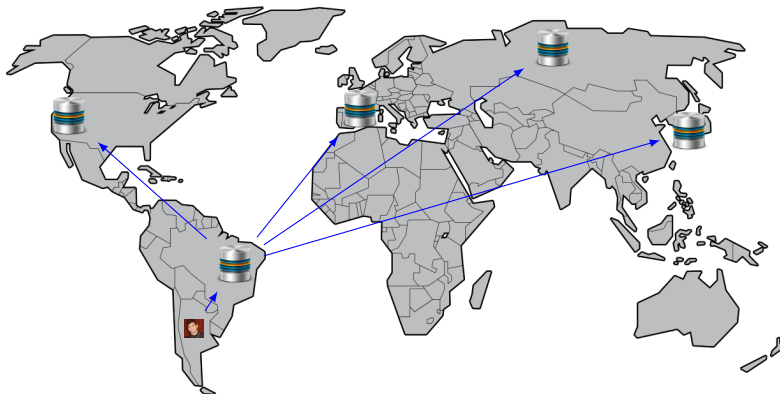


Verificación formal de protocolos distribuidos

Alejandro Naser Pastoriza

IMDEA Software Institute

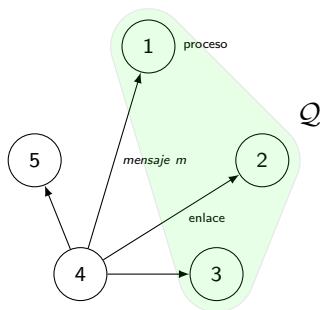




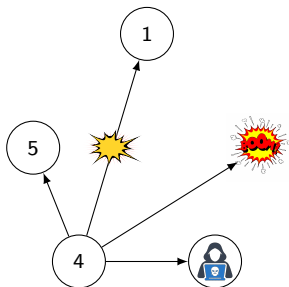
- ▶ Alta disponibilidad
- ▶ Alta confiabilidad

- ▶ Partición de datos
- ▶ etc.

Abstrayéndonos



- Objetivo: encapsular formas recurrentes de cooperación en conjunto con sus complejidades inherentes en abstracciones



- Objetivo: dotar a cada proceso con la suficiente información de manera tal que pueda continuar cooperando aún en presencia de fallos

Respecto de la correctitud...

- ▶ Las técnicas de testing, en el mejor de los casos, rasguñan la superficie de todos los posibles comportamientos
- ▶ Model checking es incompleto y no escala
- ▶ Las técnicas automáticas son limitadas por la indecidibilidad del problema subyacente
- ▶ Un enfoque más apropiado, entonces, es verificar formalmente las propiedades deseadas a través de pruebas matemáticas

Idea general

Abstracciones:

- ▶ Consenso
- ▶ Broadcast atómico

Protocolos:

- ▶ Single-Decree Paxos
- ▶ FLE Atomic Broadcast
- ▶ Vertical Paxos

Single Decree Paxos

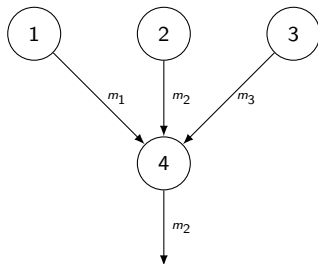
Paxos es un mecanismo para lograr acuerdo en un único valor sobre canales de comunicación no confiables

- ▶ Uno o más procesos proponen valores
- ▶ A lo más un único valor es elegido

Single Decree Paxos

Primera idea (incorrecta):

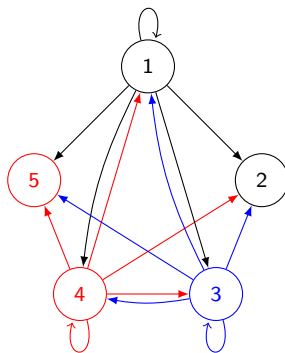
- Un único proceso acepta un valor



Single Decree Paxos

Segunda idea (incorrecta):

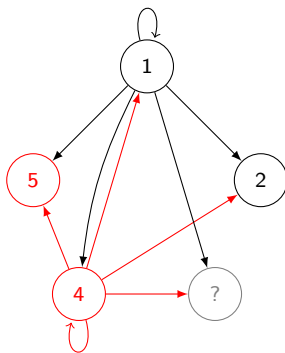
- Aceptar la primera propuesta recibida



Single Decree Paxos

Tercera idea (incorrecta):

- Aceptar múltiples propuestas

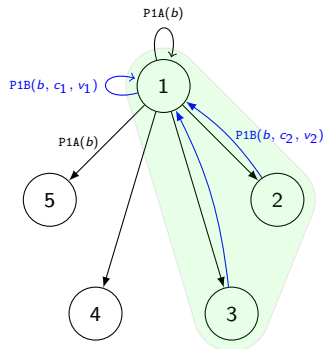


Single Decree Paxos

Fase 1:

- ▶ Proponer *ballot* b , previniendo que se acepten propuestas anteriores
- ▶ Recolectar información sobre propuestas ya aceptadas

$\{P1B(b, c_j, v_j) \mid p_j \in \mathcal{Q}\}$ de un quórum \mathcal{Q}

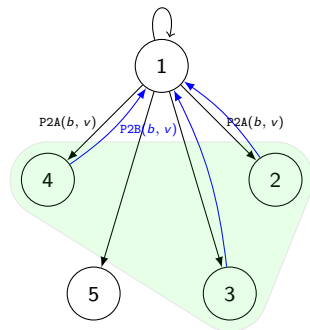


Single Decree Paxos

Fase 2:

- ▶ Si $\forall p_j. p_j \in \mathcal{Q} \Rightarrow v_j = \perp$, es libre de proponer su valor
- ▶ En caso contrario, debe proponer v_{j_0} tal que $\forall p_j. p_j \in \mathcal{Q} \Rightarrow c_j \leq c_{j_0}$

$\{P2B(b, v) \mid p_j \in \mathcal{Q}\}$ de un quórum $\mathcal{Q} \Rightarrow v$ es elegido

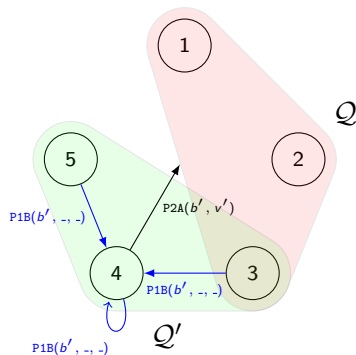


Single Decree Paxos

- ▶ Objetivo: demostrar que si v y v' son elegidos, entonces $v' = v$
- ▶ Es suficiente probar que si v es elegido en el ballot b y v' es propuesto en un ballot $b' > b$, entonces $v = v'$
- ▶ Más aún, debemos probar que si un valor v ha sido elegido o aún puede ser elegido en el ballot b y v' es propuesto en un ballot $b' > b$, entonces $v = v'$

Single Decree Paxos

- ▶ Cada $q \in \mathcal{Q}'$ se ha unido al ballot b'
- ▶ Cada $q \in \mathcal{Q}$ o bien ha recibido y aceptado $P2A(b, v)$ o su ballot es menor o igual a b
- ▶ $3 \in \mathcal{Q} \cap \mathcal{Q}'$
 - ▶ $3 \in \mathcal{Q}' \Rightarrow \text{bal} \geq b'$
 - ▶ $3 \in \mathcal{Q} \Rightarrow 3$ ha enviado $P2B(b, v)$ o tiene $\text{bal} \leq b$
 - ▶ Por lo tanto, 3 ha enviado $P2B(b, v)$ y entonces $\text{cbal} \geq b$



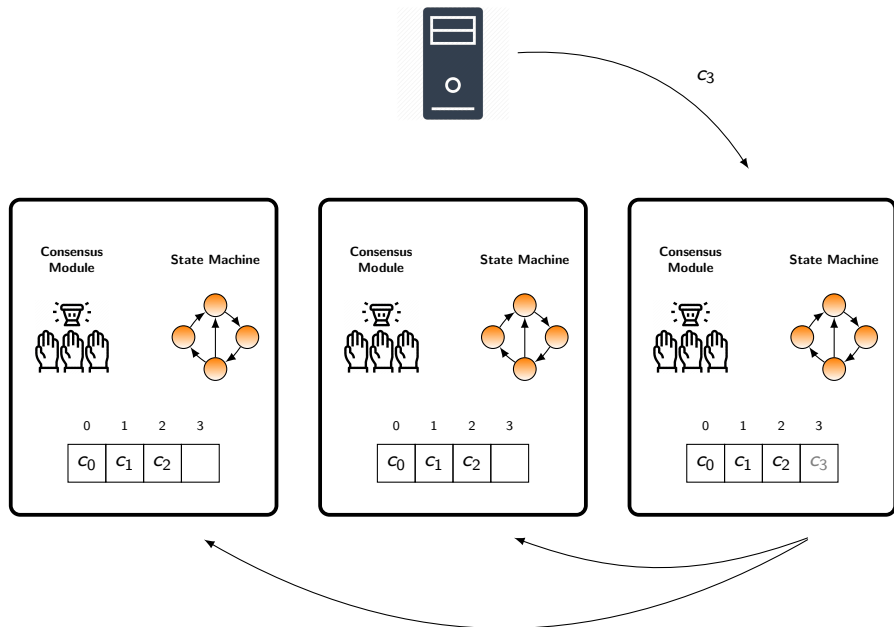
- ▶ $v' = v_{j_0}$ tal que $\forall p_j. p_j \in \mathcal{Q}' \Rightarrow c_j \leq c_{j_0}$ y por lo tanto $c_{j_0} \geq b$
- ▶ Por lo tanto $P2A(c_{j_0}, v')$ es un mensaje enviado. Si $c_{j_0} > b$, entonces (HI) $v' = v$. Si $c_{j_0} = b$, entonces $P2A(b, v)$ y $P2A(b, v')$ son enviados, por lo que $v' = v$

FLE Atomic Broadcast

Un protocolo de broadcast atómico permite a un proceso transmitir un mensaje a un grupo de procesos de manera tal que:

- ▶ Los procesos acuerdan el conjunto de mensajes entregados
- ▶ Los procesos acuerdan también el orden en que esto ocurre

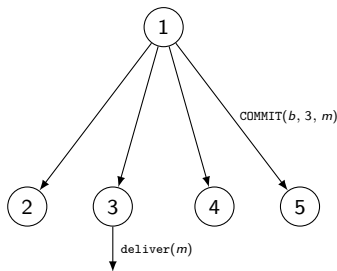
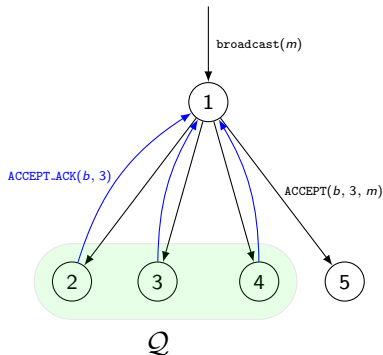
FLE Atomic Broadcast



FLE Atomic Broadcast

- Operación normal: Un único líder propone una secuencia de mensajes a sus seguidores con el objetivo de replicarlos y garantizar la durabilidad de los mensajes entregados y su orden

0	1	2	3
c_0	c_1	c_2	



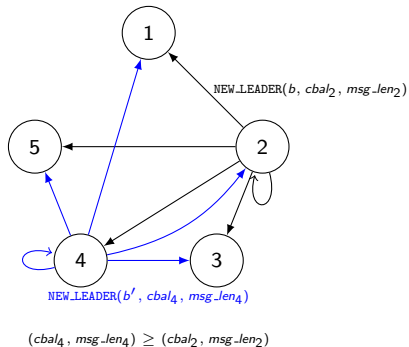
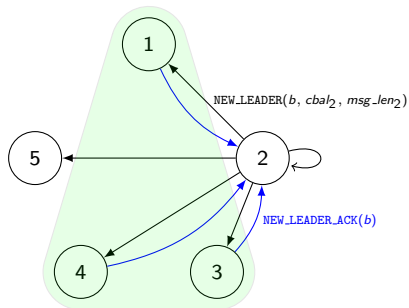
FLE Atomic Broadcast

Recuperación a fallos:

- ▶ elegir un líder quien convenza a un quórum de participar de su ballot y cuyo estado domine al quórum
- ▶ garantizar que antes de que un seguidor empiece a aceptar propuestas del nuevo líder, este ha sincronizado su estado con aquel del líder

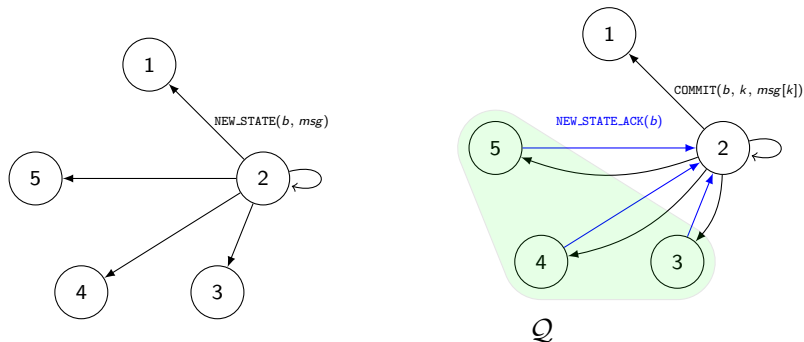
FLE Atomic Broadcast

- elegir un líder quien convenza a un quórum de participar de su ballot y cuyo estado domine al quórum



FLE Atomic Broadcast

- ▶ garantizar que antes de que un seguidor empiece a aceptar propuestas del nuevo líder, este ha sincronizado su estado con aquel del líder



FLE Atomic Broadcast

- ▶ Objetivo: demostrar que todos los procesos entregan un prefijo de una secuencia global común de mensajes
- ▶ Puesto que el k -ésimo mensaje entregado por un proceso p sólo puede ser entregado en respuesta a un mensaje $\text{COMMIT}(-, k, -)$, es suficiente demostrar que si $\text{COMMIT}(b, k, m)$ y $\text{COMMIT}(b', k, m')$ son mensajes enviados, entonces $m = m'$
- ▶ Más aún, debemos probar que si un mensaje m ha sido entregado o aún puede ser entregado en el ballot b para el slot k y $\text{COMMIT}(b', k, m')$ es un mensaje enviado con $b' > b$, entonces $m = m'$

FLE Atomic Broadcast

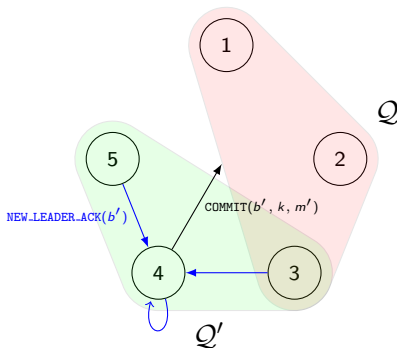
- ▶ Cada $q \in Q'$ se ha unido al ballot b'
- ▶ Cada $q \in Q$ o bien ha recibido y aceptado $\text{ACCEPT}(b, k, m)$ o $\text{NEW_STATE}(b, \text{msg})$ con $\text{msg}[k] = m$, o $q = \text{leader}(b)$, o su ballot es menor o igual a b

- ▶ $3 \in Q \cap Q'$

- ▶ $3 \in Q' \Rightarrow \text{bal} \geq b'$
- ▶ $3 \in Q \Rightarrow 3$ o bien ha recibido y aceptado $\text{ACCEPT}(b, k, m)$ o $\text{NEW_STATE}(b, \text{msg})$ con $\text{msg}[k] = m$, o $q = \text{leader}(b)$
- ▶ Por lo tanto, en 3 debemos tener $(\text{cbal}, \text{len}(\text{msg})) \geq (b, k)$

- ▶ Si $\text{cbal} > b$, entonces (HI) $\text{msg}[k] = m$. Si $\text{cbal} = b$ entonces $\text{len}(\text{msg}) \geq k$, por lo que $\text{msg}[k] = m$

- ▶ Por lo tanto si $\text{COMMIT}(b', k, m')$ es un mensaje enviado, entonces $m = m'$



Vertical Paxos

La mayoría de las soluciones por replicación requieren $2f + 1$ procesos para tolerar f fallos crash-stop. Esto es costoso: si la información se persiste en todos los procesos, sólo $f + 1$ son necesarios para garantizar durabilidad

Vertical Paxos

- ▶ El fallo de sólo un proceso bloquearía el procesamiento de mensajes, por lo que para recuperarse se requiere cambiar la membresía para reemplazar los procesos fallidos por procesos frescos
- ▶ Los procesos necesitan acordar la próxima configuración que se reduce a consenso, lo cual nuevamente requiere $2f + 1$ procesos

Vertical Paxos

Una solución consiste en utilizar un servicio de configuración independiente con $2f + 1$ procesos únicamente para realizar consenso en la configuración y sólo $f + 1$ para replicar la información relevante

Vertical Paxos

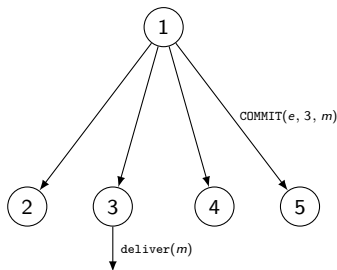
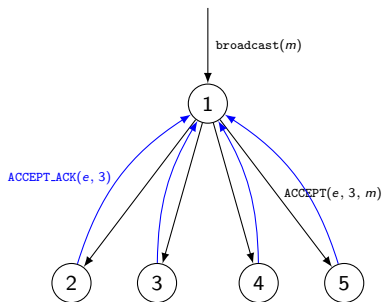
Un protocolo de broadcast atómico permite a un proceso transmitir un mensaje a un grupo de procesos de manera tal que:

- ▶ Los procesos acuerdan el conjunto de mensajes entregados
- ▶ Los procesos acuerdan también el orden en que esto ocurre

Vertical Paxos

- Operación normal: Un único líder propone una secuencia de mensajes a todos los miembros de su configuración con el objetivo de replicarlos y garantizar la durabilidad de los mensajes entregados y su orden

0	1	2	3
c_0	c_1	c_2	

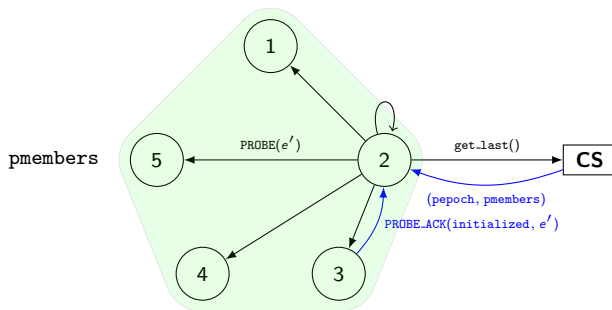


Recuperación a fallos:

- ▶ elegir un líder quien convenza a todos los miembros de una configuración de participar en su época y cuyo estado contiene todos los mensajes aceptados
- ▶ garantizar que antes de que un seguidor empiece a aceptar propuestas del nuevo líder, este ha sincronizado su estado con aquel del líder

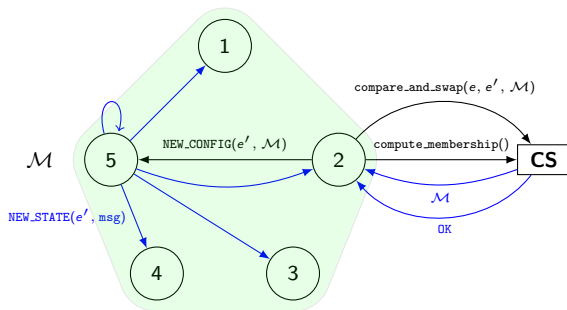
Vertical Paxos

- elegir un líder quien convenza a todos los miembros de una configuración de participar en su época y cuyo estado contiene todos los mensajes aceptados



Vertical Paxos

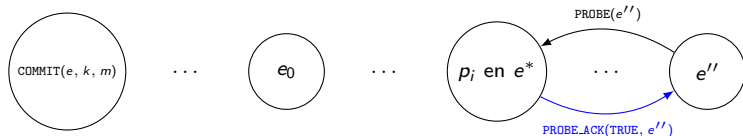
- ▶ si algún proceso responde con $\text{PROBE_ACK}(\text{FALSE}, e')$, sabemos que la época probada no puede ser operativa, por lo que probamos una época precedente
- ▶ si algún proceso (digamos 5) responde con $\text{PROBE_ACK}(\text{TRUE}, e')$, podemos garantizar que el estado de este proceso contiene todos los mensajes aceptados



Vertical Paxos

- ▶ Objetivo: demostrar que todos los procesos entregan un prefijo de una secuencia global común de mensajes
- ▶ Puesto que el k -ésimo mensaje entregado por un proceso p sólo puede ser entregado en respuesta a un mensaje $\text{COMMIT}(_, k, _)$, es suficiente demostrar que si $\text{COMMIT}(e, k, m)$ y $\text{COMMIT}(e', k, m')$ son mensajes enviados, entonces $m = m'$
- ▶ Es suficiente probar que si $\text{COMMIT}(e, k, m)$ es un mensaje enviado y en un proceso tenemos $\text{epoch} = e' > e$ entonces $\text{msg}[k] = m$

Vertical Paxos



- ▶ Todos los procesos en e se encuentran inicializados, por lo que la prueba no puede haber ido más allá de e . Por lo tanto $e \leq e^* < e''$
- ▶ Sea $e_0 = \text{epoch}$ en p_i al enviar el mensaje PROBE_ACK . Si $\text{epoch} < e$, entonces p_i no sería un miembro de e y por lo tanto no podría ser miembro de e^* . Luego $e \leq e_0$. Por otro lado, todos los miembros de la época probada e^* deben tener $\text{epoch} \leq e^* < e''$. En consecuencia $e \leq e_0 < e''$
- ▶ Si $e_0 > e$ entonces (H1) $\text{msg}[k] = m$. Si $e_0 = e$, cada proceso en e debe haber recibido y notificado su mensaje ACCEPT o NEW_STATE antes de notificar $\text{PROBE}(e'')$. Por lo tanto $\text{msg}[k] = m$

Prueba del Invariante 5. Probamos el invariante por inducción en e' . Asíumase que el invariante se cumple para cada $e' < e''$. Probamos ahora que se cumple para $e' = e''$ por inducción en la longitud de la ejecución del protocolo.

Considérese la transición en la línea [51] cuando un futuro líder p_i se incorpora a la época e'' . Mostramos que luego de esta transición tenemos $\text{msg}[k] = m$. Puesto que p_i fue elegido como el futuro líder de la época e'' , este proceso ha respondido con $\text{PROBE_ACK}(\text{TRUE}, e'')$ a $\text{PROBE}(e'')$. Por lo tanto, p_i era miembro de una configuración en una época $e' < e''$ que estaba siendo probada. La prueba finaliza cuando al menos un proceso envía un mensaje $\text{PROBE_ACK}(\text{TRUE}, e'')$. La hipótesis de que $\text{COMMIT}(e, k, m)$ es un mensaje enviado implica que en algún punto $\text{NEW_STATE}(e, \text{msg})$ es un mensaje enviado con $\text{msg}[k] = m$ y este mensaje es notificado por todos los miembros de e excepto su líder, o todos los miembros en e respondieron con ACCEPT_ACK a $\text{ACCEPT}(e, k, m)$. Esto significa que todos los miembros en e se encuentran inicializados. Esto, junto con el Invariante [2] implican que la fase de prueba no puede haber ido más allá de e . Por lo tanto, $e \leq e' < e''$.

Sea e_0 el valor de cepoch en p_i justo antes de la transición en la línea [51]. Si $e_0 < e$ entonces p_i no sería un miembro de e , pues todos los miembros de e se encuentran inicializados y tienen $\text{cepoch} = e$, y en consecuencia, por el Invariante [3] no podría haber sido elegido como el futuro líder de la época e'' . En consecuencia, $e_0 \geq e$. Por otro lado, como todos los miembros de la época e^* siendo probada tienen $\text{cepoch} \leq e^*$, debemos tener $e_0 \leq e^* < e''$. En consecuencia, $e \leq e_0 < e''$.

Si $e < e_0$, entonces por la hipótesis inductiva, tenemos $\text{msg}[k] = m$ justo luego de la transición en la línea [51] como queríamos. Asíumase ahora que $e_0 = e$. Como $\text{COMMIT}(e, k, m)$ es un mensaje enviado, entonces cada proceso en e recibe $\text{NEW_STATE}(e, \text{msg})$ con $\text{msg}[k] = m$ y responde con $\text{NEW_STATE_ACK}(e)$ o cada proceso en e recibe $\text{ACCEPT}(e, k, m)$ y responde con $\text{ACCEPT_ACK}(e, k)$. El Invariante [2] implica que p_i envió su mensaje NEW_STATE , NEW_STATE_ACK , ACCEPT o ACCEPT_ACK antes de enviar $\text{PROBE_ACK}(\text{TRUE}, e'')$. Entonces lo requerido sigue del Invariante [4].

Esto implica, en particular, que cualquier mensaje $\text{NEW_STATE}(e'', \text{msg})$ enviado con $e'' > e$ debe tener $\text{msg}[k] = m$, por lo que el resultado se cumple luego de la transición en la línea [56]. Las líneas [67] y [14] aseguran que el valor no será sobrescrito durante el caso libre de fallos del protocolo. El resultado es trivial para el resto de las transiciones. \square

Prueba del Invariante 6. Probamos el invariante por inducción en la longitud de la ejecución del protocolo. Inicialmente $\text{cbal} = 0$ para todos los procesos y no existen mensajes enviados, por lo que la afirmación es trivialmente verdadera. Para el paso inductivo, supóngase que $\text{comitable}(b, k, m)$ se cumple. El resultado sigue trivialmente de la hipótesis inductiva para todas las transiciones excepto aquellas en las líneas [37] y [17].

Considérese entonces la transición en la línea [37] por un proceso p_i la cual gestiona la recepción de mensajes $\text{NEW_LEADER_ACK}(b')$ de un quórum Q en respuesta a un mensaje $\text{NEW_LEADER}(b', \text{cbal}, \text{msgLen})$. La hipótesis implica que o bien $\text{acceptable}(b, k, m)$ o $\text{recoverable}(b, k, m)$ se cumple.

Si $\text{acceptable}(b, k, m)$ se cumple, entonces existe un quórum Q' tal que cada proceso ha recibido un mensaje $\text{ACCEPT}(b, k, m)$ y ha respondido con $\text{ACCEPT_ACK}(b, k)$ o este tiene $\text{bal} \leq b$. Sea $q \in Q \cap Q'$. Dado que $q \in Q'$ ha recibido y notificado $\text{NEW_LEADER}(b', \text{cbal}, \text{msgLen})$, este tiene $\text{bal} \geq b' > b$. Entonces, q debe haber enviado su mensaje ACCEPT_ACK antes del mensaje NEW_LEADER_ACK . Por la precondición en la línea [14] y el hecho de que si $\text{status} \in \{\text{LEADER}, \text{FOLLOWER}\}$ entonces $\text{bal} = \text{cbal}$, cuando q envió el mensaje ACCEPT_ACK , este debe haber tenido $\text{cbal} = b$. Entonces, cuando q envió el mensaje NEW_LEADER_ACK este tenía $\text{cbal} \geq b$. En consecuencia, $\text{cbal} \geq b$.

Si, por otro lado, $\text{recoverable}(b, k, m)$ se cumple, entonces existe un quórum de procesos Q' cuyos procesos p son tales que $p = \text{leader}(b)$, p ha recibido y notificado $\text{NEW_STATE}(b, \text{msg})$ con $\text{msg}[k] = m$, o este tiene $\text{bal} \leq b$. Sea $q \in Q \cap Q'$. Puesto que $q \in Q'$ ha recibido y notificado $\text{NEW_LEADER}(b', \text{cbal}, \text{msgLen})$, este tiene $\text{bal} \geq b' > b$. Entonces, q debe haber enviado su mensaje NEW_STATE o NEW_STATE_ACK antes del mensaje NEW_LEADER_ACK . Las líneas [39] y [43] implican que $\text{cbal} = b$ cuando q envió este mensaje. Entonces, cuando q envió el mensaje NEW_LEADER_ACK este tenía $\text{cbal} \geq b$. En consecuencia, $\text{cbal} \geq b$.

Puesto que p_i envió $\text{NEW_LEADER}(b', \text{cbal}, \text{msgLen})$, y por la línea [38] $\text{bal} = b'$ y $\text{status} = \text{RECOVERING}$, el Invariante [5] implica que en p_i tenemos $\text{len}(\text{msg}) = \text{msgLen}$ y $\text{cbal} \geq \text{cbal} \geq b$.

Si $\text{cbal} > b$, la hipótesis inductiva implica que $\text{msg}[k] = m$. Supóngase entonces que $\text{cbal} = b$. Por la línea [30] cuando q notificó $\text{NEW_LEADER}(b', \text{cbal}, \text{msgLen})$ este debe haber tenido $\text{len}(\text{msg}) \leq \text{msgLen}$. En el caso en que $\text{acceptable}(b, k, m)$ se cumple, q ha recibido y notificado $\text{ACCEPT}(b, k, m)$. Entonces, cuando q envió el mensaje NEW_LEADER_ACK este tenía $\text{len}(\text{msg}) \geq k$ y por lo tanto $\text{msgLen} \geq k$. Entonces, por el Invariante [4] debemos tener $\text{msg}[k] = m$. Si, por otro lado, $\text{recoverable}(b, k, m)$ se cumple, entonces existe msg tal que $\text{msg}[k] = m$ y $\text{NEW_STATE}(b, \text{msg})$ es un mensaje enviado, en cuyo caso, por el Invariante [3] msg es un prefijo de msg y por lo tanto $\text{msg}[k] = m$.

En consecuencia, el mensaje $\text{NEW_STATE}(b', \text{msg}')$ enviado durante la transición actual tiene $\text{msg}[k] = \text{msg}[k] = m$. Ningún mensaje COMMIT , ACCEPT_ACK o ACCEPT

- ▶ Estas pruebas son igualmente complejas y pueden también albergar errores sutiles
- ▶ Por lo tanto, para eliminar categóricamente errores de nuestros algoritmos, aspiramos a construir pruebas de verificación automatizadas chequeadas por máquina de sus propiedades de safety
- ▶ Utilizamos Dafny, un lenguaje de programación verifier-friendly basado en lógica de Floyd-Hoare que automatiza la verificación vía el SMT solver Z3

FIN