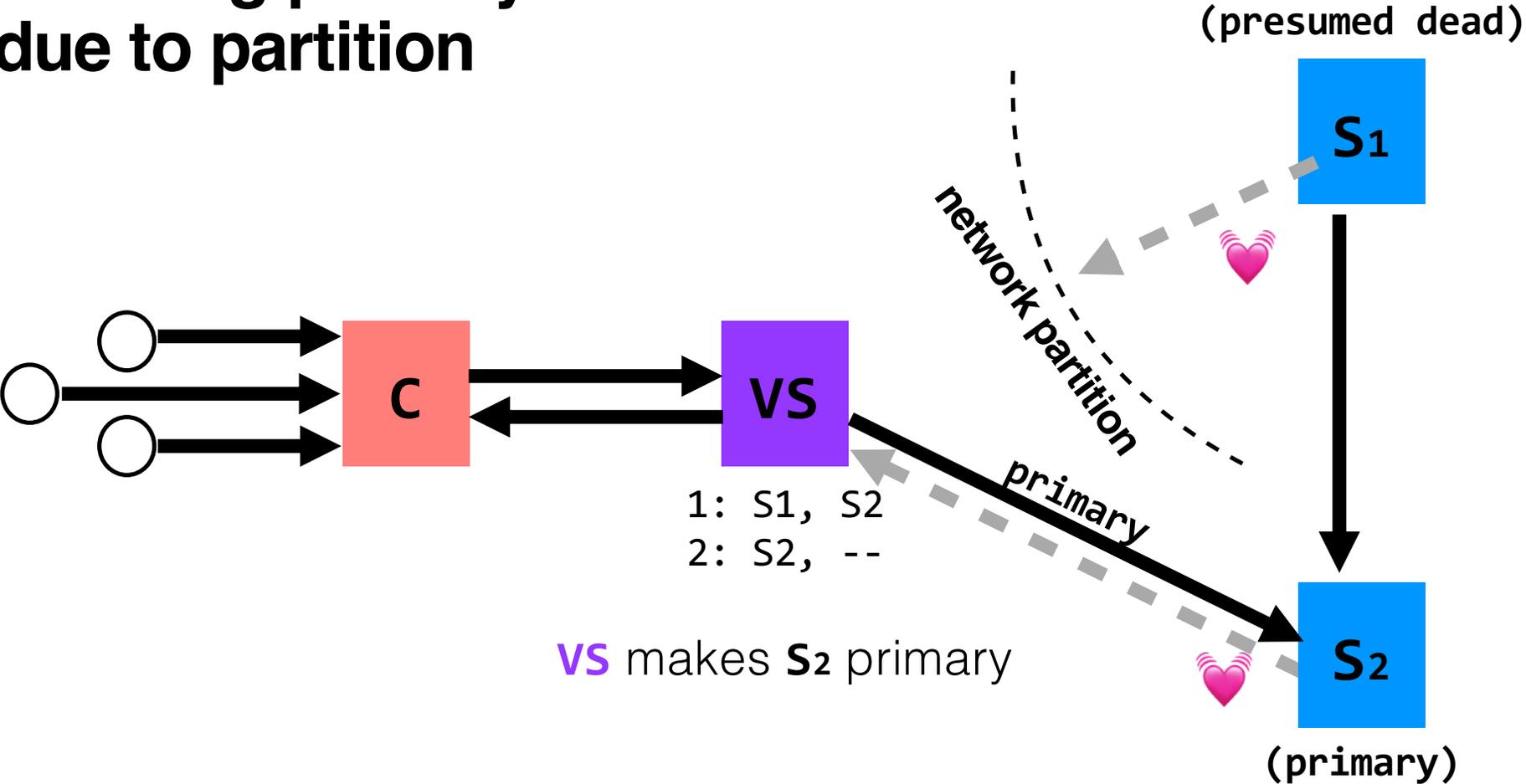


in the presence of a partition keeps **S₁** from communicating with the **view server**

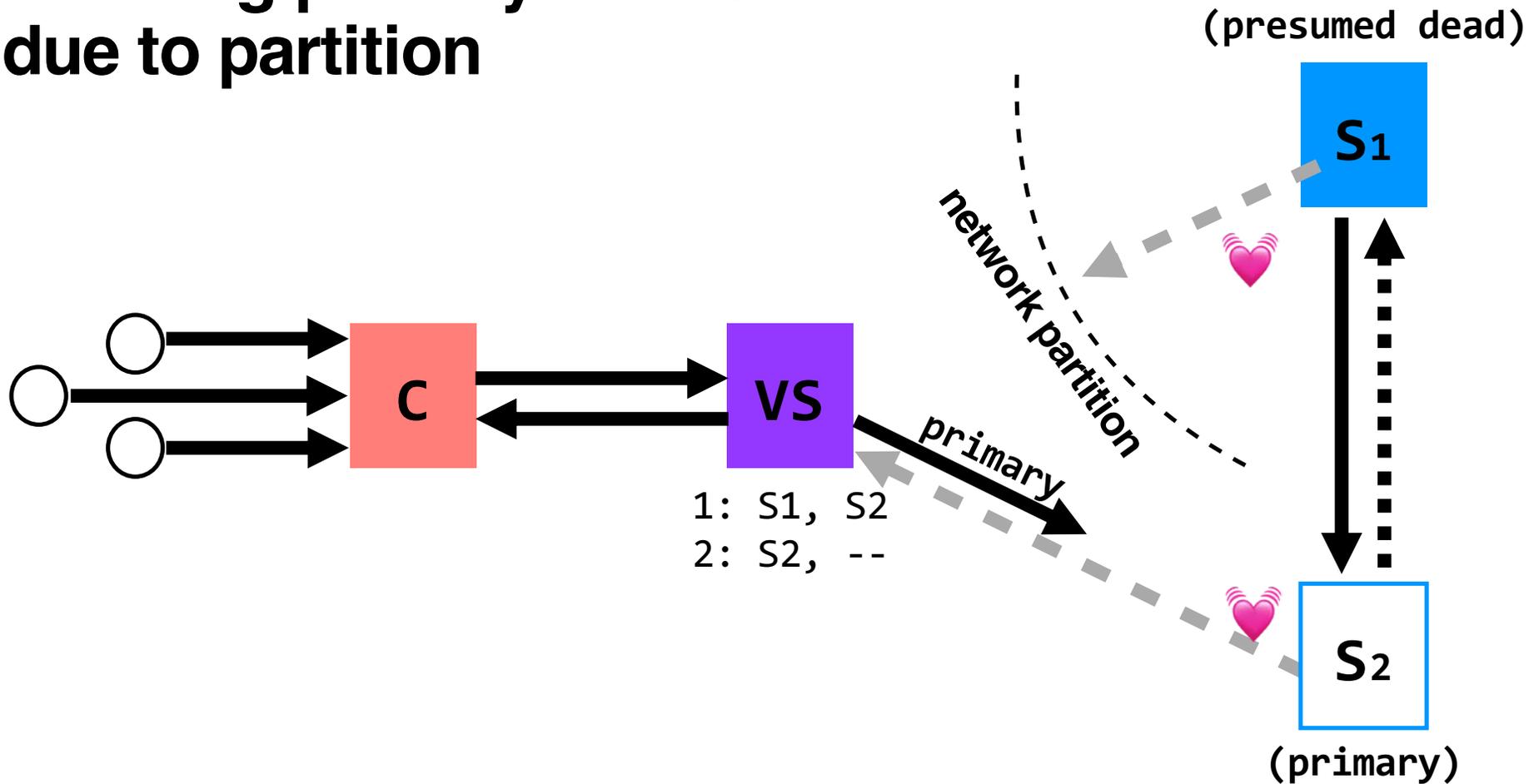
handling primary failure due to partition



handling primary failure due to partition

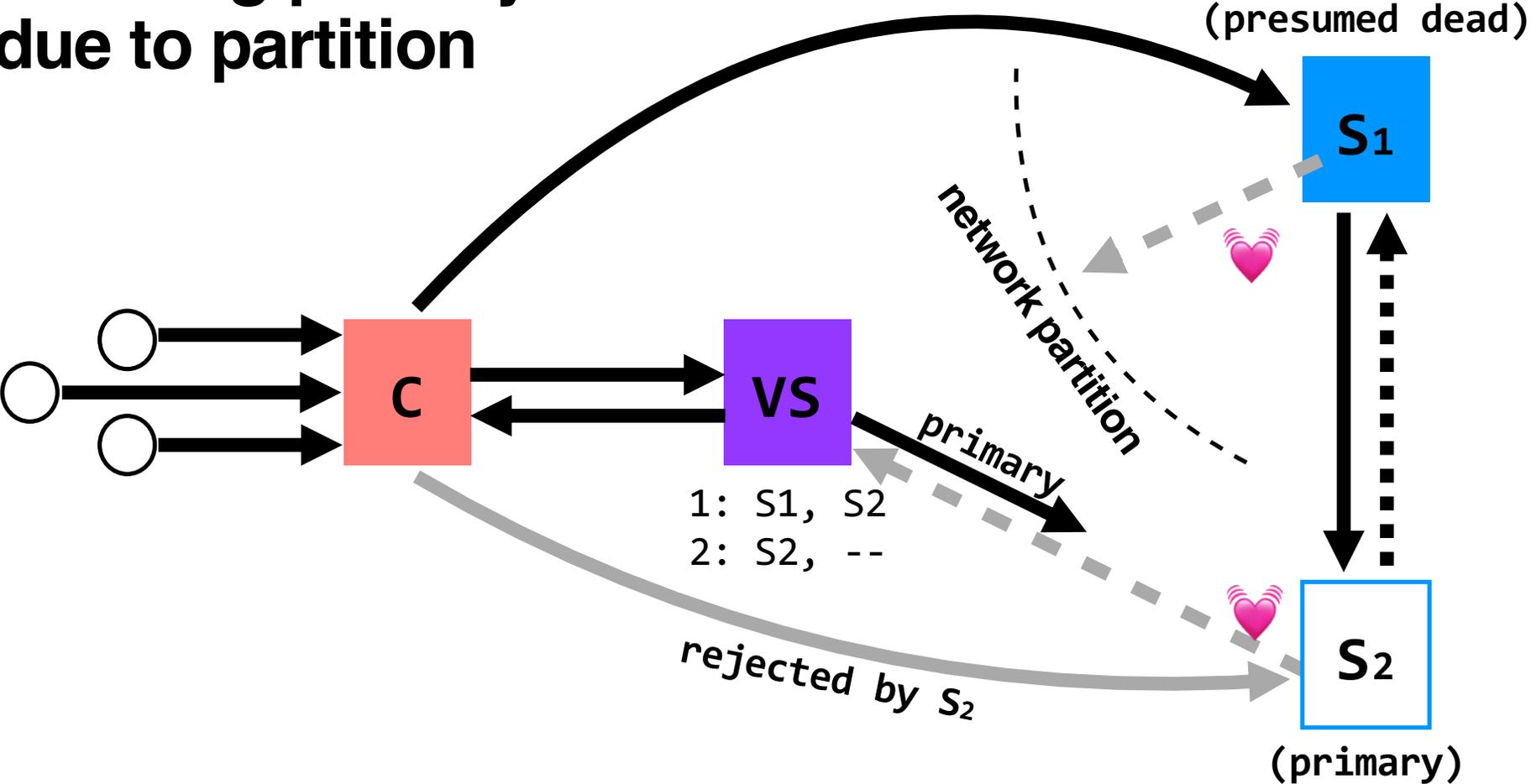


handling primary failure due to partition



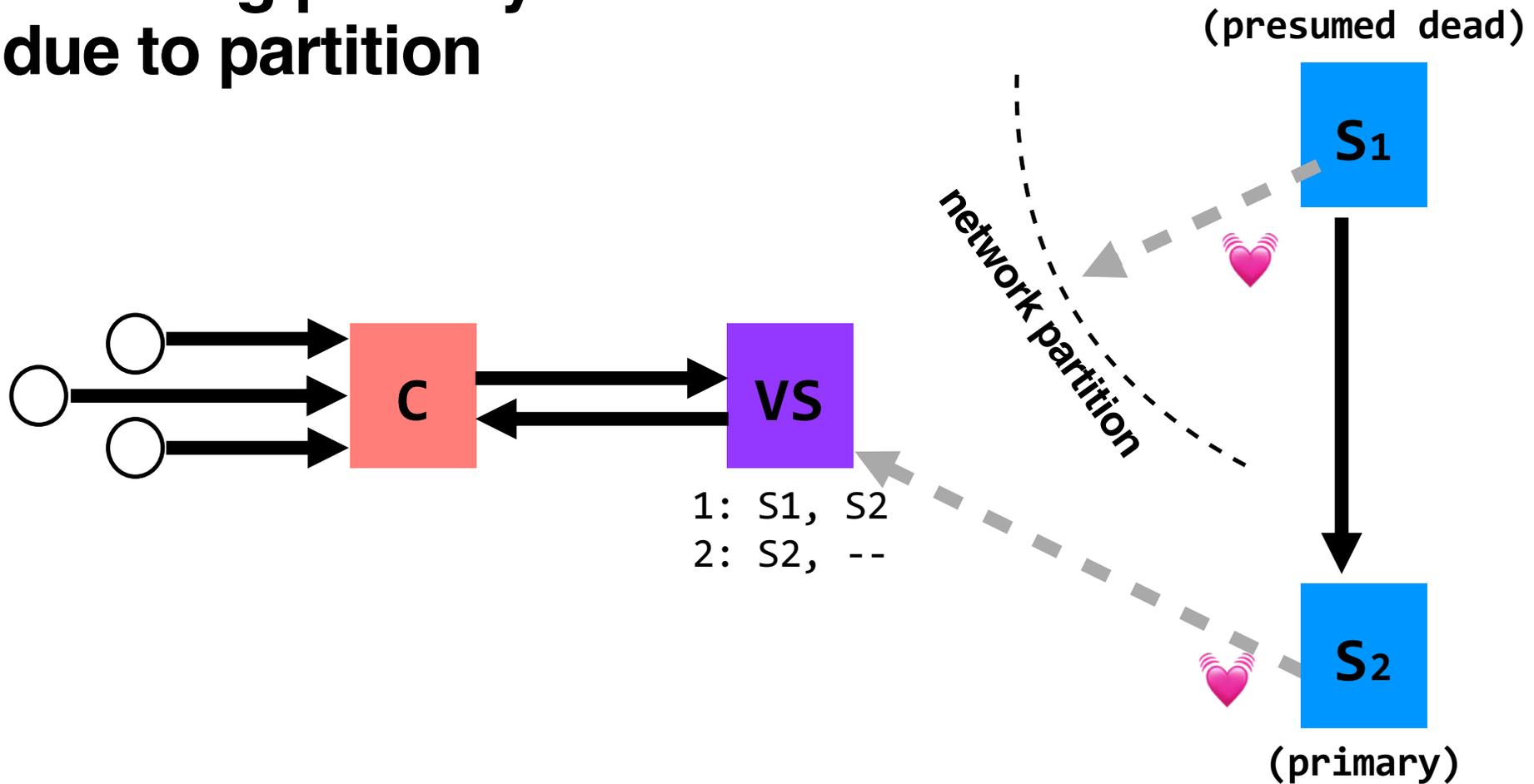
question: what happens before S_2 knows it's the primary?

handling primary failure due to partition



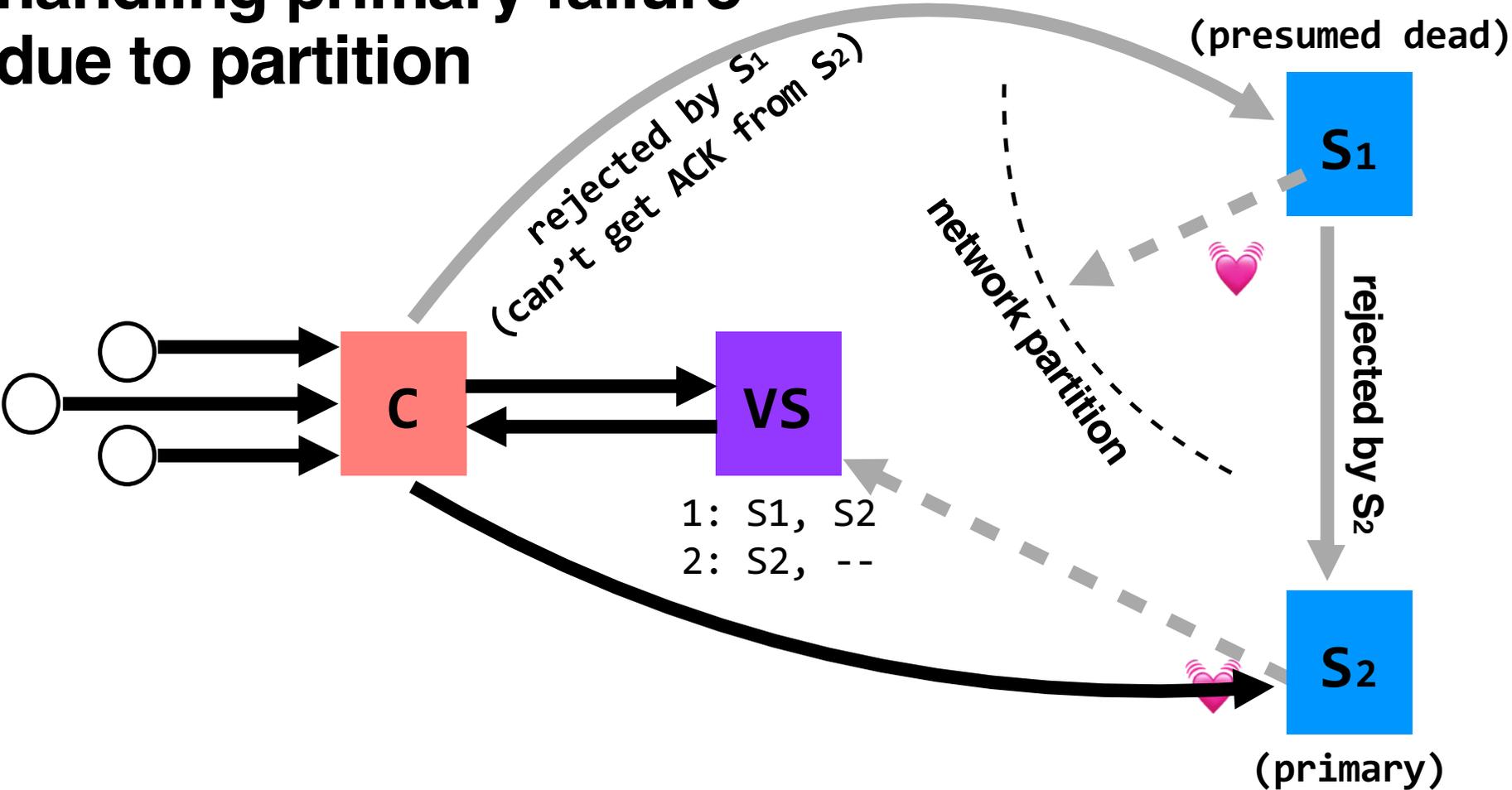
S₂ will act as backup
(accept updates from S₁, reject coordinator requests)

handling primary failure due to partition

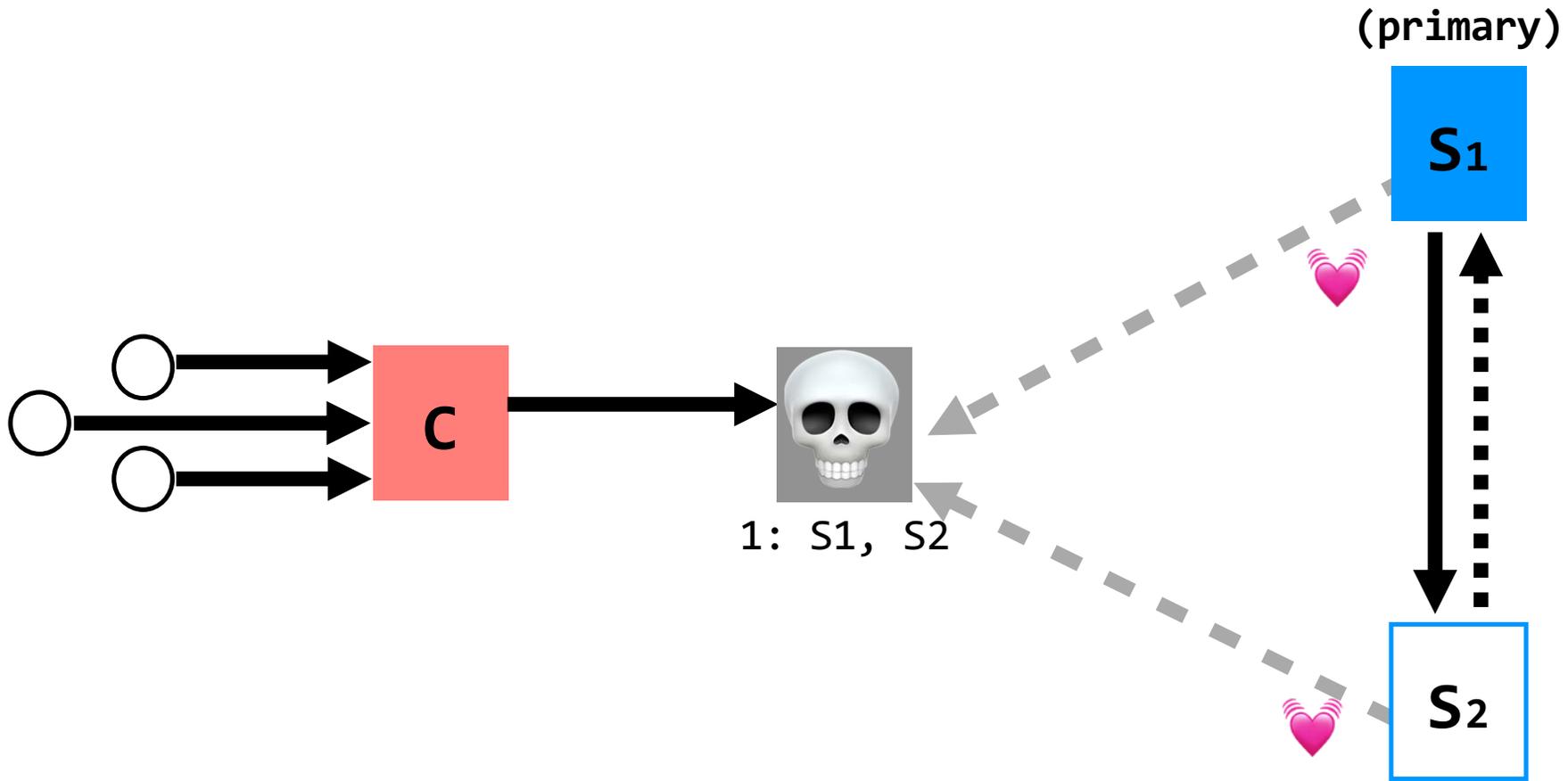


question: what happens after S₂ knows it's the primary, but S₁ also thinks it is?

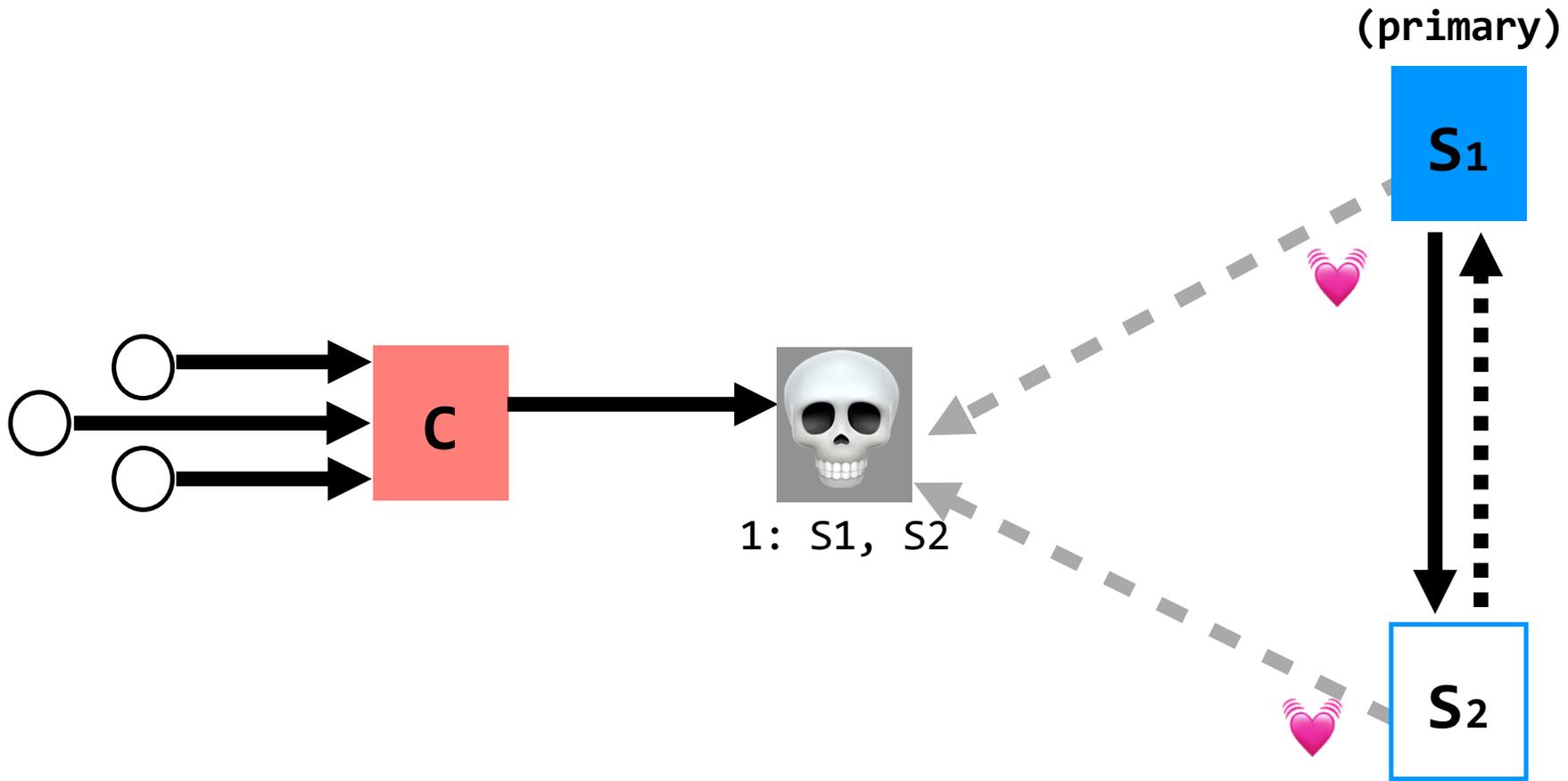
handling primary failure due to partition



S1 won't be able to act as primary
(can't accept client requests because it won't get ACKs from S2)



problem: what if view server fails?



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).

MIT OpenCourseWare
<https://ocw.mit.edu>

6.033 Computer System Engineering
Spring 2018

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.