

Tema 6 - Confiabilidad en Sistemas Distribuidos

Felipe Ortega, Gorka Guardiola

GSyC, ETSIT. URJC.

Sistemas Distribuidos (SD)

24 de noviembre, 2020



Disponibilidad

- ▶ Medida de cuanto tiempo ininterrumpido se da el servicio.
- ▶ Tiempo medio entre fallos (MTBF - Mean time between failures).
- ▶ Tiempo medio de reparación (MTTR - Mean time to repair).
- ▶ Tiempo medio hasta fallo (MTTF - Mean time to failure), si no se puede reparar.

Tiempo medio de reparación

- ▶ Tiempo de detección.
- ▶ Tiempo de respuesta - entre detección y diagnóstico.
- ▶ Tiempo de recuperación.

6.3 Seguridad

- ▶ Importante para la confiabilidad
- ▶ Uso no autorizado, acceso no autorizado, datos falsos, etc.
- ▶ No se ve en esta asignatura, pero no se puede olvidar

6.4 Tolerancia a fallos

Complejidad

6.5 Algoritmos de consenso

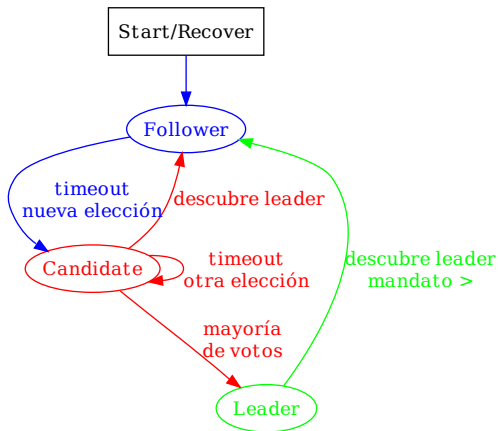
Raft

- ▶ Similar a Paxos (multi-decree).
- ▶ Diseñado alrededor de un log distribuido.
- ▶ Combina el log distribuido con un sistema de elección del líder (similar al algoritmo del dictador).
- ▶ Mucho más sencillo, bien especificado.
- ▶ Basado en timeouts y aleatorización (no totalmente asíncrono).
- ▶ Creado por Diego Ongaro y John Ousterhout (es la tesis doctoral del primero).

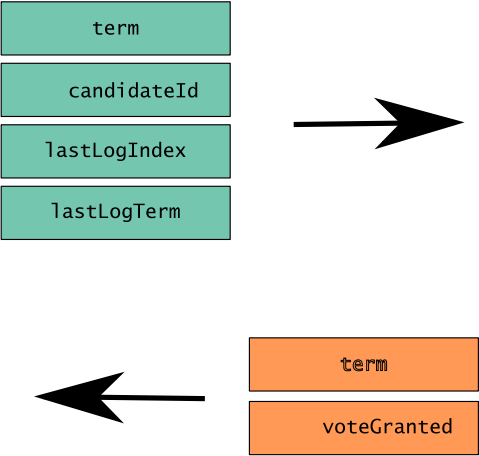
Propiedades de Raft

- ▶ No trivialidad: al log sólo van valores propuestos.
- ▶ Consistencia: Sólo va un valor para cada índice.
- ▶ Viveza: Si se propone un valor, tarde o temprano va al log de una mayoría de miembros.
- ▶ <http://thesecretlivesofdata.com/raft>.

Sistemas Distribuidos (SD)



Fase inicial de Raft: requestVotes



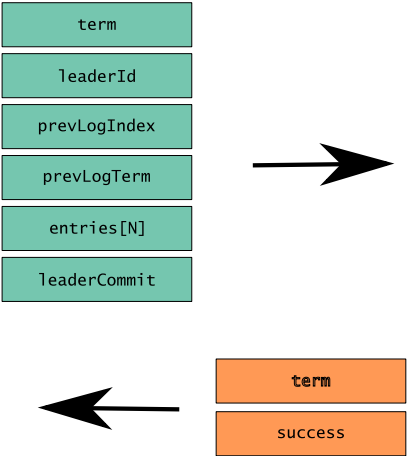
Fase normal de Raft

- ▶ Tengo un Leader.
- ▶ Periódicamente manda un AppendEntries para avisar de que está vivo.
- ▶ El AppendEntries puede tener valores para el log de los Followers.

Fase normal de Raft

- ▶ De un valor no se hace commit en el log.
 - ▶ Del servidor hasta que no han contestado una mayoría absoluta de Followers.
 - ▶ Del cliente hasta que no avanza el índice del qué ha hecho commit el servidor (un campo de `AppendEntries`).
- ▶ Cuando se hace commit, ese valor ya no debería cambiar, se avanza la máquina de estados.

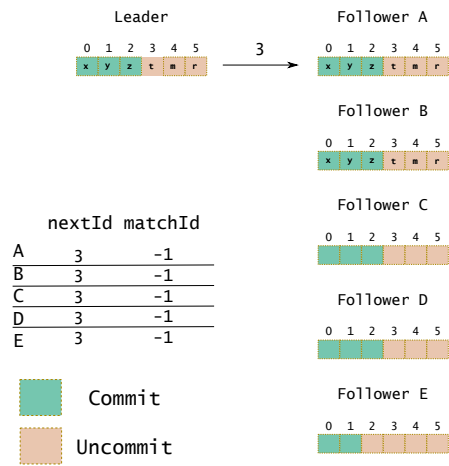
Fase normal de Raft, AppendEntries



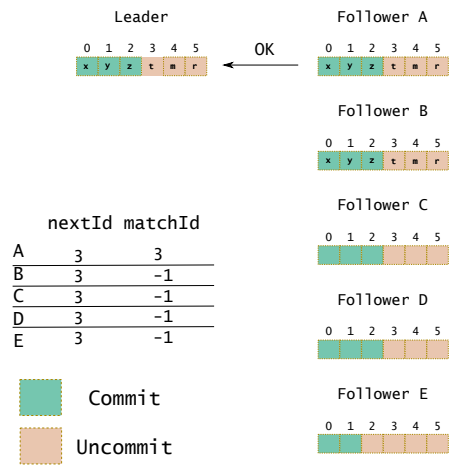
Fase normal de Raft

- ▶ Ojo, en lo que sigue está simplificado.
- ▶ El mensaje de `AppendEntries` y el `heartbeat`.
- ▶ Son lo mismo y llevan todo.
- ▶ `AppendEntries` se manda periódicamente, puede que sin entries (`heartbeat`).
- ▶ Dibujamos los mensajes a los `Follower` secuencialmente, en realidad van en paralelo.

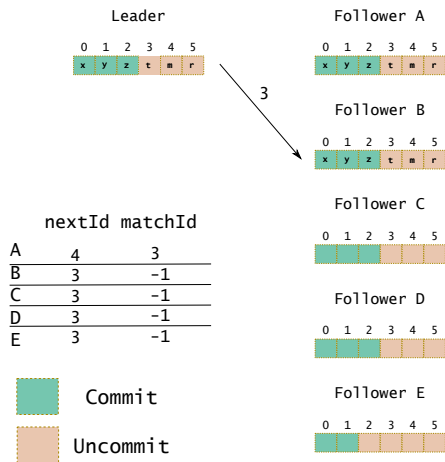
Fase normal de Raft, log



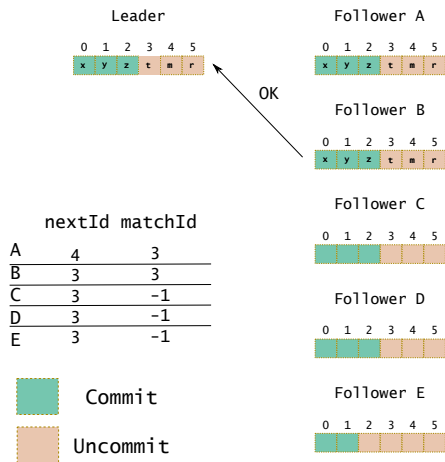
Fase normal de Raft, log



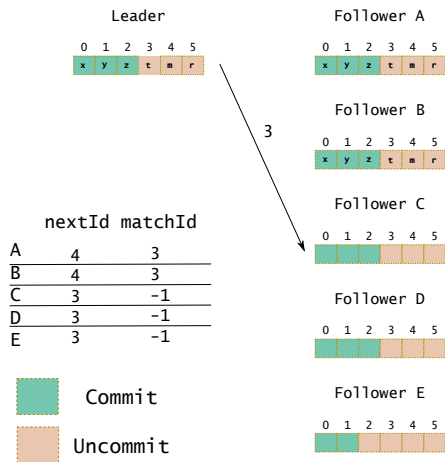
Fase normal de Raft, log



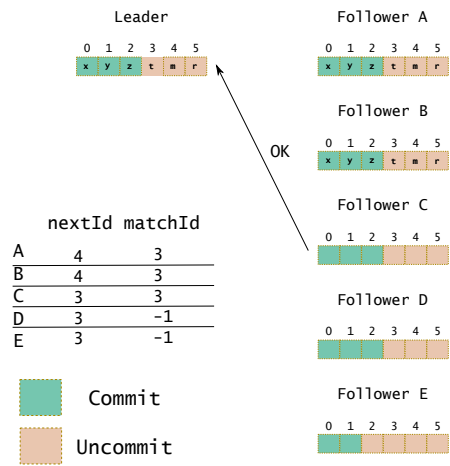
Fase normal de Raft, log



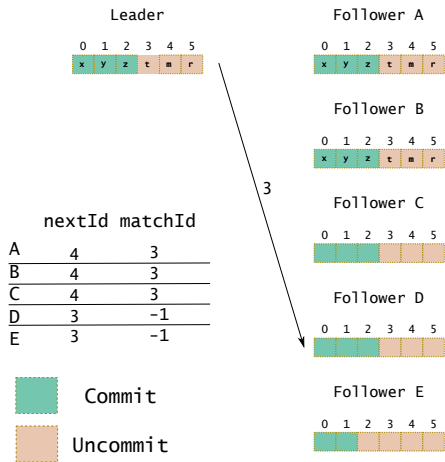
Fase normal de Raft, log



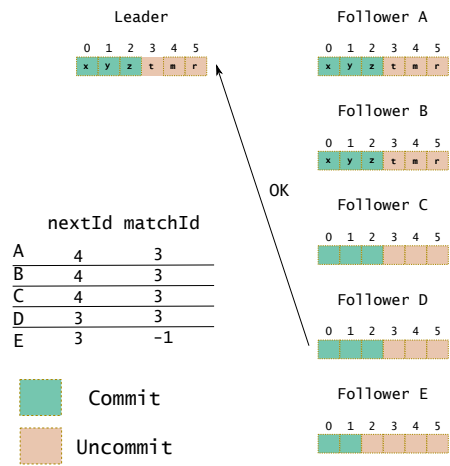
Fase normal de Raft, log



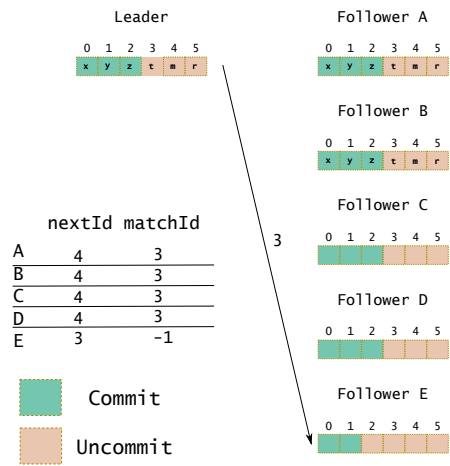
Fase normal de Raft, log



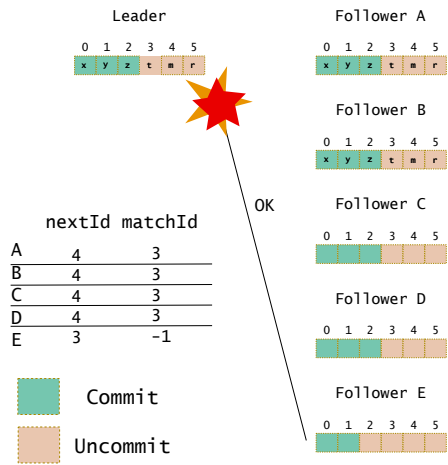
Fase normal de Raft, log



Fase normal de Raft, log



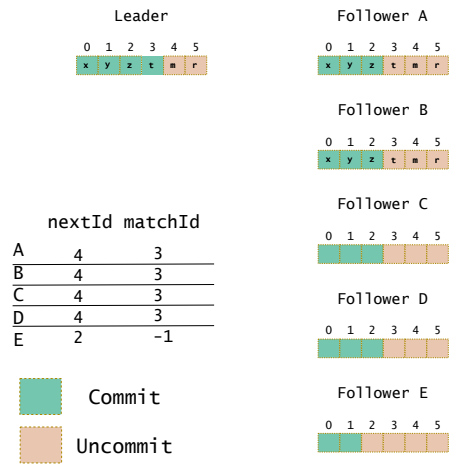
Fase normal de Raft, log



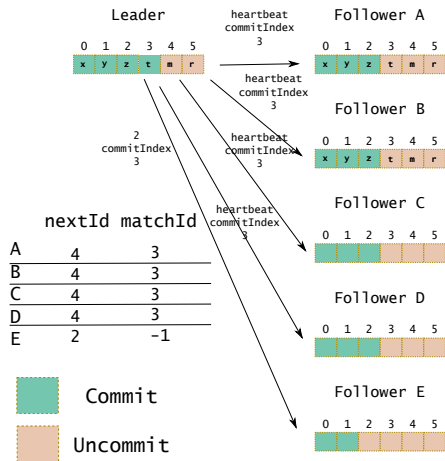
Fase normal de Raft, log

- ▶ Aunque se ha perdido uno, hay mayoría, aumento mi commit.
- ▶ De ese no me contestan así que mando con un índice menos (si no me contestan o me contestan error).
- ▶ Si me contestan error es porque el prefijo del log no coincide.

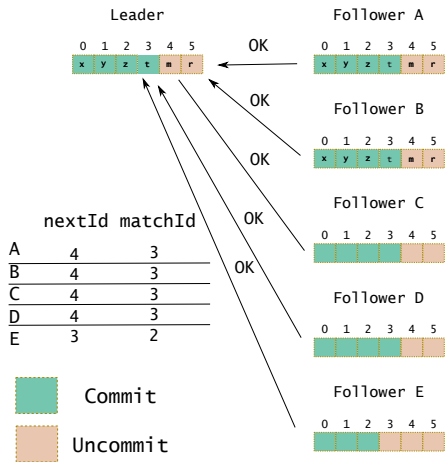
Fase normal de Raft, log



Fase normal de Raft, log



Fase normal de Raft, log



Log

- ▶ El log es lo más delicado de Raft.
- ▶ En teoría cada entrada es un cambio de estado en una máquina de estados.
- ▶ Hacen falta ciertas garantías para que funcione.

Log

- ▶ Cuando se manda `AppendEntries` a un Follower se indica la entrada anterior.
- ▶ Se intenta garantizar que hay un prefijo común entre el Leader y los Followers.
- ▶ El Leader lleva cuenta de qué le ha mandado a cada uno `nextId` y `matchId`.
- ▶ Empieza desde lo que se hizo commit en él, si hay un error, retrocede.

Raft: cuando falla

- ▶ Varios Candidate entran en un duelo.
- ▶ Piden votos a diferentes grupos.
- ▶ Entran en un bucle en el que vas teniendo muchas elecciones.
- ▶ La aleatorización hace que (finalmente) converja a un sólo líder.
- ▶ Hay que elegir bien los timeouts y depende de la red.
- ▶ Si no hay suficientes réplicas para mayoría (error duro).
 - ▶ En este caso, no avanza nunca el líder.

Flp: imposibilidad de consenso con fallos bizantinos

- ▶ En un sistema completamente asíncrono.
- ▶ Hay un teorema que dice que es imposible.
- ▶ Lo que dice es que **ningún algoritmo garantiza viveza**.
- ▶ Es imposible distinguir un fallo de un tarda mucho en contestar.
- ▶ En la práctica, no es tan importante (timeouts, como en Raft).
- ▶ Ojo, los timeouts dependen del tipo de red.

6.6 Aplicaciones en la industria

Paxos en la industria: Zookeeper

- ▶ Zookeeper, sistema de coordinación de servidores distribuidos, no implementa exactamente Paxos, sino Zab, muy similar.
- ▶ Parte de Hadoop, implementación abierta similar a GFS.
- ▶ Implementa Zab, algoritmo de consenso similar a Paxos (no exactamente igual).
- ▶ Similar a un Sistema de Ficheros (jerárquico, znodes, ficheros y directorios).
- ▶ Elección de líder.
- ▶ Propagación de estado (hay un orden total establecido).
- ▶ Eventos de cambios de estado (znodes).
- ▶ Lento.

Bibliografía I



[Birman, 2015] Birman, Kenneth P.
Reliable distributed systems: technologies, web services, and applications.
Springer, 2005.



[Lamport, 2001] Lamport. L.
Paxos made simple.
ACM SIGACT News, 32.4 (2001): 18-25.



[van Renesse, 2011] van Renesse, R.
Paxos Made Moderately Complex.
Cornell University, 2011.



[van Steen & Tanenbaum, 2017] van Steen, M., Tanenbaum, A. S.
Distributed Systems.
Third Edition, version 01. 2017.