

REPLICATION

FOR LONG-TERM

PERSISTENCE

JERRY SALTER

M.I.T. Lab. for CS

LIBRARY 2000

DESIRED TIME OF PERSISTENCE

>>

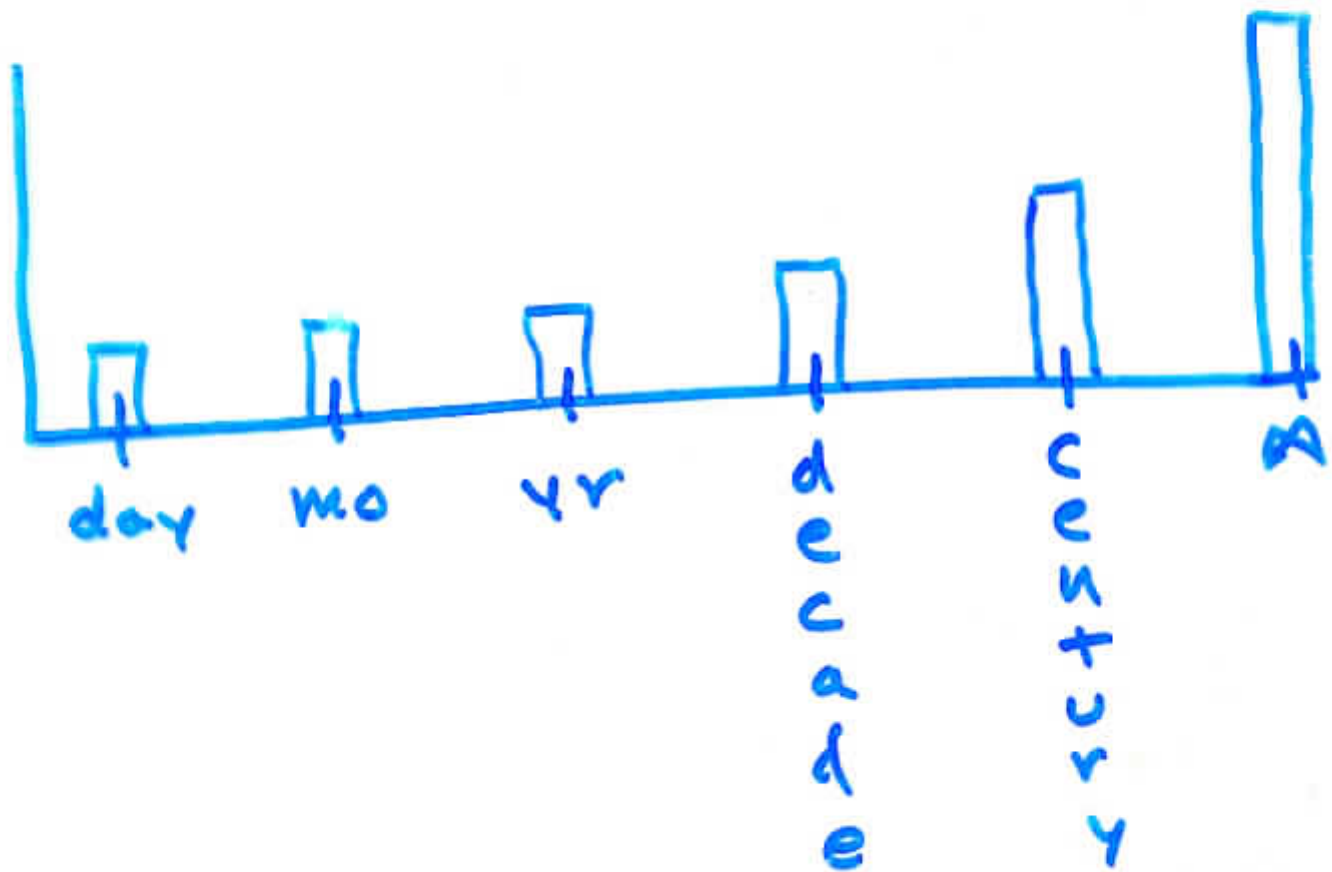
MTB SITE FAILURE

- Earthquake
- Hurricane
- Flood
- Meteor Strike
- City Fire
- Revolution

Strategy:

Full Replication, with
Geographic Diversity

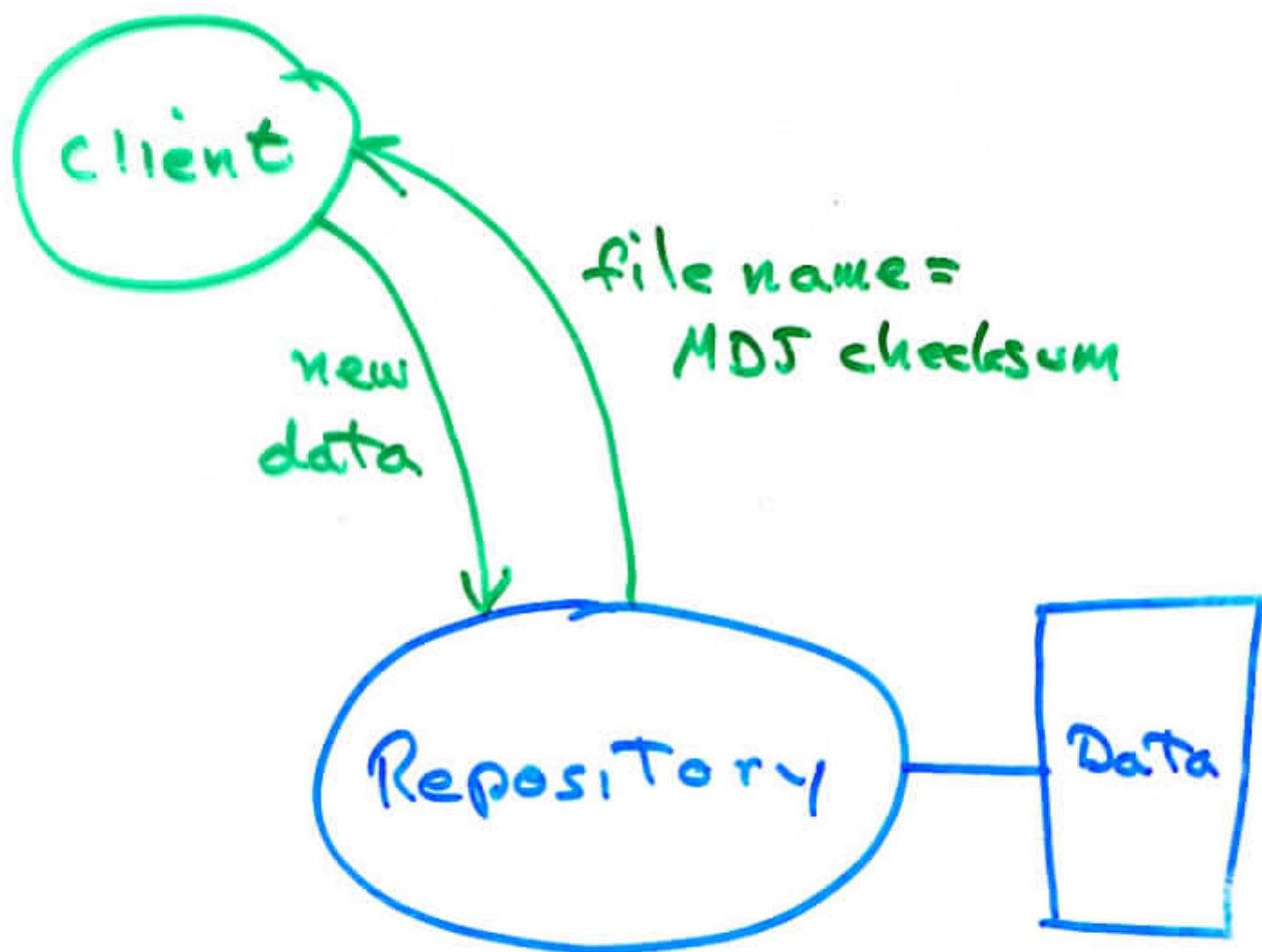
TIME BETWEEN READS



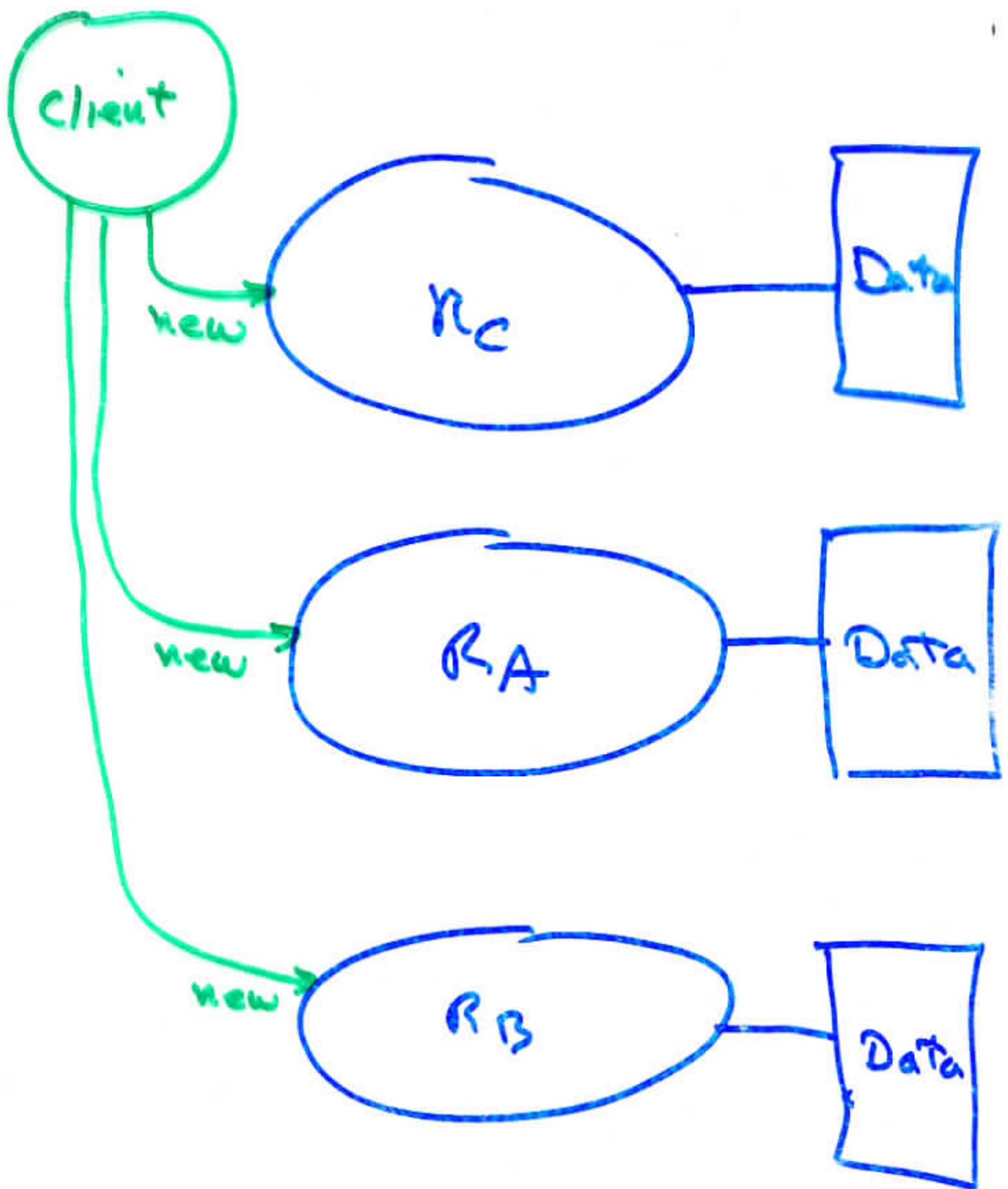
←————→
MTBF Media Failure

strategy:

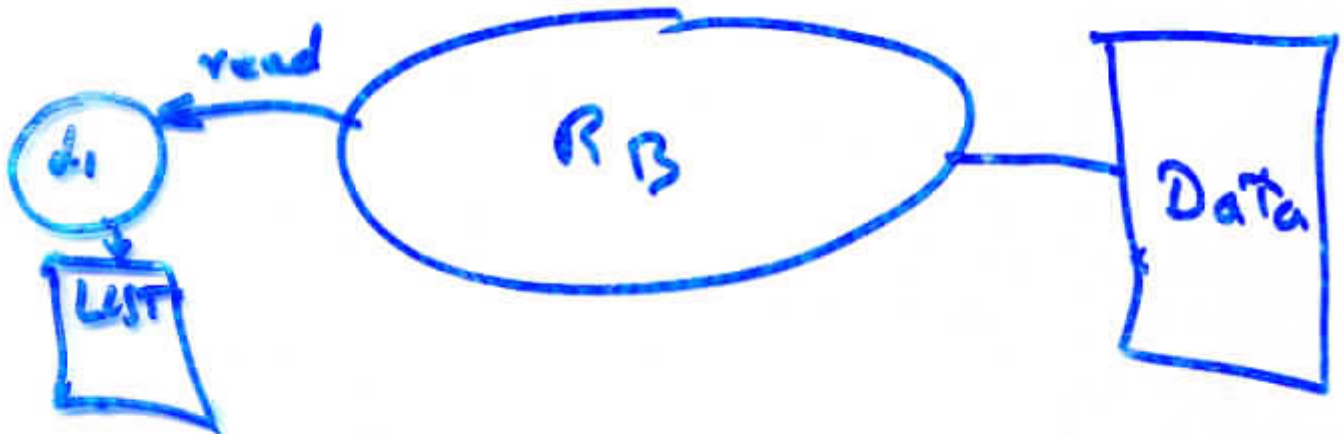
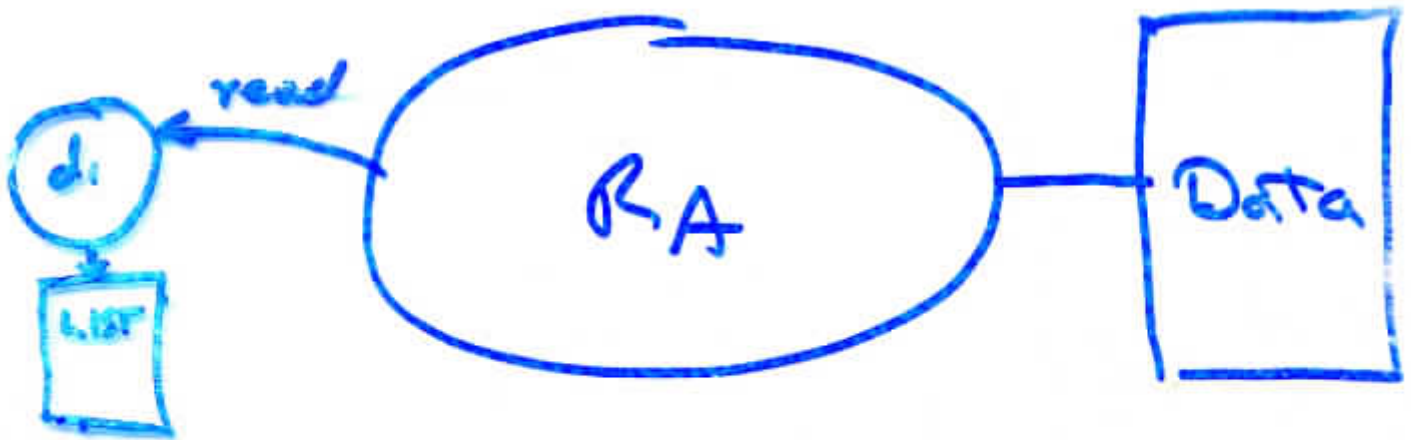
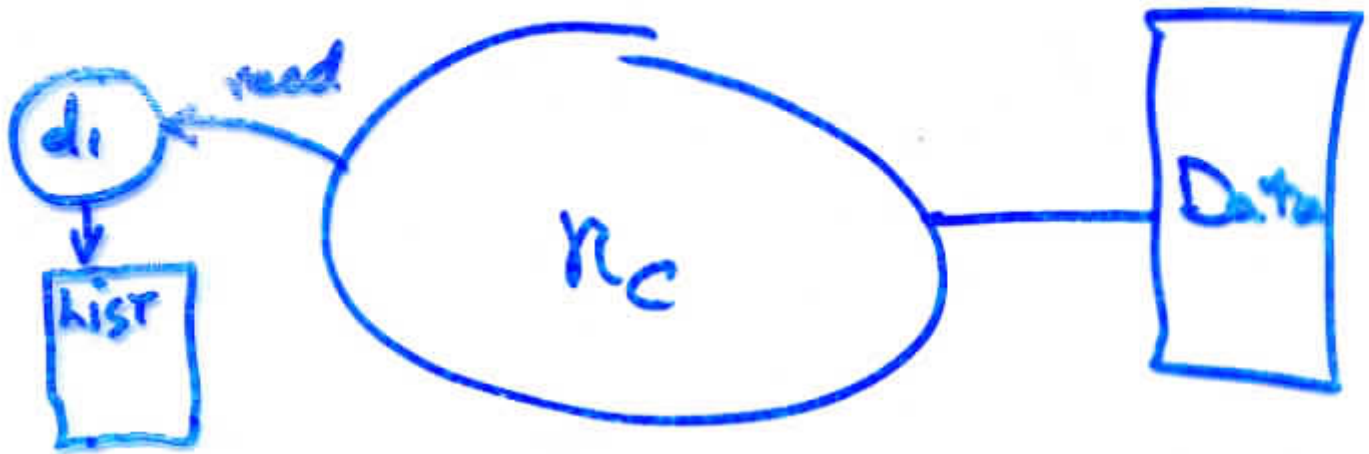
READING DAEMON



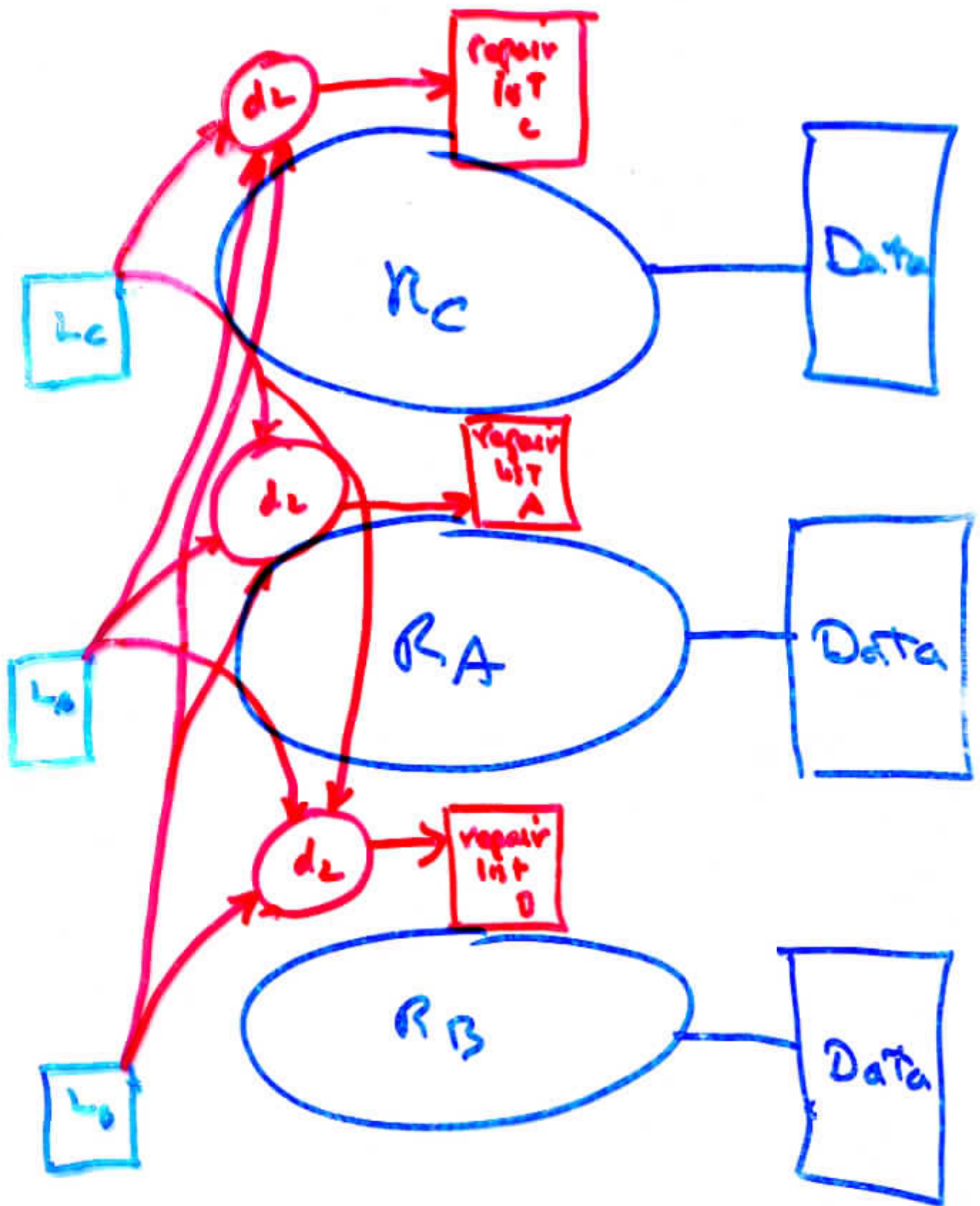
Append-Only Write Semantics



Stateless Commit: Majority
of R's accept



d_i : reading daemon



d_2 : repair order daemon

Repair List A

Delete X

Replace Y with copy from B

⋮



d3: repair execution daemon
(one/replica)

Construct repair list for replica J:

acquire list_k, k = 1...n

let U = union (list_k, k = 1...n)

for each entry in U

if entry checked in majority of k

if entry is unchecked or

missing in list_j

issue repair order:

“copy entry to j”

endif

otherwise

If entry in list_j

issue repair order:

“delete entry at j”

endif

endif

endfor

Simplicity-1

- Relaxed Consistency
- Stateless Protocols
- Append-Only Semantics
- Idempotent Repair
- Unique file name
= MD5 checksum