

Lecture #4.

Replicated Directories (Distributed Dictionary)

Fischer & Michael

- maintain a database consisting of a dictionary, i.e., a set of elements with two update operations

INSERT

DELETE

and a single query operation

LIST.

Insert(x): adds x to the set

Delete(x): deletes x from the set.

List: returns an enumeration of the elements.

Goal: database to be highly available even when nodes or the network is not operational.

Available: any operational node should be able to perform any of the basic operations at any time, regardless of the status of the system.

Each node maintains its copy or more formally VIEW of the database. - and all operations are performed initially on the local VIEW.

Periodically - node sends its VIEW to others. A node receiving such information then updates its own VIEW.

SEND(m) + receive(m)

⇒ eventually everyone converges to the global VIEW.

Correctness condition:

an element $x \in \text{View}[i]$ iff i knows of its insertion but does not know of its deletion.

2 Restrictions:

R1. At most one occurrence of $\text{insert}(x)$ for each x . (once x has been deleted it can never be reinserted).

R2. $\text{delete}(x)$ is only legal at a node j if $x \in \text{View}[j]$.

Discuss R1.

- ensures that $\text{insert}(x) \rightarrow \text{delete}(x)$
- if know about both $\text{ins}(x)$ + $\text{del}(x) \Rightarrow x \notin \text{View}$.
- associate a timestamp/unique-id
 - "foo" \rightarrow foo
 - "foo" \rightarrow foo' } control point = view
- multiple deletions of the same x are permitted.

Algorithm

$$OP = \{ \text{insert}(x), \text{delete}(x) \mid x \in D \} \cup$$

$$\{ \text{List} \} \cup$$

$$\{ \text{send}(m), \text{delete}(m) \}.$$

an event e

$$op(e) \in OP.$$

$$\text{node}(e) \in \{1, 2, \dots, N\}.$$

E is the set of all events e .

$$D(E) = \{ x \in D \mid op(e) = \text{insert}(x) \text{ for some } e \in E \}.$$

E is a partial order as defined by the logical clocks.

$x \in \text{view}(e')$ iff

V1: $\exists e$ s.t. $e \rightarrow e'$ and $\text{op}(e) = \text{insert}(x)$, and

V2: $\forall e$ s.t. if $\text{op}(e) = \text{delete}(x)$, $e \not\rightarrow e'$.

A Naive Solution.

Each node i : $\underbrace{I_i}_{\text{inserts}}$ and $\underbrace{D_i}_{\text{deletes}}$.

$$\Rightarrow V_i = I_i \setminus D_i$$

Send(m)@ j : send($\langle I_j, D_j \rangle$)

receive($\langle I_m, D_m \rangle$)@ i :

$$I_i = I_i \cup I_m$$

$$D_j = D_i \cup D_m.$$

Drawback:

- Unbounded size of I_i , + D_i

- List operation involves expensive set operation.

Ideally - maintain V_i at all times.

But \rightarrow can't replace $V_i \cup V_m$ by $V_i \cup V_m$

why?

$x \in V_i \cup V_m$ may be missing from V_k (either V_i or V_m)

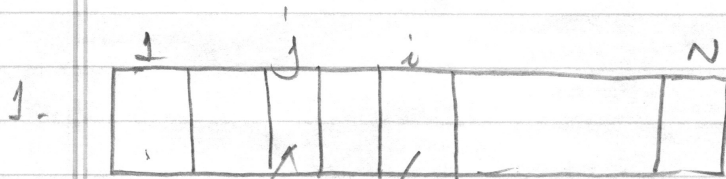
1. x used to be in V_k but it has been deleted, or

2. x was inserted so recently that node k has not yet heard about it?

Case 1 x should not be in the resultant view

Case 2 x should be in the resultant view.

To distinguish this each node maintains:



clock_i: logical clock @ i

i's view of the insertions @ j.

2. each x is tagged with (c_{rex}, T_x)

↑
id of the node
where the insertion
took place

$del(V, T, x) \text{ iff } \{x \notin V \wedge T_x \leq T[c_{rex}]\}$

↑
the timestamp

Insert(x): $T_x = \text{clock}[i]++$

$c_{rex} = i$

$V_i = V_i \cup \{ \langle x, c_{rex}, T_x \rangle \}$

Delete(x): $V_i = V_i \setminus x$;

List : return V_i

Send(m): send $\langle V_i, T_i \rangle$

receive (m) :

let m be $\langle \bar{V}, \bar{T} \rangle$.

$$V_i = \{ x \in (V_i \cup \bar{V}) \mid$$

$$\sim \text{del}(V_i, T_i, x) \wedge \sim \text{del}(\bar{V}, \bar{T}, x) \}$$

$$\forall k \quad T_i[k] = \max(T_i[k], T[k])$$

Initially

$$\forall i, j \quad T_i[j] = 0$$

$$V_i = \emptyset$$