

Systèmes de fichiers distribué : comparaison de GlusterFS, MooseFS et Ceph avec déploiement sur la grille de calcul Grid'5000.

Jean-François Garçia, Florent Lévigne,
Maxime Douheret, Vincent Claudel

Table des matières

1	Introduction	3
1.1	Contexte	3
1.2	Système de fichiers distribué	3
1.3	Le Grid'5000	3
2	NFS	4
2.1	Présentation	4
3	GlusterFs	5
3.1	Présentation	5
4	MooseFS	6
4.1	Présentation	6
4.1.1	Présentation générale	6
4.1.2	Aspect technique	6
4.2	Mise en place	7
4.2.1	Environnement logiciel	7
5	Ceph	8
5.1	Présentation	8
6	Comparaison	9
6.1	Test de performances	9
7	Conclusion	10
A	Répartition des tâches	11
A.1	Florent Lévine	11
A.2	Jean-François Garcia	11
B	Scripts	12
B.1	GlusterFs	12
B.2	MooseFs	14
B.3	Ceph	14
B.4	NFS	16
B.5	Benchmark	17

Partie 1

Introduction

1.1 Contexte

Étudiants en licence professionnelle ASRALL (Administration de Systèmes, Réseaux, et Applications à base de Logiciels libres), notre formation prévoit une période de x mois à mi-temps pour la réalisation d'un projet tuteuré.

Le projet que nous avons choisi consiste à comparer diverses solutions de systèmes de fichiers distribués.

1.2 Système de fichiers distribué

Un système de fichiers (file system en anglais) est une façon de stocker des informations et de les organiser dans des fichiers, sur des périphériques comme un disque dur, un CD-ROM, une clé USB, etc. Il existe de nombreux systèmes de fichiers (certains ayant des avantages sur d'autres), dont entre autres l'ext (Extented FS), le NTFS (New Technology FileSystem), ZFS (Zettabyte FS), FAT (File Allocation Table).

Un système de fichiers distribué est un système de fichiers permettant le partage de données à plusieurs clients au travers du réseau. Contrairement à un système de fichier local, le client n'a pas accès au système de stockage sous-jacent, et interagit avec le système de fichier via un protocole adéquat.

Un système de fichier distribué est donc utilisé par plusieurs machines en même temps (les machines peuvent ainsi avoir accès à des fichiers distants, l'espace de noms est mis en commun). Un tel système permet donc de partager des données entre plusieurs clients, et pour certains de répartir la charge entre plusieurs machines, et de gérer la sécurité des données (par réplication)

1.3 Le Grid'5000

Partie 2

NFS

2.1 Présentation

NFS (Network File System, système de fichiers en réseau en français) est un système de fichiers développé par Sun Microsystems, permettant de partager des données par le réseau.

Partie 3

GlusterFs

3.1 Présentation

Partie 4

MooseFS

4.1 Présentation

4.1.1 Présentation générale

MooseFS (Moose File System) est un système de fichiers répartis à tolérance de panne, développé par Gemius SA. Le code préalablement propriétaire a été libéré et mis à disposition publiquement le 5 mai 2008. Il permet de déployer assez facilement un espace de stockage réseau, réparti sur plusieurs serveurs.

Cette répartition permet de gérer la disponibilité des données, lors des montées en charge ou lors d'incident technique sur un serveur. L'atout principal de MooseFS, au delà du fait qu'il s'agisse d'un logiciel libre, est sa simplicité de mise en œuvre.

En effet le tutoriel de prise en main, disponible sur le site du projet, explique de manière claire comment mettre en place une architecture distribuée en quelques heures. Concernant les utilisations, elles sont multiples et surtout, après la phase de configuration, l'évolution du système est très simple. L'ajout de serveurs, d'espace disque peuvent être gérés très facilement.

4.1.2 Aspect technique

Les montages du système de fichiers par les clients se font à l'aide de FUSE.

MooseFS est constitué de trois types de serveurs :

1. Un serveur de métadonnées (MDS)
Ce serveur gère la répartition des différents fichiers
2. Un serveur métajournal (Metalogger server)
Ce serveur récupère régulièrement les métadonnées du MDS et les stocke en tant que sauvegarde.
3. Des serveurs Chunk (CSS)
Ce sont ces serveurs qui stockent les données des utilisateurs.

Le point le plus important étant de bien dimensionner le serveur Master (qui stocke les métadonnées) afin de ne pas être limité par la suite. Donc pour ceux qui ne peuvent pas mettre en place des systèmes de stockage réseaux propriétaires assez coûteux, je vous conseille d'étudier cette possibilité. Elle vous permettra de partager des données sur plusieurs machines, de manière rapide, fiable, sécurisée et surtout peu coûteuse.

4.2 Mise en place

4.2.1 Environnement logiciel

Le système utilisé est une Debian Squeeze. MosseFs ne faisant pas partie des dépôts de la distribution, nous avons compilé le paquet à partir des sources dans leurs dernières version (1.6.20).

Partie 5

Ceph

5.1 Présentation

Partie 6

Comparaison

6.1 Test de performances

Afin de ne pas avoir de différence de matériel lors nos test, ceux-ci ont tous été réalisés sur un même cluster du Grid'5000 : Graphene.

Ce cluster est composé de 144 noeuds, avec pour caractéristique :

- 1 CPU Intel de quatre cœurs cadencé à 2.53 GHz
- 16 Go de RAM
- 278 Go d'espace disque

Nous avons réalisé un benchmark mesurant les performances (débit) de quatre type d'opérations sur le système de fichier distribué :

Écriture de petits fichiers : écriture des sources du noyau linux (décompressé).

Écriture de gros fichiers : écriture d'un fichier de 3 Go¹.

Lecture de petits fichiers : lecture des fichiers du noyau linux. Pour cela, nous avons compressé le dossier contenant le noyau (impliquant la lecture des fichiers), en redirigeant la sortie vers /dev/nul (afin que les performances du disque ne rentrent pas en jeux).

Lecture de gros fichiers : lecture du fichier de 3 Go. Opération réalisé en faisant un "cat" du fichier, et en redirigeant la sortie vers /dev/nul afin de ne pas "polluer" le terminal.

1. Fichier créé avec la commande : `dd if=/dev/zero of=/lieu/voulu bs=1G count=3`

Partie 7

Conclusion

Partie A

Répartition des tâches

A.1 Florent Lévigne

- Étude sur la mise en place de GlusterFS
- Réalisation d'un script de déploiement de GlusterFS
- Étude sur la mise en place de MooseFS
- Réalisation d'un script de déploiement de MooseFS
- Réalisation d'un script de benchmark pour système de fichiers distribué

A.2 Jean-François Garcia

- Étude sur la mise en place de GlusterFS
- Étude sur la mise en place de NFS
- Réalisation d'un script de déploiement de NFS
- Étude sur la mise en place de CephFS
- Réalisation d'un script de déploiement de CephFS

Partie B

Scripts

B.1 GlusterFs

Fichier deploymentGluster.rb :

```
1  #!/usr/bin/ruby -w
2  # encoding: utf-8
3
4  # r s e r v a t i o n  d e s  n o e u d s  (a  l a n c e r  m a n u e l l e m e n t)
5  # o a r s u b  -I  -t  d e p l o y  -l  n o d e s = 8 , w a l l t i m e = 2
6  # o a r s u b  -I  -t  d e p l o y  -l  n o d e s = 8 , w a l l t i m e = 2  -p  " c l u s t e r = ' g r a p h e n e ' "
7
8  # d o i t  c o n c o r d e r  a v e c  l a  c o m m a n d e  o a r s u b
9  n u m b e r O f C l i e n t s  = 5
10 n u m b e r O f S e r v e r s  = 3
11
12 i n f i n i b a n d  = 1  # 1 : a c t i v , 0 : n o n  a c t i v   ( n e  c h a n g e  r i e n  p o u r  l ' i n s t a n t )
13
14 # c r a t i o n  d ' u n  f i c h i e r  c o n t e n a n t  l a  l i s t e  d e s  n o e u d s  r s e r v s
15 ' t o u c h  l i s t O f N o d e s '
16 F i l e . o p e n ( " l i s t O f N o d e s " , ' w ' )  d o  | f i l e |
17   f i l e  <<  ' c a t  $ O A R _ F I L E _ N O D E S  |  s o r t  - u '
18 e n d
19
20 # c r a t i o n  d e  d e u x  f i c h i e r s  c o n t e n a n t  l a  l i s t e  d e s  s e r v e u r s ,  e t  d e s  c l i e n t s
21 ' t o u c h  l i s t O f C l i e n t s  l i s t O f S e r v e r s '
22 s e r v e r W r i t e d  = 0
23 F i l e . o p e n ( " l i s t O f N o d e s " , ' r ' )  d o  | n o d e |
24   F i l e . o p e n ( " l i s t O f S e r v e r s " , ' w ' )  d o  | s e r v e r |
25     F i l e . o p e n ( " l i s t O f C l i e n t s " , ' w ' )  d o  | c l i e n t |
26       w h i l e  l i n e  =  n o d e . g e t s
27         i f  s e r v e r W r i t e d  <  n u m b e r O f S e r v e r s
28           s e r v e r  <<  l i n e
29           s e r v e r W r i t e d  += 1
30         e l s e
31           c l i e n t  <<  l i n e
32         e n d
33       e n d
34     e n d
35   e n d
36 e n d
37
38 # d e p l o i e m e n t  d e s  m a c h i n e s
39 p u t s  " M a c h i n e s _ e n _ c o u r s _ d e _ d e p l o i e m e n t ... "
```

```

40 # 'kadeploy3 -k -a ../images/mysqueezegcluster-x64-base.env -f listOfNodes' # image
    perso
41 'kadeploy3 -k -e squeeze-collective -u flevigne -f listOfNodes' # image collective
42
43
44 # Envoie d'un script de cr ation d'un r pertoire dans /tmp/sharedspace sur les
    serveurs
45 File.open("listOfServers", 'r') do |file|
46     while line = file.gets
47         machine = line.split.join("\n")
48         'ssh root@#{machine} < createFolders.sh'
49     end
50 end
51
52 # Envoie d'un script de cr ation d'un r pertoire dans /media/glusterfs sur les
    clients
53 File.open("listOfClients", 'r') do |file|
54     while line = file.gets
55         machine = line.split.join("\n")
56         'ssh root@#{machine} < createMountDirectory.sh'
57     end
58 end
59
60 masterServer = 'head -n 1 listOfServers'.split.join("\n")
61
62 # g n ration des fichiers de conf, et envoie des fichiers de conf aux machines (
    serveurs et clients)
63 puts "Configuration_des_serveurs_et_des_clients..."
64 'scp listOfServers root@#{masterServer}:'
65 'scp listOfClients root@#{masterServer}:'
66 'scp glusterfs-volgen.rb root@#{masterServer}:'
67 'ssh root@#{masterServer} ./glusterfs-volgen.rb'
68 # 'ssh root@#{masterServer} < execScript/ex-glusterfs-volgen.sh'
69
70 # d marriage des serveurs
71 puts "D marriage_des_serveurs..."
72 File.open("listOfServers", 'r') do |file|
73     while line = file.gets
74         machine = line.split.join("\n")
75         'ssh root@#{machine} < startGluster.sh'
76     end
77 end
78
79 # montage du r pertoire par les clients
80 puts "Montage_du_r_pertoire_par_les_clients..."
81 File.open("listOfClients", 'r') do |file|
82     while line = file.gets
83         machine = line.split.join("\n")
84         'ssh root@#{machine} < mountFs.sh'
85     end
86 end
87
88 # r sum des machines
89 puts "GlusterFS_op rationnel"
90 puts "\nMachines_clients:"
91 puts 'cat listOfClients'

```

```

92
93 puts "\nMachines_serveurs:"
94 puts `cat listOfServers `
95
96 puts "\nServeur_maitre: #{masterServer}"
97
98 # nettoyage
99 #`rm listOfNodes listOfClients listOfServers `

```

B.2 MooseFs

B.3 Ceph

Fichier deploymentCeph.rb :

```

1  #!/usr/bin/ruby -w
2  # encoding: utf-8
3
4  #####
5
6  # File Name : deploymentCeph.rb
7
8  # Purpose :
9
10 # Creation Date : 11-03-2011
11
12 # Last Modified : jeu. 17 mars 2011 14:54:51 CET
13
14 # Created By : Helldar
15
16 #####
17
18 # doit concorder avec la commande oarsub
19
20 if ARGV[0] != nil
21   numberOfServers = ARGV[0].to_i
22   puts "Nb_serveur: #{numberOfServers}\n"
23 else
24   puts "Veuillez relancer le script avec les bons param tres!\n"
25   exit
26 end
27
28 # cr ation d'un fichier contenant la liste des noeuds r serv s
29 `touch listOfNodes `
30 File.open("listOfNodes", 'w') do |file|
31   file << `cat $OAR_FILE_NODES | sort -u`
32 end
33 # cr ation de deux fichiers contenant la liste des serveurs, et des clients
34
35 `touch listOfClients listOfServers `
36 serverWrited = 0
37 File.open("listOfNodes", 'r') do |node|
38   File.open("listOfServers", 'w') do |server|
39     File.open("listOfClients", 'w') do |client|

```

```

40         while line = node.gets
41             if serverWrited < numberOfServers
42                 server << line
43                 serverWrited += 1
44             else
45                 client << line
46             end
47         end
48     end
49 end
50 end
51
52 # d ploiment des machines
53 #puts "Machines en cour de d ploiment..."
54 #'kadeploy3 -k -e squeeze-collective -u flevigne -f listOfNodes' # image collective
55
56 # configuration du serveur
57 serveur_1 = 'head -1 listOfServers | cut -d "." -f1'.strip
58 ip_serveur = 'ssh root@#{serveur_1} hostname -i'.strip
59
60 # g n ration du fichier de ceph.conf
61
62 'touch ceph.conf'
63 File.open("ceph.conf", 'w') do |file|
64     file << "[global]
65     .....pid_file=_/var/run/ceph/$name.pid
66     .....debug_ms=_1
67     .....keyring=_/etc/ceph/keyring.bin
68     [mon]
69     .....mon_data=_/tmp/partage/mon$id
70     [mon0]
71     .....host=_#{serveur_1}
72     .....mon_addr=_#{ip_serveur}:6789
73     [mds]
74     .....debug_mds=_1
75     .....keyring=_/etc/ceph/keyring.$name"
76     if numberOfServers > 3
77         1.upto(3) { |i| file << "
78     [mds#{i-1}]"
79         host = 'sed -n #{i + 1}p listOfServers | cut -d '.' -f1'.strip
80         file << "
81     .....#{host}" }
82     else
83         file << "[mds0]"
84         host = 'sed -n 2p listOfServers | cut -d '.' -f1'.strip
85         file << "
86     .....#{host}"
87     end
88     file << "
89     [osd]
90     .....sudo=_true
91     .....osd_data=_/tmp/partage/osd$id
92     .....keyring=_/etc/ceph/keyring.$name
93     .....debug_osd=_1
94     .....debug_filstore=_1
95     .....osd_journal=_/tmp/partage/osd$id/journal

```

```

96  _____osd_journal_size__=1000"
97      1.upto(numberOfServers - 1) { |i| file << "
98  [osd#{i__1}]"}
99      host = `sed -n #{i + 1}p listOfServers | cut -d '.' -f1 '.strip
100      file << "
101  _____#{host}" }
102  end
103  # copie du fichier ceph.conf vers le serveur
104  `scp ceph.conf root@#{serveur_1}:/etc/ceph`
105  puts "Envoy !"
106  # g n ration du fichier keyring.bin
107  `ssh root@#{serveur_1} cauthtool --create-keyring -n client.admin --gen-key keyring
    .bin`
108  `ssh root@#{serveur_1} cauthtool -n client.admin --cap mds 'allow' --cap osd 'allow
    *' --cap mon 'allow rwx' keyring.bin`
109  `ssh root@#{serveur_1} mv keyring.bin /etc/ceph/`
110  puts "Keyring_g n r !"
111  # montage
112  `ssh root@#{serveur_1} mount -o remount,user_xattr /tmp`
113  1.upto(numberOfServers - 1) { |i| serveurs = `sed -n #{i + 1}p listOfServers | cut
    -d "." -f1 '.strip
114  `ssh root@#{serveurs} mount -o remount,user_xattr /tmp` }
115  puts "Montage_fait!"
116  # d mariage du serveur
117  `ssh root@#{serveur_1} mkcephfs -c /etc/ceph/ceph.conf --allhosts -v -k /etc/ceph/
    keyring.bin`
118  `ssh root@#{serveur_1} /etc/init.d/ceph -a start`
119  puts "Serveur_ceph_d marr !"
120  # configuration des clients
121  0.upto(`wc -l listOfClients` - 1) { |i| clients = `sed -n #{i + 1}p listOfClients |
    cut -d "." -f1 '.strip
122  `ssh root@#{clients} mkdir /ceph`
123  `ssh root@#{clients} cfuse -m #{ip_serveur} /ceph` }
124  puts "Clients_mont s!"

```

B.4 NFS

Fichier deploimentNFS.rb :

```

1  #!/usr/bin/ruby -w
2  # encoding: utf-8
3
4  #####
5
6  # File Name : deploimentNFS.rb
7
8  # Purpose :
9
10 # Creation Date : 17-03-2011
11
12 # Last Modified : jeu. 17 mars 2011 16:29:07 CET
13
14 # Created By : Helldar
15
16 #####

```



```

17
18 'cat $OAR_FILE_NODES | sort -u > listOfNodes '
19
20 # D ploiment des machines
21 #puts "Machines en cour de d ploiment...\n"
22 #'kadeploy3 -k -e squeeze-collective -u flevigne -f listOfNodes # image collective '
23
24 serveur = 'head -l listOfNodes '.strip
25 puts "Le_serveur_:_#{serveur}!\n"
26 # Suppression du serveur de la liste
27 'sed -i 1d listOfNodes '
28
29 puts "Configuration_du_serveur...\n"
30 'scp exports root@#{serveur}:/etc/'
31 'ssh root@#{serveur} /etc/init.d/nfs-kernel-server restart '
32
33 puts "Configuration_des_clients...\n"
34 line = 'wc -l listOfNodes | cut -d ' ' ' -f1 '.strip.to_i
35 puts "Il_y_a_#{line}_nodes"
36 1.upto(line) { |i| clients = 'sed -n #{i}p listOfNodes | cut -d "." -f1 '.strip
37 'ssh root@#{clients} mkdir /tmp/partage '
38 'ssh root@#{clients} mount #{serveur}:/tmp -t nfs /tmp/partage ' }

```

B.5 Benchmark

Fichier benchmark.rb :

```

1 #!/usr/bin/ruby -w
2 # encoding: utf-8
3
4 if "#{ARGV[0]}" == ""
5   puts "Doit_prendre_en_param tre_le_nombre_de_clients_participant_au_bench."
6   exit(1)
7 end
8
9 $clientsOfBench = "#{ARGV[0]}"
10
11 # chemin du fichier contenant la liste des clients
12 listOfClients = "/home/flevigne/glusterFs/listOfClients"
13
14 # chemin ou ecrire les donn es du benchmark
15 whereToWrite = "/media/glusterfs"
16
17 # chemin du fichier contenant les r sultats
18 $outputRes = "./resOfBench"
19
20 # le client doit avoir dans /home/flevigne :
21 # - linux-2.6.37.tar.bz2 : noyau linux compress
22 # - bigFile : un fichier de 3 Go
23
24 # fichier contenant la liste des clients participant au benchmark
25 'touch clientOfBench '
26 'head -#{ $clientsOfBench } #{listOfClients} > clientOfBench '
27
28 # si le fichier $outputRes n'existe pas, on le cr e .

```

```

29 if !File.exist?($outputRes)
30   'touch #{ $outputRes }'
31 end
32
33 'echo "\nBenchmark sur #{ $clientsOfBench } clients" >> #{ $outputRes }'
34
35 $numberOfClients = open("clientOfBench").read.count("\n").to_i
36
37 puts "Lancement du benchmarck sur #{ $numberOfClients } clients."
38
39
40 # lance un travail
41 # parametres :
42 # - name : nom du travail (str)
43 # - work : chemin du script de travail (str)
44 # - whereToWrite : chemin ou cree les donnees du benchmark (str)
45 # - size : taille (en Mo) du/des fichier(s) a ecrire/lire (float)
46 def startBench(name, work, whereToWrite, size)
47   puts "bench : #{name} en cours..."
48
49   totalSize = size.to_i * $clientsOfBench.to_i
50   workFinished = 0
51   startOfBench = Time.now
52
53   # execution du sript pour tous les clients
54   File.open("clientOfBench", 'r') do |file|
55     while line = file.gets
56       fork do
57         machine = line.split.join("\n")
58         'scp #{work} root@#{machine}:/root'
59         'ssh root@#{machine} ./#{work} #{whereToWrite}'
60         exit(0)
61       end
62     end
63   end
64
65   # on attend que tous les clients aient fini leur travail
66   1.upto($numberOfClients) do
67     pid = Process.wait
68     workFinished += 1
69     puts "Machine(s) ayant termin leur travail : #{workFinished}"
70   end
71
72   endOfBench = Time.now
73   duration = endOfBench - startOfBench
74
75   puts "Toute les machines ont termin leur travail."
76
77   puts "——> Le benchmark \"#{name}\" a dur #{duration} secondes. (debit : #{
     totalSize / duration } Mo/s)"
78
79   'echo " #{name} : #{duration} sec : #{totalSize / duration } Mo/s" >> #{ $outputRes
     }'
80 end
81
82

```

```

83 # lancement du benchmark
84 startBench("écriture_de_petits_fichiers", "writingSmallFiles.sh", whereToWrite,
    479)
85 startBench("écriture_de_gros_fichiers", "writingBigFiles.sh", whereToWrite, 3076)
86 startBench("lecture_de_petits_fichiers", "readingSmallFiles.sh", whereToWrite, 479)
87 startBench("lecture_de_gros_fichiers", "readingBigFile.sh", whereToWrite, 3076)
88
89 # nettoyage du syst me de fichier distribue (necessaire pour enchaîner les
    benchmark)
90 puts "Nettoyage_de_l'espace_de_travail..."
91 oneClient = 'head -1 clientOfBench'.strip
92 'ssh root@#{oneClient} rm -r #{whereToWrite}/*'
93
94 puts "\nBenchmark_termine"

```

Fichier writingSmallFiles.sh :

```

1 #!/ bin / bash
2
3 whereToWrite=$1
4
5 nameOfMachine='uname -n'
6
7 # creation du repertoire de travail de la machine
8 mkdir "$whereToWrite/$nameOfMachine"
9
10 # decompression dans ce repertoire
11 cd "$whereToWrite/$nameOfMachine"
12 tar -xf /home/flevigne/linux-2.6.37.tar.bz2
13
14 # copie dans le rep partage
15 #cp -r /home/flevigne/linux-2.6.37 "$whereToWrite/$nameOfMachine"

```

Fichier writingBigFiles.sh :

```

1 #!/ bin / bash
2
3 whereToWrite=$1
4
5 nameOfMachine='uname -n'
6
7 # on copie le gros fichier au lieu voulu
8 cp /home/flevigne/bigFile "$whereToWrite/$nameOfMachine"

```

Fichier readingSmallFiles.sh :

```

1 #!/ bin / bash
2
3 whereToWrite=$1
4
5 nameOfMachine='uname -n'
6
7 cd "$whereToWrite/$nameOfMachine"
8
9 # lecture des fichiers du noyau linux (compression (donc lecture) redirig vers /
    dev/nul)
10 tar -cf /dev/null linux-2.6.37

```

Fichier readingBigFile.sh :

```
1 #!/bin/bash
2
3 whereToWrite=$1
4
5 nameOfMachine=$(uname -n)
6
7 cd "$whereToWrite/$nameOfMachine"
8
9 # lecture du gros fichier
10 cat bigFile > /dev/nul
```