



Joint work with Hernan Melgratti & Christian Roldan

On the semantics of replicated data types

Fabio Gadducci
University of Pisa

The quickest background, I

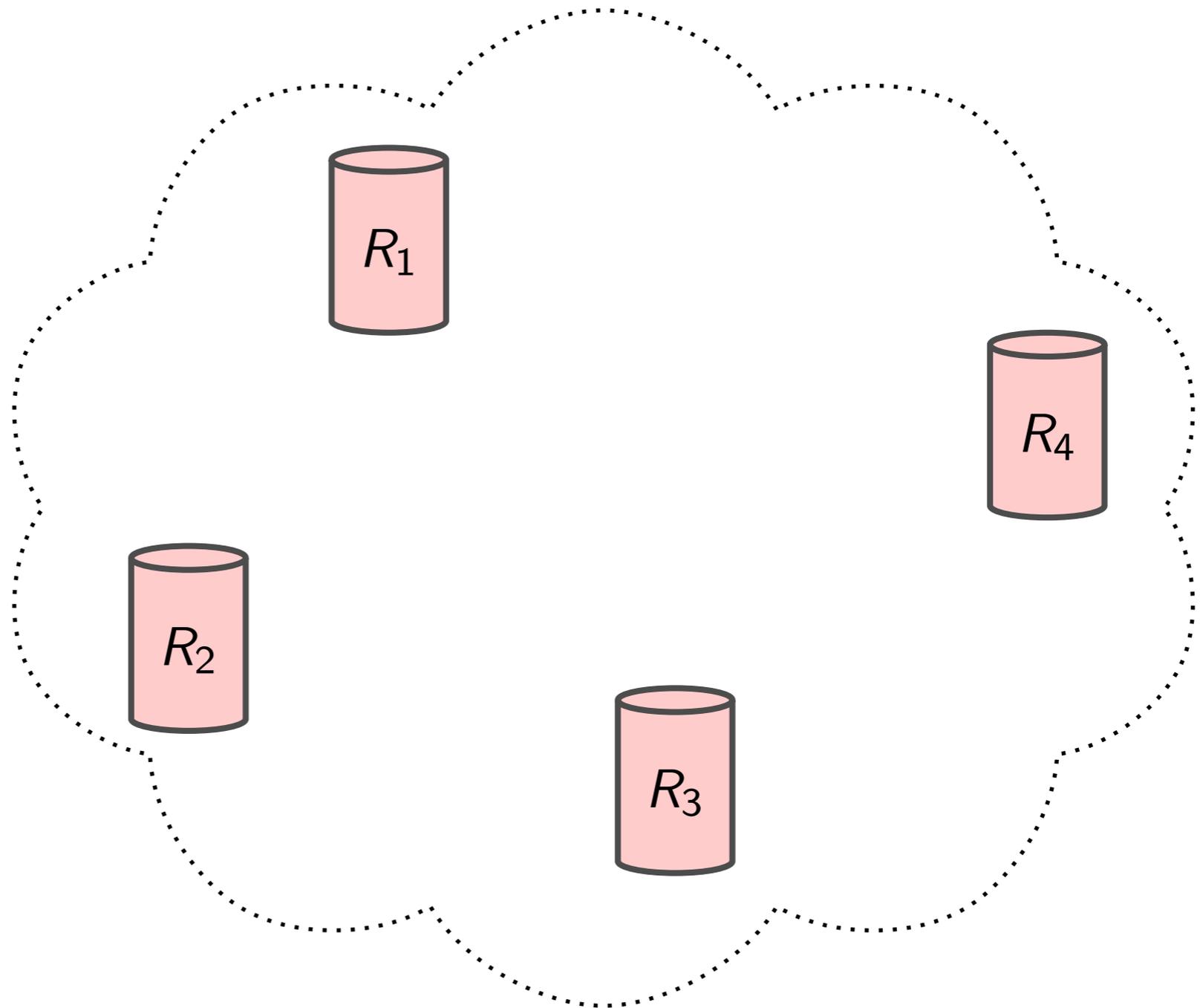
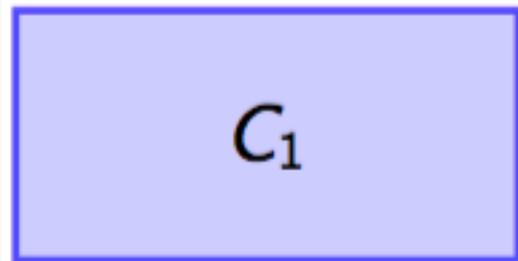
- ❖ Distributed systems replicate their state over different nodes in order to satisfy non-functional requirements.
- ❖ *Strong consistency* (every request receives the most recent update) of replicated data is in conflict with *availability* (every request is eventually executed) and tolerance to network *partitions* (the system operates even in the presence of failures that temporarily prevent communication among components).
- ❖ CAP theorem: it is impossible to simultaneously achieve strong Consistency, Availability and Partition tolerance [GL2002].

The quickest background, II

- ❖ *Weak* consistency: replicas may (temporarily) exhibit discrepancies (every request receives a correct update).
- ❖ How are the data specified? States, state transitions and returned values should account for the different views that a data item may simultaneously have.
- ❖ In the end, consistency has to be *eventually* guaranteed (if no new updates are made to a data item, eventually all accesses to that item will return the most recent update).

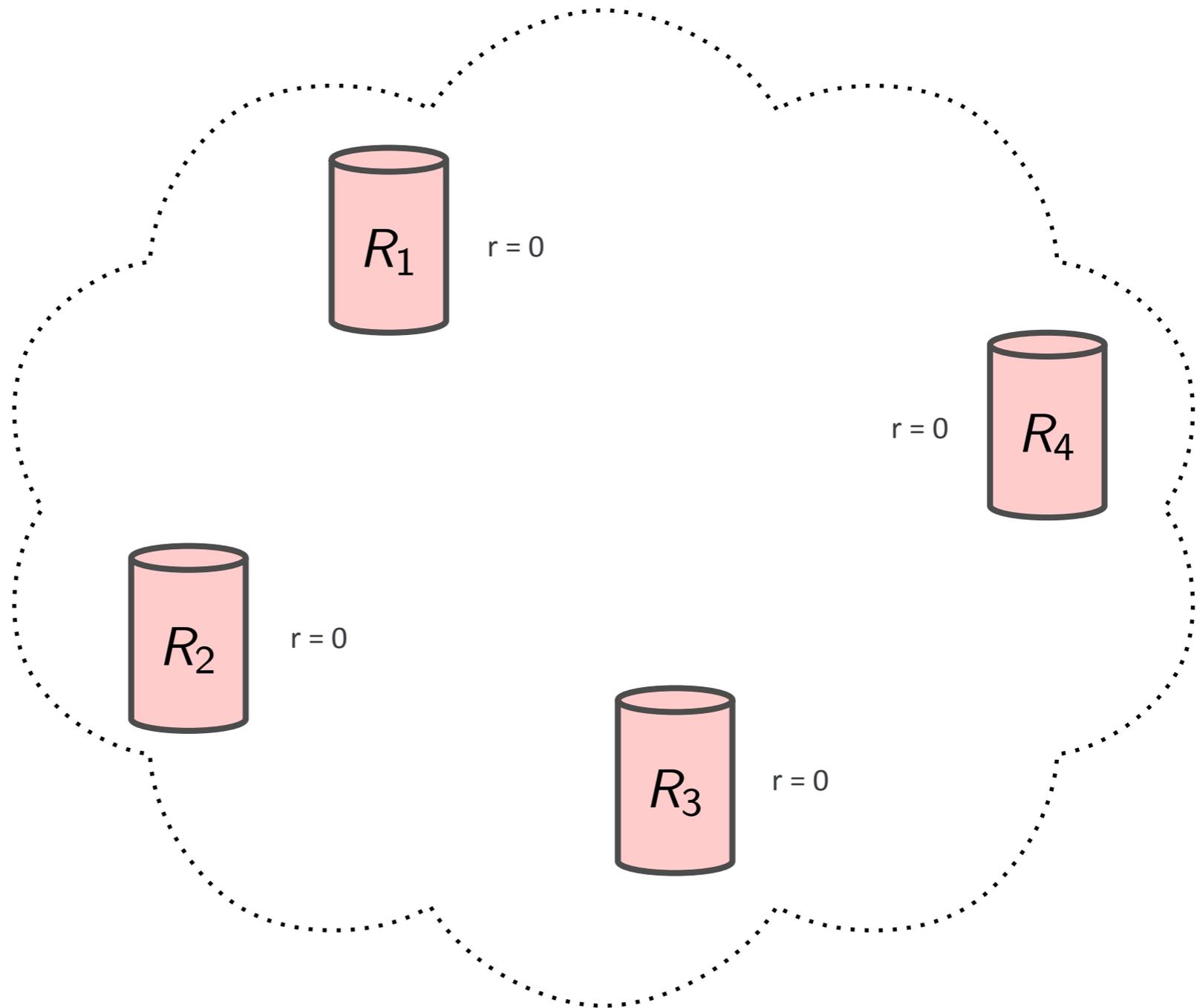
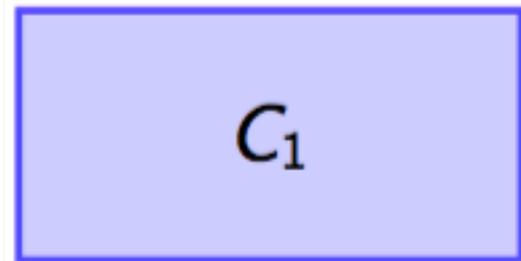
A register

register $r = 0$



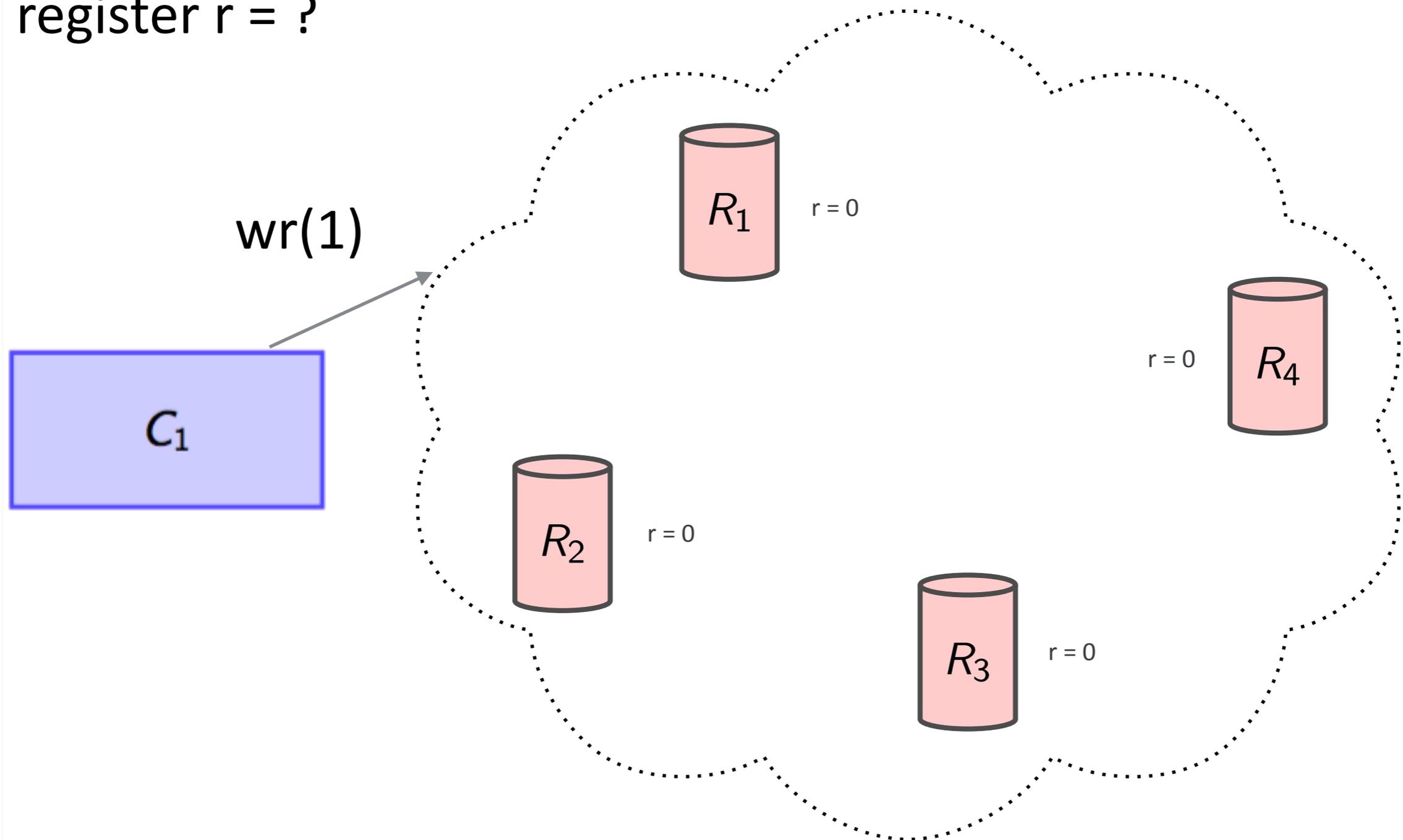
A register

register $r = 0$



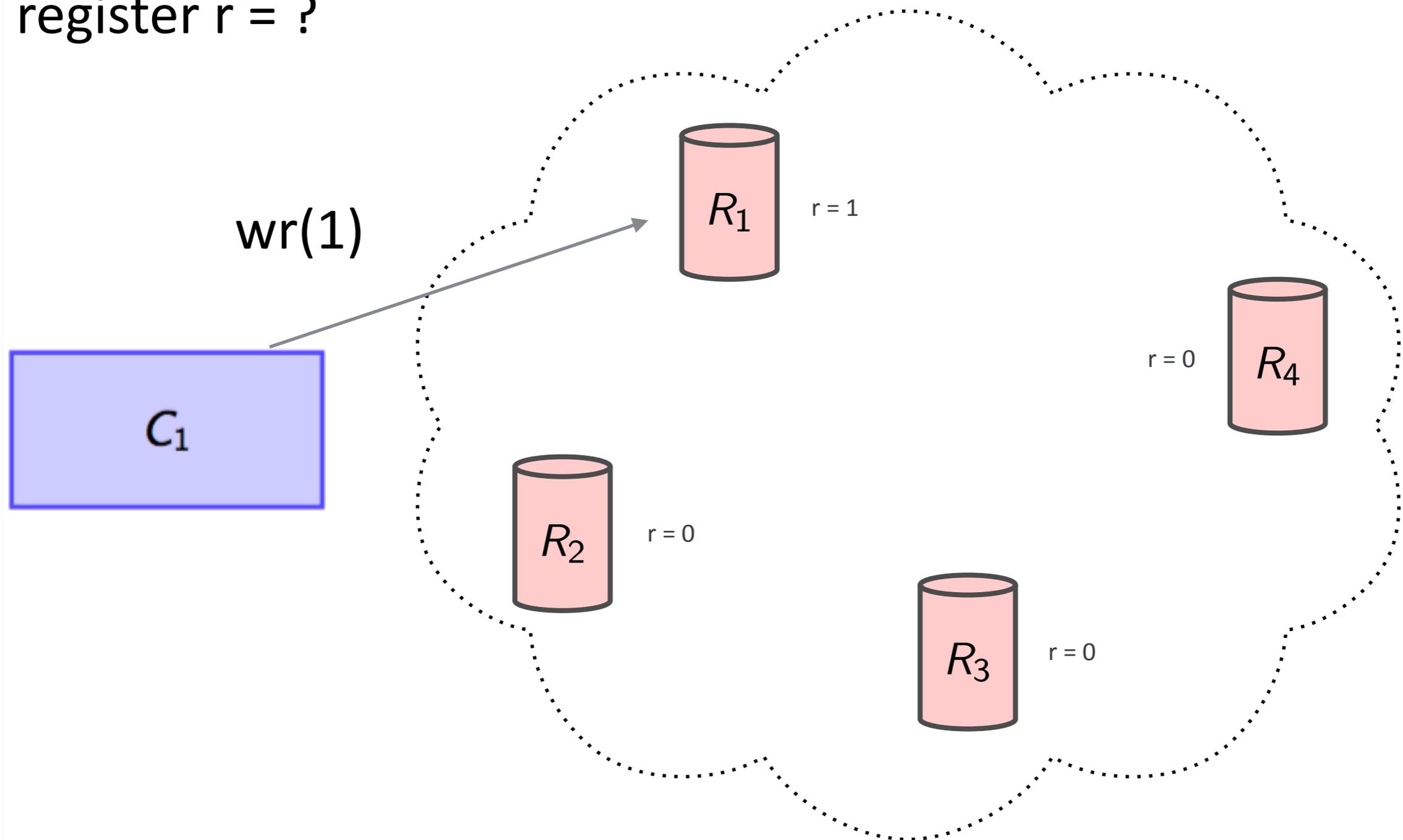
A register

register $r = ?$



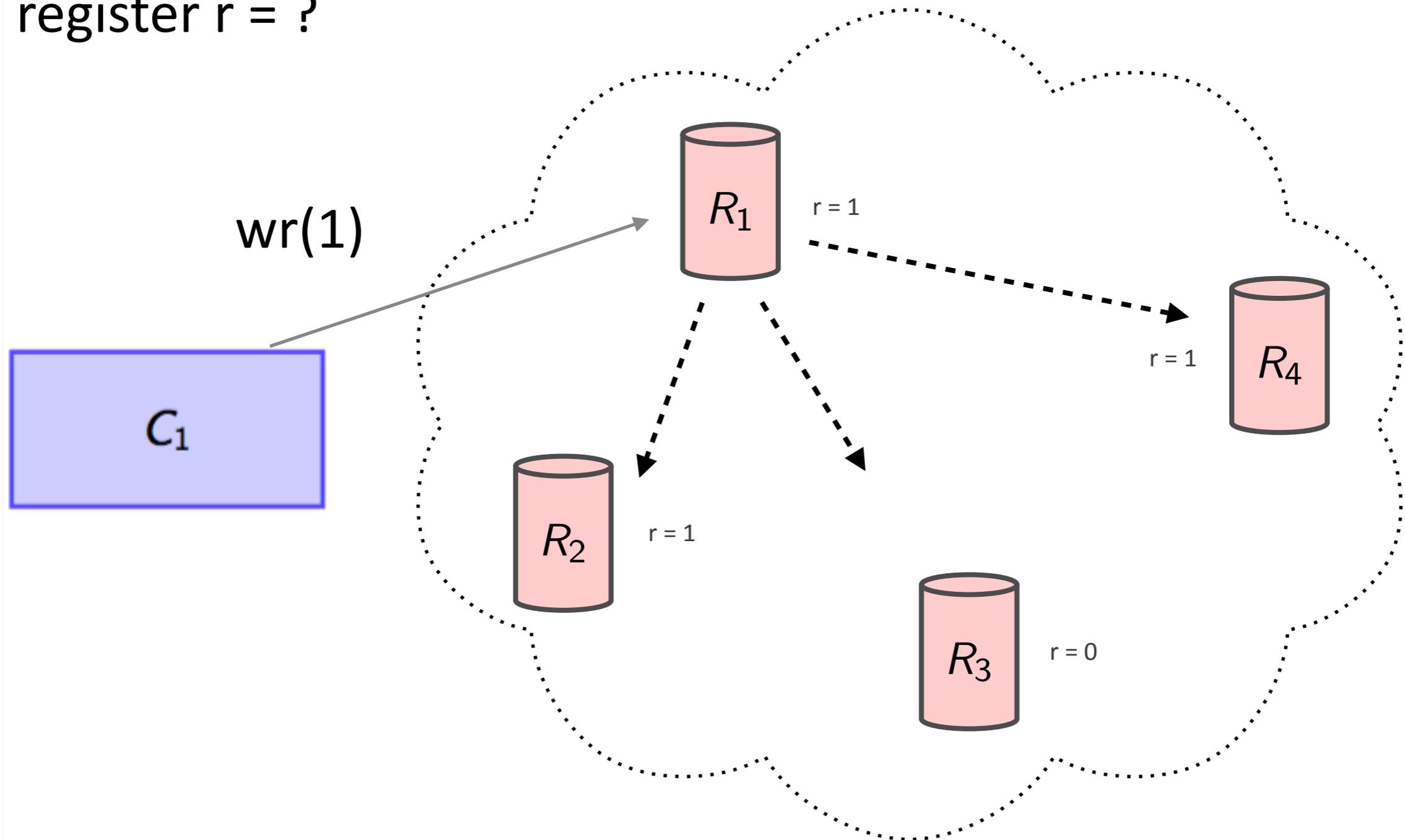
A register

register $r = ?$



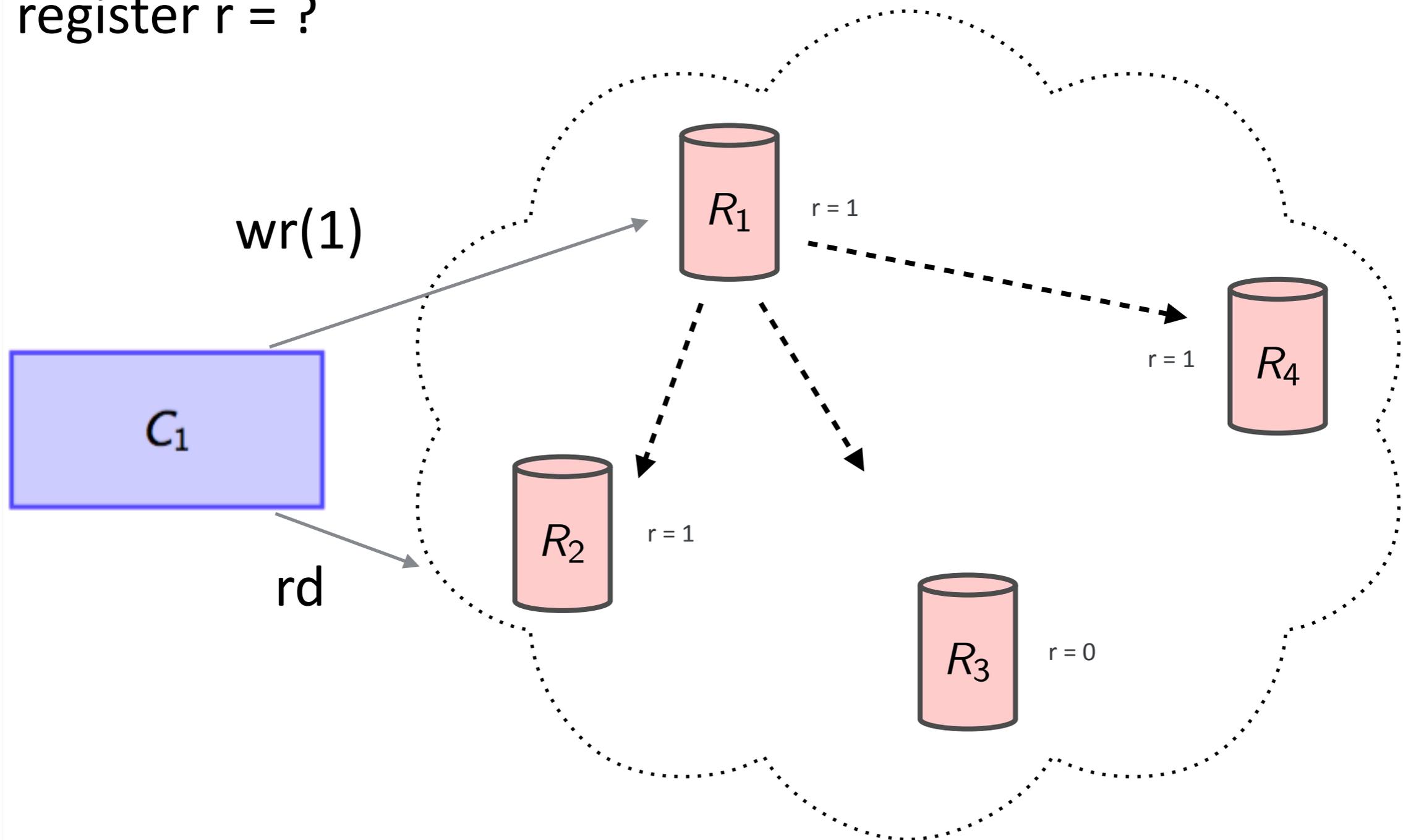
A register

register $r = ?$



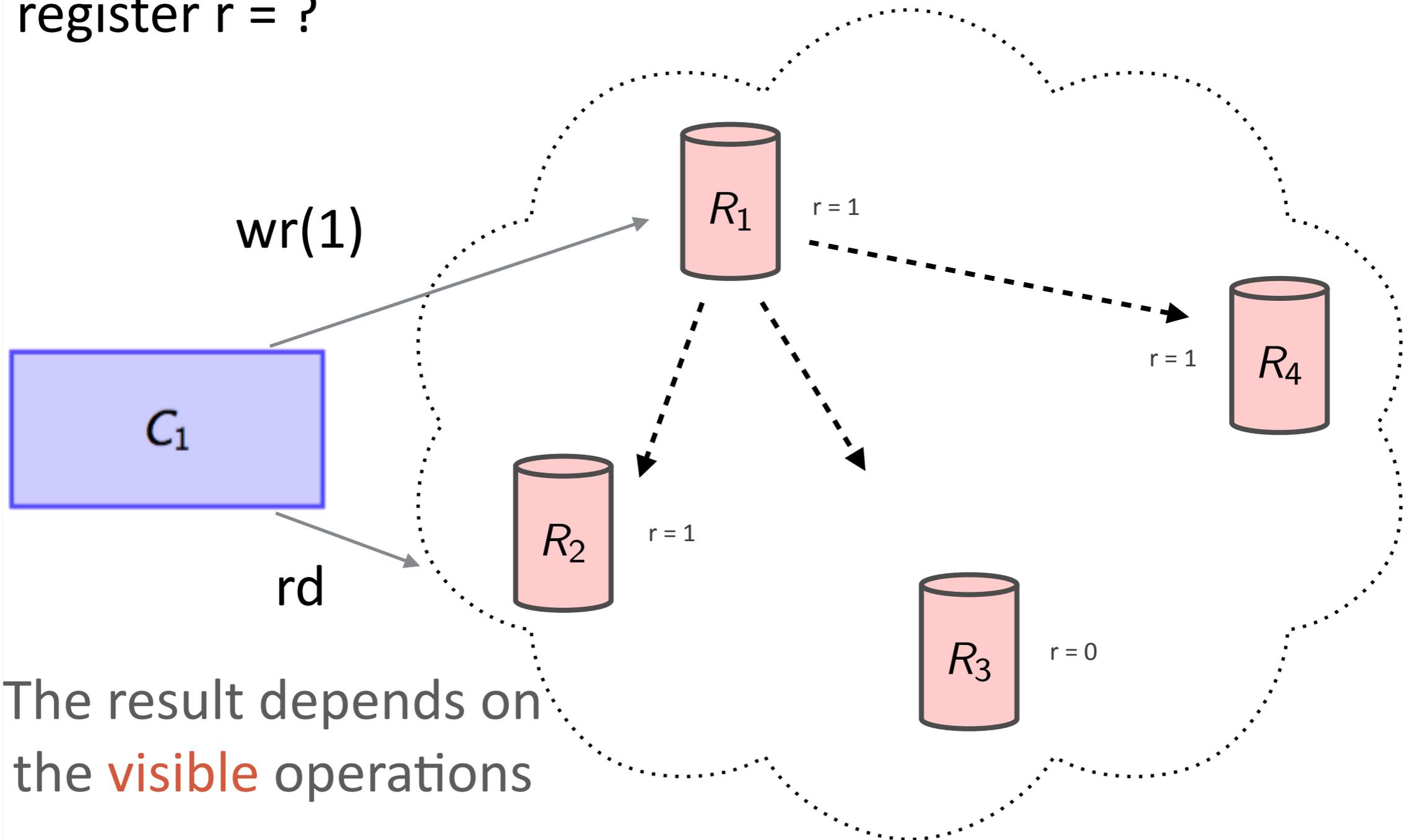
A register

register $r = ?$



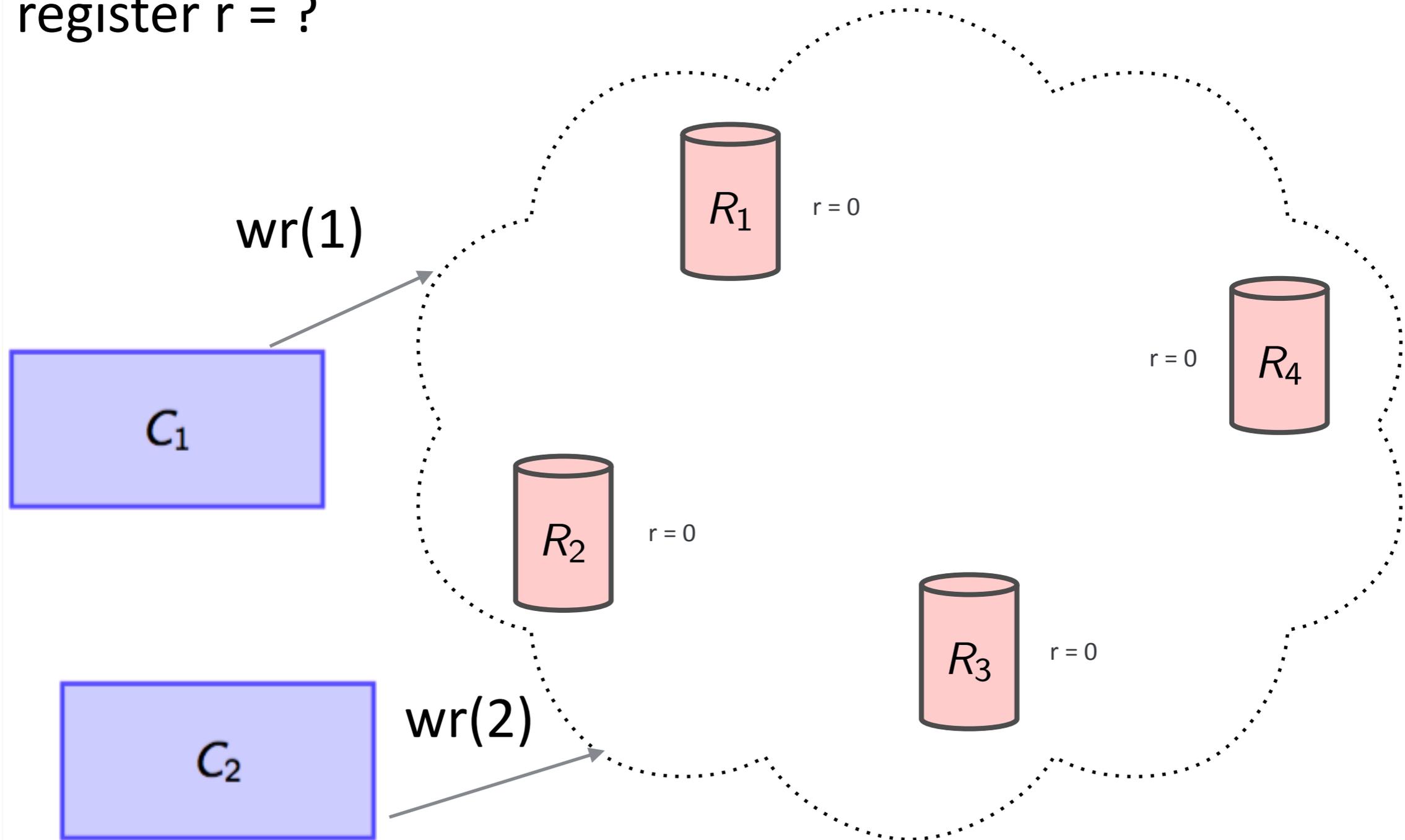
A register

register $r = ?$



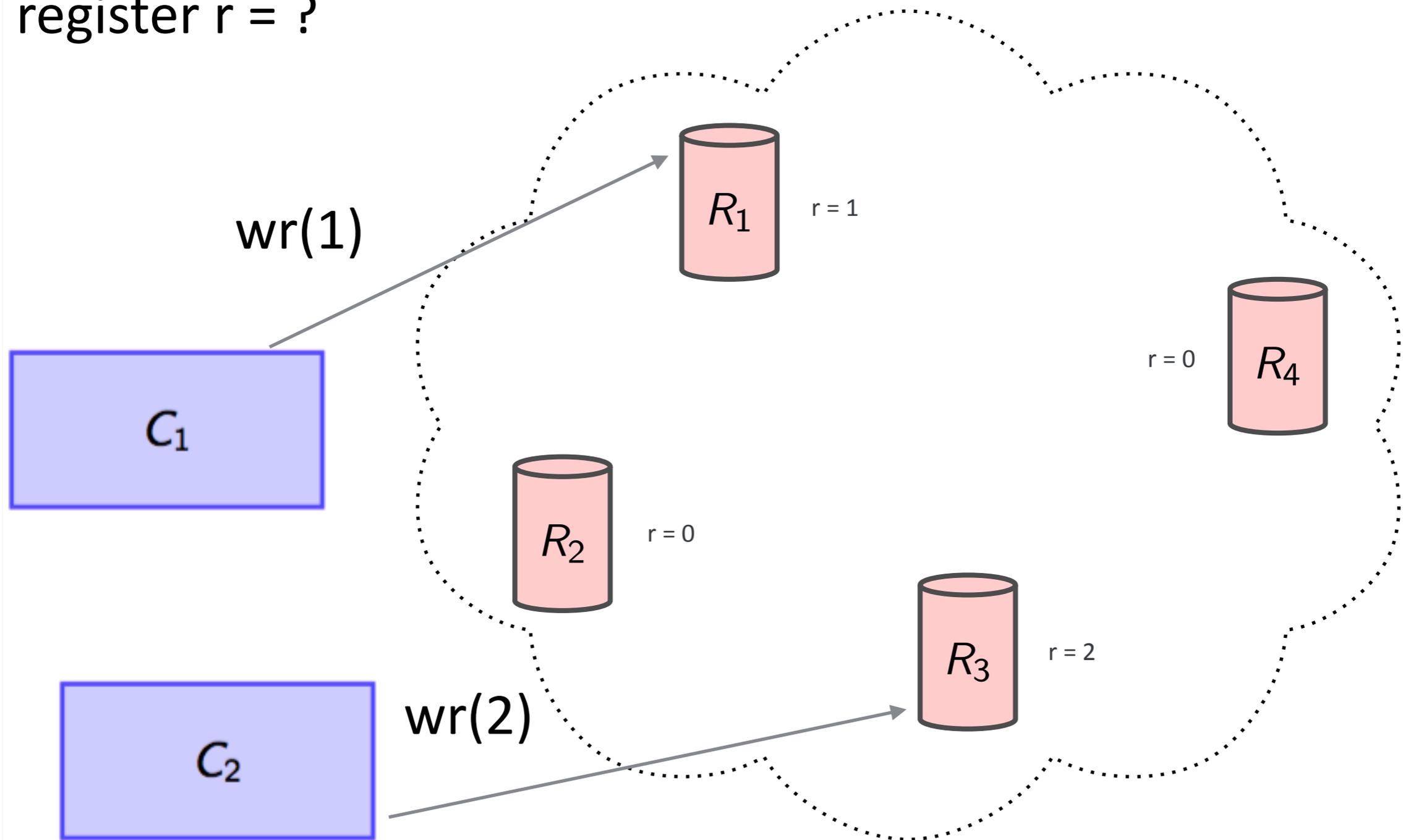
A register

register $r = ?$



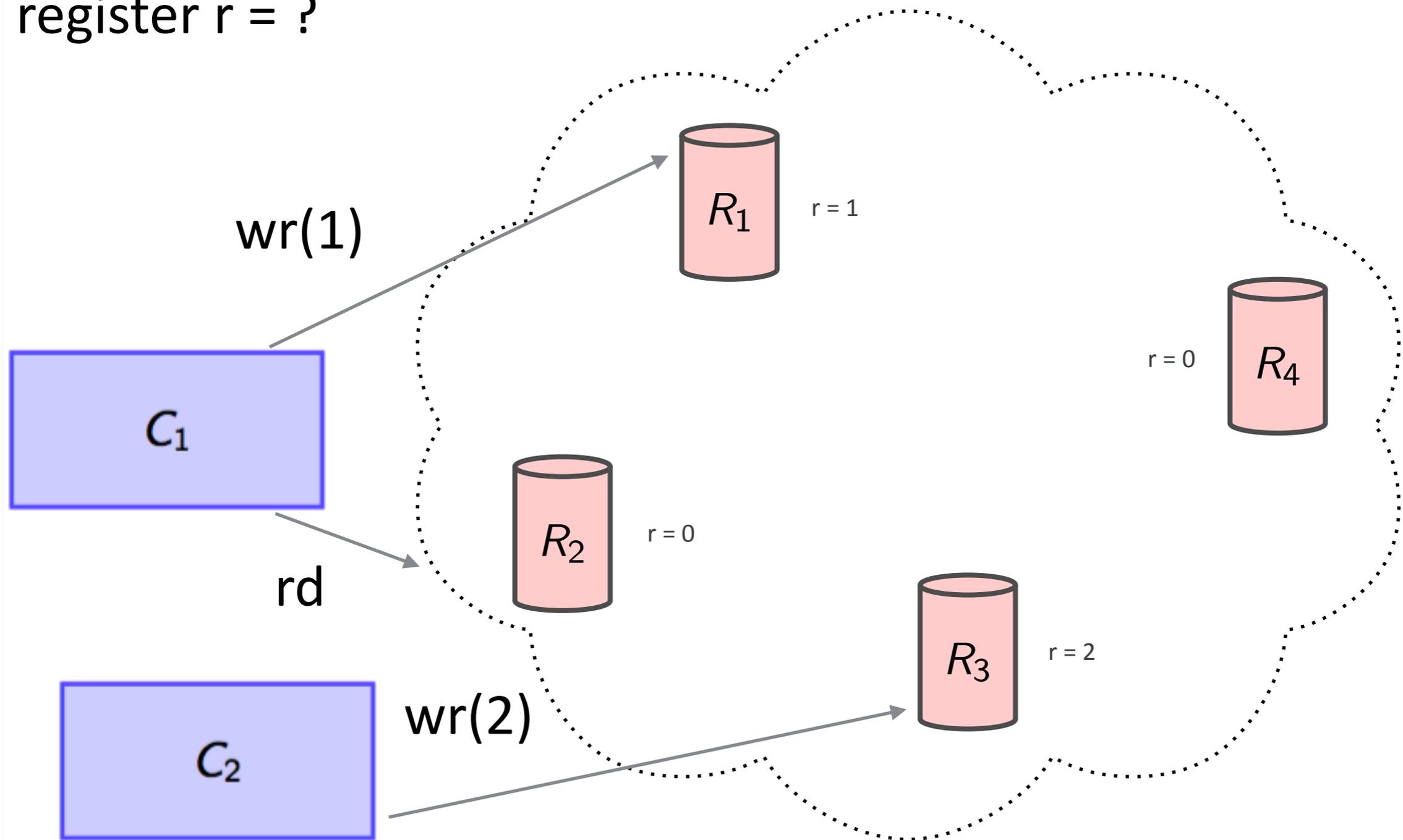
A register

register $r = ?$



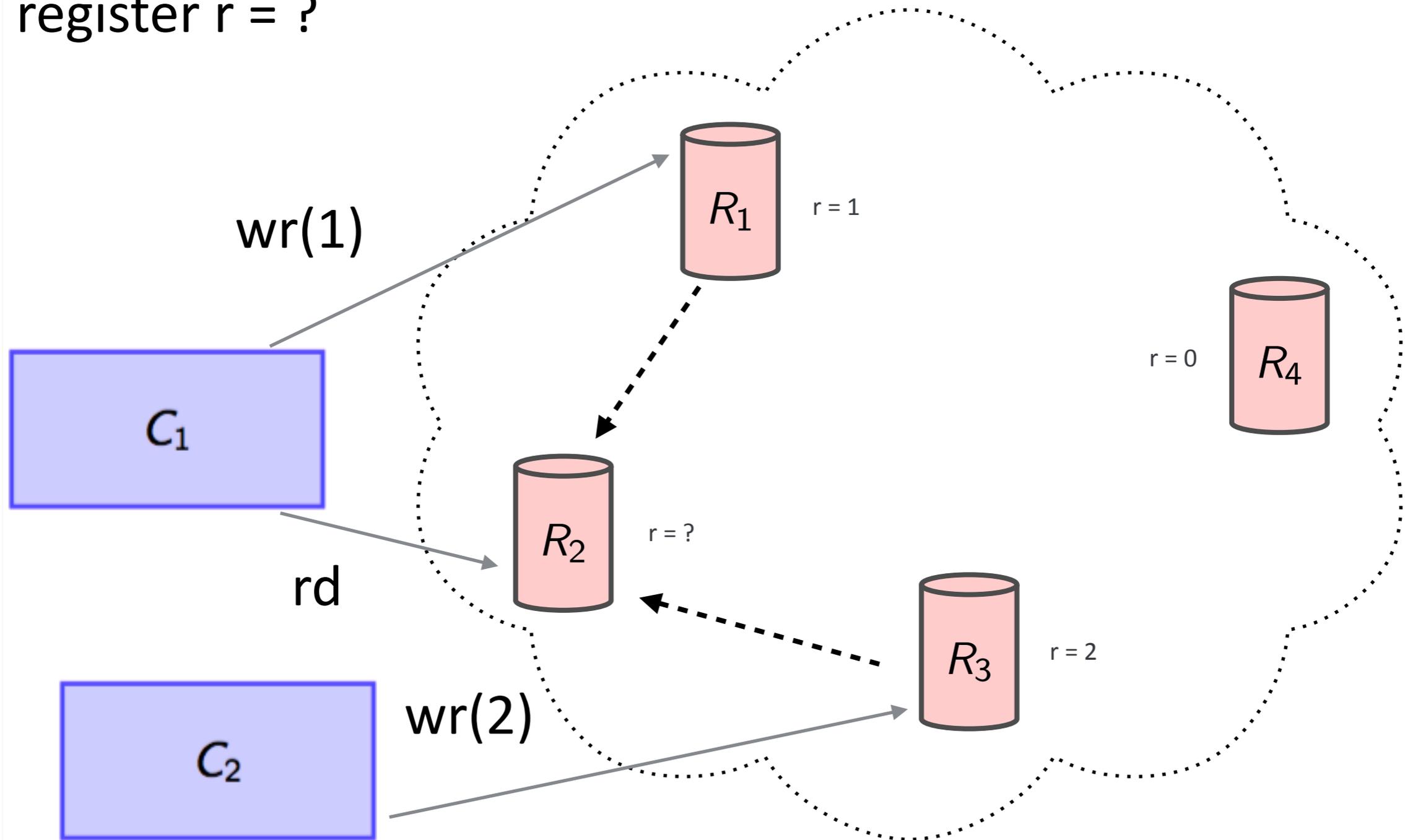
A register

register $r = ?$



A register

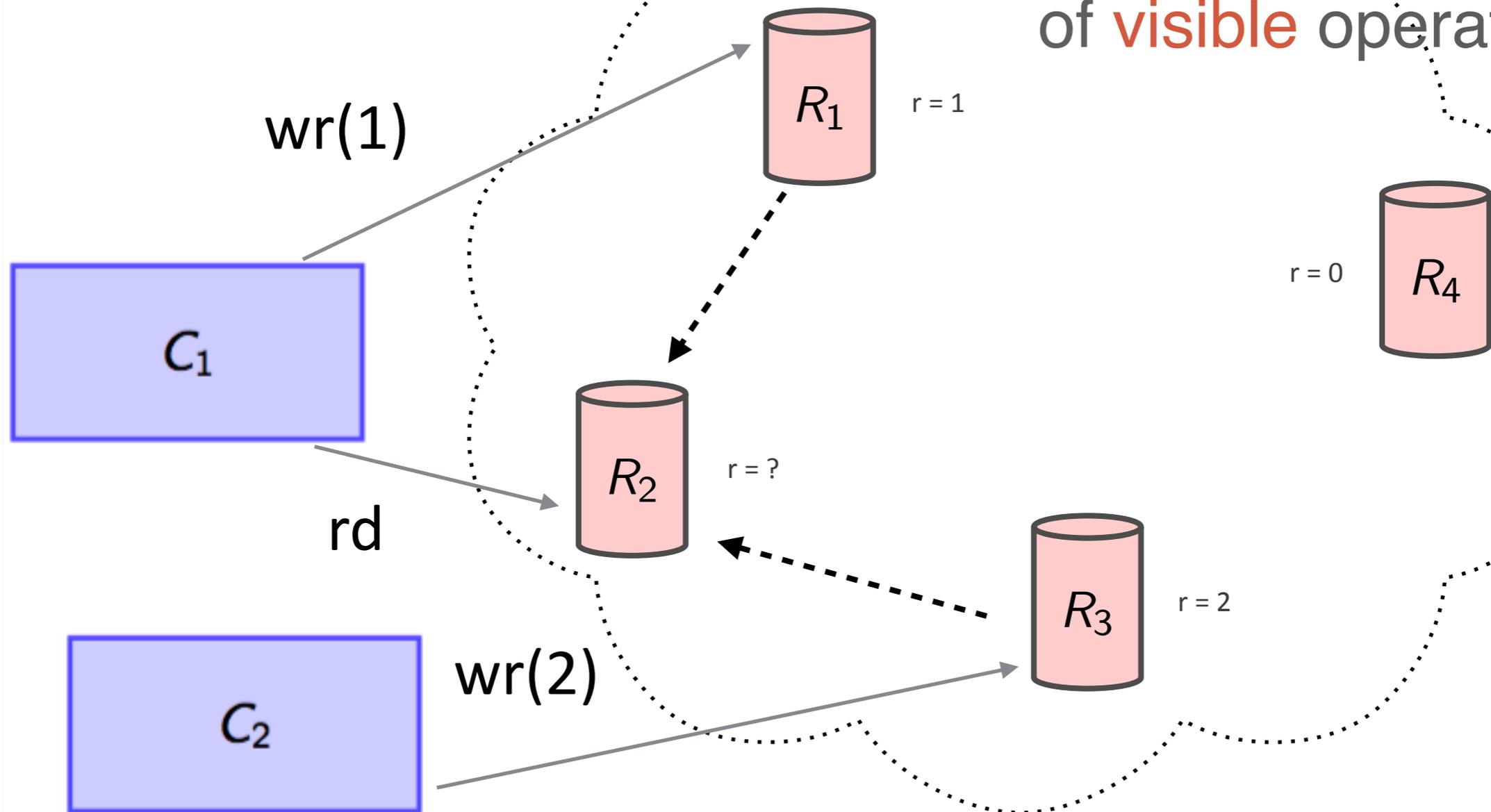
register $r = ?$



A register

register $r = ?$

The result depends on the **relative order** of **visible** operations



Replicated Data Types

op : VIS × ARB → RVAL

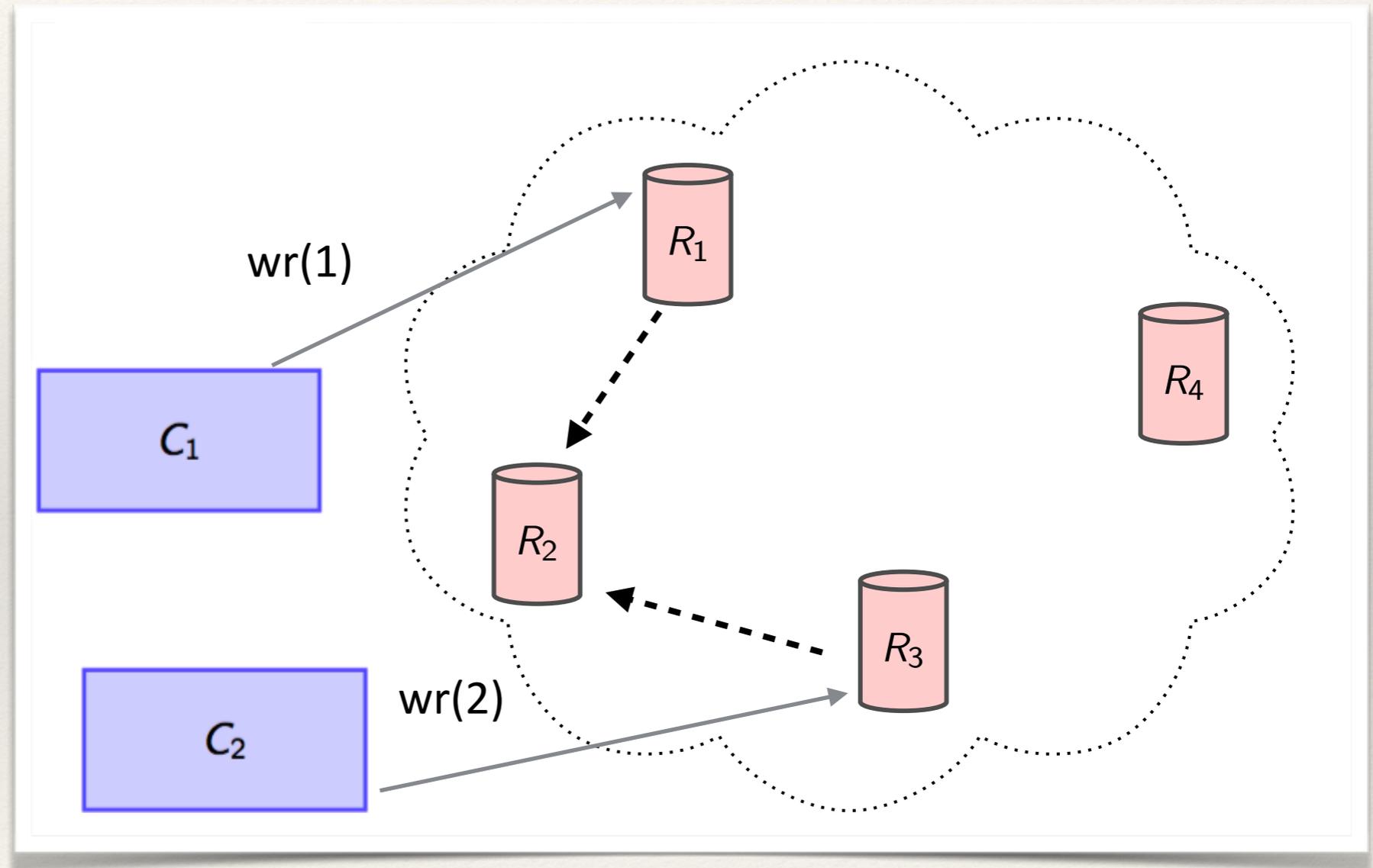
- ❖ **VIS**ibility: A partial order of operations over a replica
- ❖ **ARB**itration: A total order of such operations
- ❖ **Return VAL**ue: The value returned by the last operation

A register

- ❖ Two operations
 - ❖ $\text{rd}(_, _) = ?$
 - ❖ $\text{wr}(k)(_, _) = \text{ok}$

A register

- ❖ Two operations
 - ❖ $rd(_, _) = ?$
 - ❖ $wr(k)(_, _) = ok$

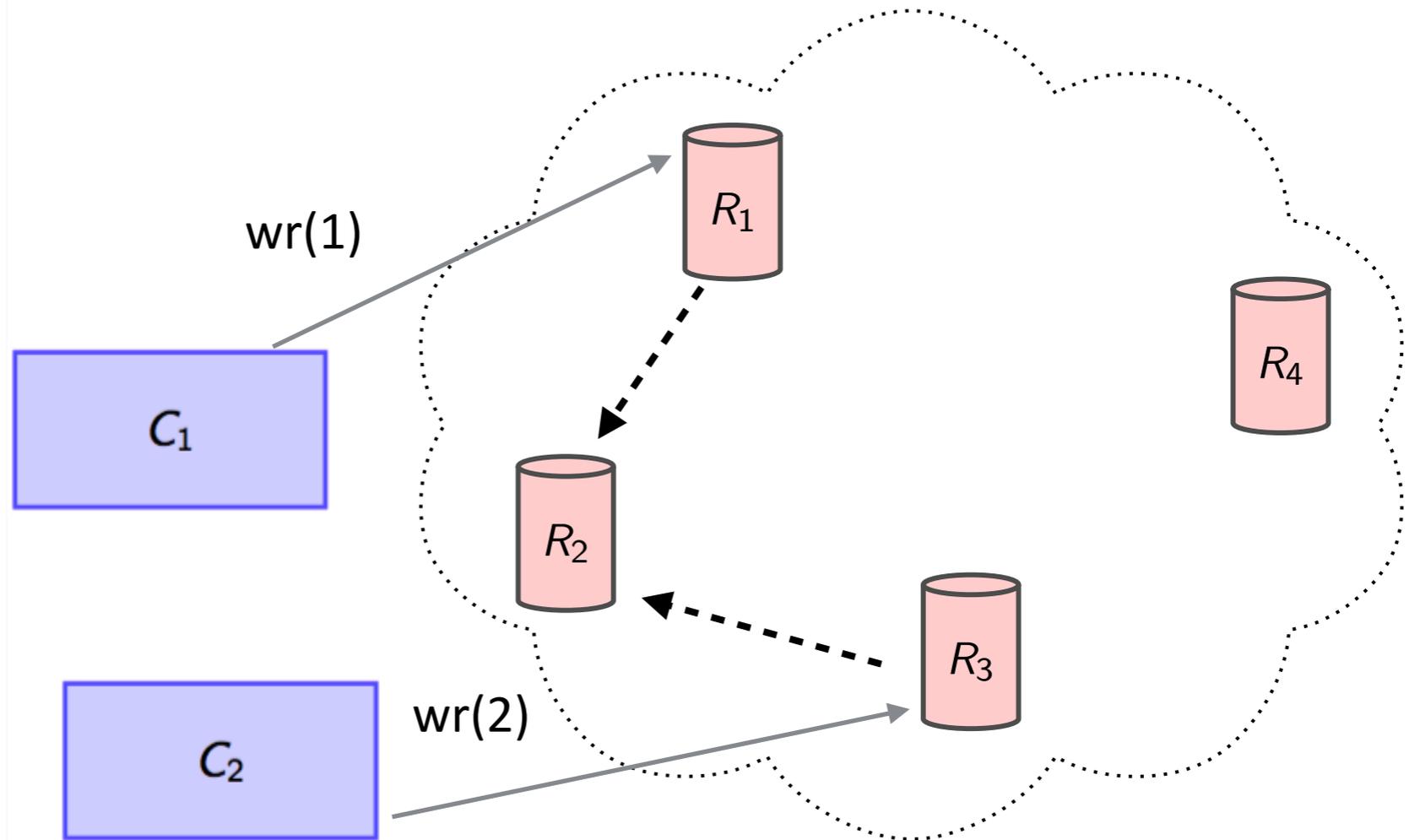


A register

`wr(1)`

`wr(2)`

VISibility



A register

wr(1)

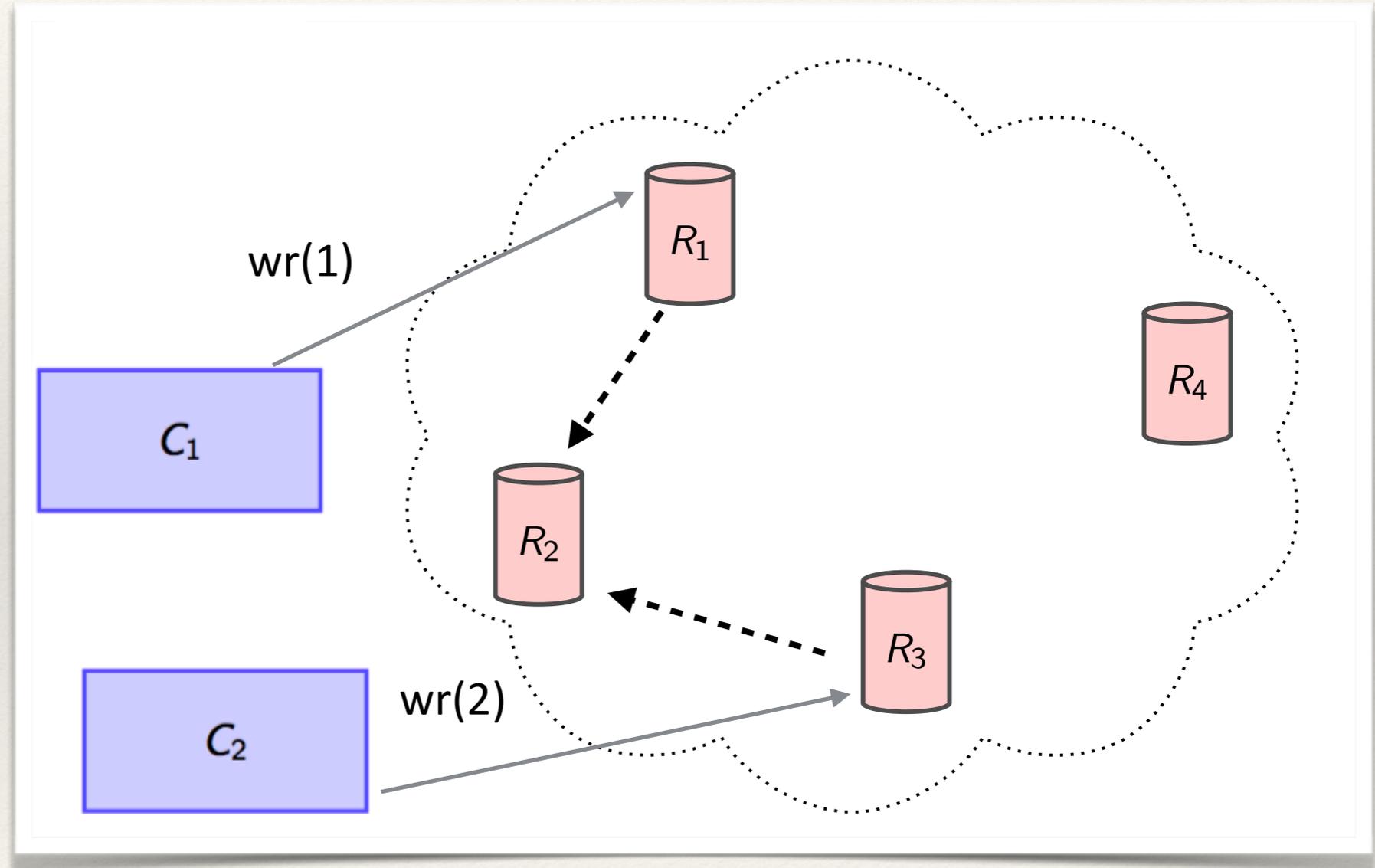
wr(2)

VISibility

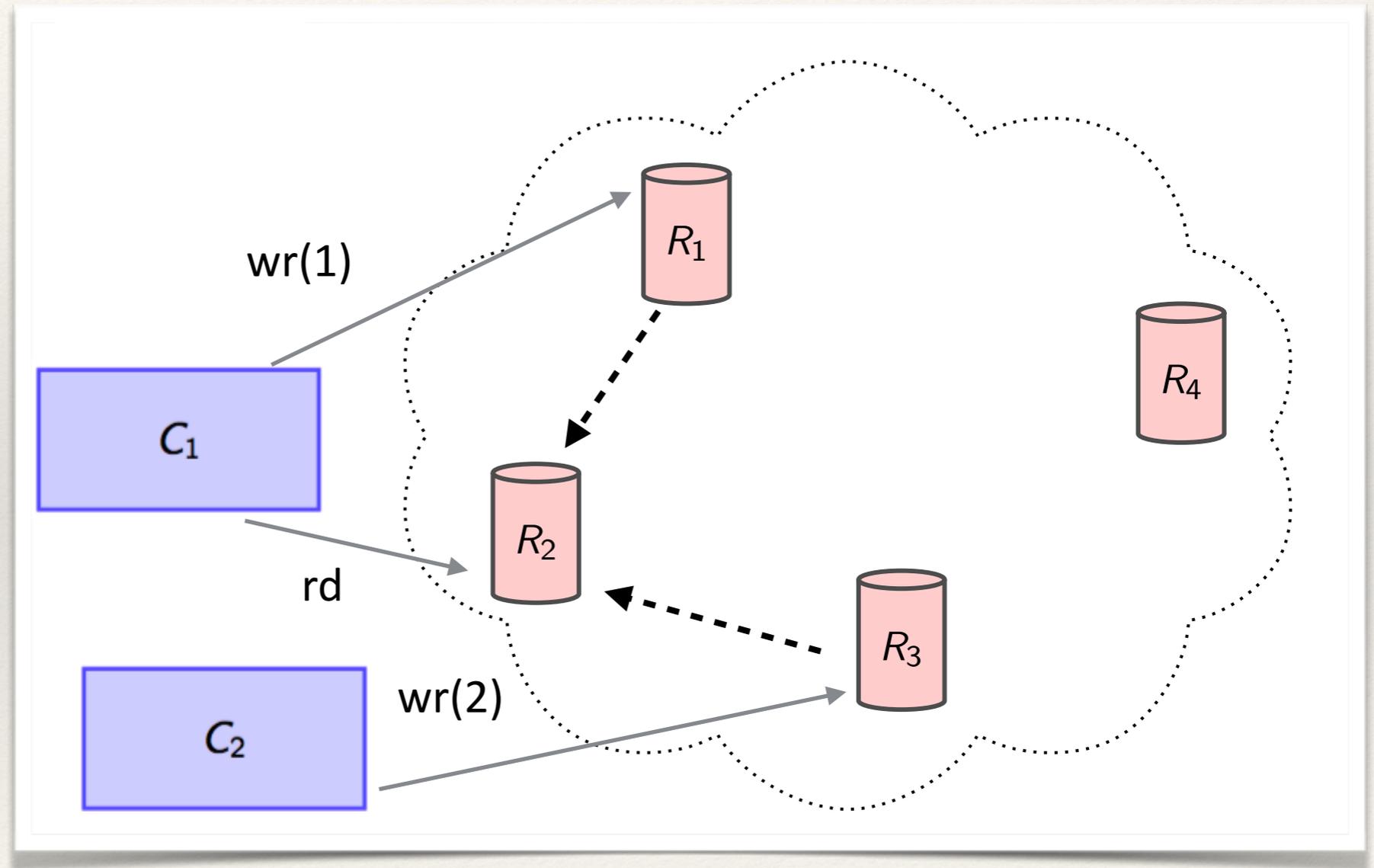
wr(1)

wr(2)

ARBitration

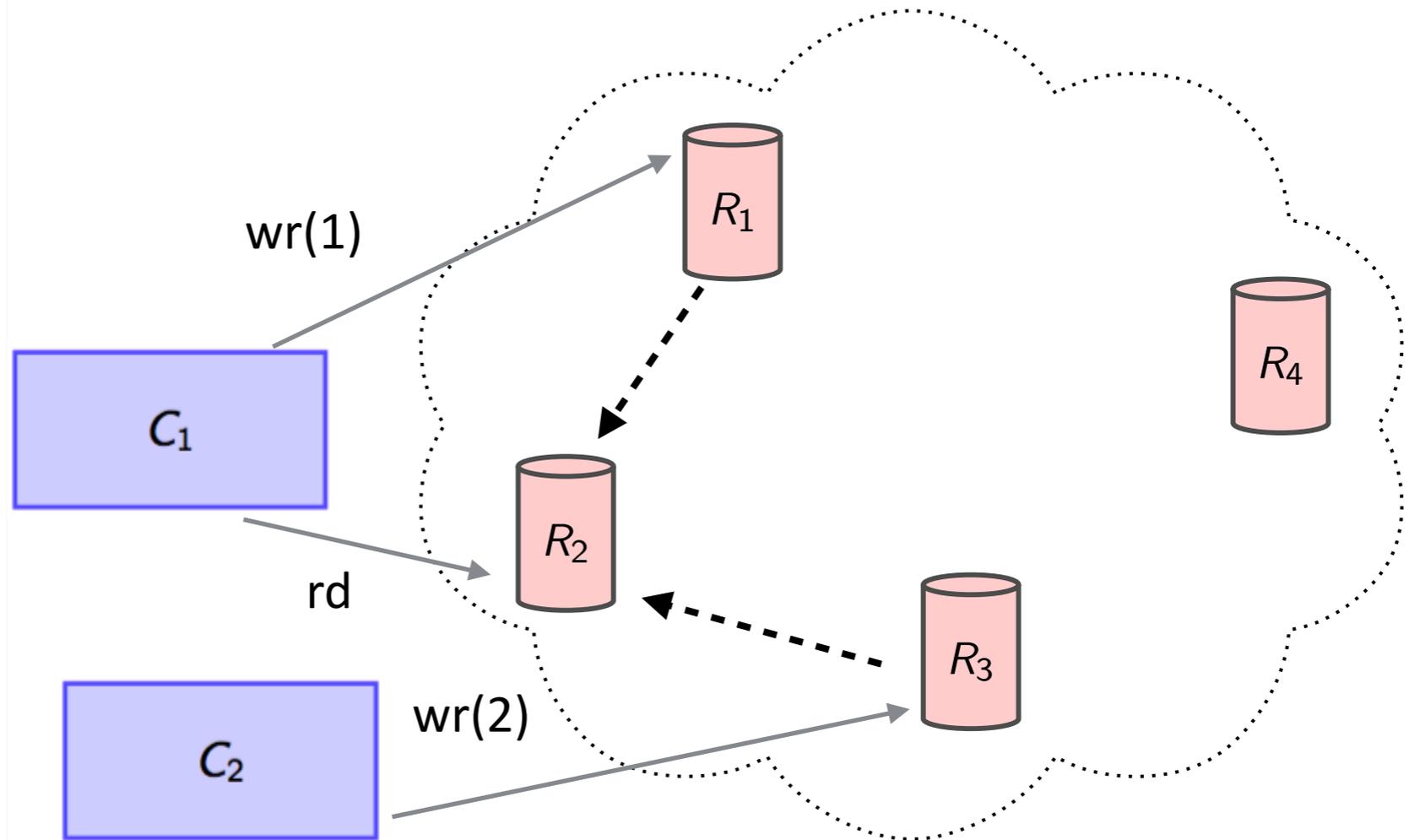


A register



A register

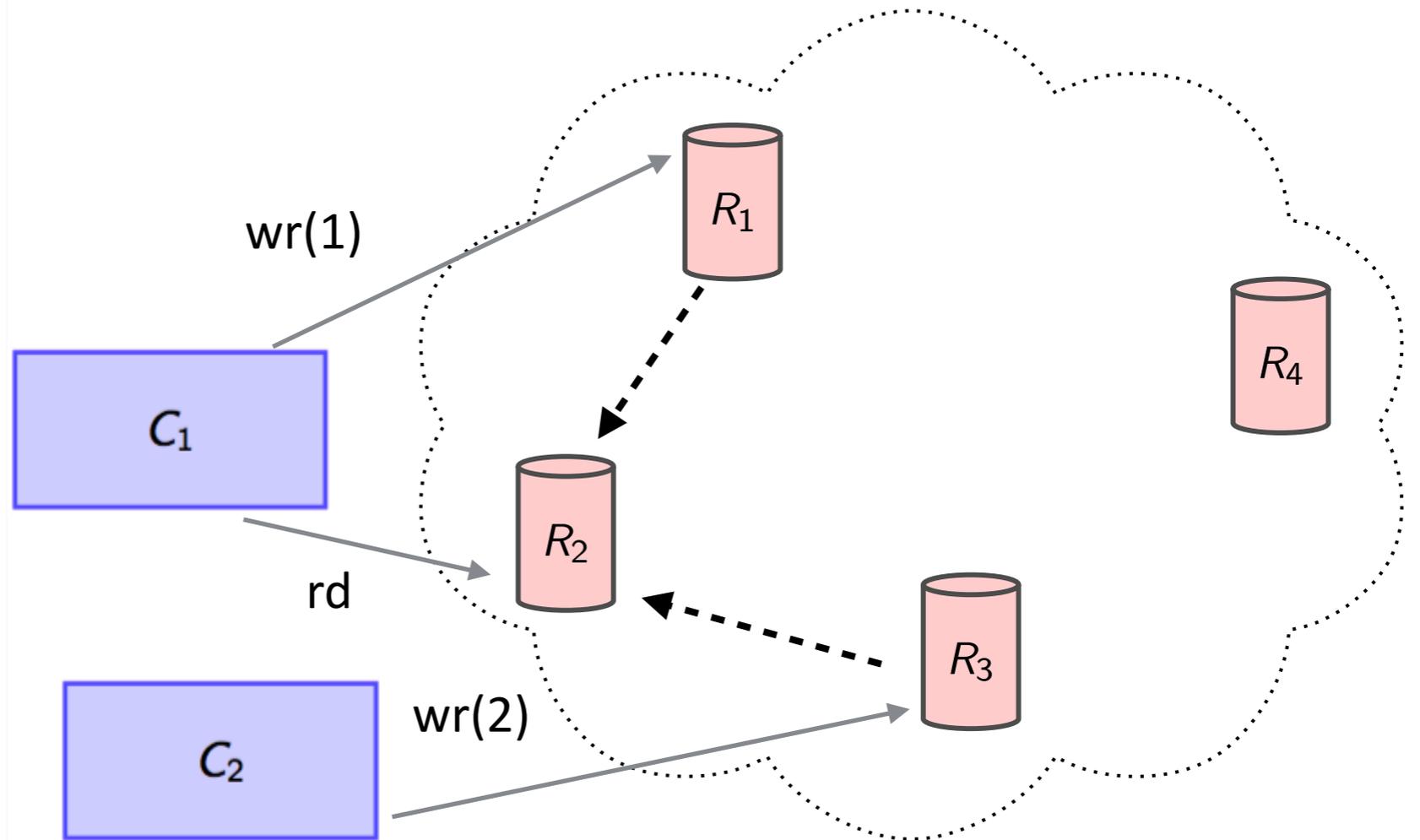
$$\text{rd} \left(\begin{array}{cc} \text{wr}(1) & \text{wr}(2) \\ \text{wr}(1) & \text{wr}(2) \end{array} \right) = 2$$



A register

$$\text{rd} \left(\begin{array}{cc} \text{wr}(1) & \text{wr}(2) \\ \text{wr}(1) & \text{wr}(2) \end{array} \right) = 2$$

$$\text{rd} \left(\begin{array}{cc} \text{wr}(1) & \text{wr}(2) \\ \text{wr}(2) & \text{wr}(1) \end{array} \right) = 1$$

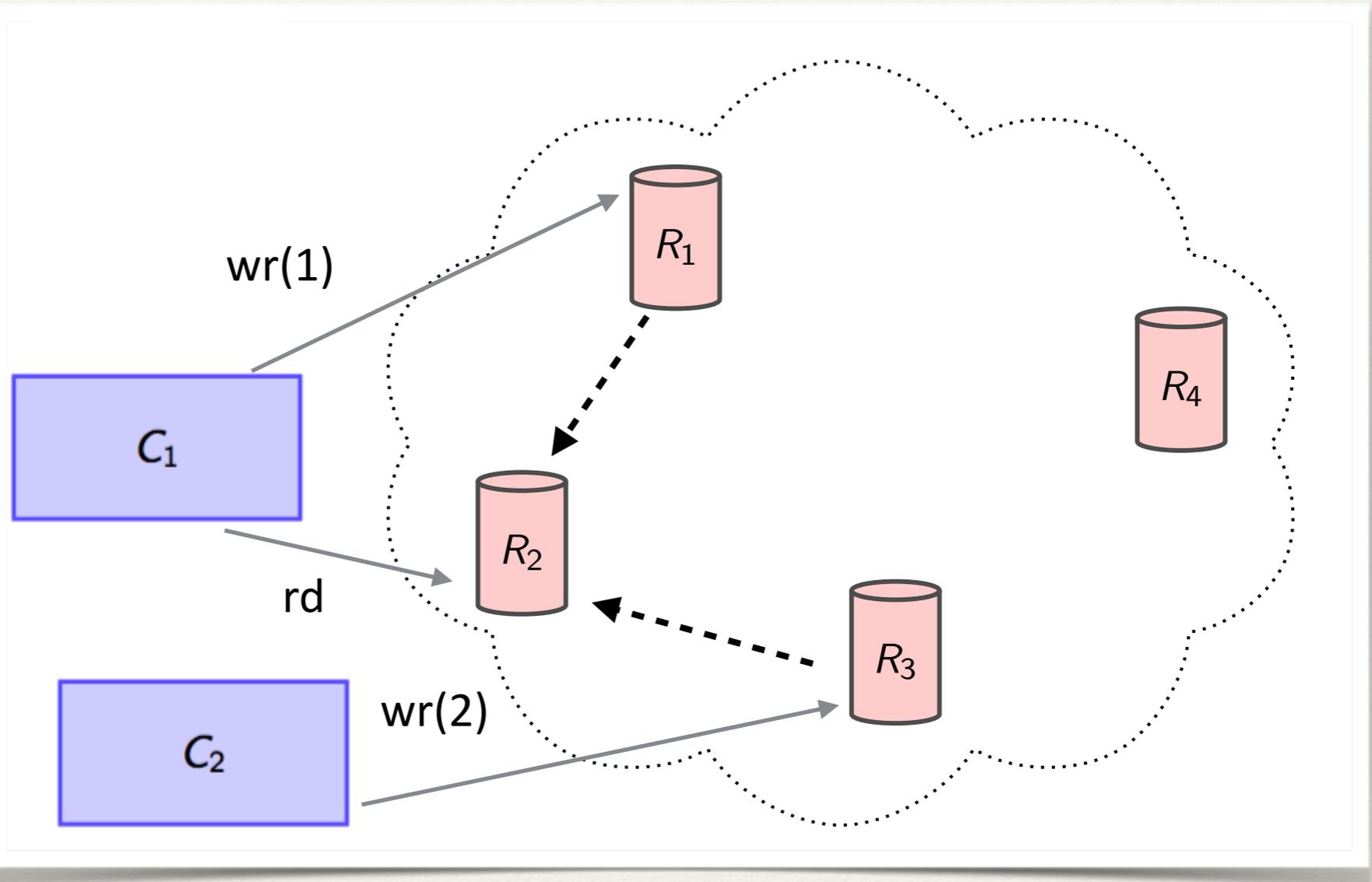


A register

$$\text{rd} \left(\begin{array}{cc} \text{wr}(1) & \text{wr}(2) \\ \text{wr}(1) & \text{wr}(2) \end{array} \right) = 2$$

$$\text{rd} \left(\begin{array}{cc} \text{wr}(1) & \text{wr}(2) \\ \text{wr}(2) & \text{wr}(1) \end{array} \right) = 1$$

Last-write-wins



Understanding RDTs

- ❖ Implementing RDTs means to provide a communication mechanism among replicas, to ensure its compatibility wrt. the behaviour of the operations and to guarantee that some global properties (e.g. eventual convergence of replicas) are preserved.
- ❖ But first...
 - ❖ Is it possible to get a traditional presentation of RDTs?
 - ❖ Is there any implicit assumption on the arbitrations?
 - ❖ Are RDTs compositional? I.e., are arbitrations of larger visibility orders explained in terms of smaller ones?

Internalising values

$\langle \text{wr}(1), \text{ok} \rangle$ $\langle \text{wr}(2), \text{ok} \rangle$

,

Internalising values

$\langle \text{wr}(1), \text{ok} \rangle \quad \langle \text{wr}(2), \text{ok} \rangle$

,

$$\text{rd} \left(\begin{array}{cc|c} & & \text{wr}(1) \\ \text{wr}(1) & \text{wr}(2) & | \\ & & \text{wr}(2) \end{array} \right) = 2$$

Internalising values

$$\text{rd} \left(\begin{array}{cc} & \text{wr}(1) \\ \text{wr}(1) & \text{wr}(2) \\ & | \\ & \text{wr}(2) \end{array} \right) = 2$$

$$\mathcal{S} \left(\begin{array}{cc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \\ & \swarrow \quad \searrow \\ & \langle \text{rd}, 2 \rangle \end{array} \right) = \left\{ \begin{array}{c} \langle \text{wr}(1), \text{ok} \rangle \\ | \\ \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{rd}, 2 \rangle \end{array} \right\}$$

A **specification** goes
from **configurations**
to **sets of arbitrations**

Internalising values

$$\text{rd} \left(\begin{array}{cc} \text{wr}(1) & \text{wr}(2) \\ \text{wr}(1) & \text{wr}(2) \\ \text{wr}(2) & \end{array} \right) = 2$$

$$\text{rd} \left(\begin{array}{cc} \text{wr}(1) & \text{wr}(2) \\ \text{wr}(2) & \text{wr}(1) \\ \text{wr}(1) & \end{array} \right) = 1$$

$$S \left(\begin{array}{cc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \\ & \downarrow \quad \downarrow \\ & \langle \text{rd}, \bar{0} \rangle \end{array} \right) = \left\{ \begin{array}{c} \dots \\ \dots \\ \dots \end{array} \right\}$$

A **specification** goes
from **configurations**
to **sets of arbitrations**

Recovering RDTs: saturation

$$\text{rd} \left(\begin{array}{cc} & \text{wr}(1) \\ \text{wr}(1) & \text{wr}(2) \\ & | \\ & \text{wr}(2) \end{array} \right) = 2$$

$$S \left(\begin{array}{cc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \\ & \swarrow \quad \searrow \\ & \langle \text{rd}, 2 \rangle \end{array} \right) = \left\{ \begin{array}{c} \langle \text{wr}(1), \text{ok} \rangle \\ | \\ \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{rd}, 2 \rangle \end{array} \right\}$$

A **specification** goes
from **configurations**
to **sets of arbitrations**

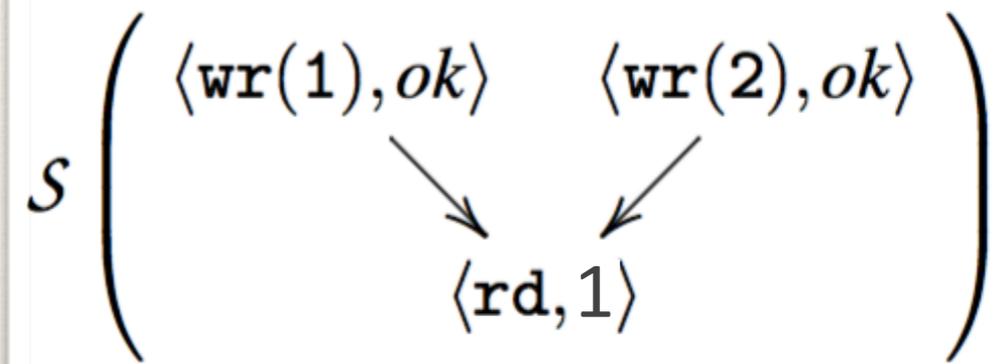
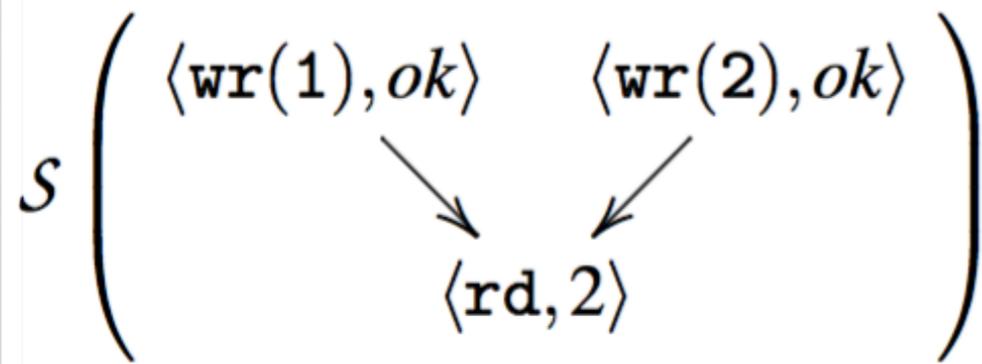
Recovering RDTs: saturation

$$\text{rd} \left(\begin{array}{cc} & \text{wr}(1) \\ \text{wr}(1) & \text{wr}(2) \\ & | \\ & \text{wr}(2) \end{array} \right) = 2$$

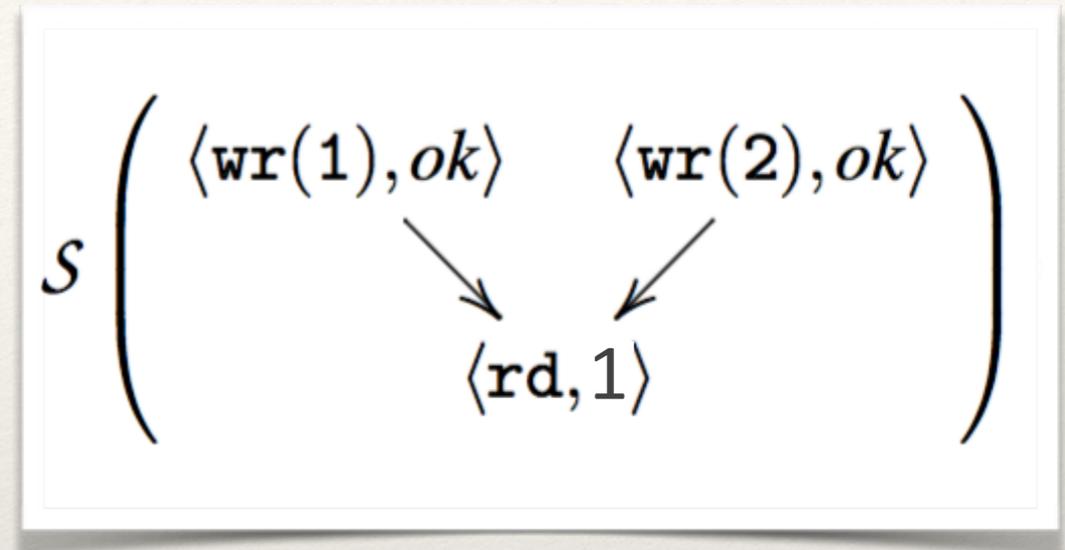
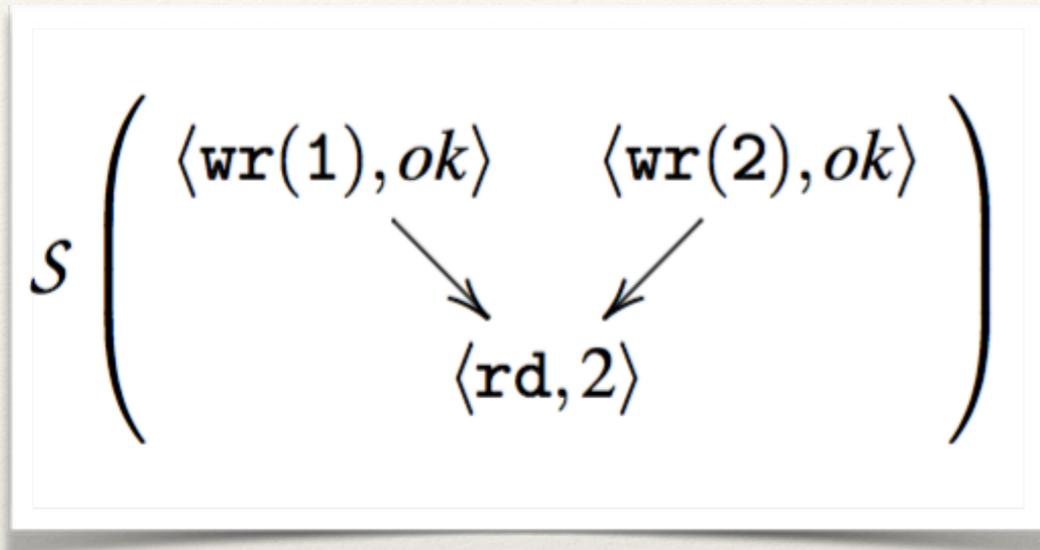
$$S \left(\begin{array}{cc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \\ & \swarrow \quad \searrow \\ & \langle \text{rd}, 2 \rangle \end{array} \right) = \left\{ \begin{array}{ccc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(1), \text{ok} \rangle & \langle \text{rd}, 2 \rangle \\ | & | & | \\ \langle \text{wr}(2), \text{ok} \rangle & \langle \text{rd}, 2 \rangle & \langle \text{wr}(1), \text{ok} \rangle \\ | & | & | \\ \langle \text{rd}, 2 \rangle & \langle \text{wr}(2), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \end{array} \right\}$$

A **specification** goes
from **configurations**
to **sets of arbitrations**

Recovering RDTs: determinism

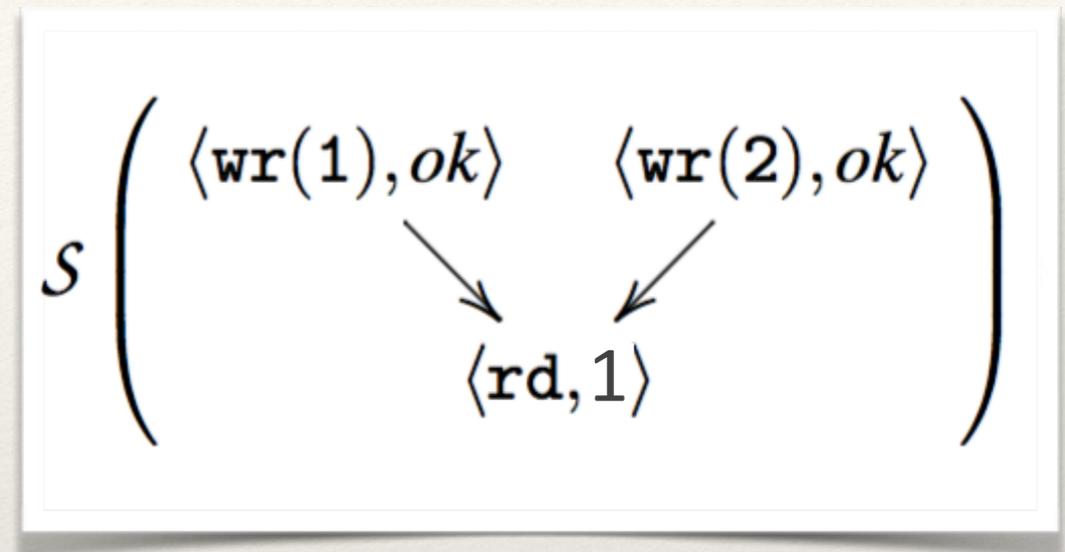
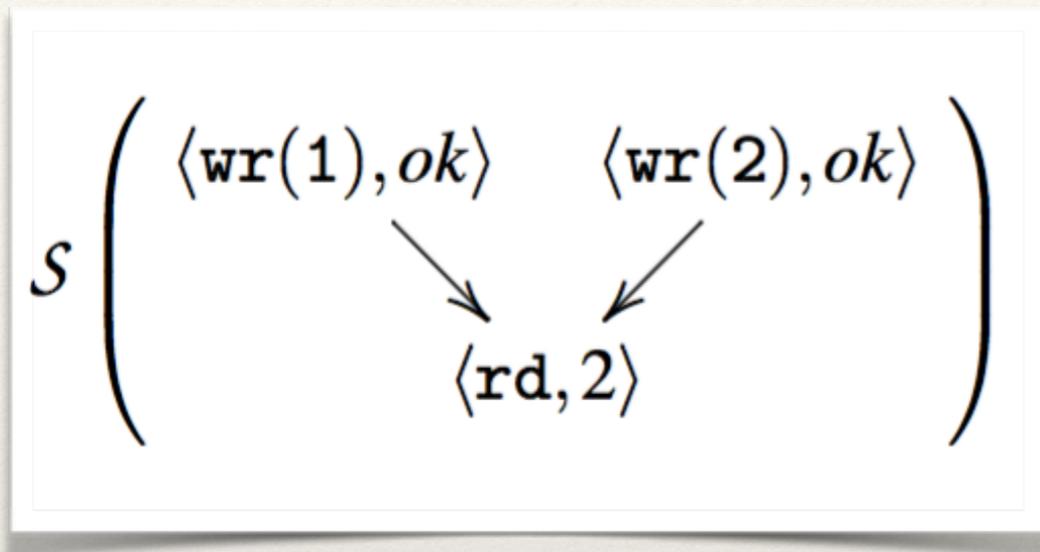


Recovering RDTs: determinism



value-deterministic: empty intersection (after removing the last event)

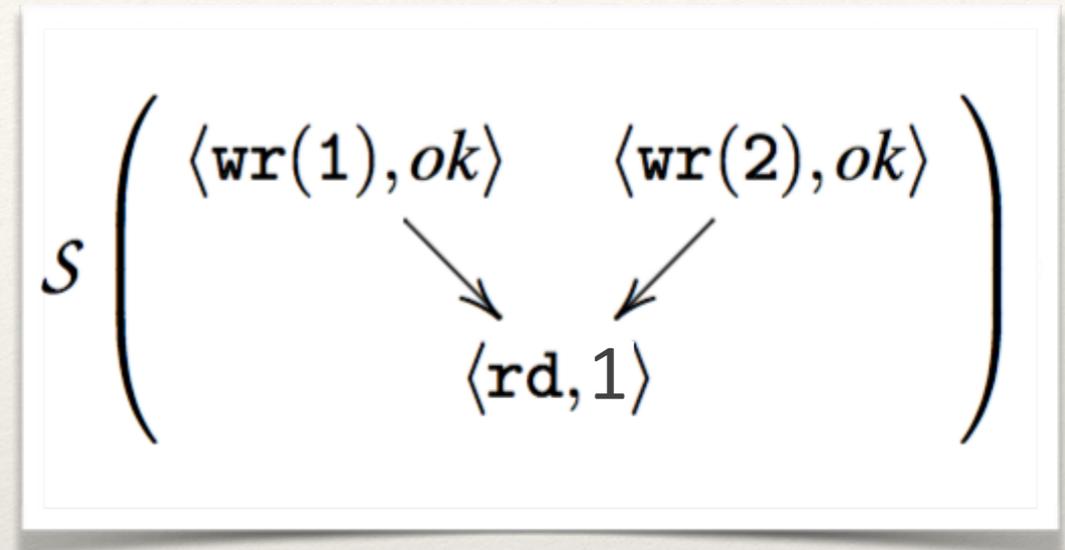
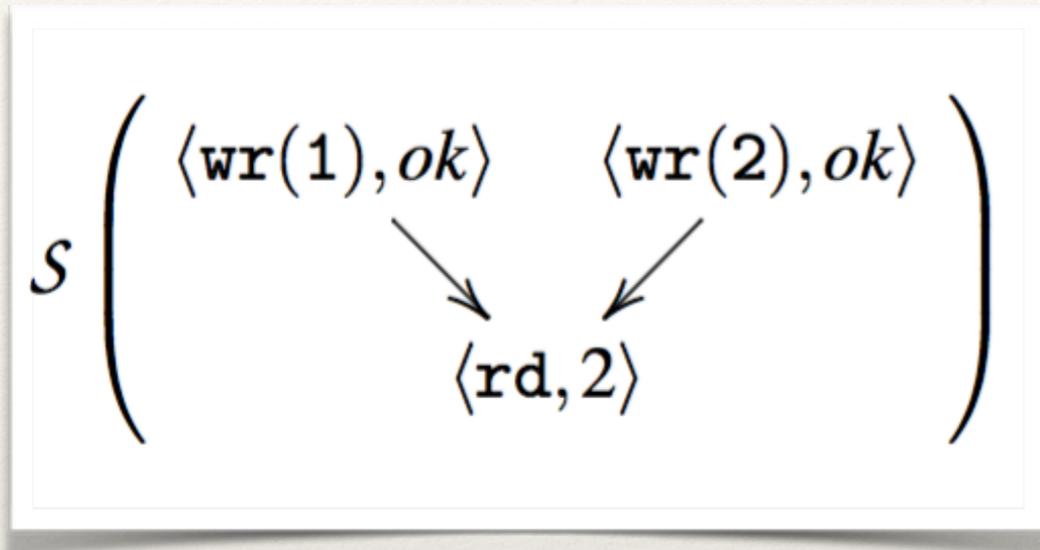
Recovering RDTs: determinism



value-deterministic: empty intersection (after removing the last event)

deterministic: empty intersection *even* forgetting the value component

Recovering RDTs: determinism



value-deterministic: empty intersection (after removing the last event)

deterministic: empty intersection *even* forgetting the value component

RTDs have chosen the second path, thus e.g. forbidding write failures

Recovering RDTs: coherence

$$\begin{array}{c} \langle \text{wr}(1), \text{ok} \rangle \\ | \\ \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{rd}, 2 \rangle \end{array} \otimes \begin{array}{c} \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{wr}(3), \text{ok} \rangle \end{array} = \left\{ \begin{array}{cc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(1), \text{ok} \rangle \\ | & | \\ \langle \text{wr}(2), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \\ | & | \\ \langle \text{rd}, 2 \rangle & \langle \text{wr}(3), \text{ok} \rangle \\ | & | \\ \langle \text{wr}(3), \text{ok} \rangle & \langle \text{rd}, 2 \rangle \end{array} \right\}$$

Recovering RDTs: coherence

$$\begin{array}{c} \langle \text{wr}(1), \text{ok} \rangle \\ | \\ \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{rd}, 2 \rangle \end{array} \otimes \begin{array}{c} \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{wr}(3), \text{ok} \rangle \end{array} = \left\{ \begin{array}{cc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(1), \text{ok} \rangle \\ | & | \\ \langle \text{wr}(2), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \\ | & | \\ \langle \text{rd}, 2 \rangle & \langle \text{wr}(3), \text{ok} \rangle \\ | & | \\ \langle \text{wr}(3), \text{ok} \rangle & \langle \text{rd}, 2 \rangle \end{array} \right\}$$

$$\forall G. \mathcal{S}(G) = \bigotimes_{e \in \mathcal{E}_G} \mathcal{S}(G|_{-\prec^* e})$$

Recovering RDTs: coherence

$$\begin{array}{c} \langle \text{wr}(1), \text{ok} \rangle \\ | \\ \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{rd}, 2 \rangle \end{array} \otimes \begin{array}{c} \langle \text{wr}(2), \text{ok} \rangle \\ | \\ \langle \text{wr}(3), \text{ok} \rangle \end{array} = \left\{ \begin{array}{cc} \langle \text{wr}(1), \text{ok} \rangle & \langle \text{wr}(1), \text{ok} \rangle \\ | & | \\ \langle \text{wr}(2), \text{ok} \rangle & \langle \text{wr}(2), \text{ok} \rangle \\ | & | \\ \langle \text{rd}, 2 \rangle & \langle \text{wr}(3), \text{ok} \rangle \\ | & | \\ \langle \text{wr}(3), \text{ok} \rangle & \langle \text{rd}, 2 \rangle \end{array} \right\}$$

$$\forall G. \mathcal{S}(G) = \bigotimes_{e \in \mathcal{E}_G} \mathcal{S}(G|_{-\prec^* e})$$

Admissible arbitrations
never increases when
extending the visibility

Recovering RDTs: the theorem

- ❖ There is a one-to-one correspondence between RDTs and saturated, deterministic, and coherent specifications

Recovering RDTs: the theorem

- ❖ There is a one-to-one correspondence between RTDs and saturated, deterministic, and coherent specifications

....which is bad for RDTs

Recovering RDTs: the theorem

- ❖ There is a one-to-one correspondence between RTDs and saturated, deterministic, and coherent specifications

....which is bad for RDTs

both saturation and determinism are bad!!

Recovering RDTs: the theorem

- ❖ There is a one-to-one correspondence between RTDs and saturated, deterministic, and coherent specifications

....which is bad for RDTs

both saturation and determinism are bad!!

$$\mathcal{S} \left(\begin{array}{c} \langle \text{inc, ok} \rangle \\ \downarrow \\ \langle \text{rd, 1} \rangle \end{array} \right) = \left\{ \begin{array}{c} \langle \text{inc, ok} \rangle \\ | \\ \langle \text{rd, 1} \rangle \end{array} \right\} \quad \mathcal{S} \left(\begin{array}{c} \langle \text{inc, fail} \rangle \\ \downarrow \\ \langle \text{rd, } \perp \rangle \end{array} \right) = \left\{ \begin{array}{c} \langle \text{inc, fail} \rangle \\ | \\ \langle \text{rd, } \perp \rangle \end{array} \right\}$$

value-deterministic, yet not deterministic

From specifications to transition systems

$\langle G, P \rangle$

$P \in \mathcal{S}(G)$

states

From specifications to transition systems

 $\langle G, P \rangle$ $P \in \mathcal{S}(G)$

states

 $\langle G, P \rangle \xrightarrow{\ell} \langle G', P' \rangle$

transitions

 $G' = G^\ell$ $P'|_{\mathcal{E}_G} = P$

From specifications to transition systems

(COMP)

$$\frac{\langle G_1, P|_{\mathcal{E}_{G_1}} \rangle \xrightarrow{\ell} \langle G'_1, P'_1 \rangle \quad P' \in P \otimes P'_1}{\langle G_1 \sqcup G_2, P \rangle \xrightarrow{\ell} \langle G'_1 \sqcup G_2, P' \rangle}$$

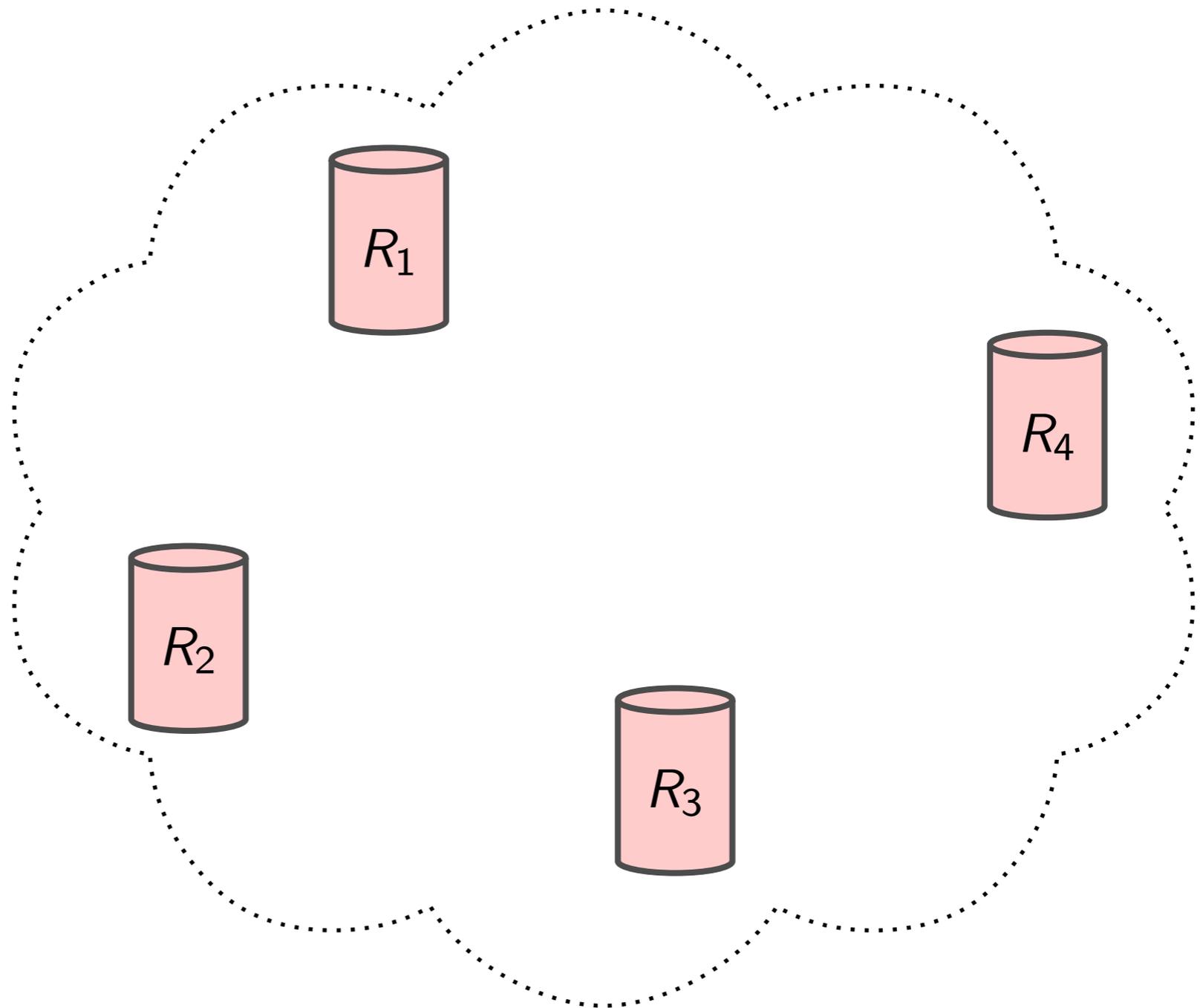
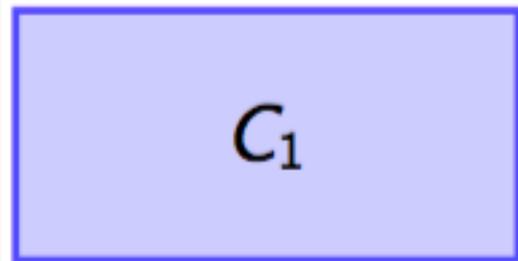
an abstract transition system against which
to compare (by asynchronous simulation)
those of actual implementations...

Implementing a specification

- ❖ Proposed implementation of RDTs [Burckhardt et al.].
 - ❖ Each replica propagates its state to the other replicas.
 - ❖ It is assumed that all replicas have the same behaviour.

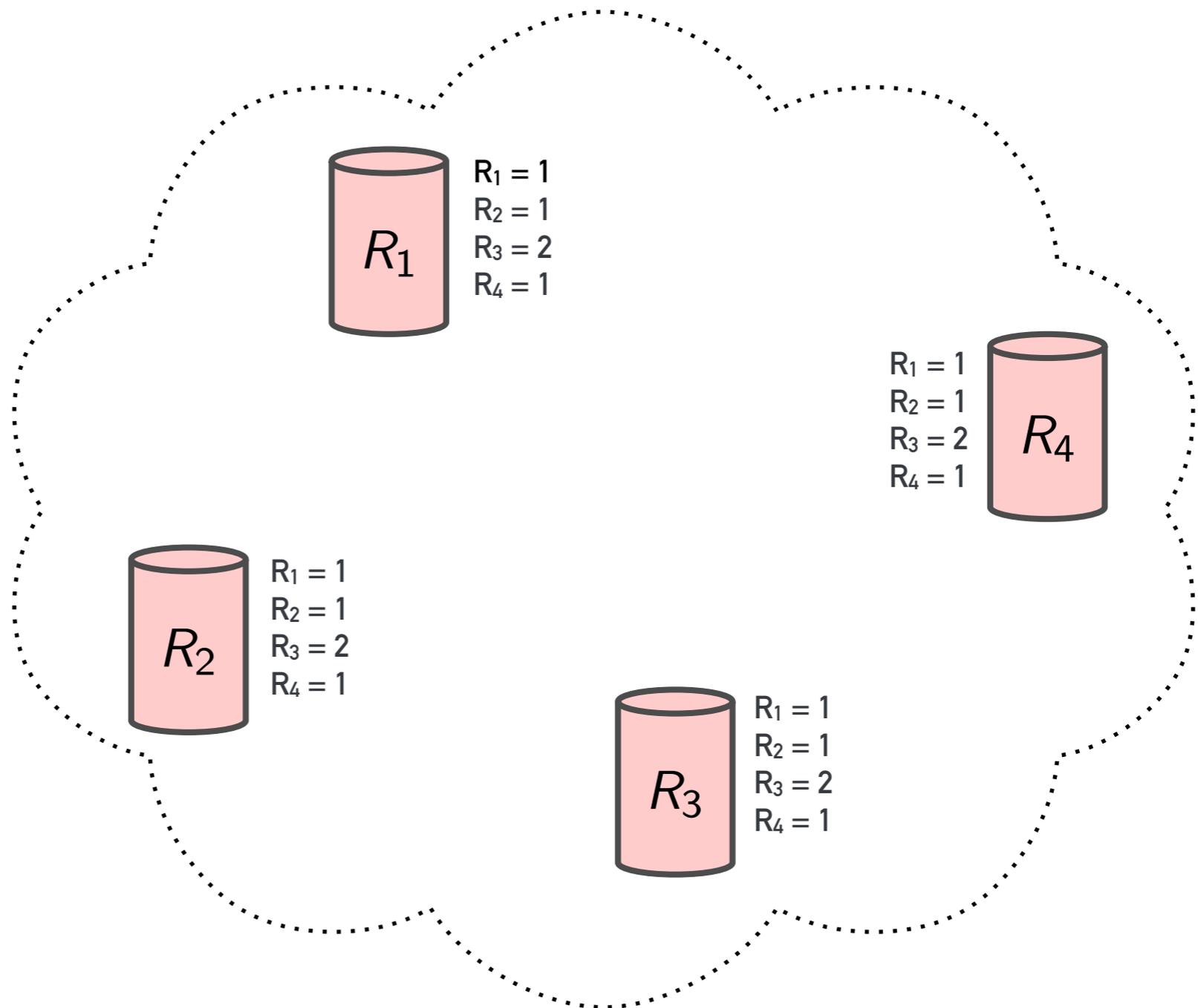
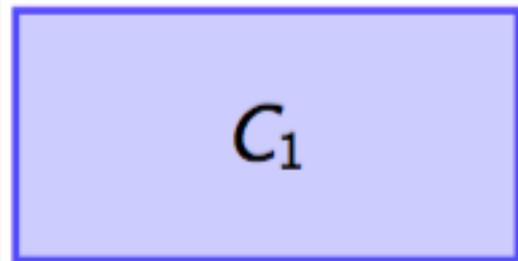
A counter

counter $c = 5$



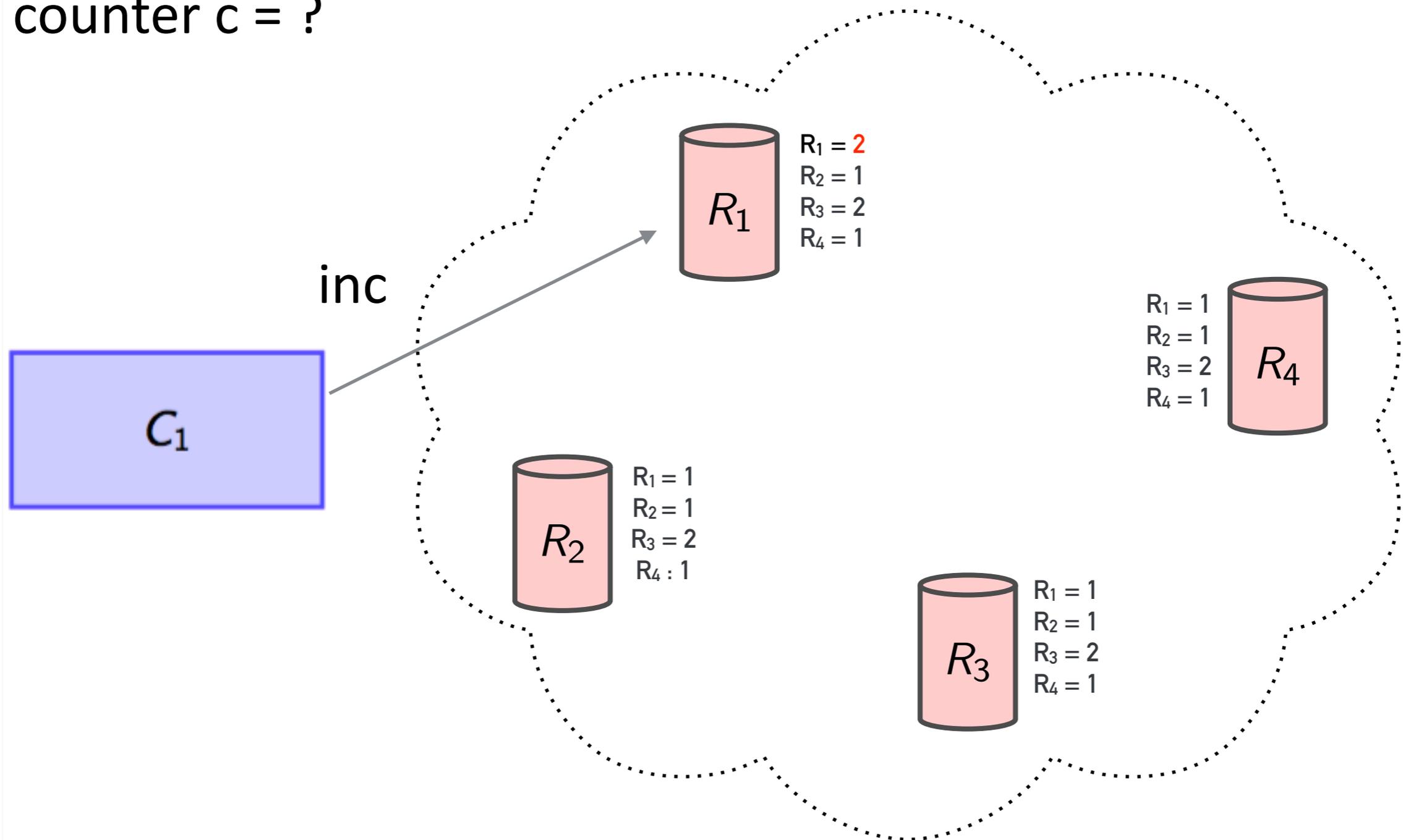
A counter

counter $c = 5$



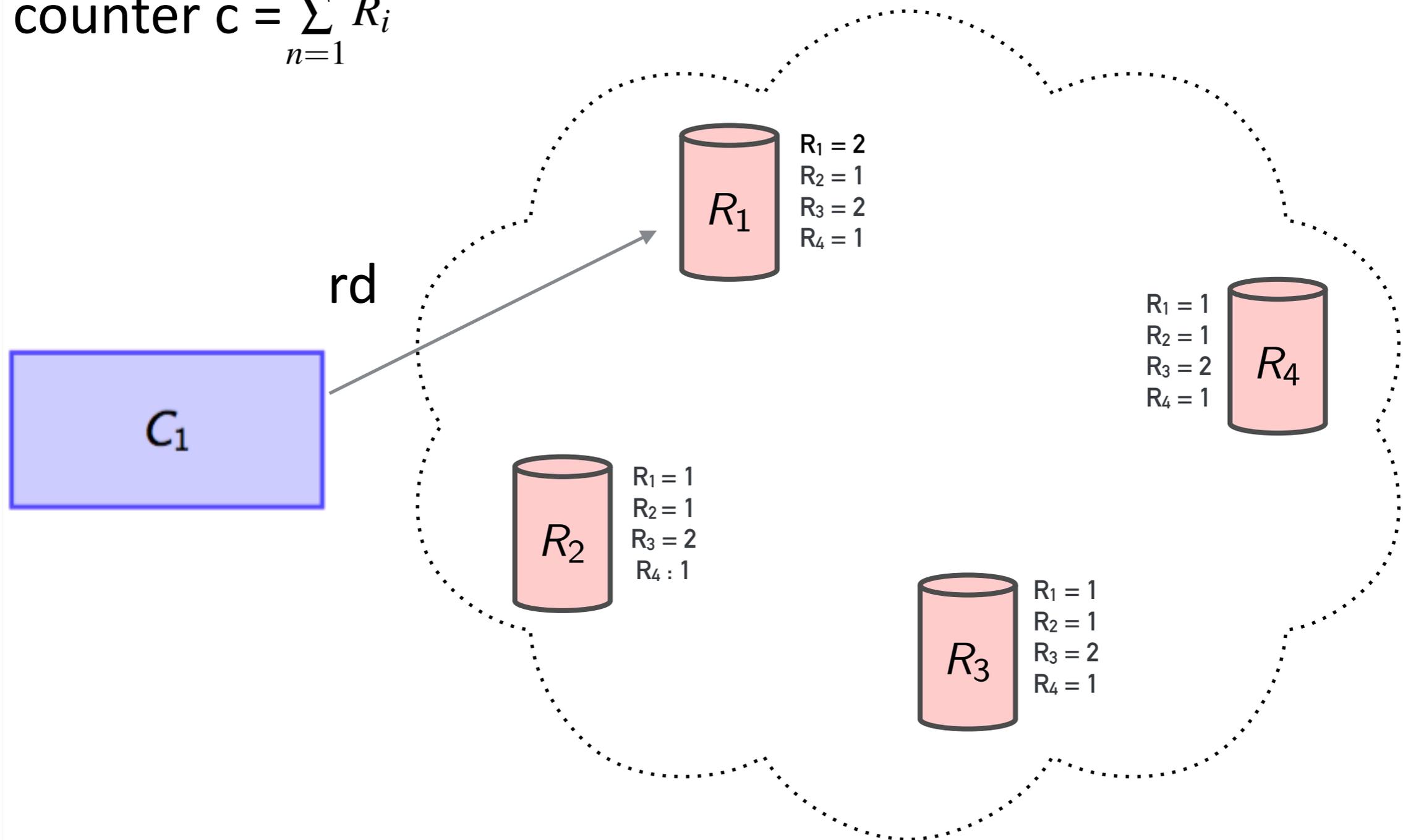
A counter

counter $c = ?$



A counter

$$\text{counter } c = \sum_{n=1}^4 R_n$$



an LTS for a counter

$$\Sigma = \mathcal{R} \times (\mathcal{R} \mapsto \mathbb{N})$$

$$\mathcal{L} = \{\langle \text{inc}, \text{ok} \rangle\} \cup (\{\text{rd}\} \times \mathbb{N})$$

(READ)

$$\frac{k = \sum_{s \in \text{dom}(v)} v(s)}{\langle r, v \rangle \xrightarrow{\text{rd}, k} \langle r, v \rangle}$$

(INC)

$$\langle r, v \rangle \xrightarrow{\text{inc}, \text{ok}} \langle r, v[r \mapsto v(r) + 1] \rangle$$

an LTS for a counter

(SEND)

$$\langle r, v \rangle \xrightarrow{\text{send}, \langle r, v \rangle} \langle r, v \rangle$$

(RCV)

$$\langle r, v \rangle \xrightarrow{\text{rcv}, \langle r_k, v_k \rangle} \langle r, \max\{v, v_k\} \rangle$$

$$\forall s. \max\{v, v_k\}(s) = \max\{v(s), v_k(s)\}$$

the resulting LTS is correct wrt. the abstract LTS
(via a suitable simulation)

an LTS for multiple counters

(SEND)

$$\langle r, v \rangle \xrightarrow{\text{send}, \langle r, v \rangle} \langle r, v \rangle$$

(RCV)

$$\langle r, v \rangle \xrightarrow{\text{rcv}, \langle r_k, v_k \rangle} \langle r, \max\{v, v_k\} \rangle$$

(PARL)

$$\frac{\sigma_1 \xrightarrow{\ell} \sigma'_1}{\sigma_1 \parallel \sigma_2 \xrightarrow{\ell} \sigma'_1 \parallel \sigma_2}$$

(COMM)

$$\frac{\sigma_1 \xrightarrow{\text{send}, \sigma} \sigma'_1 \quad \sigma_2 \xrightarrow{\text{rcv}, \sigma} \sigma'_2}{\sigma_1 \parallel \sigma_2 \xrightarrow{\tau} \sigma'_1 \parallel \sigma'_2}$$

the resulting LTS is still correct wrt. the abstract LTS

Conclusions & future works

Conclusions & future works

- ❖ We provided a denotationally-flavoured characterisation of a well-accepted definition of replicated data types
- ❖ thus making explicit some implicit assumptions

Conclusions & future works

- ❖ We provided a denotationally-flavoured characterisation of a well-accepted definition of replicated data types
 - ❖ thus making explicit some implicit assumptions
- ❖ ...& a mechanism for proving correctness of implementations

Conclusions & future works

- ❖ We provided a denotationally-flavoured characterisation of a well-accepted definition of replicated data types
 - ❖ thus making explicit some implicit assumptions
- ❖ ...& a mechanism for proving correctness of implementations
- ❖ We are looking for a categorical presentation
 - ❖ ...in order to get operators for composing specifications

Conclusions & future works

- ❖ We provided a denotationally-flavoured characterisation of a well-accepted definition of replicated data types
 - ❖ thus making explicit some implicit assumptions
- ❖ ...& a mechanism for proving correctness of implementations
- ❖ We are looking for a categorical presentation
 - ❖ ...in order to get operators for composing specifications
- ❖ We plan to recast guarantee properties via the abstract LTS