

Rapport de TP - Cassandra DB

loïc divad

Mai 2015

1 Introduction

...

2 Installation

Pour réaliser l'installation du logiciel avec la dernière version en date: *2.1.7*,

on procède d'une manière différente en indiquant toute les étapes.

```
lmdadm $ wget ftp://mirrors.ircam.fr/.../apache-cassandra-2.1.7-bin.tar.gz
```

```
lmdadm $ tar -zxvf apache-cassandra-2.1.7-bin.tar.gz
```

On redirige en suite les logs en éditant le fichier *logback.xml*.

...

```
<file> /home/lmdadm/log/cassandra/system.log </file>
```

...

Puis on édite le fichier *cassandra.yaml*. La liste suivante présente tous la paramètres personnalisés dans le cadre du TP:

- *cluster_name*: Cassandra Cluster
- *data_file_directories*: /home/lmdadm/data/cassandra
- *commitlog_directory*: /home/lmdadm/log/cassandra/commitlog
- *saved_caches_directory*: /home/lmdadm/tmp/cassandra/saved_caches

On lance ensuite la base de donnée. On note que l'option indiquée dans le TP -f signifie: *force foreground*. Conclusion, nous n'utiliserons sur tout pas cette option. Après l'avoir lancé en mode démon, on ouvre la console:

```
lmdadm $ ./bin/cassandra
```

```
lmdadm $ ./bin/cqlsh
```

2.1 Logiciel d'administration

Il existe un bon nombre de logiciel d'administration pour cassandra, la plupart propulsés par des communautés opensources. Voici une petite liste des les logiciels trouvée après un brève recherche.

Le logiciel retenu pour réaliser ce TP est la dernière version (*1.3.1*) de DevCenter édité par *DataStax*. Contrairement à d'autres logiciels essayés (OpsCenter, helenos) DevCenter est un client lourd, on n'y accède pas par un navigateur web. Il ne sera donc pas installé sur le serveur mais sur notre ordinateur personnel. On ouvre en suite une connection vers les noeuds du Cleusteur.

téléchargement: <http://www.datastax.com/download-ops-dev>

Une fois le logiciel installé et la connection enregistrée (figure 1.) on demande à Cassandra d'accepter les connections à distance, les 'remote connections '. Pour cela on édite de nouveau le fichier *cassandra.yaml*.

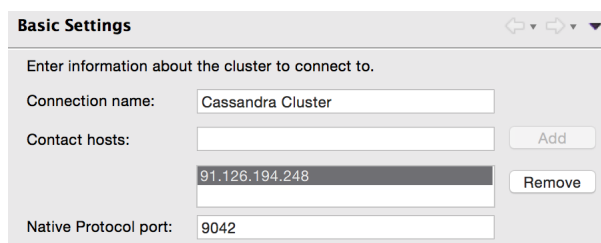


Figure 1: Ajout du serveur sur le quel tourne Cassandra.

- `rpc_address`: 0.0.0.0 (Anciennement localhost)
- `broadcast_rpc_address`: 1.2.3.4

Une fois le tout configuré il est possible d'accéder à l'ensemble des keyspaces. Et de réaliser des requêtes sur la base. Nous reviendrons sur cette interface pour illustrer des différents points du TP.

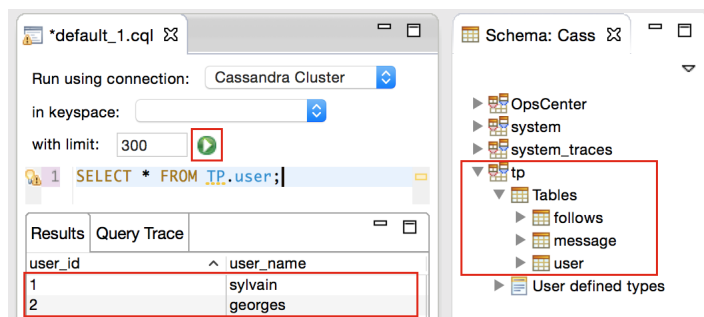


Figure 2: Exemple d'interface DevCenter par DataStax.

3 Base de donnée

3.1 Description de l'espace de clef

Pour mieux comprendre la structure de donnée dans cassandra cette partie propose un petit point sur ce qu'est un espace de clef, une brève définition et une comparaison avec d'autre systemes connus. On propose également de décrire la commande de création proposé dans le TP:

```
cqlsh> CREATE KEYSPACE TP WITH REPLICATION = {
'class' : 'SimpleStrategy',
'replication_factor' : '1'
};
```

Definition(Keyspace): Il s'agit de la plus grande structure de données de cassandra. Elle est comparable aux bases (ou schéma) dans les systemes classiques ou dans les systemes document (mongodb, couchdb). Elle contient des familles de colonnes, qui elles, représentent la séparation entre les tables (mysql) ou les collections (mongodb). Chaque ligne insérée présente des informations pour certaines colonnes d'un famille de colonnes.

Propriétés:

- **Replication factor**: définit le nombre de noeud possédant un copie de la donnée.

- **Replica placement strategy:** il s'agit du mode de répartition des réplicats. Il en existe trois. simple strategy, old network topology strategy, network topology strategy
- **Column families:** Liste des familles des colonnes associées. Dans notre cas sur le keyspace TP il y a user, follows et message.

La commande de création du Keyspace TP signifie donc: créer l'espace de clefs TP dont les lignes ne sont présentes que sur un noeud. Si il doit y avoir de la réplication le *partitioner* redistribuera automatiquement sur le noeud suivant.

La commande DESCRIBE nous permet d'avoir la liste de keyspaces.

```
cqlsh> DESCRIBE KEYSPACES;
system_traces OpsCenter system tp
```

On peut ensuite accéder aux familles de colonnes sur un espace de clefs particulier.

```
cqlsh> use TP;
cqlsh:tp> DESCRIBE TABLES;
follows message user
```

On peut ensuite retrouver les colonnes d'une famille en particulier.

```
cqlsh:tp> DESCRIBE TABLE user;
CREATE TABLE tp.user (
  user_id bigint PRIMARY KEY,
  user_name text
)
```

```
~/products/cassandra/bin • lmdadm $>./cqlsh
Connected to Cassandra Cluster at 127.0.0.1:9042.
[cqlsh 5.0.1 | Cassandra 2.1.7 | CQL spec 3.2.0 | Native protocol v3]
Use HELP for help.
cqlsh> DESCRIBE KEYSPACES;

system_traces "OpsCenter" system tp

cqlsh> use TP;
cqlsh:tp> DESCRIBE TABLES;

follows message user

cqlsh:tp> SELECT * FROM TP.user;

 user_id | user_name
-----|-----
      2 | georges
      1 | sylvain

(2 rows)
cqlsh:tp>
```

Figure 3: Quelques commandes cql.

3.2 le CQL, langage de requête Cassandra

remarque: On note que l'utilisation de double quotes (") dans la ligne de commande provoque l'erreur:

```
SyntaxException: <ErrorMessage code=2000 [Syntax error in CQL query] message="line
1:58 extraneous input hello world expecting ) (...) values (1,1 hello world[])...>
```

3.3 Clefs composites

Via Devcenter on insère les lignes demandées comme sur les figures suivantes.

```

1 INSERT INTO user (user_id, user_name) VALUES (3, 'patrick');
2
3 INSERT INTO follows (follower_id, followed_id) VALUES (2, 3);

```

Figure 4: Ajout de l'utilisateur patrick suivi par georges.

Puis on vérifie la table des relations d'abonnement 'follows'. La première relation d'abonnement (georges suit sylvain) a été supprimée. En effet, lors de la définition de la famille de clef follows, la colonne follower_id est déclarée comme une clef primaire, elle est donc unique.

CREATE TABLE follows(follower_id bigint PRIMARY KEY, followed_id bigint);
Autrement dit, on ne peut suivre qu'une seule personne à la fois.

```

1 SELECT followed_id FROM follows WHERE follower_id = 2;
2

```

Results	Query Trace
follower_id	followed_id
2	3

Figure 5: Unique relation de la table follows après l'ajout de patrick.

A l'aide de la commande indiquée dans l'énoncé du tp on définit un couple de clef. Désormais, c'est la relation suivi-abonné qui est unique. On ne peut suivre quelqu'un qu'une seule fois. On note que seules 3 lignes ont été insérées car les deux dernières INSERT sont des duplicats.

```

INSERT INTO follows (follower_id, followed_id) VALUES (2,3);
INSERT INTO follows (follower_id, followed_id) VALUES (2,1);
INSERT INTO follows (follower_id, followed_id) VALUES (3,1);
INSERT INTO follows (follower_id, followed_id) VALUES (3,1);

```

Figure 6: Série d'insertion d'abonnement.

```

1 SELECT * FROM follows;
2

```

Results	Query Trace
follower_id	followed_id
2	1
2	3
3	1

Figure 7: Relevé des différents abonnements.