# 6.033 Spring 2019
## 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** $\longrightarrow$ **shadow copies** (simple, poor performance) or **logs** (better performance, a bit more complex)

**isolation** $\longrightarrow$ **two-phase locking**

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

$C_1$ write$_1$(X)

$S_1$

$C_2$ write$_2$(X)

$S_2$
(replica of $S_1$)

| | |
|---|---|
| $C_1$ | $S_1$   $write_1(X)$ <br> $write_2(X)$ |
| $C_2$ | $S_2$   $write_2(X)$ <br> $write_1(X)$ <br> (replica of $S_1$) |

**problem:** replica servers can become inconsistent

(primary)

**S₁**

primary chooses order
of operations, decides
all non-deterministic
values

**C**

(backup)

**S₂**

primary ACKs coordinator
only after it's sure that
backup has all updates

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

if primary fails, **C** switches to backup **(primary)**

(**C** knows how to contact backup servers)

**primary chooses order of operations, decides all non-deterministic values**

S₁

C

**primary ACKs coordinator only after it's sure that backup has all updates**
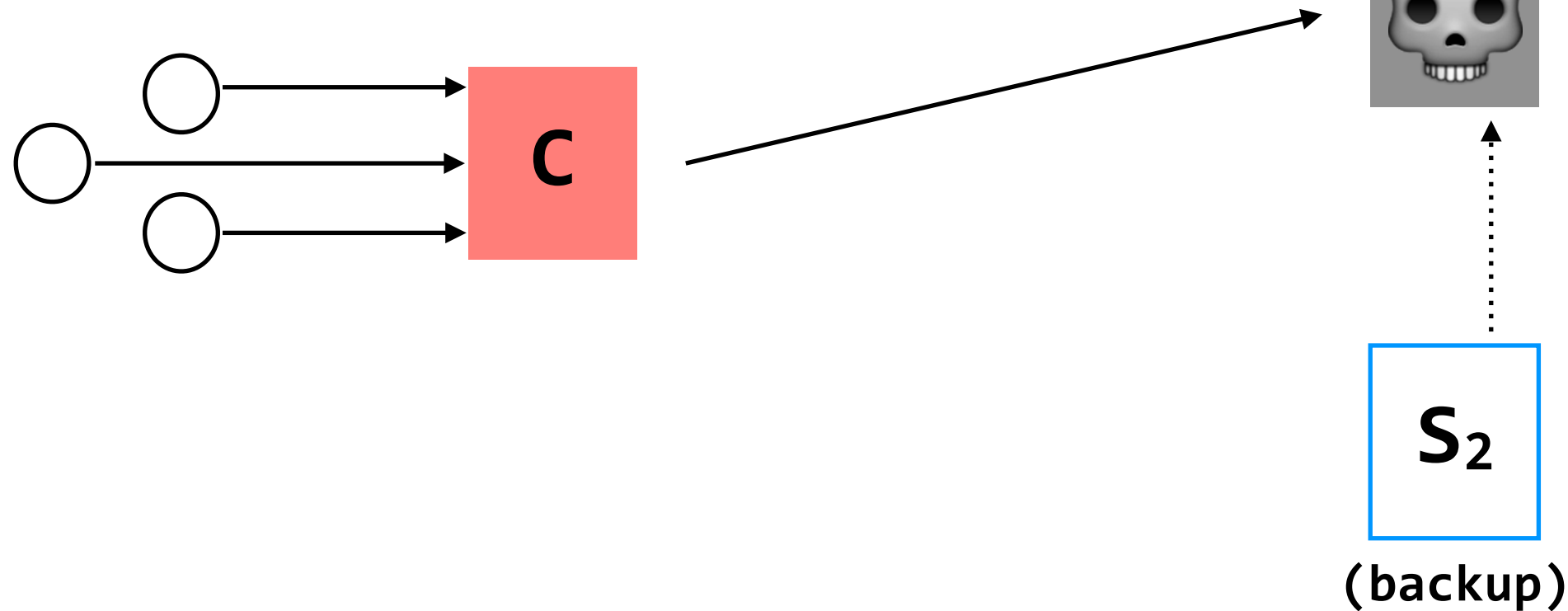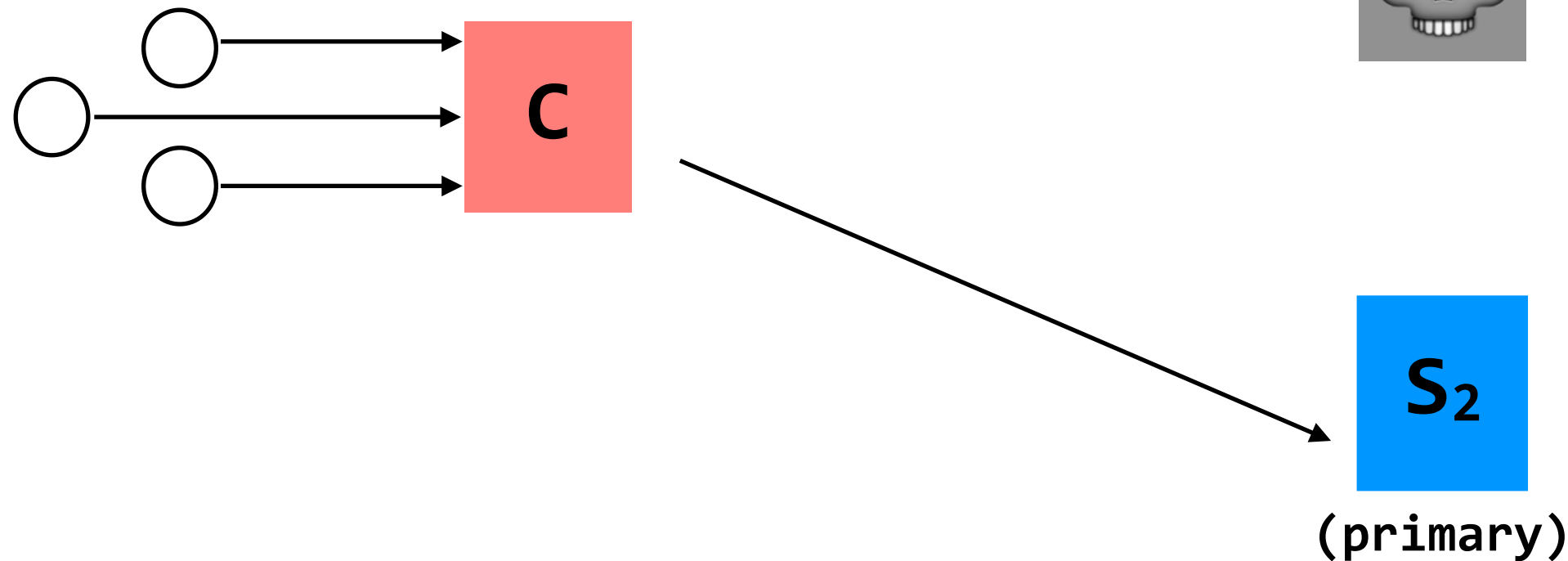
S₂

**(backup)**

**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)

**(dead)**



**C**

**S₂**

**(backup)**

**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)

**(dead)**



**C**

**S$_2$**

**(primary)**

**attempt:** coordinators communicate with primary servers, who communicate with backup servers
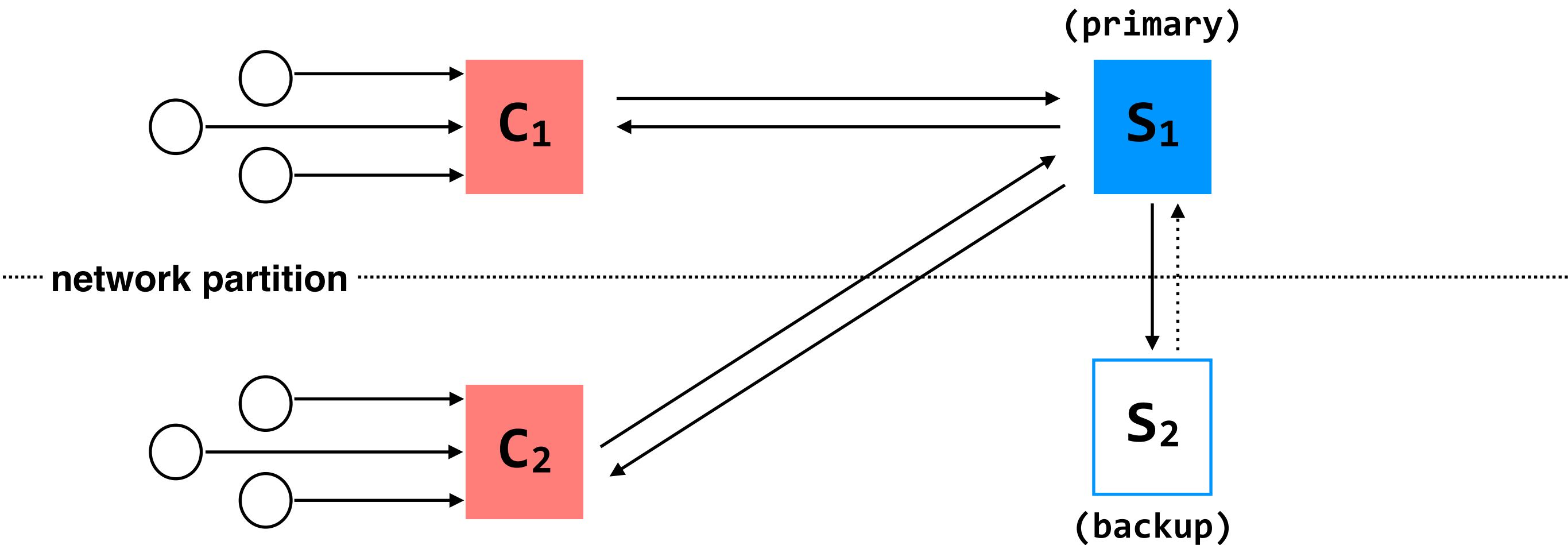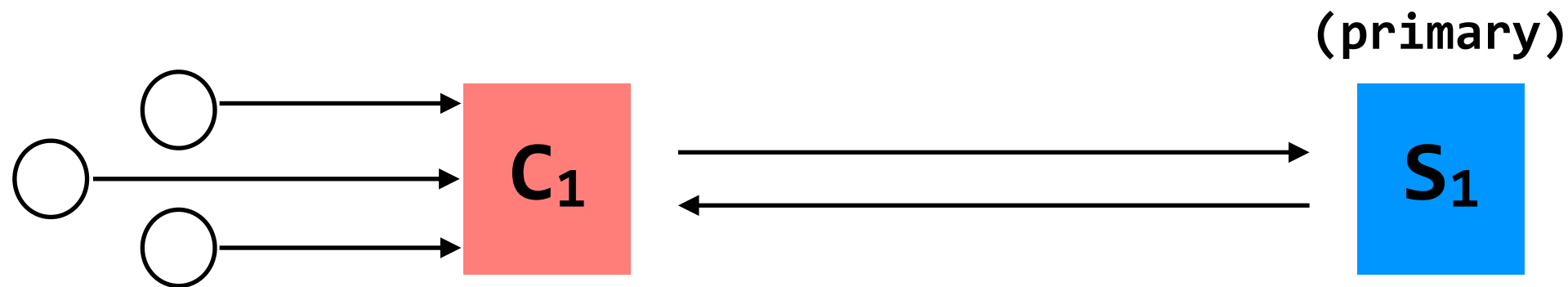
# multiple coordinators + the network = problems



**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

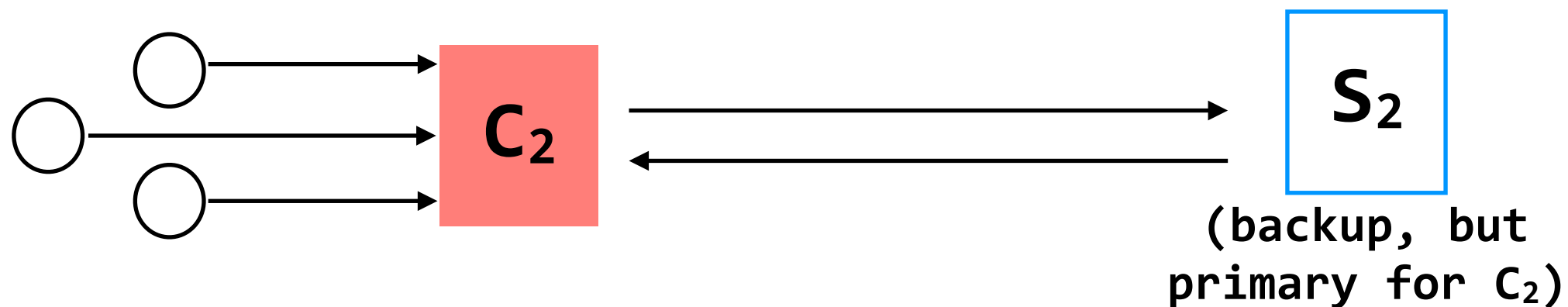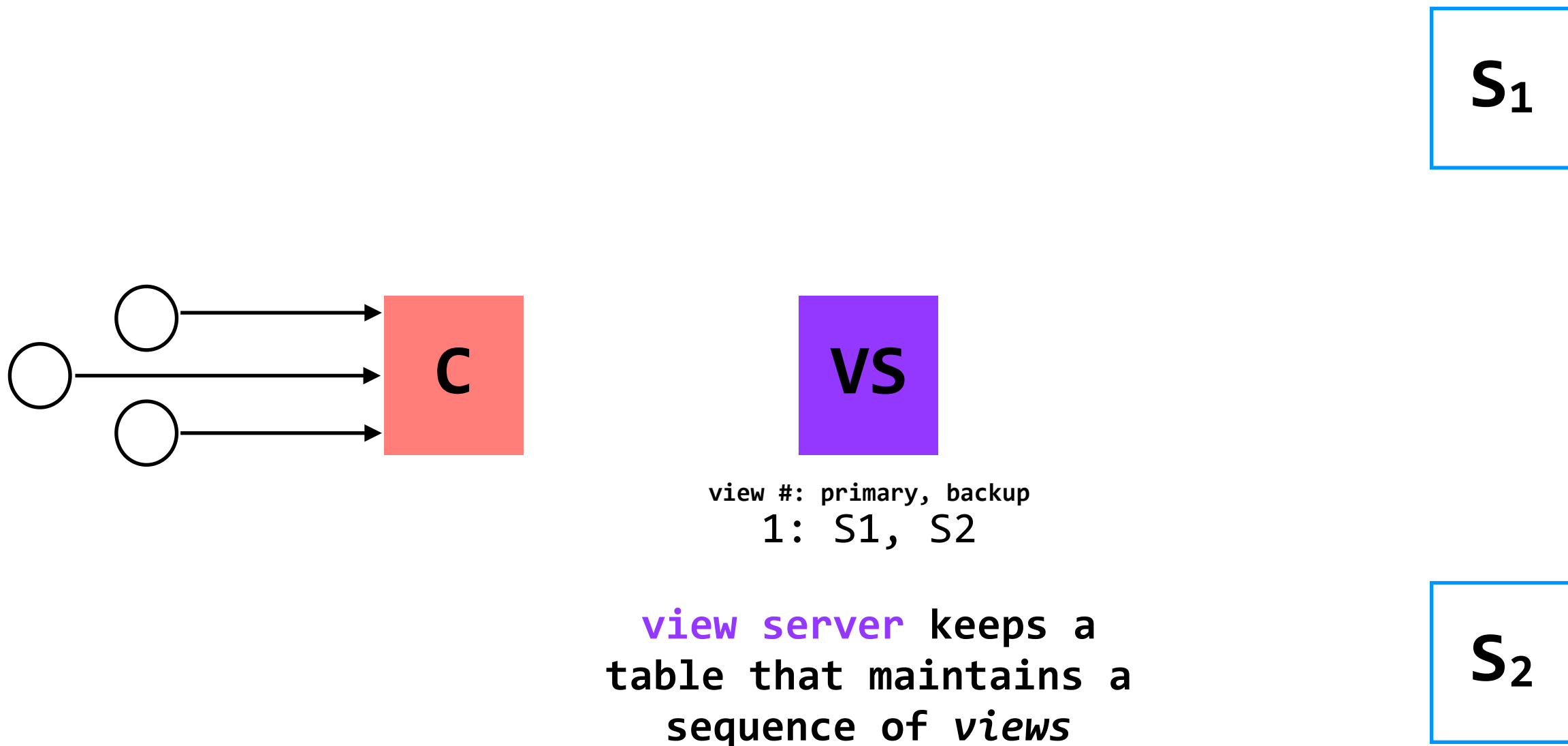# multiple coordinators + the network = problems



**(primary)**

**C₁** → **S₁**

**network partition**

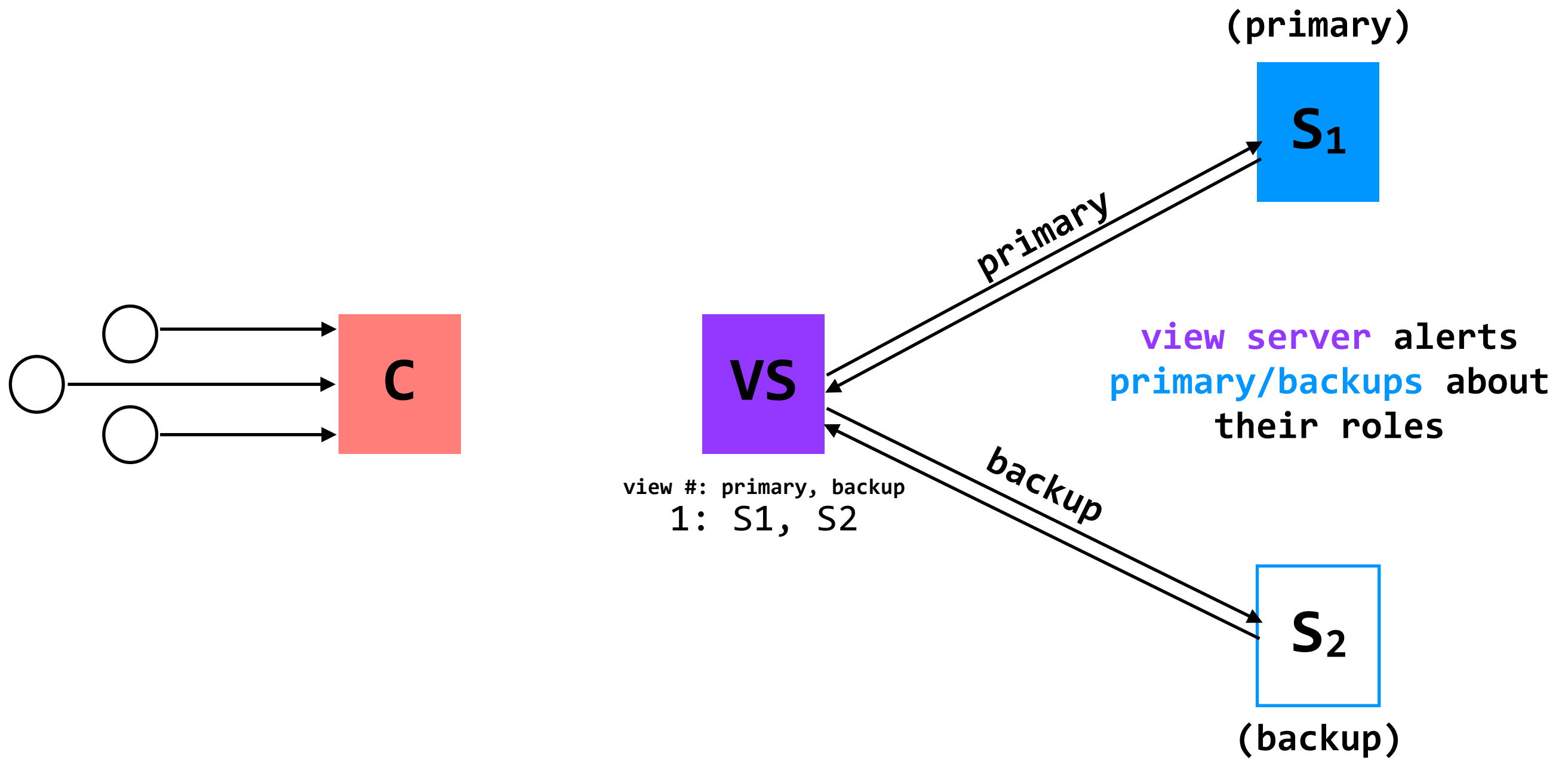**C₂** → **S₂**

**(backup, but primary for C₂)**

$C_1$ 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
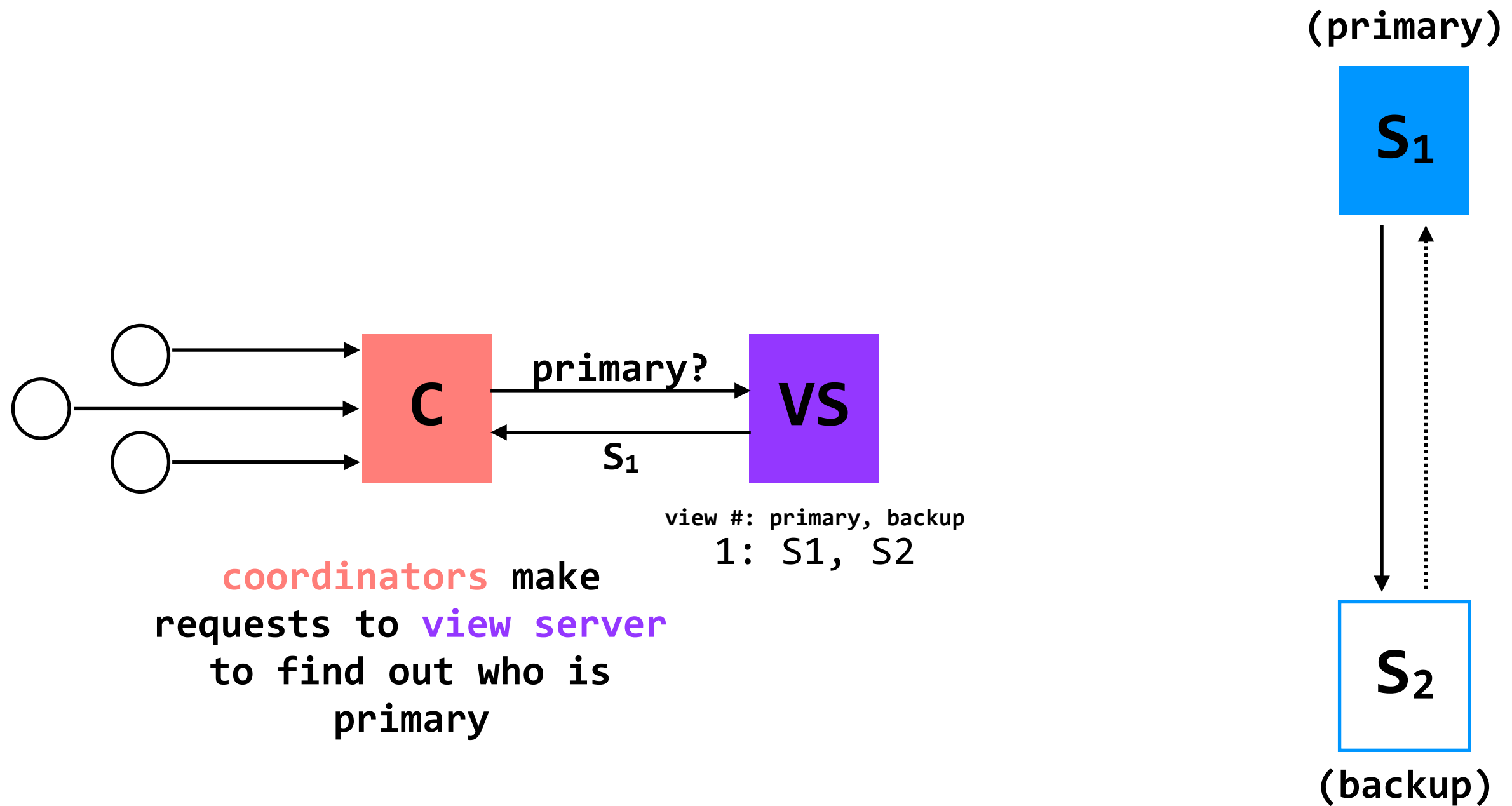
**S₁**

**C**

**VS**

```
view #: primary, backup
    1: S1, S2
```

**view server** keeps a
table that maintains a
sequence of *views*

**S₂**

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

(primary)

**S₁**

primary

**VS**

**view server** alerts
**primary/backups** about
their roles

view #: primary, backup
1: S1, S2

backup

**S₂**

(backup)

**C**

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

**primary** sends updates to, gets ACKs from backup (as before)

(primary)

**S₁**

primary

**C**

**VS**

view #: primary, backup
1: S1, S2

backup

**S₂**

(backup)

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

(primary)

**S₁**

(backup)

**S₂**

**C**

primary? → **VS**

← S₁

view #: primary, backup
1: S1, S2

coordinators make
requests to view server
to find out who is
primary

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

coordinators contact
primary (as before)

(primary)

S₁

C primary? VS

S₁

view #: primary, backup
1: S1, S2

S₂

(backup)

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

(primary)

**S₁**

**C**  primary? → **VS**
← S₁

view #: primary, backup
1: S1, S2

**primary/backup(s)** ping
**view server** so that it
can discover failures

**S₂**

(backup)

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

# handling primary failure

(dead)

C

VS

view #: primary, backup
1: S1, S2

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

S₂

(backup)

# handling primary failure

(dead)

**C**

**VS**

view #: primary, backup
1: S1, S2
2: S2, --

primary

**S₂**

(primary)

# handling primary failure

(dead)

C —— primary? ——> VS

VS —— S₂ ——> C

view #: primary, backup
1: S1, S2
2: S2, --

S₂
(primary)

# handling primary failure

(dead)



**C**

**VS**

```
view #: primary, backup
1: S1, S2
2: S2, --
```

**S₂**

(primary)

# handling primary failure
# due to partition

(primary)

**S₁**

network partition

**C**

**VS**

view #: primary, backup
1: S1, S2

**S₂**

(backup)

suppose a partition keeps **S₁** from communicating with the
**view server**

# handling primary failure due to partition

(presumed dead)

**S₁**

network partition

**C**

**VS**

view #: primary, backup
1: S1, S2

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

**S₂**

(backup)

# handling primary failure due to partition



(presumed dead)
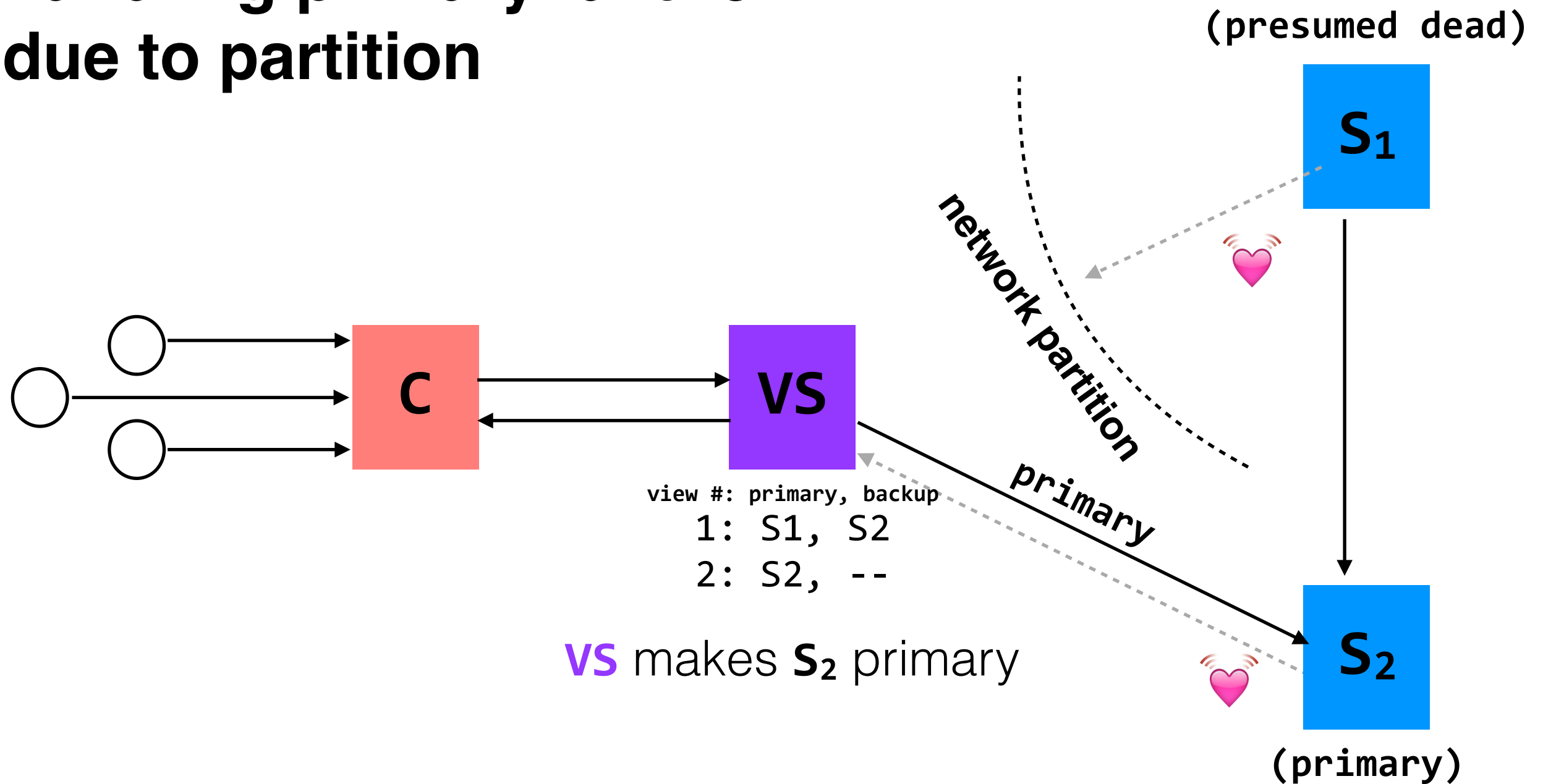
S₁

network partition

C

VS

view #: primary, backup
1: S1, S2
2: S2, --

VS makes S₂ primary
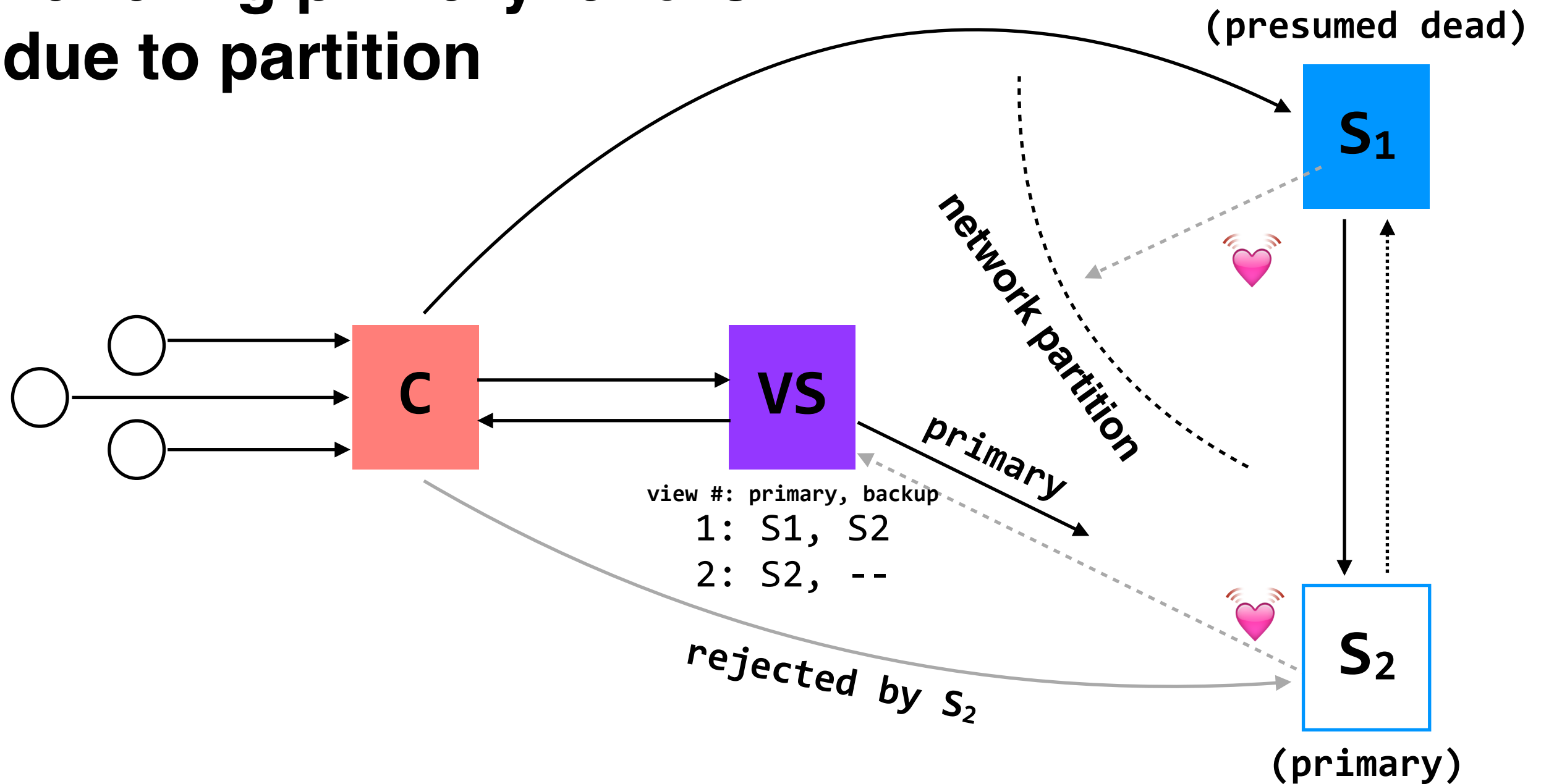
primary

S₂

(primary)

# handling primary failure
# due to partition



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

# handling primary failure due to partition



(presumed dead)

**S₁**

network partition

(primary)

**S₂**

**C**

**VS**

```
view #: primary, backup
1: S1, S2
2: S2, --
```

primary

rejected by S₂

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

# handling primary failure due to partition

(presumed dead)

**S₁**

network partition

**C**          **VS**
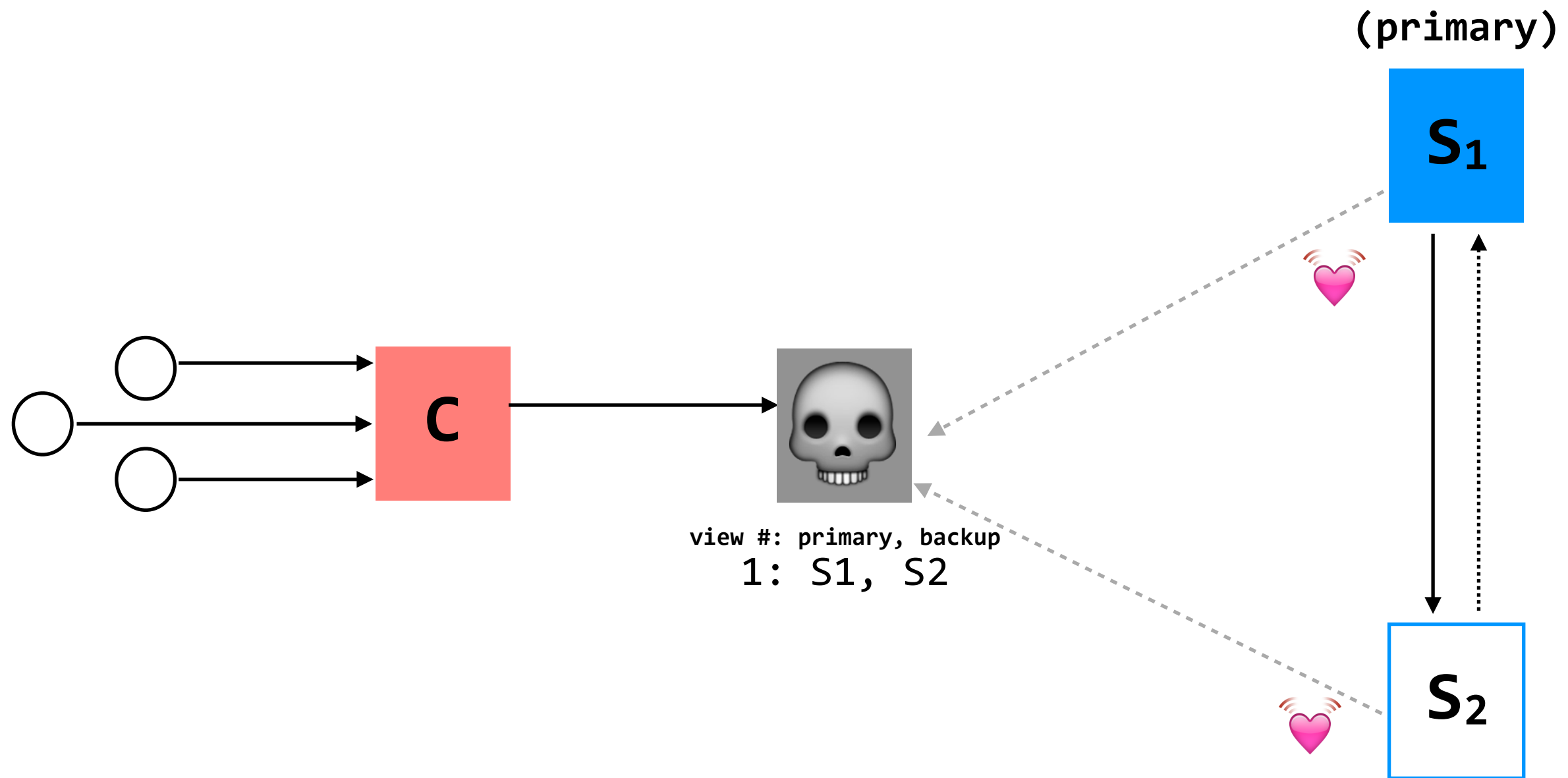
view #: primary, backup
    1: S1, S2
    2: S2, --

**S₂**

(primary)

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

# handling primary failure
# due to partition

rejected by $S_1$
(can't get ACK from $S_2$)

(presumed dead)

**S₁**

network partition

rejected by $S_2$

**C**

**VS**

```
view #: primary, backup
1: S1, S2
2: S2, --
```

**S₂**

(primary)

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

(primary)

**S₁**

**C**

view #: primary, backup
   1: S1, S2

**S₂**

**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).