

PADI-DSTM

Abstract

O PADI-DSTM é um sistema distribuído que permite gerir objectos que residem em memória e são partilhados por programas transaccionais que correm em máquinas diferentes.

1. Introdução

Inicialmente iremos abordar a nossa solução e compará-la com algumas alternativas que considerámos, apresentando uma vista global da sua arquitectura. De seguida falaremos das estruturas de dados e algoritmos escolhidos. Concluiremos abordando a tolerância a faltas.

2 Solução

2.1 Optimista vs Pessimista

Inicialmente considerámos seguir uma abordagem optimista. No entanto, não conhecendo a relação existente entre o número de leituras e escritas, admitimos que existe igual número de leituras e escritas, logo o número de conflitos poderá vir a ser grande.

Constatando esse potencial bottleneck de performance no nosso sistema, optámos por abandonar as soluções optimistas. Outro dos argumentos contra as soluções optimistas é o facto das transacções poderem abortar (se surgir algum conflito) depois de já terem realizado trabalho, o que levaria a refazer esse trabalho.

Assim, escolhemos utilizar *Strict Two Phase Locking* ou *S2PL*

2.2 Replicação activa vs replicação passiva

Depois de escolhermos qual o protocolo a usar, deparamo-nos com a escolha entre replicação activa e passiva. Inicialmente considerámos usar replicação activa com um protocolo de *Quorum Consensus*. No entanto, por esta precisar de mais servidores do que os necessários para a replicação passiva, e tendo também em conta que teríamos

que enviar e receber mais mensagens do que as necessárias ao usar replicação passiva, optámos por esta última opção.

3 Estruturas de Dados

3.1 Master

Esta classe gere o sistema de memória distribuída e é responsável por armazenar os dados globais do sistema.

A lista *registered servers* tem como finalidade registar os servidores primários existentes, assim como, permitir identificar quais os *PadInts* atribuídos a cada servidor. Garantindo, que no caso em que o servidor primário é substituído, o endereço registado é actualizado e os clientes continuam a aceder aos *PadInts* sem perturbações. Para além disso, os *Transaction ID (TID)* são atribuídos de forma sequencial e única.

Variável	Descrição
Last transaction identifier	Último <i>Transaction ID</i> atribuído
Registered servers	Estrutura que armazena <i>Server Registry</i> referente a cada servidor primário registado no master

Table 1. Atributos da classe Master

3.1.1 Server Registry

A classe *Server Registry* é responsável por armazenar o identificador e endereço de cada servidor primário registado no master. Para além disto, guarda também uma lista com os identificadores dos *PadInts* (UID) que foram atribuídos ao servidor.

Na nossa solução optámos por guardar apenas os servidores primários no *Master* em detrimento de guardar também os servidores secundários. Estes últimos precisam de ser conhecidos apenas pelos servidores primários, sendo que o master necessita apenas de lhes atribuir o seu papel (de servidor secundário) quando se tentam registar.

3.1.2 LoadBalancer

Com o objectivo de distribuir a carga pelos vários servidores criamos a classe *Load Balancer*. Esta classe é usada, pelo *Master*, para decidir qual o servidor primário a que deve ser atribuído um *PadInt* na sua criação. Para esta distribuição o critério escolhido foi atribuir o novo *PadInt* ao servidor primário que actualmente tem menos *PadInts* atribuídos.

Esta classe é também utilizada na redistribuição dos *PadInts* quando um novo par (primário, secundário) de servidores é criado. Nessa altura os *PadInts* são distribuídos de forma a que cada par de servidores fique com um número de *PadInts* aproximadamente igual. Este número resulta da média entre o número total de *PadInts* criados e o número de pares de servidores existentes.

3.2 Server

O conjunto das instâncias da classe *Server* representam a memória distribuída onde são armazenados os *PadInts*.

Um dos atributos desta classe é uma instância da classe *Server machine*. Esta a classe encapsula o *Server* de forma a permitir simular que este deixou de funcionar em resposta ao pedido de *Fail*. Permite ainda re-iniciá-lo como se tratasse de um novo *Server* acabado de criar.

Quando o *Server* precisa de realizar uma leitura ou escrita num dado *PadInt* usa os métodos descritos nas secções 4.1.1 e 4.1.2 para obter o lock respectivo e poder depois efectuar a leitura ou escrita, respectivamente.

Variável	Descrição
Identifier	Identificador do servidor
Address	Endereço do servidor
Server state	Estado do servidor
Server machine	Classe que contém o servidor

Table 2. Atributos da classe Servidor

3.2.1 Server State

A classe servidor contém um dos seguintes estados: *primary*, *backup*, *failed*, *froze*. Indicando respectivamente que o servidor é o servidor primário, secundário, recebeu pedido *Fail* ou *Freeze*. Cada estado modela o comportamento do servidor, dependendo do papel que lhe foi atribuído.

Variável	Descrição
padInt Dictionary	Estrutura que mapeia <i>uid</i> em <i>PadInt</i>
Im alive timer	Timer usado para o envio e recepção de mensagens <i>Im alive</i>
Pair server reference	Referência para o outro servidor do par
Pair server address	Endereço do outro servidor do par

Table 3. Atributos da classe Server State

3.3 PadInt

Esta classe representa o objecto gerido pelo *PADI-DSTM* que guarda um inteiro, onde o *lock* é apenas o TID. Esta classe é composta por:

Variável	Descrição
UID	Identificador do <i>PadInt</i> que representa
Actual value	Valor no momento actual da transacção
Original value	Valor no início da transacção
Lock type	Tipo do lock (leitura ou escrita) atribuído actualmente
Readers	Lista de transacções com locks de leitura atribuídos
Writer	Transacção com lock de escrita atribuído

Table 4. Atributos da classe PadInt

3.4 Stub do PadInt

A *Library* envia ao cliente stubs da classe *PadInt*. Esta classe tem a seguinte estrutura:

Variável	Descrição
UID	Identificador do <i>PadInt</i> que representa
TID	Identificador da transacção
ServerID	Identificador <i>Server</i> onde o <i>PadInt</i> está guardado
Address	Endereço do <i>Server</i> onde o <i>PadInt</i> está guardado
Cache	Referência para a <i>Cache</i> da <i>Library</i>

Table 5. Atributos da classe Stub do PadInt

Esta classe disponibiliza os seguintes métodos:

- *int Read()*: lê o valor guardado na *Cache* se este estiver disponível. Isto é, já foi feito anteriormente um pedido de leitura ou escrita ao *Server* e foi obtido o *lock* de leitura ou escrita, respectivamente. Caso contrário, efectua a leitura no *Server* e guarda o valor em *Cache*.

- *void Write(int value)*: escreve o valor na *Cache* se este estiver disponível. Isto é, já foi feito anteriormente um pedido de escrita ao *Server* e foi obtido o *lock* de escrita. Caso contrário, efectua a escrita no *Server* e guarda o valor em *Cache*.

3.5 Biblioteca

A biblioteca usada pelos clientes para comunicar com o sistema de memória distribuída encontra-se representada na classe *Library*. Esta classe tem a seguinte estrutura:

Variável	Descrição
Master server	Referência para o <i>Master</i>
Actual tid	Identificador da transacção atribuído pelo <i>Master</i>
Cache	Cache usada para guardar os valores temporários dos <i>PadInts</i>
Channel	Canal usado pela biblioteca

Table 6. Atributos da classe Library

De seguida apresentam-se alguns métodos da Biblioteca:

- *bool init()*: cria o canal a ser usado pela *Library* e inicializa a variável master server. Quando o canal é criado é também definido um valor de *timeout* ao fim do qual a *Library* deixa de estar bloqueada à espera de resposta a pedido que efectuou;
- *bool TxBegin()*: a Biblioteca pede ao Master para criar um novo *TID* para a transacção e regista-o. É criada uma nova *Cache*;
- *bool TxCommit()*: descrito na secção 4.2;
- *bool TxAbort()*: descrito na secção 4.3;
- *PadInt CreatePadInt(int UID)*: a *Library* pede ao *Master* para registar o *PadInt* e se este não tiver já sido registado anteriormente o *Master* retorna um tuplo contendo o identificador e endereço do *Server* onde esse *PadInt* deverá ser criado. De seguida, o *Server* primário, em resposta ao pedido da *Library*, cria um *PadInt* inicializado a zero, sem *locks* e pede ao secundário para fazer o mesmo, só respondendo à *Library*, com um *ack*, depois de ter recebido o *ack* do secundário. Por fim, a *Library* insere o *PadIntRegistry*, relativo ao *PadInt*, na cache, criando de seguida o *Stub* do *PadInt* para retornar ao cliente;
- *PadInt AccessPadInt(int UID)*: a *Library* pergunta ao *Master* qual é o *Server* onde está o *PadInt*. Caso o *PadInt* exista, o *Master* retorna um tuplo contendo o identificador e endereço do *Server* onde esse *PadInt* se encontra. De seguida a *Library* pergunta ao *Server* se tem o *PadInt*. Caso a resposta seja afirmativa, a *Library* insere o *PadIntRegistry*, relativo ao *PadInt*, na cache, e retorna ao cliente uma nova instância do *Stub* do *PadInt*. Caso contrário é lançada uma excepção;

3.5.1 Cache

Com o objectivo de melhorar o desempenho relativo a leituras e escritas criámos a classe *Cache*. Esta classe permite guardar valores temporários relativos a cada *PadInt* usado pelo cliente. A ideia base do funcionamento das leituras e escritas usando a *Cache* foi já descrita anteriormente na secção 3.4. Para tornar definitivos os valores escritos na *Cache* antes de ser efectuado o *commit*, descrito na secção 4.2, todos os *PadInt* acedidos para escrita são escritos no *Server* que os armazena.

Uma das vantagens obtidas com o uso da *Cache* é que se reduz drasticamente o número de pedidos efectuados ao *Server*, reduzindo assim a carga a que este está sujeito. Outra vantagem é reduzir o tempo consumido nas transacções, pois as leituras e escritas são mais rápidas devido a serem realizadas localmente.

4 Algoritmos propostos

4.1 Locking

Para a obtenção de locks a classe *PadInt*, referente a um dado *PadInt* identificado por *UID*, disponibiliza um conjunto de métodos cujo funcionamento será explicado de seguida. Estes métodos são utilizados pela classe *Server* para responder a pedidos de leitura ou escrita.

4.1.1 GetReadLock(TID)

Quando este método é invocado começa-se por verificar se o lock de leitura pretendido pela transacção, identificada por *TID*, já lhe foi atribuído ou se já lhe foi atribuído um lock de escrita. Nesse caso, retorna-se *true*. Caso contrário, é chamado o método *AcquireLock*, descrito na secção 4.1.3, com o tipo de lock pretendido, neste caso leitura.

4.1.2 GetWriteLock(TID)

Primeiro verifica-se se o lock de escrita pretendido pela transacção, identificada por *TID*, já lhe foi atribuído. Nesse caso, retorna-se *true*. Caso contrário, é chamado o método *AcquireLock*, descrito na secção 4.1.3, com o tipo de lock pretendido, neste caso escrita.

4.1.3 AcquireLock(TID, requiredLockType)

Este método é invocado para obter locks de leitura ou escrita, consoante o valor do argumento *requiredLockType*. O argumento *TID* identifica a transacção que está a tentar obter o lock.

Se for possível obter o lock do tipo pedido, a variável *lockType* da classe *PadInt* é actualizada para o novo tipo de lock, *requiredLockType*. De seguida, caso o tipo de lock pedido seja leitura, o *TID* é adicionado à variável *readers*. Caso contrário, isto é o tipo de lock pedido é de escrita, é verificado se se trata de uma promoção, isto é, a transacção já tinha o lock de leitura e pretende obter o lock de escrita e nesse caso, o *TID* é removido da variável *readers*. Por fim, é atribuído à variável *writer* o valor *TID*.

Quando não é possível obter o lock o pedido é posto em espera no máximo durante um determinado intervalo de tempo, usando o método *Wait* da classe *Monitor* disponibilizado pela linguagem C#. Assim que é efectuado um *commit* ou *abort*, descrito nas secções 4.2 e 4.3, todos os pedidos são retirados da fila de espera e aos que for possível tentar adquirir o lock este ser-lhes-á atribuído. Os restantes voltarão a ficar em espera. Se ao fim do intervalo de tempo máximo estipulado um pedido ainda se encontra em espera então foi detectado um *deadlock* e esse pedido recebe uma excepção a indicar que tem que abortar.

Um pedido fica em espera num de dois casos. O primeiro é quando o tipo de lock actualmente atribuído e o lock pedido são ambos do tipo escrita. O segundo ocorre quando o tipo de lock actualmente atribuído e o pedido são diferentes e a transacção não está a tentar realizar uma promoção possível, isto é, se o tipo de lock atribuído é de leitura e a transacção está a tentar obter um lock de escrita e não é possível realizar a promoção porque não é só esta a transacção que possui o lock de leitura.

4.1.4 FreeWriteLock(TID)

O servidor remove o lock de escrita da transacção identificada por *TID*, associado ao *PadInt*. Neste passo basta apenas inutilizar a variável *writer*, atribuindo-lhe um valor que não identifique numa transacção.

4.1.5 FreeReadLock(TID)

O servidor remove o lock de leitura, da transacção identificada pelo *TID*, associado ao *PadInt* identificado por *UID*. Neste passo basta apenas remover o *TID* da variável *readers*.

4.2 Commit

Quando é invocado o método *TxCommit* da *Library*, caso tenham sido criados ou acedidos *PadInts*, são efectua-

dos dois passos. O primeiro passo corresponde à escrita de todos os valores escritos em *Cache* e encontra-se descrito na secção 3.5.1. No segundo passo, é enviado a cada *Server* um pedido de commit contendo todos os *UID* dos *PadInts* acedidos nesse *Server* e a cache é re-iniciada.

O *Server* ao receber este pedido verifica se guarda todos os *PadInts* identificados pelos *UID* recebidos. Caso não guarde algum deles lança uma excepção a indicar que não guarda o *PadInt*. Caso contrário, usando os métodos da classe *PadInt* descritos nas secções 4.1.5 e 4.1.4 são libertados os locks de leitura ou escrita atribuídos à transacção que está a efectuar o commit.

4.3 Abort

O método *TxAbort* da *Library* é em tudo semelhante ao método *TxCommit* excepto em dois pontos. A primeira diferença é que antes ser libertado cada lock de escrita associado a cada *UID* referenciado pela transacção, o valor actual do *PadInt* é substituído pelo valor registado como sendo o valor original antes da transacção o ter alterado, isto é, é re-posto o valor do último commit realizado com sucesso. A segunda diferença está relacionada com a *Cache*, pois não existe a escrita dos valores nos *Servers*, devido à primeira diferença enunciada.

5 Tolerância a Faltas - Fail, Freeze e Recover

Sempre que um método do servidor é invocado é verificado se este se encontra no estado freeze, fail ou normal de forma a ser realizado o comportamento adequado.

O servidor primário envia uma mensagem de *i'm alive* ao respectivo servidor secundário a cada *x* segundos. Caso o secundário não receba a mensagem após o tempo limite, este regista-se no Master como primário e cria uma nova instância de secundário.

Quando o antigo servidor primário, que recebeu o pedido de freeze ou fail e não enviou um *i'm alive* ao respectivo secundário, volta ao estado normal(recebeu recover) verifica se o secundário já assumiu o papel de servidor primário:

- Se não aconteceu, é enviado um *i'm alive* para o secundário. No caso de ter estado em freeze executa os pedidos que registou e continua o seu funcionamento normal, no caso de ter estado em fail continua a operar de forma normal enquanto primário.
- Se aconteceu, independentemente do estado anterior ser fail ou freeze, o servidor termina a sua execução.

Na situação inversa, ou seja, o servidor secundário não enviou a resposta a um pedido no intervalo de tempo máximo, o servidor primário cria um novo servidor secundário e

quando o antigo servidor secundário receber recover, termina a sua execução independentemente do estado anterior.

No caso em que recebeu recover, vindo do estado de freeze e o servidor primário ainda está a espera da resposta, executa os pedidos pendentes, responde ao primário e continua a operar de forma normal.

O caso em que o servidor secundário retorna do estado fail, por opção nossa não acontece. A alternativa seria esperar que o servidor voltasse e fazê-lo terminar de seguida, no entanto, como à partida sabemos que se o secundário falhou, assim que o primário atender um pedido o seu estado vai ficar inconsistente, então este termina assim que receber o pedido de fail, parecendo ao servidor primário que não respondeu no intervalo de tempo máximo previsto.

6 Conclusão

A nossa solução garante as propriedades ACID para transações, consistência sequencial e tolera a falha de um servidor.