

SISTEMAS DISTRIBUIDOS

PRÁCTICA 3

874055 - Ariana Porroche Llorén

871627 - Irene Pascual Albericio

ÍNDICE

Introducción.....	1
Diseño de raft.....	1
Diagramas de secuencia:.....	5

Introducción

En esta práctica, construiremos un sistema de almacenamiento de datos que opere de manera distribuida y sea resistente a fallos. Para lograrlo, utilizaremos el algoritmo Raft, que nos garantiza que todos los nodos en la red estén siempre alineados respecto al estado de los datos, incluso si algunos de ellos fallan.

Raft organiza los nodos de tal manera que uno de ellos actúa como líder, quien toma decisiones y se asegura de que los demás sigan el mismo "plan". En esta primera parte de la práctica, implementaremos el proceso de elección de líder: cada nodo puede proponerse como líder y, si cuenta con el apoyo de la mayoría, asume el mando. Una vez elegido, el líder es responsable de mantener a todos sincronizados mediante el envío de latidos de confirmación (heartbeat) para que los demás nodos sepan que sigue activo.

Además comprobaremos que en caso de que caiga un líder, los demás nodos proponen una nueva elección y eligen a otro líder para continuar con la ejecución.

Esta práctica es el primer paso hacia un sistema robusto, capaz de operar incluso ante problemas. En la siguiente práctica, nos centraremos en someter las operaciones y en mejorar aún más la tolerancia a fallos.

Diseño de raft

A continuación vamos a explicar cómo funciona un nodo en Raft, gracias a explicar los tipos de datos y funciones clave que hemos implementado. Estos datos y operaciones las llevan a cabo cada nodo para acordar un líder, comunicarse entre ellos, informar de desactualizaciones a otros nodos...

Tipos de Datos Principales

- **TipoOperacion:** este tipo define las operaciones básicas que el nodo puede realizar, como "leer" o "escribir". También incluye una clave y un valor asociados, dependiendo de si se trata de una lectura o escritura de datos.
- **AplicaOperacion:** indica una operación específica que debe llevarse a cabo en el sistema y registra su posición en el registro.
- **TipoEstadoNodo:** Este tipo determina el rol del nodo, es decir, si es líder, candidato o seguidor.
- **TipoEstadoVital:** este tipo determina si el nodo se encuentra Vivo o Caído en dicho momento.

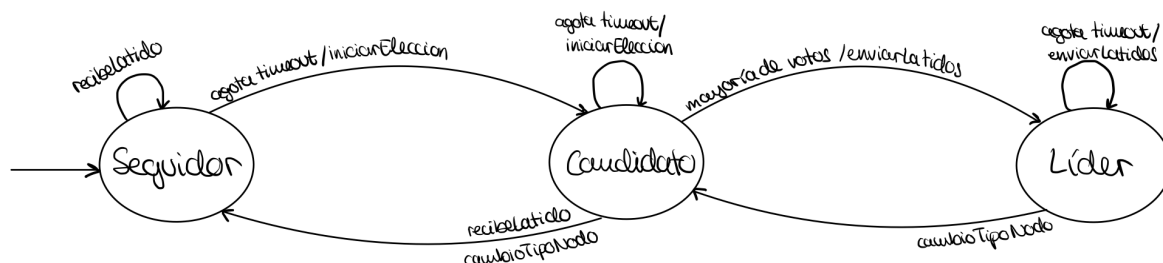
- **TipoReg:** es el registro de cada entrada del nodo, que almacena el mandato y la operación.
- **NodoRaft:** es el componente nodo en sí, y almacena las siguientes variables:
 - un mutex Mux para garantizar el acceso a las variables en exclusión mutua
 - un vector Nodos con las direcciones ip:puerto de cada réplica del algoritmo raft
 - una componente Yo que guarda el identificador del propio nodo
 - un IdLider que almacena el identificador del líder en curso
 - el MandatoActual
 - el identificador del nodo al que ha votado en la elección VotadoA
 - el EstadoActual del nodo (si es Seguidor, Candidato o Líder)
 - el EstadoVital del nodo (si está Vivo o Caído)
 - el registro Reg que almacena las entradas con las operaciones y el mandato
 - un IndiceReg que indexa la última componente del registro Reg
 - finalmente, tenemos 4 canales para los distintos aspectos de sincronización y cambio de estados del nodo: CanalAplicarOperacion, CanalVotos, CanalCambioEstadoActual y CanalLatidos
- **ArgsPeticionVoto:** es la estructura que se manda a los nodos para solicitar su voto. Enviamos el mandato y el identificador del candidato que será el de nuestro nodo.
- **RespuestaPeticionVoto:** es la estructura que recibimos tras enviar la petición de voto. Recibimos si nos han concedido o no el voto, y el mandato del nodo al que hemos enviado la petición, para el caso en el que estuviésemos desactualizados.
- **ArgAppendEntries:** es la estructura que enviamos para mandar un latido o para comprometer una entrada del registro. Tiene las componentes mandato, identificador del líder y las entradas a someter. En el caso de enviar un latido, estas entradas estarán vacías. En la próxima práctica tendremos más componentes en este registro.
- **Results:** es la estructura recibida tras haber enviado ArgAppendEntries. Nos dice si ha tenido éxito la acción y el mandato del nodo al que ha sido enviada.

Operaciones:

- **NuevoNodo:** crea y configura un nuevo nodo en raft, iniciandizándolo y poniendo en marcha la máquina de estados de los nodos con una goroutine.
- **para:** sirve para desactivar un nodo raft (es decir, la simulación de la caída de un nodo).
- **obtenerEstado:** devuelve el estado actual de un nodo (índice del nodo, mandato actual, si es líder o no, índice del líder).
- **someterOperacion:** si el nodo no es el líder, devolverá un error ya que sólo el líder puede someter una operación. En caso de que sea el líder, añadirá la operación a su registro, actualizará el índice y devolverá los argumentos necesarios.
- **PedirVoto:** es la llamada RPC en la que le llega una solicitud de otro nodo para que le vote. Compara los mandatos y concede el voto o no según si el nodo solicitante está actualizado o si somos nosotros los que nos tenemos que actualizar.
- **AppendEntries:** esta llamada RPC, de momento, si recibe un latido del líder y está todo correcto, envía una señal a su canal de latidos para reiniciar el timeout. Si está desactualizado, se actualiza, y si el líder es el desactualizado, le manda actualizarse. En la siguiente práctica, gestionaremos el hecho de recibir una entrada a aplicar.
- **enviarPeticionVoto:** es la llamada RCP que envía PedirVoto al nodo indicado en los argumentos. En función de la respuesta, si le han concedido el voto, enviará una señal por el canal de votos. Si no le han concedido el voto y además está desactualizado, se actualiza y pasa a ser seguidor.
- **maquinaEstados:** es la máquina de estados de los nodos del algoritmo, muy similar a la vista en clase. Se encarga de comprobar el estado del nodo y redirigir a su máquina de estados correspondiente.
- **maquinaEstadosSeguidor:** el seguidor programa su timeout y se queda esperando a que le llegue un latido o se agote el tiempo. Si le llega un latido, resetea el timer. Pero si se agota el tiempo, quiere decir que no hay un líder y pasa a ser candidato.
- **maquinaEstadosCandidato:** el candidato programa el timer y comienza una elección. Si se agota el timer y no ha recibido respuesta o la mayoría de los votos, vuelve a iniciar otra elección, incrementando el mandato y reseteando los votos y el timer. Si recibe un latido, es que hay otro líder en curso, entonces pasaría a ser seguidor. Si recibe un voto, comprueba si tiene la mayoría, y en ese caso pasaría a ser líder. Si recibe que tiene que cambiar de estado, acabaría y volvería a la máquina de estados global para ser redirigido.
- **maquinaEstadosLider:** el líder envía cada cierto tiempo un latido a todos los nodos menos a él, para que el resto sepan que sigue vivo. Cuando se agota este tiempo, vuelve a enviar los latidos. Si le llega que debe cambiar de estado, acaba su ejecución y vuelve a la máquina de estados global para ser redirigido.

- **iniciarEleccion:** el candidato inicia una elección y envía la petición a todos los nodos vivos excepto a él mismo, lanzando una goroutine de enviarPeticonVoto.
- **enviarLatidos:** manda a todos los nodos menos a él mismo un latido mediante la goroutine enviar1Latido (ArgAppendEntries con las entradas vacías).
- **enviar1Latido:** envía un latido a un nodo en concreto mediante la llamada CallTimeout para que ejecute AppendEntries. En caso de que haya errores, o el mandato del líder esté desactualizado, el líder pasará a ser seguidor.

Primeramente, queremos presentar la **máquina de estados** seguida en nuestro programa. Más concretamente, las acciones realizadas por los nodos en función de su estado actual, las cuales están explicadas anteriormente.



```

sequenceDiagram
    participant test1 as test1-test.go
    participant ssh as ssh...go
    participant rpc as rpcclient.go
    participant main as main.go
    participant raft as raft.go

    Note over test1: solitaire.ParallelTest
    Note over test1: startDistributedProcesses()

    test1->>ssh: (ExecMultipleTest(cwd, wdir, cwd))
    ssh->>main: go execTest(hostName, cwd, cwd)
    ssh->>main: conn := Dial(hostName: 27)
    ssh->>main: session := conn.NewSession()
    ssh->>main: executeCmd(hostName, cwd, session)
    ssh->>main: session.Run(cwd)

    main->>raft: nr := NewNode(nodes, me, cwd, plier)
    raft->>raft: nr := initialize attributes
    raft->>raft: go raftProtoCall(nr)
    raft->>main: rpc.Register(nr)
    raft->>main: net.Listen(42.468.3.9:25730)
    raft->>main: rpc.Accept()

    test1->>ssh: 3 replicas
    ssh->>main: compId := ExtractPowers(0,0, false, -1)
    ssh->>main: id, wdir, eId, idId := ExtractPowers(0)
    ssh->>main: callTimedOut("nodeTest okNewBenchNode", 11, 1, reply, 2ms)

    main->>raft: client := Dial(42.468.3.9:25730)
    main->>raft: done := client.Go("nodeRaft okNewBenchNode", 11, 1, reply, cancel)

    raft->>main: reply, 11ms, reply: NewNodeNotId
    raft->>main: reply, 11ms, reply: 11ms, okNewBenchNode()
    raft->>main: return, 11ms, reply: NewNodeNotId, 11ms, 11ms, 11ms

    main->>ssh: return reply, 11ms, reply: NewNodeNotId, reply: 11ms, reply: 11ms
    ssh->>main: return call Error
    ssh->>main: stopDistributedProcesses()
  
```

The diagram illustrates the workflow of a distributed system test. It begins with a test runner (test1-test.go) initiating a parallel test and starting distributed processes. This leads to a series of calls between the test runner, a shell script (ssh...go), and a client (rpcclient.go). The client then interacts with a server (main.go), which in turn communicates with a raft node (raft.go). The raft node performs internal operations like initializing attributes and making protocol calls. The process then moves to a replication phase (3 replicas), where the client connects to the raft node and sends a request. The raft node returns a reply, which is then processed by the client and the test runner. Finally, the distributed processes are stopped.

test1.test.go	ssh...go	raftwast.go	main.go	raft.go
<pre> eleggirPrimarLiderTest2 startDistributedProcesses() ... primeiraLider(3) </pre>				
<pre> to intentos 3vozes - , mandando, elider, ... := obtenerEstadoRauchto(i) callTimeout("hacerTest obtenerEstadoRauchto", 10, &reply, 20ms) </pre>				
		<pre> client := Dial("192.168.3.9:29730") done := client.Go("underRaft obtenerEstadoRauchto", 10, &reply, cancelCall) </pre>		
				<pre> reply Estado, reply MandadoRauchto, reply Elider, reply IdLider = m. obtenerEstadoRauchto() return m. Estado m MandadoRauchto, m Elider, m IdLider </pre>
		return call.Error		return reply
<pre> ate {elider} mapalideres[guamulato]=append(...,i) </pre>				
<pre> qto [hay más de 1 lider] Fatal error [hay 1 lider] return mapalideres[ultimoMandadoRauchto][5] </pre>				
<pre> stop DistributedProcesses ... </pre>				

test1_test.go	ssh...go	raft-ws.go	main.go	raft.go
<ul style="list-style-type: none"> ↳ followUpVorBegründenWiederTest3 ↳ startDistributedProcesses() ... ↳ probeAllLider(s) ... ↳ disconnectLider() 				
<p>3 replicas</p> <p>akt [i = csg.idLider]</p> <p>ep CallTimeWait("nodeRaft-Parallel", {i, 3, 20ms})</p> <ul style="list-style-type: none"> ↳ csg.comebacked[i] = false ↳ csg.idLider = -1 		<ul style="list-style-type: none"> ↳ client := Dial("192.168.39.255:0") ↳ done := client.Go("nodeRaft-Parallel", {i, 3, 20ms}) 		<ul style="list-style-type: none"> ↳ nr.para() ↳ nr.EstadoVital = Caido ↳ Sleep(5ms)
<ul style="list-style-type: none"> ↳ probeAllLider(s) ... ↳ stopDistributedProcesses ... 				