

Seminar Report: Chordy

June 5, 2019

Carles Tornel
Jesus Alfaro
Ricard Abril

1 Introduction

En aquesta pràctica implementarem una hash table distribuïda d'acord amb l'esquema Chord. Primerament podrem afegir nodes i posteriorment afegirem l'opció de tenir storage i controlar la caiguda d'un node així com recuperar storage.

2 Experiments

2.1 Building a Ring

- 2.1.1 Do some experiments to test the ring integrity, initially in one Erlang instance but then in a network of several instances. When connecting nodes on different instances remember to start Erlang in distributed mode (giving a -name argument) and make sure that you use the same cookie (-setcookie).**

En aquest punt, el nostre codi, és capaç de crear un ring i afegir-ne nodes. Així que per comprovar que funciona correctament hem creat 4 nodes i hem enviat missatges ha aquest node, per tal de veure si es reben correctament i el temps que trigen a arribar. Tal com podem veure en la imatge següent:

```
(Ni@127.0.0.1)> N1 = node1:start(1).
<0.84.0>
(Ni@127.0.0.1)> N3 = node1:start(3,N1).
<0.86.0>
(Ni@127.0.0.1)> N5 = node1:start(5,N1).
<0.88.0>
(Ni@127.0.0.1)> N7 = node1:start(7,N1).
<0.90.0>
(Ni@127.0.0.1)> N1 ! probe.
Create probe 1!
Forward probe 1!
Forward probe 1!
probe
Received probe 1 in 0 ms Ring: [7,5,3,1]
(Ni@127.0.0.1)> N3 ! probe
(Ni@127.0.0.1)> N3 ! probe.
* 2: syntax error before: N3
(Ni@127.0.0.1)> N3 ! probe.
Create probe 3!
Forward probe 3!
Forward probe 3!
Forward probe 3!
probe
Received probe 3 in 0 ms Ring: [1,7,5,3]
(Ni@127.0.0.1)> N5 ! probe.
Create probe 5!
Forward probe 5!
Forward probe 5!
Forward probe 5!
probe
Received probe 5 in 0 ms Ring: [3,1,7,5]
(Ni@127.0.0.1)> N7 ! probe.
* 1: variable 'N0' is unbound
(Ni@127.0.0.1)> N7 ! probe.
Create probe 7!
Forward probe 7!
Forward probe 7!
Forward probe 7!
probe
Received probe 7 in 0 ms Ring: [5,3,1,7]
(Ni@127.0.0.1)> █
```

2.2 Adding Store

2.2.1 If we now have a distributed store that can handle new nodes that are added to the ring we might try some performance testing. You need several machines to do this. Assume that we have eight machines and that we will use four in building the ring and four in testing the performance.

Ara el nostre codi ja és capaç de gestionar un "Storage" en els nodes, de forma que podrem emmagatzemar i buscar elements en forma de clau - valor en ells. Per comprovar que funciona correctament, hem afegit diferents elements als nodes. Tal com podem veure a continuació:

```
(N1@127.0.0.1)5> N1=node2:start(1).
<0.89.0>
(N1@127.0.0.1)6> P=chordy:connect(N1).
<0.91.0>
(N1@127.0.0.1)7> P ! {add,0,0}.
[Add] Key added correctly
{add,0,0}
(N1@127.0.0.1)8> P ! {add,1,1}.
[Add] Key added correctly
{add,1,1}
(N1@127.0.0.1)9> P ! {add,2,2}.
[Add] Key added correctly
{add,2,2}
(N1@127.0.0.1)10> P ! {add,3,3}.
[Add] Key added correctly
{add,3,3}
(N1@127.0.0.1)11> N1 ! probe
(N1@127.0.0.1)11> N1 ! probe.
* 2: syntax error before: N1
(N1@127.0.0.1)11> N1 ! probe.
Create probe 1!
probe
Received probe 1 in 0 ms Ring: [1]
(N1@127.0.0.1)12> N0 = node2:start(0,N1).
(N1@127.0.0.1)12> N0 = node2:start(0,N1).
** exception error: no match of right hand side value <0.100.0>
(N1@127.0.0.1)13> N0 = node2:start(0,N1).
<0.103.0>
(N1@127.0.0.1)14> N1 ! probe
(N1@127.0.0.1)14> N1 ! probe.
* 2: syntax error before: N1
(N1@127.0.0.1)14> N1 ! probe.
Create probe 1!
Forward probe 1!
probe
Received probe 1 in 0 ms Ring: [0,1]
(N1@127.0.0.1)15>
```

2.2.2 As a first test, build a ring with one node only and let each of the four test machines add a number (e.g., 1000) of random elements (key-value pairs) to the ring and then do a lookup of the elements, measuring the time it takes. You can use the test procedure given in 'Appendix A', which should be given with the ID of the node to contact. You should use the function `key:generate()` to assign the keys of the nodes in the ring, so that elements will get uniformly distributed among them.

En aquest test, només tindrem un node al ring, al qual afegirem un nombre aleatori d'entrades i comprovarem el seu funcionament mitjançant el codi de prova facilitat en l'apèndix de la pràctica.

```
Erlang/OTP 20 [erts-9.2] [source] [64-bit] [smp:2:2] [ds:2:2:10] [async-threads:10] [kernel-poll:false]
Eshell V9.2 (abort with ^C)
(N1@127.0.0.1)1> c(node2).
{ok,node2}
(N1@127.0.0.1)2> c(key).
{ok,key}
(N1@127.0.0.1)3> c(chordy).
{ok,chordy}
(N1@127.0.0.1)4> c(storage).
{ok,storage}
(N1@127.0.0.1)5> N1 = node2:start(1).
<0.89.0>
(N1@127.0.0.1)6> chordy:test(N1,2300).
Requests done. Waiting for answers...
Answers received. Making lookups...
Lookups done. Waiting for values...
Values received. Checking correctness...
Everything correct!
Elapsed Time: 103 ms
ok
(N1@127.0.0.1)7> █
```

2.2.3 For the next test, you should add another node to the ring and check how the performance changes. Then, add two more nodes to the ring, any changes? Check also the impact on the performance if all the test machines access the same node or different nodes in the ring.

En aquest experiment, repetirem el procediment de l'anterior, però afegint més nodes al ring, per tal de veure el canvi de rendiment del sistema.

Amb 2 nodes:

```
(N1@127.0.0.1)5> N1 = node2:start(1).
<0.89.0>
(N1@127.0.0.1)6> P = chordy:connect(N1).
<0.91.0>
(N1@127.0.0.1)7> N2 = node2:start(2,N1).
<0.93.0>
(N1@127.0.0.1)8> chordy:test(N1,2300).
Requests done. Waiting for answers...
Answers received. Making lookups...
Lookups done. Waiting for values...
Values received. Checking correctness...
Everything correct!
Elapsed Time: 112 ms
ok
(N1@127.0.0.1)9> chordy:test(N1,2300).
Requests done. Waiting for answers...
Answers received. Making lookups...
Lookups done. Waiting for values...
Values received. Checking correctness...
Everything correct!
Elapsed Time: 129 ms
ok
(N1@127.0.0.1)10> █
```

Amb 4 nodes:

```
(N1@127.0.0.1)5> N1 = node2:start(1).
<0.89.0>
(N1@127.0.0.1)6> P = chordy:connect(N1).
<0.91.0>
(N1@127.0.0.1)7> N2 = node2:start(2,N1).
<0.93.0>
(N1@127.0.0.1)8> chordy:test(N1,2300).
Requests done. Waiting for answers...
Answers received. Making lookups...
Lookups done. Waiting for values...
Values received. Checking correctness...
Everything correct!
Elapsed Time: 112 ms
ok
(N1@127.0.0.1)9> chordy:test(N1,2300).
Requests done. Waiting for answers...
Answers received. Making lookups...
Lookups done. Waiting for values...
Values received. Checking correctness...
Everything correct!
Elapsed Time: 129 ms
ok
(N1@127.0.0.1)10> N3 = node2:start(3,N1).
<0.97.0>
(N1@127.0.0.1)11> N4 = node2:start(4,N1).
<0.99.0>
(N1@127.0.0.1)12> chordy:test(N1,2300).
Requests done. Waiting for answers...
Answers received. Making lookups...
Lookups done. Waiting for values...
Values received. Checking correctness...
Everything correct!
Elapsed Time: 125 ms
ok
(N1@127.0.0.1)13> chordy:test(N1,2300).
Requests done. Waiting for answers...
Answers received. Making lookups...
Lookups done. Waiting for values...
Values received. Checking correctness...
Everything correct!
Elapsed Time: 140 ms
ok
(N1@127.0.0.1)14> █
```

2.3 Handling Failure

2.3.1 Experiments. Do some experiments to evaluate the behavior of your implementation when nodes can fail.

Ara el nostre codi, hauria de ser tolerant a fallades, és a dir que encara que un node caigués el programa que avalua el codi no hauria de donar errors.

```
(N1@127.0.0.1)5> N2 = node3:start(2).
<0.90.0>
(N1@127.0.0.1)6> P = chordy:connect(N2).
<0.92.0>
(N1@127.0.0.1)7> P ! {add,0,0}.
[Add] Key added correctly
[add,0,0]
(N1@127.0.0.1)8> P ! {add,1,1}.
[Add] Key added correctly
[add,1,1]
(N1@127.0.0.1)9> P ! {add,2,2}.
[Add] Key added correctly
[add,2,2]
(N1@127.0.0.1)10> N0 = node3:start(0,N2).
<0.97.0>
(N1@127.0.0.1)11> N1 = node3:start(1,N2).
<0.99.0>
(N1@127.0.0.1)12> N2 ! probe.
Create probe 2!
Forward probe 2!
probe
Received probe 2 in 0 ms Ring: [1,0,2]
(N1@127.0.0.1)13> N0 ! stop.
(N1@127.0.0.1)13> N0 ! stop.
^ Z: syntax error before: N0
(N1@127.0.0.1)13> N0 ! stop.
stop
(N1@127.0.0.1)14> N2 ! probe.
Create probe 2!
Forward probe 2!
probe
Received probe 2 in 0 ms Ring: [1,2]
(N1@127.0.0.1)15> <
```

2.4 Replication

2.4.1 Do some experiments to demonstrate that your replication strategy works. You can for instance execute the following commands, and check the contents of the Store and the Replica.

En aquesta última implementació, el nostre codi hauria de permetre replicar les entrades d'un node en un altre, per comprovar-ho, utilitzarem el codi donat. Tal com podem veure a continuació.

```
(N1@127.0.0.1)> N2 = nodes:start(N2).
* 1: variable 'N2' is unbound
(N1@127.0.0.1)> N2 = node4:start(2).
=> 94.0s
(N1@127.0.0.1)> P = chordy:connect(N2).
=> 96.0s
(N1@127.0.0.1)> P ! {add,0,0}.
[add] Key added correctly
[add,0,0]
(N1@127.0.0.1)> P ! {add,1,1}.
[add] Key added correctly
[add,1,1]
(N1@127.0.0.1)> P ! {add,2,2}.
[add] Key added correctly
[add,2,2]
(N1@127.0.0.1)> P ! {add,3,3}.
[add] Key added correctly
[add,3,3]
(N1@127.0.0.1)> N2 ! probe.
Create probe 2! Replica: [[3,3],[2,2],[1,1],[0,0]]
probe
Received probe 2 in 0 ms Ring: [2]
(N1@127.0.0.1)> N0 = node4:start(0,N2).
=> 103.0s
(N1@127.0.0.1)> N2 ! pr
pr_in_val  pr_in_file  pr_in_inet  pr_in_ip  proc_lbb  proplists

(N1@127.0.0.1)> N2 ! probe.
Create probe 2! Replica: [[3,3],[0,0]]
Forward probe 2! Replica: [[2,2],[1,1]]
probe
Received probe 2 in 0 ms Ring: [0,2]
(N1@127.0.0.1)> N1 = node4:start(1,N2).
=> 106.0s
(N1@127.0.0.1)> N2 ! probe.
Create probe 2! Replica: [[1,1]]
Forward probe 2! Replica: [[2,2],[1,1]]
Forward probe 2! Replica: [[3,3],[0,0]]
probe
Received probe 2 in 0 ms Ring: [1,0,2]
(N1@127.0.0.1)> N1 ! stop.
stop
(N1@127.0.0.1)> N2 ! probe.
Create probe 2! Replica: [[3,3],[0,0]]
Forward probe 2! Replica: [[1,1],[2,2]]
probe
Received probe 2 in 0 ms Ring: [0,2]
(N1@127.0.0.1)> N2
```

3 Open Questions

3.1 Building a Ring

3.1.1 What are the pros and cons of a more frequent stabilizing procedure?

Amb una freqüència d'estabilització alta podrem determinar més rapidament quins nodes encara no són estables a l'anell i reduïrem l'estat d'inestabilitat a l'anell. Encara que hauríem de tenir en compte si el overhead creat pels missatges d'estabilització és òptim respecte la velocitat d'estabilització dels nodes, ja que no sempre amb un procediment d'estabilització més freqüent s'aconsegueix una millora en la velocitat d'enllaçament dels nodes.

3.1.2 Do we have to inform the new node about our decision? How will it know if we have discarded its friendly proposal?

L'hauréu d'informar sempre, ja que no tindrà informació suficient per determinar el seu lloc a l'anell per si mateix. Si no acceptem la seva proposició significa que no pot ser el nostre predecessor, per tant, deleguem la feina al nostre predecessor.

3.1.3 What would happen if we didn't schedule the stabilize procedure? Would things still work?

No del tot. Els nodes en els quals hi haurà el predecessor dels nous nodes entrants no se n'adonaràn d'aquests nous nodes i no actualitzaran la seva informació, fent que el ring esdevingue malformat.

3.2 Handling Failure

3.2.1 What will happen if a node is falsely detected of being dead (e.g. it has only been temporally unavailable)?

Serà automàticament descartat i el ring canviarà temporalment fins que el node aparentment mort estigui novament disponible, on retornarà a la seva posició inicial.

4 Personal opinion

El Seminari ha d'inclou-re's als pròxims cursos ja que expandeix informació sobre els "name servers" i dóna l'oportunitat d'implementar un sistema de noms distribuïts similar als DNS, clarificant i exemplificant com aquests sistemes es porten a terme.