

Francisco Torres-Rojas, Mustaque Ahamad y Michel Raynal

Timed Consistency for Shared Distributed Objects

CARLOS MARTÍN FLORES GONZÁLEZ, Carné: 2015183528

Instituto Tecnológico de Costa Rica
Maestría en Computación
Sistemas Operativos Avanzados
Profesor: Francisco Torres Rojas, Ph.D

Ordenamiento y tiempo son dos aspectos diferentes de consistencia de objetos compartidos en sistemas distribuidos. Uno evita conflictos entre operaciones, el otro dice que tan rápido los efectos de una operación son percibidos por el resto del sistema.

Consistency Criteria: La historia global H de un sistema distribuido es un conjunto parcialmente ordenado de todas las operaciones que ocurren en todos los sitios del sistema. H_i es la secuencia de operaciones que son ejecutadas en el sitio i . Si a ocurre antes que b en H_i decimos que a precede a b en orden de programa. Se asume que todas las operaciones en H son **read** o **write**. A estas operaciones les toma un tiempo finito ejecutarse, por lo tanto, hay un intervalo que va desde el momento cuando un **read** o un **write** “inicia” al momento cuando alguna de estas operaciones “finaliza”. Sin embargo, para los propósitos de la *timed consistency*, se asocia un instante a cada operación llamada el tiempo efectivo (*effective time*) de la operación. El *effective time* de una operación corresponde a algún instante entre su inicio y su fin. Si el tiempo efectivo de a es t , se dice que a fue “ejecutado” en el tiempo t . Si D es un conjunto de operaciones, entonces la serialización de D es una secuencia lineal S conteniendo exactamente todas las operaciones de D tal que cada operación **read** para un objeto particular retorna el valor que fue escrito por la operación **write** más reciente. Sea \sim una relación de orden parcial arbitraria definida sobre D , decimos que la serialización S “respeta \sim ” si $\forall a, b \in D$ tal que $a \sim b$, pasa que a precede a b en S . La historia H satisface linealización (LIN) si hay una serialización de H que respeta el orden inducido por los tiempos efectivos de las operaciones. Serialización estricta se define sobre historias formadas por transacciones y requiere de la existencia de una serialización de H que respete el orden de tiempo real de las transacciones. Normalidad es equivalente a LIN cuando operaciones en objetos son unarias, pero es estrictamente más débil que LIN cuando las operaciones pueden abarcar varios objetos. Un nivel más débil de consistencia es ofrecido por consistencia secuencial (SC). Este modelo, no garantiza que una operación **read** retorna el valor más reciente con respecto al tiempo real sino solo que el resultado de cualquier ejecución sea el mismo como si las operaciones en todos los sitios fuera ejecutados en orden secuencial, y las operaciones de cada sitio individual aparecen en esta secuencia en el orden especificado por su programa. La historia H satisface SC si hay una serialización de H que respete el orden del programa para cada sitio en el sistema. La relación de causalidad “ \rightarrow ” para sistemas de paso de mensajes puede ser modificada para ordenar las operaciones de H . Sea a, b y $c \in H$, se dice que $a \rightarrow b$ (a precede causalmente a b) si se da una de las siguientes: (i) a y b son ejecutados en el mismo sitio y a es ejecutado antes que b , (ii) b lee el valor de un objeto escrito por a , (iii) $a \rightarrow c$ y $c \rightarrow b$. Dos operaciones distintas a y b son concurrentes si ninguna de estas condiciones se dan entre ellas. Sea H_{i+w} el conjunto de todas las operaciones en H_i más todas las operaciones **write** en H . La historia H satisface consistencia causal (CC) si para cada sitio i existe una serialización de H_{i+w} que respete el orden causal “ \rightarrow ”. Así, si $a, b, c \in H$ son tales que a escribe un valor v en el objeto X , c lee el mismo valor v del objeto X , y b escribe el valor v' en el objeto X , nunca se da el caso que $a \rightarrow b \rightarrow c$. CC requiere que todas las relaciones relacionadas con causalidad sean vistas en el mismo orden por todos los sitios, mientras que diferentes sitios pueden percibir operaciones concurrentes en órdenes diferentes. CC es un modelo de consistencia más débil que SC pero puede ser implementado eficientemente.

Timed Consistency. Ni en SC ni en CC, el tiempo real es explícitamente capturado. En SC, operaciones puede aparecer fuera de orden en relación con sus tiempos efectivos. En CC, cada sitio puede ver operaciones de escritura concurrente en diferente orden. Por otro lado, LIN requiere que las operaciones sean observadas en un orden que respete su ordenamiento de tiempo real. Timed consistency requiere que si el tiempo efectivo de un **write** es t , el valor escrito por esta operación tiene que estar visible a todos los sitios en el sistema distribuido en el tiempo $t + \Delta$, donde Δ es un parámetro de ejecución. Se puede ver que cuando Δ es 0, *timed consistency* se convierte en LIN¹. Reading on Time: en modelos basados en tiempo, el conjunto de los valores que un **read** puede retornar está restringido por la cantidad de tiempo que ha transcurrido

¹ *Timed Consistency* es una generalización (o un debilitamiento) de LIN

a partir de los **write** que le preceden. Un **read** ocurre a tiempo (*on time*) si no devuelve datos obsoletos cuando hay valores más recientes que han estado disponibles por más de Δ unidades de tiempo. Esta definición depende las propiedades del reloj usando para assignar marcas de tiempo (*timestamps*) a las operaciones en ejecución. Primero, se asume relojes físicos perfectamente sincronizados en donde $T(a)$ es el instante en tiempo real correspondiente al tiempo efectivo de la operación a . **Definición 1:** Sea $D \subseteq H$ un conjunto de operaciones y S una serialización de D . Sea $w, r \in D$ tal que w escribe un valor dentro de un objeto X que luego es leído por r (w es la operación **write** más cercana dentro del objeto X que aparece a la izquierda de r en la serialización S). Se define el conjunto W_r asociado con r como: $W_r = \{w' \in D : (w' \text{ escribe un valor dentro del objeto } X) \wedge (T(w) < T(w') < (T(r) - \Delta))\}$. Se dice que la operación r ocurre o lee a tiempo en la serialización S si $W_r = \emptyset$. S es *timed* si cada operación **read** en S ocurre a tiempo. Approximately-synchronized real-time clocks: En relojes de tiempo real aproximadamente sincronizados, resincronizaciones periódicas garantizan que no hay dos relojes difieran en más de ε unidades de tiempo. Típicamente también se garantiza que la diferencia entre cualquier reloj y el tiempo “real”² nunca es mayor que $\varepsilon/2$. De esta forma, si el tiempo efectivo de una operación a fue reportada por un sitio en particular como $T(a)$, entonces, desde el punto de vista del servidor de tiempo, este tiempo efectivo corresponde a algún instance en el intervalo $[T(a) - \varepsilon/2, T(a) + \varepsilon/2]$. Se dice que $\forall a, b \in H$, a definitivamente ocurrió antes de b si $T(a) + \varepsilon/2 < T(b) - \varepsilon/2$, o, de forma equivalente si $T(a) + \varepsilon < T(b)$. Si ni a ni b definitivamente ocurrieron antes que el otro, se se dice que sus marcas de tiempo son no comparables o concurrentes. Sea $w, r \in D \subseteq H$ tal que w escribe un valor dentro del objeto X que luego es leído por r . Dejar que $w' \in D$ actualice también el valor del objeto X . Si $T(w')$ definitivamente ocurrió antes que $T(w)$, o si $(T(r) - \Delta)$ ocurrió definitivamente antes que $T(w')$, entonces el claro que w' no afecta el hecho que r ocurra a tiempo. Ahora, si ambos $T(w')$ y $T(w)$ son concurrentes, o si $(T(r) - \Delta)$ y $T(w)$ son concurrentes, entonces no hay evidencia que w' sea más reciente que w , o de w' ocurriendo en más de Δ unidades de tiempo real antes que r respectivamente. En estos casos, aun se puede argumentar que r ocurre a tiempo. **Definición 2:** sea $D \subseteq H$ un conjunto de operaciones y S una serialización de D . Sea $w, r \in D$ como se presenta en la definición 1. Se asume un conjunto W_r , asociado con r como: $W_r = \{w' \in D : (w' \text{ escribe un valor dentro del objeto } X) \wedge (T(w) + \varepsilon < T(w')) \wedge (T((w')) \wedge (T(w') + \varepsilon < (T(r) - \Delta))\}$. Se dice que una operación r ocurre o lee a tiempo en la serialización S si $W_r = \emptyset$. S es *timed* si cada operación de **read** en S ocurre a tiempo. TSC and TCC: se combinan los requerimientos de los modelos de consistencia bien conocidos como **SC** y **CC** con requerimientos de lectura a tiempo. **Definición 3:** la historia H satisface *timed serial consistency* (**TSC**) si hay una *timed serialization* S de H que respete el orden del programa para cada sitio en el sistema. **Definición 4:** la historia H satisface *timed causal consistency* (**TCC**) si para cada sitio i hay una *timed serialization* S_i de H_{i+w} que respete el orden causal “ \rightarrow ”.

Applications of Timed Consistency: A pesar de la utilidad y la importancia de modelos de consistencia como **SC** y **CC**, su definición formal y su implementación práctica permite ejecuciones que ignoran la ocurrencia en tiempo real de ciertos eventos. En muchos casos eso podría considerarse una ventaja que permite implementaciones eficientes y convenientes. Por ejemplo, **CC** se ajusta bien a aplicaciones móviles y tiene la habilidad de manejar desconexiones sin problemas. Los modelos de *timed consistency* son útiles en aplicaciones en donde la observancia de paso del tiempo real y sus efectos es parte de la denificación de una definición de correctitud/exactitud. En la misma forma que ocurre en la naturaliza, se puede permitir el paso de un período finito antes que los efectos de una operación san conocidos por todos los sitios en el sistema. Este período es precisamente el parámetro Δ .

Implementación: Lifetime Based Consistency Protocol: esta técnica provee consistencia a través de objetos diferentes pero relacionados. Se asume una arquitectura donde cada objeto tiene un conjunto de sitios-servidor que proveen almacenamiento a largo plazo para el objeto y donde sitios-cliente ponen en caché objetos antes de accederlos. Intentos fallidos de caché (*cache misses*) son manejados por cada sitio-servidor, que tiene una copia del objeto solicitado o puede obtenerlo. Sea el conjunto C_i denotar el caché del sitio i , en donde se almacena copias de objetos que han sido accesados recientemente. Si un *cache miss* ocurre cuando se accede el objeto X , algún servidor provee una copia de su versión actual de X . Una vez que esta copia es almacenada en C_i , se denota como X_i . El tiempo de inicio de X_i , denotado como X_i^α , es el momento cuando el valor X_i fue escrito. El último momento en que se conoce el valor almacenado en X_i se conoce como su tiempo de finalización y se denota con X_i^ω . El intervalo $[X_i^\alpha, X_i^\omega]$ es conocido como el tiempo de vida (*lifetime*) del valor almacenado en X_i . Los valores X_i y Y_i (en el caché C_i) son mutuamente consistentes si $\max(X_i^\alpha, Y_i^\alpha) \leq \min(X_i^\omega, Y_i^\omega)$ - sus tiempos de vida se sobreponen y por lo tanto, coexisten en algún instance. C_i es consistente si el tiempo máximo de inicio de cualquier valor de objeto en C_i es menor o igual que el tiempo mínimo de finalización de cualquier valor de objeto en C_i - cada par de objetos en C_i es mutuamente consistente. En general, el tiempo de vida de valores de objetos arbitrarios no se conoce. Cuando el sitio i actualiza la versión del objeto X_i en el tiempo t , la marca de tiempo t se asigna a ambos X_i^α y X_i^ω . Se tiene que descubrir conforme se avanza que ninguna copia $X_j (i \neq j)$ ha sido sobrescrita y usa esta información para avanzar X_i^ω . Una variable de la marca de tiempo local llamada **Context** _{i} es asociada con C_i . Su valor inicial es 0, y es actualizada con las siguiente reglas: 1. Cuando una copia del objeto X se introduce en C_i (convirtiéndose en X_i): **Context** _{i} := $\max(X_i^\alpha, \text{Context}_i)$. 2. Cuando la copia del objeto X_j se actualiza en el tiempo t : **Context** _{i} := $X_i^\alpha := t$. **Context** _{i} mantiene el último tiempo de inicio de cualquier valor de objeto que está o ha sido almacenado en C_i . Cuando

²Mantenido por un servidor de tiempo

una copia del objeto X se introduce en C_i , su valor final no debe ser menor que $\mathbf{Context}_i$, si es necesario, otros sitios-servidor o sitios-cliente con contactados hasta que una versión de X que satisfice la condición se encuentra. Además, cualquier objeto $Y_i \in C_i$ tal que $Y_i^\omega < \mathbf{Context}_i$ es invalidado. Se ha probado que este protocolo induce **SC** en la ejecución. TSC implementation: en el protocolo *lifetime*, el sitio i accede un estado de los objetos que fue consistente en el tiempo $\mathbf{Context}_i$, pero el tiempo actual puede ser mucho más tarde. Al controlar que la diferencia entre el tiempo actual y $\mathbf{Context}_i$ es mejor que o igual a Δ , se induce **TSC**. Los valores de los objetos cuyo tiempo de finalización es mayor/más viejo que Δ unidades de tiempo real en el pasao son invalidados localmente. Esto se lleva a cabo agregando una regla extra para actualizar el $\mathbf{Context}_i$ (Sea t_i el tiempo actual en el sitio i): 3. $\mathbf{Context}_i := \max(t_i - \Delta, \mathbf{Context}_i)$. Esto trabaja muy bien con objetos cuyos tiempos de finalización son conocidos pero puede llegar a generar invalidaciones innecesarias para objetos arbitrarios cuyos tiempos de vida no son conocidos de forma precisa. TCC Implementation: todas las marcas de tiempo usadas en el sistema son tomadas de relojes vectoriales o plausibles. Se hacen variates al protocolo *lifetime* para no requerir de relojes físicos y generar menor *overhead* en las comunicaciones que **SC**. $\mathbf{Context}_i$ se actualiza con versiones de las reglas 1 y 2³ adaptadas a relojes lógicos. Entre otras cosas, esto requiere el cómputo del máximo y el mínimo de dos marcas de tiempo lógicas que fueron tomadas ya sea de relojes vectoriales o plausibles. Cuando una copia de un objeto es traído desde el servidor a C_i , se verifica que su tiempo de finalización no sea causalmente antes de $\mathbf{Context}_i$. De ser necesario, otros sitios-servidor o sitios-cliente son contactados hasta que una vesion del objeto que satisfaga esta condición se encuentre. Cualquier objeto $Y_i \in C_i$ cuyo tiempo de finalización es causalmente antes de $\mathbf{Context}_i$ ($Y_i^\omega \rightarrow \mathbf{Context}_i$) es invalidado. Copias locales de objetos podrían ser invalidadas cuando un nuevo objeto es traído hacia el caché pero nunca son invalidadas como consecuencia de la actualización de un objeto local poque el tiempos de finalización lógicos de las copias locales se avanzan/adelfantan juntos con el tiempo lógico del sitio. Para garantizar que los efectos de una operación de **write** sea observada luego de Δ unidades de tiempo real, se agrega una nueva marca de tiempo tomada del reloj físico: el *checking time* de X_i , denotado como X_i^β . Esta marca de tiempo es el último instante de tiempo real en que el valor almacenado en X_i se conoce como válido. Sea t_i el tiempo real del sitio i . Cuando una copia del objeto X se trae desde el servidor s hacia C_i , ahora también se requiere que $X_i^\beta \geq t_i - \Delta$. De forma similar, cualquier objeto local $Y_i \in C_i$ tal que $Y_i^\beta < t_i - \Delta$ es invalidado o marcado como antiguo. Estas reglas garantizan que el sitio siempre lee a tiempo, porque los valores de los objetos que se sospechan que son antiguos son o validados o reemplazados por versiones nuevas. TCC using logical clocks: **TCC** puede ser aproximado en un sistema distribuido cuyos sitios solo comparten un reloj lógico. Se define un mapa ξ del conjunto de marcas de tiempo lógicas al conjunto de números reales. El parámetro Δ no será expresado en unidades de tiempo real, sino como un número real que define la máxima diferencia entre los valores que ξ produce para ciertas marcas de tiempo. **Definición 5**: la función ξ mapea marcas de tiempo tomadas de un reloj lógico al conjunto de números reales. Sea t y u dos marcas de tiempo lógicas, ξ tiene las siguientes propiedades: (1) $t = u \Rightarrow \xi(t) = \xi(u)$, (2) $t \rightarrow u \Rightarrow \xi(t) < \xi(u)$. Informalmente, $\xi(t)$ representa la “cantidad” de actividad global del sistema que es conocido cuanod el evento asociado con la marca de tiempo t es generada. Incluso si las marcas de tiempo t y u son concurrentes, se desea que si en el tiempo u el sistema está conciente de mayor actividad global que en el tiempo t , entonces $\xi(u) > \xi(t)$. **Definición 6**: Sea $D \subseteq H$ un conjunto de operaciones y S una serialización de D . Sea $w, r \in D$ como se presenta en la definición 1. Se define el conjunto W_r , asociado a r como: $W_r = \{w' \in D: (w' \text{ escribe un valor dentro del objeto } X \wedge (\xi(L(w)) < \xi(L(w')) < (\xi(L(r)) - \Delta))\}$. Se dice que la operación r ocurre o lee a tiempo en la serialización S , si $W_r = \emptyset$. S es *timed* si cada operación **read** en S ocurre a tiempo. Por lo tanto, *timed consistency* requer que si una operación **write** se ejecuta en el tiempo lógico t , tiene que estar visible en el sitio i antes de $\xi(t_i) - \xi(t) > \Delta$, donde t , es el tiempo lógico del sitio i .

1. ¿CUÁL ES EL PROBLEMA QUE PLANTEA EL PAPER?

Los modelos de consistencia como consistencia secuencial y la consistencia causal no consideran el tiempo/momento particular en que una operación es ejecutada para establecer un orden válido enter todas las operaciones de una computación.

2. ¿POR QUÉ EL PROBLEMA ES INTERESANTE O IMPORTANTE?

Accesos no locales a los objetos en sistemas distribuidos es caro debido a los costos elevados en la comunicación. Con el fin de aliviar esta ineficiencia y para mejorar la confianza, disponibilidad y escalabilidad del sistema, técnicas como *caching* y replicación son normalmente usadas. Sin embargo, estos métodos conducen a la coexistencia de muchos y posiblemente diferentes versiones del mismo objeto.

En muchos modelos de consistencia, el tiempo real no es explícitamente capturado. Por ejemplo, consistencia secuencial solo requiere que sea posible serializar todas las operaciones ejecutadas en el sistema distribuido. Consistencia causal, un modelo más débil, solo requiere que las operaciones relacionadas con causalidad sean vistas en el mismo orden por todos los sitios del sistema, mientras operaciones concurrentes puede ser percibidas en orden diferente por diferentes sitios. Por otro lado, serialización estricta y linealización requieren que las serializaciones respeten el ordenamiento en tiempo real entre las transacciones y operaciones respectivamente.

³Del protocolo *lifetime*

3. ¿QUÉ OTRAS SOLUCIONES SE HAN INTENTADO PARA RESOLVER ESTE PROBLEMA?

La noción de tiempo ha sido explorada también en sistemas de comunicación como *delta causal broadcast*, en sistemas de memoria que tienen un modelo de consistencia temporal y en varios protocolos de caché para la Web.

La simulación de ambientes virtuales interactivos se ha explorado por Singla, donde el concepto de “consistencia delta” es definida como un criterio de corrección para controlar accesos a una memoria compartida.

La noción de causalidad- Δ en redes no confiables la presenta Baldoni. Este protocolo soporta aplicaciones colaborativas multimedia en tiempo real. Con el fin de entregar un mensaje a un proceso, se verifica que el orden causal de los mensajes se cumple.

4. ¿CUÁL ES LA SOLUCIÓN PROPUESTA POR LOS AUTORES?

Este problema se aborda definiendo modelos de consistencia convenientes para el estado de los objetos distribuidos. Un modelo de consistencia establece una serie de garantías acerca de las relaciones válidas entre las operaciones ejecutadas por los sitios en el sistema, que a su vez limitan los valores posibles que objetos compartidos pueden retornar cuando operaciones de **read** son ejecutadas.

Se propone un modelo de **timed consistency** que define un umbral aceptable máximo de tiempo después de lo cual los efectos de una operación **read** deben ser observados por todos los sitios de un sistema distribuido. *Timed consistency* no solo cumple las necesidades de muchas aplicaciones que deben observar actualizaciones para cambiar dinámicamente objetos en el momento oportuno, sino que también unifica modelos de consistencia existentes como consistencia secuencial y linealización. Se concentra en modelos de *timed serial consistency* (**TSC**) y *timed causal consistency* (**TCC**).

5. ¿QUÉ TAN EXITOSA ES ESTA SOLUCIÓN?

Al combinar los requerimientos de *timed consistency* con criterios con **SC** y **CC** se propone **TSC** y **TCC**, los cuales son muy adecuados para cumplir las necesidades de muchas aplicaciones. Además, este concepto también unifica modelos existentes de consistencia. Se exploró una posible implementación de **TSC** y **TCC** basado en el concepto de tiempos de vida en los valores de los objetos. Finalmente, una posible definición de **TCC** usando solo relojes lógicos fue presentado.

Algunos problemas de esta investigación serán abordados como parte de un trabajo futuro. Con el fin de tener un mejor entendimiento la relación entre el valor de Δ y el costo de lograr un nivel particular de oportunidad (*timeliness*), se están completando detalladas simulaciones de sistemas basadas en los criterios de consistencia descritos en el artículo. Otras posibles implementaciones de **TSC** y **TCC** tienen que ser consideradas. Para el caso de **TCC** usando solo relojes lógicos, es necesario explorar diferentes mapeos de marcas de tiempo lógicas a números reales, y proveer las con semántica apropiada para la selección del parámetro Δ .