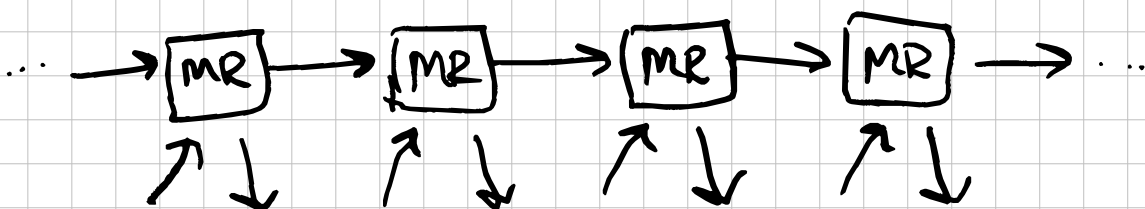


CSE138 Lecture 18

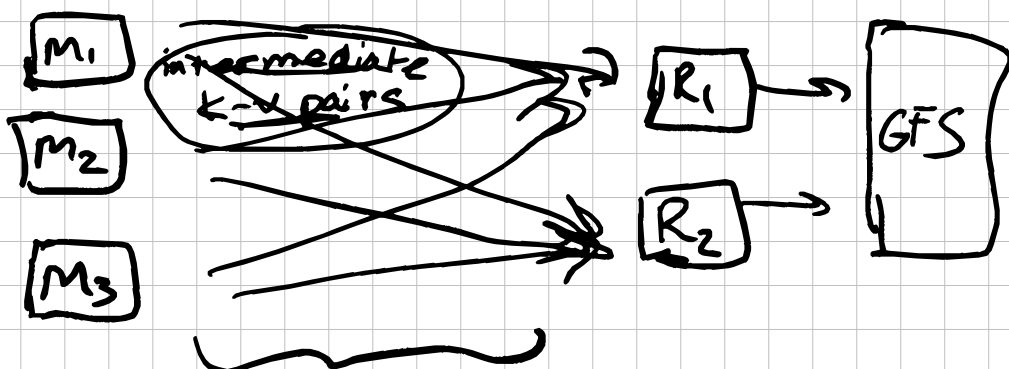
this time:

- MapReduce wrap-up
- the math behind replica conflict resolution



it would be slow to read/write data to disk for each MR job!

a system for orchestrating MR jobs:
Flume Java (2010)



$$\text{shuffle} \quad \text{hash}(\text{key}) \bmod N \leftarrow \# \text{ of reduce workers}$$

With MapReduce, this is not expected to change.

Another MapReduce example:

word count

doc ID	contents
doc1	"my", "dog", "spot"
doc2	
doc3	

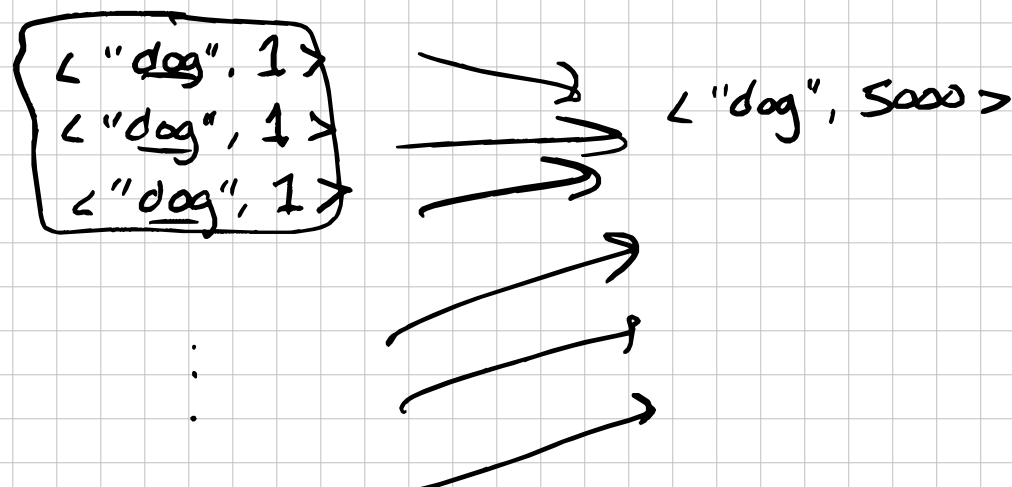
$\langle \text{"dog"}, 1 \rangle$ map function

For document in documents:

for word w in document:

emit $\langle \text{word}, 1 \rangle$

intermediate k-v pairs



could have combined these intermediate k-v pairs locally on each map worker
 $\langle \text{"dog"}, 3 \rangle$

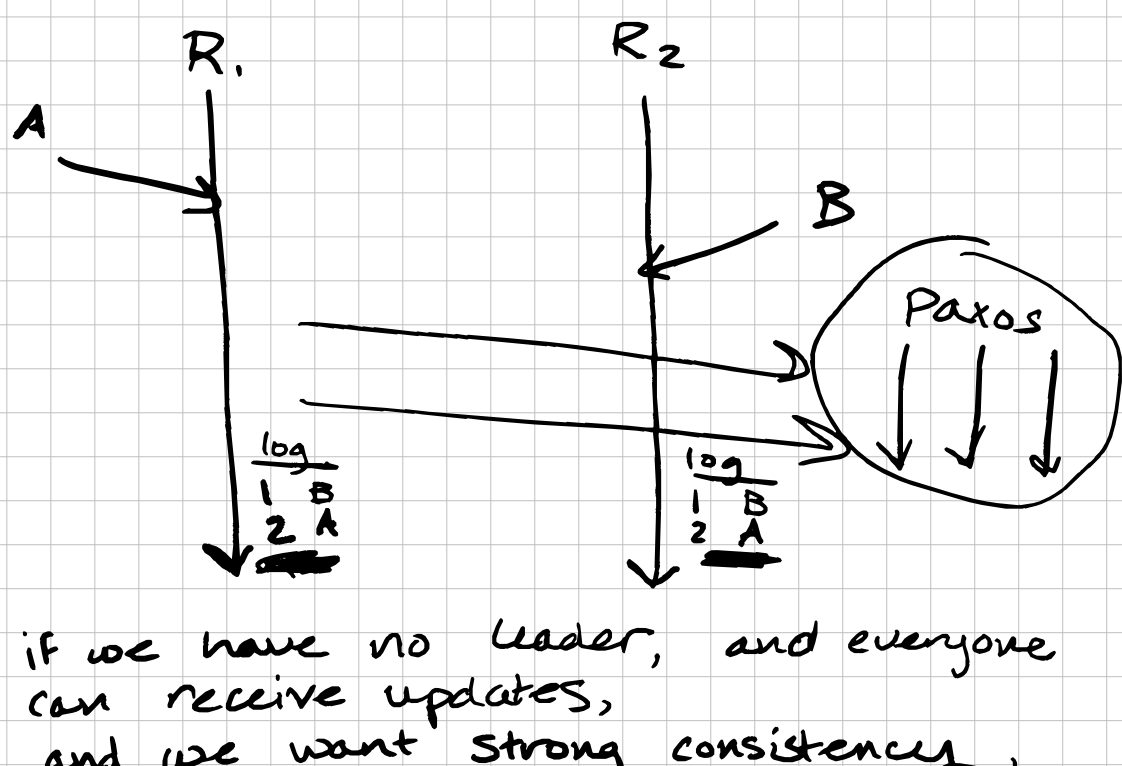
This is using a combine function.

$$(a+b) + c = a + (b+c)$$

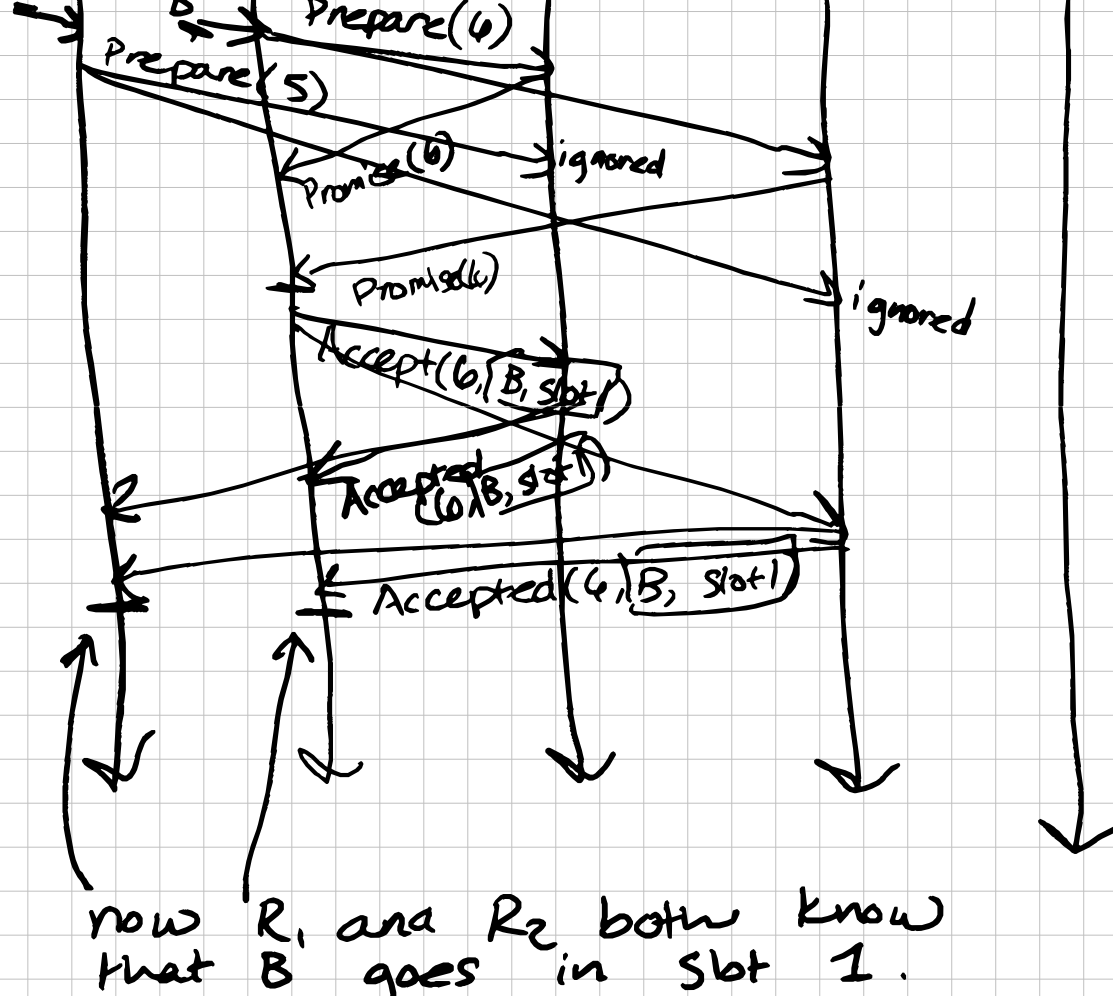
← associativity of addition

As long as the operation that reduce workers would do is associative, it's OK to do some of it in advance on the map workers, using a combine function.

The math behind replica conflict resolution



if we have no leader, and everyone can receive updates, and we want strong consistency,



now R1 and R2 both know that B goes in slot 1.

Consensus protocols are expensive!

often, you just need

Strong convergence:

replicas that have delivered the same set of updates have equivalent state.

$\{\square, \square\}$ $\{\square, \square\}$

Remember partially ordered sets?

A partially ordered set (poset) is

a set S ,

together with a binary relation \subseteq

relating elements of S .

$S =$ subsets of a set $\{\Delta, \emptyset, \square\}$

with \subseteq (set inclusion) as the relation

$\{\Delta\} \subseteq \{\Delta, \emptyset\}$

some elements of S are comparable

$\{\Delta\} \not\subseteq \{\emptyset\}$

$\{\emptyset\} \not\subseteq \{\Delta\}$

some elements of S aren't comparable.

Also, some axioms have to hold for our relation \subseteq :

- reflexivity: $\forall a \in S, a \subseteq a$.

- antisymmetry: $\forall a, b \in S,$

if $a \subseteq b$

and $b \subseteq a,$

then $a = b$.

- transitivity: $\forall a, b, c \in S,$

if $a \subseteq b$ and $b \subseteq c,$

then $a \subseteq c$.

$S = \{\emptyset, \{\emptyset\}, \{\Delta\}, \{\square\}, \{\emptyset, \Delta\}, \{\Delta, \square\}, \{\emptyset, \square\}, \{\emptyset, \Delta, \square\}\}$

All that was review.

But, something new: upper bounds.

pick 2 elements of S :

e.g. $\{\emptyset\}, \{\Delta\}$

and ask: which elements of S are at least as "big" as both of those, according to the ordering relation?

in this case: it's $\{\emptyset, \Delta\}$ and $\{\emptyset, \Delta, \square\}$. those are the upper bounds of $\{\emptyset\}$ and $\{\Delta\}$.

\emptyset and $\{\Delta\}$ have

$\{\{\Delta\}, \{\emptyset, \Delta\}, \{\emptyset, \Delta, \square\}, \{\Delta, \emptyset, \square\}\}$ as upper bounds.

Given a poset $(S, \subseteq),$

an upper bound of $a, b \in S$

is an element $u \in S$

such that $a \subseteq u$ and $b \subseteq u$.

An upper bound u of $a, b \in S$

is the least upper bound (lub),

or join,

if $u \subseteq v$ for each upper bound v of a and b .

$(u, v \in S)$

The lub of $\emptyset, \{\Delta\}$ is $\{\Delta\}$.

The lub of $\{\square\}$ and $\{\square\}$

is $\{\square, \square\}$.

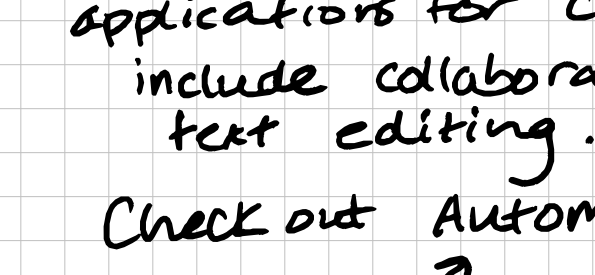
A poset where every two

elements have a lub is

called a join-semilattice.

Here's one that isn't a join-semilattice

a Boolean register



$S = \{\text{empty}, T, F\}$

$\text{empty} \subseteq T,$

$\text{empty} \subseteq F,$

$T \subseteq T$

$F \subseteq F$

$\text{empty} \subseteq \text{empty}.$

But there's no element here

that is at least as "big" as

both true and false!

so this is not a join-semilattice.

What does any of this have

to do with CSE138?!

informal claim:

if the states that a replica can take on are elements of a join-semilattice, then there is a natural way of resolving conflicts between replicas, without consensus. Just use the lub!

Conflict-free replicated

data types (CRDTs)

are distributed data structures

that exploit this idea.

* Marc Shapiro, et al. (2011)

applications for CRDTs

include collaborative

text editing.

Check out Automerge!

↗ Martin Kleppmann's

(and others)

work.

local-First software

* "merge what you can, fork what you can't"