



Les sauvegardes de Zeste de Savoir

20 janvier 2019

Table des matières

1. Introduction	1
2. Sauvegarde de la base via Xtrabackup	1
3. rsync et snapshots ZFS	4
4. La cerise sur le gâteau : compression à la volée des volumes	5
5. Conclusion	6

% LES SAUVEGARDES DE ZESTE DE SAVOIR % Sandhose % 11 mars 2018

1. Introduction

J'ai dans mon appart un petit serveur tournant sous FreeBSD, me servant notamment de routeur et de NAS. J'ai 8 To de stockage dessus en RAID¹, et me sert depuis peu de ce serveur comme espace de sauvegarde pour la production de Zeste de Savoir.

Les données de Zeste de Savoir se trouvent à deux endroits en production. D'une part dans le dossier `/opt/zds/data`, contenant les dépôts git des contenus, les médias uploadés, et autres PDFs générés. D'autre part dans une base de donnée MySQL.

2. Sauvegarde de la base via Xtrabackup

Auparavant, pour sauvegarder la base, `mysqldump` était utilisé pour générer un dump complet de la base sous forme de fichier `sql`. Ces dumps pèsent un peu moins de 500 Mo non-compressés, et 100 Mo `gzippés`. Aussi, chaque dump contient l'ensemble de la base de données, alors qu'en général entre deux sauvegardes, seul une partie de la base a changé.

Aussi, `mysqldump` a tendance à vouloir verrouiller les tables pendant le dump, et ainsi ralentir le site.

Mon objectif était d'effectuer des backups plus régulièrement, pour éviter des [retours dans le passés trop importants](#) en cas d'incident majeur. Les dumps étaient effectués une fois tous les jours, et je ne pouvais pas me permettre d'en faire un toutes les heures.

Xtrabackup est un outil développé par Percona permettant de faire des sauvegardes incrémentales et sans `LOCK` d'une base MySQL. Il se base sur le fait qu'InnoDB² maintient sur le disque un historique des transactions lui permettant de récupérer la base en cas d'incident majeur, tant que le disque est intacte. Il utilise donc directement les fichiers InnoDB sur le disque, plutôt que de récupérer toute la base via une connexion `mysql` au serveur de BDD.

Il offre donc

1. J'ai en réalité ~16 To de disques durs répartis sur 4 disques, qui sont en miroirs deux par deux.

2. Sauvegarde de la base via Xtrabackup

- des sauvegardes sans lock des tables³
- des sauvegardes incrémentales, puisqu'il lui suffit de sauvegarder les nouvelles entrées du journal de transactions d'InnoDB
- une restauration beaucoup plus rapide, puisqu'il s'agit d'une simple copie de fichier, et non un énorme script SQL à exécuter

J'ai donc écrit un petit script permettant de faire des backups complètes et incrémentales. Les backups complètes sont effectuées une fois par jours, et les backups incrémentales toutes les heures.

```
1  #!/bin/sh
2
3  ## Bonne pratique: arrêter le script si une commande retourne un
4  # code d'erreur autre que 0,
5  # et si une variable non-existante est utilisée
6  set -eu
7
8  ## L'endroit où sont stockées les sauvegardes
9  WD=/var/backups/mysql
10 ## `latest` est un lien symbolique vers la dernière sauvegarde
11 LATEST=$WD/latest
12
13 ## On récupère le nom complet de l'ancienne backup
14 PREVIOUS=`readlink -f $LATEST`
15 ## La nouvelle backup aura comme nom AAMMJJ-HHSS
16 NEXT=$WD/`date '+%Y%m%d-%H%M'`
17
18 if [ -d "$NEXT" ]; then
19     echo "\`$NEXT' already exists."
20     exit 1
21 fi
22
23 if [ "$#" -ge 1 ] && [ "$1" = "full" ]; then
24     # Les backups "full" ont en plus `-full` dans le nom
25     NEXT=$NEXT-full
26     xtrabackup --backup --target-dir=$NEXT 2> $NEXT.log
27 else
28     if ! [ -L "$LATEST" ]; then
29         echo
30         "\`$LATEST' does not exists. Consider doing a full backup first."
31         exit 1
32     fi
33
34     xtrabackup --backup --target-dir=$NEXT
35     --incremental-basedir=$PREVIOUS 2> $NEXT.log
36 fi
37
38 ## On supprime les anciens liens symboliques
```

2. Sauvegarde de la base via Xtrabackup

```
36 rm -f $LATEST $LATEST.log
37
38 ## Et on les recréé avec la sauvegarde fraîche
39 ln -s $NEXT $LATEST
40 ln -s $NEXT.log $LATEST.log
```

Plus qu'à exécuter ce script régulièrement le script via cron :

```
1 root@zdsprod1:~# crontab -l
2 ## min hour dom month dow command
3 0 * * * * /var/backups/mysql/backup.sh
4 15 3 * * * /var/backups/mysql/backup.sh full
```

Et voilà le résultat !

```
1 root@zdsprod1:~# ls -l /var/backups/mysql
2 ...
3 drwxr-x---  5 root root 4.0K Mar 10 02:00 20180310-0200
4 -rw-r--r--  1 root root  47K Mar 10 02:00 20180310-0200.log
5 drwxr-x---  5 root root 4.0K Mar 10 03:00 20180310-0300
6 -rw-r--r--  1 root root  47K Mar 10 03:00 20180310-0300.log
7 drwxr-x---  5 root root 4.0K Mar 10 03:15 20180310-0315-full
8 -rw-r--r--  1 root root  46K Mar 10 03:15 20180310-0315-full.log
9 drwxr-x---  5 root root 4.0K Mar 10 04:00 20180310-0400
10 -rw-r--r--  1 root root  46K Mar 10 04:00 20180310-0400.log
11 drwxr-x---  5 root root 4.0K Mar 10 05:00 20180310-0500
12 -rw-r--r--  1 root root  46K Mar 10 05:00 20180310-0500.log
13 ...
```

À titre indicatif, une sauvegarde complète pèse 1.2 Go non-compressée, 200 Mo compressé, et une incrémentale pèse en moyenne 15 Mo.

Ces sauvegardes sont ensuite récupérées sur mon serveur toutes les heures via `rsync` :

```
1 rsync -avr zdsprod1:/var/backups/mysql/
   /tupperware/backups/zestedesavoir/db
```

Histoire de pas exploser le stockage de la prod, j'ai écrit un petit script pour garder que les deux dernières backups complètes (+ les incrémentales entre). Les anciennes sauvegardes restent bien sûr sur mon serveur.

3. rsync et snapshots ZFS

```
1  ##!/bin/sh
2
3  set -eu
4
5  WD=/var/backups/mysql
6
7  BACKUPS="`echo $WD/*-*/ | tr ' ' '\n' | sort -nr`"
8
9  TO_DELETE=""
10     echo "$BACKUPS" | awk '
11     BEGIN { full=0 }
12     { if (full > 1) { print $0 } }
13     /full/ { full++ }
14     '
15 `"`
16
17 [ -z "$TO_DELETE" ] || rm -r "$TO_DELETE"
```

3. rsync et snapshots ZFS

J'ai donc mentionné plus tôt que j'utilise ZFS sur mon serveur comme système de fichier.

ZFS a été développé à l'origine par Oracle pour Solaris. Il a été ensuite porté notamment sur FreeBSD, et plus récemment sur Linux (via ZFSOnLinux).

ZFS a la particularité de faire à la fois gestionnaire de volumes logiques (comme LVM) & RAID (mdadm sur linux) et du système de fichier. Il permet de gérer un ensemble de disques durs comme un seul ensemble (éventuellement avec de la redondance pour être tolérant aux pannes), que l'on peut subdiviser en multiples sous-volumes (**datasets**, dans le jargon ZFS) avec différentes propriétés. On pourra donc attribuer différents droits aux utilisateurs en fonction du volume, limiter la taille maximum de chaque volume, et activer certaines fonctionnalités type compression et deduplication à la volée par volume.

ZFS permet également de faire des snapshots instantanément et sans coût d'un volume. Une snapshot ne va prendre comme place sur le disque que la différence entre la snapshot et le volume actif. Idéal donc, pour des backups.

```
1  spaetzle/root ~ > zfs list -t snap
2  NAME                                     USED  AVAIL
3  REFER  MOUNTPOINT
4  ...
```

2. Le moteur de stockage par défaut sur MySQL

3. Je sais, **mysqldump** peut ne pas locker les tables, mais il lock quand même les transactions pendant la backup, ce qui "met en pause" toute opération pendant le dump.

4. La cerise sur le gâteau : compression à la volée des volumes

4	tupperware/backups/zestedesavoir/data@20180310-1602	412K	-
	5.56G -		
5	tupperware/backups/zestedesavoir/data@20180310-1702	884K	-
	5.56G -		
6	tupperware/backups/zestedesavoir/data@20180310-1802	876K	-
	5.56G -		
7	tupperware/backups/zestedesavoir/data@20180310-1902	844K	-
	5.57G -		
8	tupperware/backups/zestedesavoir/data@20180310-2303	0	-
	5.57G -		
9	tupperware/backups/zestedesavoir/data@20180311-0002	0	-
	5.57G -		
10	tupperware/backups/zestedesavoir/data@20180311-0102	0	-
	5.57G -		
11	...		

Les snapshots sont donc très légères, avec en général moins d'1 Mo par snapshot.

Le script pour synchroniser le dossier `/opt/zds/data` et les backups de la base, et faire la snapshot est relativement simple :

```
1  ##!/bin/sh
2
3  BASE=tupperware/backups/zestedesavoir
4
5  echo "Syncing data"
6  rsync -azvr --delete zdsprod1.zestedesavoir.com:/opt/zds/data/
   /$BASE/data
7  echo "Creating snapshot"
8  zfs snapshot $BASE/data@`date +%Y%m%d-%H%M`
9  echo "Syncing database"
10 rsync -azvr zdsprod1.zestedesavoir.com:/var/backups/mysql/
    /$BASE/db
```

Ce script est lancé 1 fois par heure, 2 min après la backup incrémentale de la base.

4. La cerise sur le gâteau : compression à la volée des volumes

ZFS permet de compresser des volumes entiers de manière transparente.

5. Conclusion

```
1 spaetzle/root ~ > zfs get
  used,compressratio,compression,logicalused
  tupperware/backups/zestedesavoir/{db,data}
2 NAME          PROPERTY      VALUE
  SOURCE
3 tupperware/backups/zestedesavoir/data  used          5.57G      -
4 tupperware/backups/zestedesavoir/data  compressratio 1.17x       -
5 tupperware/backups/zestedesavoir/data  compression   lz4
  local
6 tupperware/backups/zestedesavoir/data  logicalused   5.89G      -
7 tupperware/backups/zestedesavoir/db    used          666M       -
8 tupperware/backups/zestedesavoir/db    compressratio 3.11x       -
9 tupperware/backups/zestedesavoir/db    compression   lz4
  local
10 tupperware/backups/zestedesavoir/db    logicalused   1.88G      -
```

Ainsi, si les backups de la base semblent peser près de 2 Go, sur le disque ils ne prennent que 700 Mo.

Ça fait partie des fonctionnalités qui rendent ZFS *awesome*. Sachez que je n'ai couvert ici qu'une petite partie des fonctionnalités de ZFS.

5. Conclusion

Conclusion, j'aime ZFS.

À plus ou moins long terme j'aimerais mettre en place un autre serveur à un autre endroit, pour notamment avoir de la redondance à deux endroits physiques distincts sur mes backups. Oh, est-ce que j'ai mentionné que ZFS permet d'envoyer et recevoir des datasets entiers ou incrémentalement via de simples pipes ?

Et les backups de ZdS sont à nouveau un peu plus sérieuse, et on devrait pouvoir réagir plus rapidement en cas d'incident.