

Déduplication extrême d'OS avec SIDUS : un petit pas pour la reproductibilité ?

Emmanuel Quemener

Ecole Normale Supérieure de Lyon,
Centre Blaise Pascal - 46, allée d'Italie
69007 Lyon - France
Emmanuel.Quemener@ens-lyon.fr

Résumé

Lors d'une expérience numérique menée sur un système informatique complet, comment s'assurer que le composant «*Operating System*» reste bien identique, d'un redémarrage à l'autre sur une même machine, ou au même instant sur un ensemble de nœuds ? Initialement, ce n'est pas en réponse à cette épineuse difficulté que SIDUS a été développé : SIDUS l'a été essentiellement pour simplifier l'administration de la majorité des plateaux techniques du Centre Blaise Pascal à l'ENS-Lyon.

En effet, SIDUS simplifie la gestion de nœuds de calcul, de stations de travail ou de nouveaux matériels, limite l'empreinte du stockage de l'OS sur ses machines, voire offre à un utilisateur un environnement scientifique complet en quelques secondes. Scalable et polyvalent, SIDUS a comme principales propriétés : l'unicité de son socle sur un serveur et l'usage des ressources locales de chaque client.

Dès lors, si seules les ressources locales sont exploitées et si ce socle technique unique demeure inchangé, alors la comparaison entre les expériences numériques que nous menons gagnent en pertinence. A contrario, si la reproductibilité fait défaut, les origines sont à rechercher ailleurs que dans l'OS.

Dans un cadre opérationnel, nous illustrerons que le simple examen de durées d'exécution suffit à découvrir cette variabilité temporelle : reste à en déterminer la nature. Étonnement, c'est essentiellement dans les BIOS des matériels ou dans leur vieillissement que ces variabilités trouvent leur source.

Mots-clés : déduplication, variabilité temporelle, reproductibilité.

1. Introduction

Nous aborderons, en premier lieu, la variabilité inhérente à celle de toute expérience numérique. Puis nous chercherons des solutions permettant d'extraire la partie «*Operating System*» du système informatique dans son ensemble. Ensuite, SIDUS, la solution que nous proposons, sera décrite au travers de questions fondamentales, notamment comment l'installer voire l'administrer simplement. Enfin, au travers d'expériences dans un cadre opérationnel, nous illustrerons la nature de cette variabilité, ses sources et comment les combattre pour gagner en pertinence.

2. La reproductibilité sous le regard d'un expérimentateur

2.1. Rapide état des lieux de l'informatique scientifique

Chaque traitement, chaque simulation et même chaque visualisation lancés sur un système informatique sont des «expériences numériques». Depuis plus de 30 ans, depuis l'extinction massive des calculateurs analogiques, nous pensions vivre dans le fabuleux monde du déterminisme numérique. En fait, une rapide introspection sur le traitement des opérations en virgule flottante modère ce point de vue[1] : des erreurs d'arrondi affectent une majorité de nos calculs et les différentes approches de parallélisme les rendent de moins en moins prédictibles. A cela s'ajoute désormais un nouveau type de variabilité : les processeurs disposent de plus en plus de cœurs, les systèmes d'exploitation nécessitent de plus en plus de processus pour conserver actifs, en totalité, leurs services courants ; les réseaux utilisent des méthodes d'accès non déterministes à leurs médias de communication, les processeurs embarquent trois niveaux de cache et des douzaines d'ALU (unités arithmétiques et logiques) ou d'UC (unités de contrôle) ; les fréquences des processeurs voire de la mémoire changent continuellement. Pour «couronner le tout», les processeurs récents n'ont plus qu'une seule contrainte : respecter leur enveloppe thermique maximale (la TDP ou *Thermal Design Power*) ou des stratégies d'économie d'énergie au détriment de tout le reste !

«Cerise sur le gâteau», les applications tournent parfois sur des machines différentes (de la douzaine à des milliers) pour lesquelles il est très difficile de s'assurer qu'elles partagent exactement le même système d'exploitation.

2.2. Qu'espérer d'une infrastructure devenue complexe ?

Ce panorama établi, comment pouvons-nous réduire autant que possible la variabilité de notre système dans son ensemble, la variabilité dans les résultats eux-mêmes, ou, c'est le cas de notre étude, la variabilité sur les durées d'exécution sur nos systèmes ? Comment s'assurer que le système reste identique d'un redémarrage à l'autre (sur une même machine) ou que toutes les machines disposent, au même instant, du même système au bit près ?

Dans le même ordre d'idée, comment pouvons-nous réduire l'effort de portage de notre station de travail ou de petits clusters départementaux vers les grosses machines ?

2.3. De l'émergence d'un système unique

2.3.1. Le pragmatisme d'une distribution complète

La première étape relève du choix de «la» distribution offrant le plus large panorama possible de logiciels scientifiques précompilés et préconfigurés afin d'éviter à tout prix toutes inutiles recompilation et réinstallation du socle de nos codes : la distribution Debian (ou toute autre distribution aussi complète que la Debian) est une approche pragmatique. L'assurance qualité de Debian est aussi une des plus tatillonnes qui soit : ce serait dommage de ne pas en profiter et de réinventer la roue ! A titre de petite expérience, nous invitons les sceptiques à installer *from scratch* le logiciel de chimie quantique Abinit[2].

2.3.2. De l'instantané à l'instance unique

La seconde étape consiste alors à fournir ce système d'exploitation issu de cette distribution à nos clients. Tout d'abord, il nous faut figer cet OS pour le retrouver exactement dans le même état à une date ultérieure. Ces solutions de «photographie» de système sont légion, ce sont essentiellement des solutions de sauvegarde (MondeRescue, SystemImager, CloneZilla) ou sont proches du monde HPC comme Kadeploy. Ensuite, une fois les systèmes sauvegardés, ces outils se proposent de restaurer un système sous forme d'une ou plusieurs partitions. Enfin, le

redémarrage des toutes les machines ainsi fraîchement installées : toute variabilité est absente d'un système à l'autre.

Toutefois, la gestion de ce type d'infrastructure devient rapidement très lourde dès lors que le système devient imposant. A titre d'illustration, une distribution *Debian Science* complète occupe un espace de quarantaine de giga-octets. Même avec d'improbables taux de compression (de 90%), démarrer une centaine de machines va générer 400 Go de trafic pour le serveur, avec à la clé vingt bonnes minutes d'écriture sur chaque disque dur, donc une demi-heure d'immobilisation, sans même évoquer la partie création de la partition «maître»...

Une autre approche vise à s'affranchir de la partie «stockage local» en offrant un système directement par le réseau. Là le protocole iSCSI est parfait pour offrir un système en «mode bloc». Le boot PXE permet, par TFTP, la récupération d'un noyau Linux, d'une séquence de démarrage InitRD. Ensuite, le démarrage permet le montage distant par un client iSCSI au InitRD. Les «cibles» iSCSI existent suivant deux implémentations sous Linux. Pour leur duplication, les mécanismes d'instantanés sont à trouver auprès des gestionnaires de volumes logiques comme LVM, ZFSonLinux voire BTRFS. Ainsi, chaque instantané, chaque «image système» est clonée puis va servir chaque nœud individuellement par iSCSI. Cette méthode est efficace, mais exige la maîtrise d'un environnement *chrooté* pour la création du système de base. Généralisée sur une centaine de nœuds, elle devient difficile à dimensionner : ce n'est pas tant le réseau mais le serveur qui se retrouve victime d'un fabuleux *boot storm*. Tout est donc à recommencer à chaque redémarrage pour s'assurer de la cohérence avec l'image originelle.

La solution que nous proposons dépasse toutes ces limitations en exploitant le partage réseau le plus simple qui soit : NFS. C'est un environnement configuré une seule et unique fois et déployé sur tous les clients, où chaque modification est instantanément propagée sur toutes les autres : SIDUS pour *Single Instance Distributing Universal System*, n'est pas qu'un projet. Il est pleinement opérationnel depuis 4 années au Centre Blaise Pascal (maison de la simulation lyonnaise) sur une centaine de machines permanentes et au Pôle Scientifique de Modélisation Numérique depuis 2 ans (un des méso-centres lyonnais) sur plus de 330 nœuds.

2.4. Qu'allons nous présenter ?

Dans une première partie, nous détaillerons SIDUS par le truchement de réponses à des questions fondamentales. Qu'est-ce que SIDUS ? Où & quand est-il exploité ? A qui peut-il servir ? Comment s'installe-t-il ? Comment s'administre-t-il ?

3. SIDUS : un couteau suisse ?

SIDUS est l'acronyme de *Single Instance Distributing Universal System*.

Son origine latine «d'ensemble de corps stellaires» est une allégorie. Ainsi SIDUS partage le même système d'exploitation avec des machines aux ressources matérielles différentes. Tout comme les étoiles d'une constellation, aux observables physiques différentes, appliquent les mêmes mécanismes de fusion nucléaire. SIDUS a comme principales propriétés :

- son unicité de configuration : deux machines démarrant sous SIDUS ont exactement le même système d'exploitation ;
- son exploitation des ressources locales : les processeurs et mémoire vive sollicités sont ceux de la machine locale.

SIDUS n'est donc :

- ni LTSP pour *Linux Terminal Server Project* : LTSP propose une gestion simplifiée de terminaux légers en offrant un accès X11 ou RDP à un serveur. Ce dernier supporte ainsi toute la charge de traitement. A contrario, SIDUS exploite entièrement (ou à discrétion de l'utilisateur) toute

- la machine qui s'y raccroche. Seul le stockage du système d'exploitation est déporté sur des machines tierces ;
- ni FAI pour *Fully Automatic Installation* : FAI et Kickstart proposent une installation complète simplifiée permettant de limiter voire d'éliminer toute action de l'administrateur. A contrario, SIDUS propose un système unique dans un arbre intégrant à la fois le système de base et toutes les applications installées manuellement ;
 - ni un LiveCD sur réseau : un LiveCD démarre un système minimaliste, nécessairement figé. Il est toujours possible de créer son propre LiveCD mais c'est une opération lourde. Avec SIDUS, il est possible d'installer à la volée sur tous ses clients un nouveau composant instantanément ou de reconfigurer l'instance ;
 - ni monolithique : dans le cadre de formations en informatique, trois solutions s'offrent aux utilisateurs : exploiter les machines mises à disposition, utiliser leur équipement personnel, installer un environnement virtuel complet figé au téléchargement et donc difficilement modifiable. A contrario, SIDUS offre un environnement unique aisément configurable ;
 - ni original : SIDUS exploite des services disponibles sur n'importe quelle distribution : DHCP, PXE, TFTP, NFSroot, DebootStrap, AUFS. Ces quelques mots clés permettant d'installer SIDUS. Il utilise en outre des astuces de distributions de LiveCD et fonctionne sur la distribution Debian depuis sa version Etch.

En définitive, SIDUS est :

- universel : toutes les plates-formes x86 ou x86-64 fonctionnent instantanément ;
- efficace : installation en quelques dizaines de minutes, démarrage en quelques secondes ;
- économe : à l'origine, 1 cœur, 1Go de RAM, 40Go d'espace disque, et un réseau (GigaBit) Ethernet
- scalable : éprouvé sans difficulté sur une centaine de nœuds, en production maintenant sur plus de 330 nœuds ;
- robuste : avec un réseau standard routé, des *uptime* de plusieurs mois dans sa version COMOD (pour *Compute On My Own Device*) ;
- polyvalent : avec la Debian et tout «science», un excellent socle pour toutes les sciences.

3.1. Où & Quand , ou quelles *success stories* ?

Sur quels types d'équipements SIDUS peut-il se déployer ? C'est là que s'illustre sa polyvalence :

- **Sur les postes utilisateurs** : machines mutualisées ou stations de travail individuelles ? Tout a commencé avec une douzaine de clients légers Neoware gonflés en mémoire et overclockés. Ce sont maintenant plus de 20 machines équipées de cartes graphiques différentes et offertes à 300 utilisateurs de l'ENS-Lyon ;
- **Sur les nœuds de cluster** :
 - après un démonstrateur de 24 nœuds en mars 2010 au Centre Blaise Pascal, SIDUS sert actuellement 76 nœuds permanents, sur 3 architectures matérielles différentes,
 - après une période de qualification d'une année au Pôle Scientifique de Modélisation Numérique sur quelques dizaines de nœuds, SIDUS équipe maintenant plus de 330 nœuds, dont le nouvel équipement Equip@Meso,
 - après seulement quelques heures de configuration à l'Institut de Génomique Fonctionnelle de Lyon, sur 6 nœuds, comme puissance de traitement d'un serveur Galaxy ;
- **Sur les postes virtuels** : depuis 2011, l'Université Joseph Fourier organise chaque année une école d'été sur le calcul numérique en physique. Au programme, 10 jours intenses ponctués de travaux pratiques : offrir un environnement homogène quasi-instantanément est indispensable. Co-organisateur de ces écoles, le CBP met en place deux images virtuelles

- de systèmes : l'une autonome utilisable après l'école d'été, l'autre par SIDUS nécessitant seulement une connexion réseau filaire. Ainsi, les professeurs peuvent-ils quotidiennement adapter leurs TP. C'est une évolution de cette version qui est utilisée, depuis l'été 2012, par le laboratoire de chimie de l'ENS-Lyon et proposée aux laboratoires de biologie LBMC et IGFL ;
- **Sur les machines suspectes** : le démarrage par le réseau offre une investigation de la mémoire de masse système éteint : inutile d'utiliser un LiveCD sur lequel manque toujours son outil *forensics* préféré ;
 - **Sur les machines de prêt** : les fabricants de matériels proposent souvent des équipements d'évaluation. La phase d'installation peut être pénible sur des matériels très (voire trop) récents. Avec SIDUS, le système démarre comme sur les autres équipements déjà en service : quelques minutes pour 20 nœuds.

3.2. Pour qui : Quels avantages ?

- **Côté utilisateur** : la machine démarre avec seulement les ressources associées. La version VirtualBox fonctionne au moins sur Linux, Windows et MacOSX : accélération 3D et partage avec l'hôte via un dossier partagé sont disponibles. L'utilisateur retrouve exactement le même environnement que sur les nœuds : l'intégration des codes est donc grandement facilitée. Côté performances, les pertes liées à la virtualisation oscillent entre 10 et 20% (pour VirtualBox) et autour de 5% pour KVM ;
- **Côté administrateur** : une opération impacte l'ensemble de l'infrastructure, de l'ordre du simple sync sur l'arbre SIDUS. L'installation se déroule en quelques dizaines de minutes pour un système complet. Si des différences entre les systèmes sont minimes, un usage simple de scripts ou de Puppet suffit. Si la différence entre les systèmes est importante, un autre arbre SIDUS est construit, voire cloné instantanément avec des mécanismes de *snapshots* : LVM, ZFSonLinux ou Btrfs ;
- **Côté expérimentateur** : ingénieur système ou scientifique, l'environnement SIDUS lui offre la reproductibilité. Deux nœuds démarrant sur le même socle SIDUS disposent exactement du même système au bit près. Une même machine redémarrant SIDUS va retrouver le même système que celui précédemment démarré, quelles que soient les modifications entreprises. Cela permet ainsi, que les machines soient identiques ou pas, de mener des tests vraiment pertinents.

3.3. Combien : Quelles ressources ?

A titre d'exemple, le serveur des clusters du CBP, également passerelle, héberge les services DHCP, DNS, TFTP, NFS et le serveur de batch OAR. Au démarrage de toute l'infrastructure (76 nœuds), le serveur NFS encaisse sans broncher jusqu'à 900 Mb/s.

Le serveur d'instance SIDUS du PSMN, durant la phase de test, était une machine archaïque (Sunfire v40z) : elle remplissait sans souci le service de plus de 330 nœuds avant son remplacement par une machine plus récente.

3.4. Comment installer le système ?

3.4.1. En quelques lignes...

Le pré-requis est réduit, à commencer par un réseau idéalement comparable au débit d'un disque dur. Côté services, sont nécessaires : serveurs DHCP, DNS, TFTP et NFS. Les deux derniers vont «porter» SIDUS. La version opérationnelle du CBP exploite en outre des serveurs tiers LDAP (identification/authentification) et NFSv4 (espaces utilisateurs). Côté client, seul suffit un démarrage PXE opérationnel (par la carte, ou par GPXE sur CDROM ou clé USB).

L'installation comporte 8 phases :

1. préparation du système
2. installation de base (socle Debian, debootstrap)
3. installation des paquets complémentaires (TOUT Debian-Science)
4. purge des paquets non désirés
5. adaptation du système à l'environnement local
6. pointage du système vers les serveurs tiers : authentification et partages utilisateurs
7. création de la séquence de démarrage
8. détachement de SIDUS du système hôte

Durant l'installation, les phases coûteuses sont le téléchargement des paquets et le paramétrage de quelques composants (Perl et LaTeX). Elle dure au mieux 45 minutes pour un arbre complet de 32 Go. Quelques précautions sont cependant nécessaires, liées au montage des dossiers systèmes et à l'inhibition du démarrage des services à leur installation dans SIDUS.

Comment maintenant l'offrir sans le dupliquer ? Une première approche a été d'utiliser un corège de montages volatils (à base de TMPFS) : elle n'est pas viable. La préférence s'est portée sur mécanisme de LiveCD très répandu : la séquence de démarrage intègre ainsi la superposition de deux couches par le liant AUFS (évolution de UnionFS), l'une lecture seule (le CDRROM pour le LiveCD et NFS chez nous), l'autre lecture/écriture (en TMPFS).

Mais comment bénéficier de SIDUS et disposer d'un paramétrage conservé d'un démarrage à l'autre ? La première approche avec un montage NFS exclusif pour chaque nœud a été abandonnée, remplacée par un montage iSCSI associé à chaque nœud. Actuellement, au CBP, les machines SIDUS nécessitant une persistance (comme les nœuds Distonet) utilisent le mécanisme NFSroot+iSCSI=AUFS, les autres NFSroot+TMPFS=AUFS.

Pour conclure, les machines mises à disposition sont assez hétérogènes : les nœuds de clusters (disposant d'équipements réseau rapides), les stations de travail (embarquant des cartes graphiques) ou les machines virtuelles (exigeant un partage des données et une accélération graphique) demandent quelques adaptations. Une première solution serait la persistance, mais trop lourde pour les grands parcs de machines : seront alors préférées l'utilisation de scripts de démarrage, l'exploitation d'un arbre SIDUS séparé ou l'installation de composants tierces.

Dans la suite de l'article, nous ne détaillerons que les points nous semblant indispensables, les autres étant disponibles dans la documentation officielle[4] ou dans la presse spécialisée[5].

3.4.2. Préparation du système

Nous devons préparer un peu notre système afin d'accueillir SIDUS. Nous avons la main sur plusieurs services pour déployer nos clients : serveurs DHCP, TFTP, NFS. Nous entretenons de très bonnes relations avec notre service IT ou nous sommes assez libres pour accéder sans contraintes aux serveurs LDAP et DNS bien définis :

- le service DHCP fournit à notre client une adresse IP mais diffuse deux informations complémentaires : l'adresse du serveur TFTP via la variable `next-server` et le nom du binaire PXE, souvent nommé `pxelinux.0`.

- le service TFTP entre alors en scène. Il offre par TFTP tout le nécessaire permettant le démarrage du système : le binaire `pxelinux.0`, le noyau et le démarrage du système du client. Si nous avons besoin d'offrir des paramètres à tel ou tel client, nous construisons un document spécifique dont le nom sera construit à partir de son adresse MAC (préfixé de 01 et dont les «:» sont remplacés par des «-»).
- le serveur NFS s'invite alors dans la boucle : il va offrir la racine du système par son protocole (donc NFSroot). C'est donc dans cette racine, par exemple `/srv/nfsroot/sidus` que nous allons installer notre système client.

Sur nos configurations nous utilisons respectivement `isc-dhcp-server`, `tftpd-hpa` et `nfs-kernel-server` pour les serveurs DHCP, TFTP et NFS.

3.4.3. Installation de base par Debootstrap

Debootstrap permet l'installation d'un système dans une racine. Il exige plusieurs paramètres comme la racine d'installation, l'architecture matérielle, la distribution et l'archive FTP ou HTTP Debian à solliciter pour le téléchargement. Là commence la «spécialisation» de notre installation, à l'origine construite autour d'une distribution Debian. Cet outil est familier de toutes les distributions *Debian-like* : il sera donc disponible chez les dérivées du système à la spirale (à commencer par Ubuntu). Il sera cependant assez facile de réaliser cela sur les Redhat-like, Fedora intégrant par exemple un clone, Febootstrap, mais que nous n'avons pas souhaité tester.

Debootstrap accepte aussi en entrée une liste d'archives (vous savez que Debian est très tatillon sur les licences en séparant les archives en main, contrib et non-free), une liste de paquets à inclure et une liste de paquets à exclure. Nous aurions été ravis de pouvoir, ici, préciser la liste complète des paquets à inclure ou à exclure, mais, malheureusement, cette approche est une voie sans issue : nous installerons donc, dès cette commande `debootstrap`, un ensemble d'outils indispensables (par exemple le noyau, des micro-logiciels pour un support étendu de tous les matériels et un ensemble d'outils d'audit).

3.4.4. Précautions & création d'un «cordon ombilical» avec l'hôte

A la suite de cette commande, nous devons prendre quelques précautions : normalement, si le paquet Debian est un service, ce dernier démarre après son installation. Nous devons donc inhiber le lancement de ce service par la définition d'un hook. Ce hook sera supprimé à la fin de l'installation ; des paquets exigent l'accès à la liste des processus, du système, des périphériques, de la mémoire virtuelle, des pointeurs de périphériques. Nous devons donc « bind » le montage de ces dossiers du système hôte au système SIDUS. Les dossiers concernés sont : `/proc`, `/sysfs`, `/dev/shm`, `/dev/pts`.

3.4.5. Création de la séquence de démarrage

Comment partager SIDUS sans le dupliquer ? Nous allons nous inspirer d'un mécanisme utilisé dans certains LiveCD : le montage de la racine du système consiste en la superposition de deux couches, l'une en lecture seule (le système NFSroot) et l'autre en lecture/écriture (un TMPFS dans le cas le plus simple). Les deux couches sont liées par la glue AUFS, le projet successeur de UnionFS.

Tout réside dans un seul et unique hook : `rootaufs`, placé très tôt dans le démarrage `initrd`. Son principe repose sur cinq étapes :

1. création de dossiers temporaires `/ro`, `/rw` et `/aufs` ;

2. déplacement de la racine NFSroot du point de montage originel vers dans `/ro` ;
3. montage d'un partage distant, d'une partition locale ou distante, ou d'un volume TMPFS dans un autre point de montage `/rw` ;
4. superposition des deux dossiers `/ro` et `/rw` dans le dossier `/aufs` ;
5. déplacement de `/aufs` vers le point de montage originel.

Ce script `rootaufs` se place dans `$$SIDUS/etc/initramfs-tools/scripts/init-bottom`. Le script original a été inspiré par le projet `rootaufs`[11] de Nicholas A. Schembri. Il a été profondément modifié pour l'adapter à notre infrastructure : une version est disponible en ligne[10]. Un dernier petit effort avant de disposer d'un système fonctionnel : nous allons créer par `update-initramfs` un InitRD spécifique pour notre boot NFS contenant les modifications suivantes :

- `aufs` dans `/etc/initramfs-tools/modules`
 - `eth0` comme `DEVICE` dans `/etc/initramfs-tools/initramfs.conf`
- Il suffit ensuite de copier le noyau et le InitRD dans la définition du serveur TFTP.

A l'origine, nous avons exploré la possibilité d'offrir un second partage NFS en lecture/écriture pour la persistance des modifications associées aux clients d'un redémarrage à l'autre. Cette version, bien que fonctionnelle, exigeait l'ouverture d'un partage NFS atomique pour chaque client : imaginons la charge supplémentaire pour le serveur !

Nous avons donc préféré une autre approche de la persistance, sous la forme d'un disque réseau de technologie iSCSI. Ainsi, nous créons un volume partagé iSCSI par client. Pour des raisons de simplicité, le volume offert porte l'IP du client et nous ne fournissons que le serveur de volume iSCSI. Les login et mot de passe d'accès par défaut sont dans le fichier `rootaufs`.

3.4.6. Rupture du «cordon ombilical» avec le système hôte

De manière à installer correctement SIDUS, nous avons été contraints de lier étroitement système hôte et chroot. Il est donc nécessaire de :

- démonter les dossiers système : `/proc`, `/sysfs`, `/dev/shm`, `/dev/pts` ;
- effacer les dossiers temporaires ;
- purger des processus liés au dossier d'installation de l'instance SIDUS au besoin.

3.5. Comment administrer le système ?

Si l'administration est plus lourde que son installation, le bénéfice du premier efface la perte récurrente du second. Finalement, avec SIDUS, chaque phase d'administration intègre les mêmes mécanismes qu'à l'installation : protection contre le démarrage des nouveaux services et montages de dossiers systèmes. Le reste est identique.

En définitive, une instance SIDUS s'administre comme tout système chrooté : il y a cependant des précautions à prendre, toutes les fonctions d'administration n'exigeant pas les mêmes contraintes. Souvent, un chroot sur la racine de SIDUS suivi de la commande `systeme` suffit. Une nouvelle approche, basée sur l'ouverture d'un SSH directement dans l'instance SIDUS, est en préparation.

4. Expérimentation

4.1. De GlusterFS sur InfiniBand aux effets néfastes du C-states

4.1.1. Du choix d'un espace temporaire partagé au protocole expérimental

Le Pôle Scientifique de Modélisation Numérique, centre de calcul de l'ENS-Lyon, a choisi depuis longtemps l'internalisation de la gestion de son infrastructure : tout nouvel équipement se borne généralement à l'achat de matériel.

Début 2013, l'arrivée du nouvel équipement Equip@Meso (dans le cadre des investissements d'avenir) exigeait l'installation et donc le choix d'un espace de partage haute performance entre les nœuds. Le Centre Blaise Pascal avait déjà, dans le cadre du projet DistoNet[7], étudié plusieurs systèmes de fichiers répartis, notamment GlusterFS : il était donc intéressant que le centre de production (le PSMN) demande au centre d'essais (le CBP) d'étudier GlusterFS sur des nœuds équivalents aux nœuds du futur équipement, disposant d'une interconnexion haute passante et basse latence sous Infiniband. Vingt nœuds Hewlett Packard SL230, disposant de 16 cœurs physiques et 64 Go de RAM, interconnectés en Infiniband, ont donc été mis à disposition pour cette étude.

4.1.2. Le choix d'un RAMdisk

Tout d'abord, la première étape fut le démarrage de ces machines dans l'environnement du CBP : cela ne prit que quelques minutes avec SIDUS, tout au plus une demi-journée avec la récupération de l'adresse MAC et la configuration du contrôle IPMI. Puis, d'un point de vue expérimental, il était indispensable de se libérer de tout goulot d'étranglement matériel de stockage en supprimant toute latence disque. Le choix de système de fichiers en mémoire vive s'imposait de lui-même : TMPFS en est une implémentation courante. Une seconde possibilité se basait sur le module BRD créant un «mode bloc» lequel devait être formaté pour ensuite être offert. L'analyse de presque une dizaine de systèmes de fichiers ont montré que le système le plus efficace pour cet usage était, sans surprise, le plus simple : ext2. Ensuite, il fallait offrir le volume GlusterFS pointant vers cet espace de *Ramdisk* et le monter à l'aide d'un ou plusieurs clients. Enfin, venait l'évaluation de la performance d'accès exploitant massivement l'outil IOzone3[3], sur ces espaces partagés. De manière à disposer d'échantillon représentatif, la première expérience exploitait dix clients se connectant sur dix serveurs différents. Les machines (clients comme serveurs) étaient strictement identiques, simplement sorties de leurs cartons. Pour chacun des couples client-serveur, vingt expériences furent menées.

Les premiers résultats ont été très surprenants et sont illustrés sur les deux schémas «radar» de la figure 1, *BIOS Initial, Noeud 1 sur 11* et *BIOS Initial, Noeud 3 sur 13*. Dans le second cas, des performances autour de 500 Mo/s très stables et homogènes. Dans le premier, de bien meilleures performances, mais une variabilité très surprenante atteignant des valeurs du simple au double ! Etant donné que nous étions dans un environnement SIDUS, les configurations ne pouvaient pas ne pas être identiques ! Les machines démarraient le même noyau, lançaient les mêmes séquences de démarrage. A noter que, sur les dix couples, deux seulement présentaient ces variabilités, les huit autres restant d'une intéressante stabilité.

Après des heures d'investigation et de tests multiples, l'examen des BIOS respectifs des nœuds montra des différences : dans le cas où les machines présentaient des résultats stables mais peu performants, la gestion de la consommation était celle par défaut. Dans l'autre cas, teinté de performance mais de forte variabilité, un des deux nœuds était en *Max Performance*. Le passage de ces réglages en *Max Performance* a permis de retrouver les deux schémas «radar» de la figure 1, *BIOS corrigé, Noeud 1 sur 11* et *BIOS corrigé, Noeud 3 sur 13*. Tous les couples présentaient des taux de transfert d'une remarquable stabilité autour de 1 Go/s.

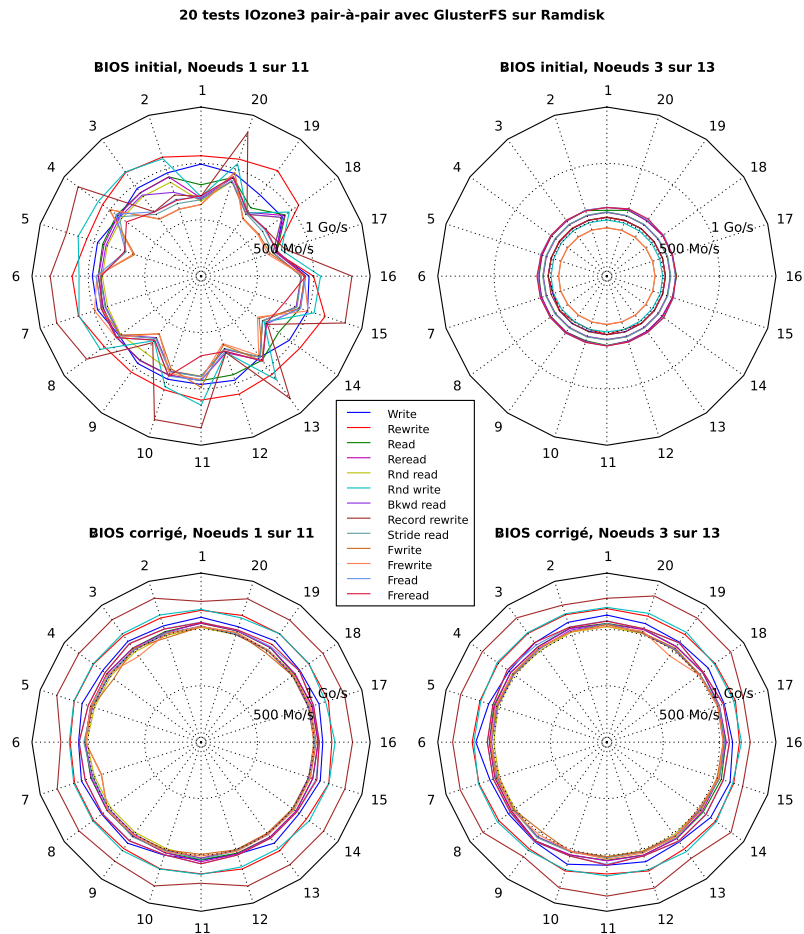


FIGURE 1 – Série de 20 tests IOzone3 entre un client et un serveur, pour deux couples différents de machines.

L'accélération générale (rapport entre la médiane des durées d'exécution totale pour les 20+20 expériences) et la diminution de la variabilité (rapport entre la médiane et l'écart-type sur les 20+20 expériences) sont illustrées sur la figure 2 pour les dix couples clients serveurs ayant participé à la campagne d'essais. La variabilité a été divisée d'un facteur 30 pour les deux premiers couples et la performance d'un facteur 2 pour les 8 derniers.

4.2. Lorsque la non-reproductibilité devient ... reproductible

Mi-2013, lors de l'intégration opérationnelle de GlusterFS sur l'Equip@Meso de l'ENS-Lyon, des variabilités d'un ordre de grandeur affectaient la communication entre deux nœuds via InfiniBand en IP, même avec des réglages identiques à ceux de la première expérience (mais sur des matériels différents). Trouver le bon paramétrage de BIOS (différent de celui évoqué ci-dessus) a été long et fastidieux mais a permis de limiter la variabilité à une valeur de quelques pour-cent. Ainsi, les mécanismes d'autoconfiguration énergétiques, trop «laissés» au processeur, génèrent une variabilité dont l'écart-type laisse perplexe : les intégrer dans le processus même de l'expérience numérique devient une exigence ! Pour parodier Clémenceau, «*La gestion de l'énergie est une chose trop sérieuse pour la confier au processeur lui-même !*»

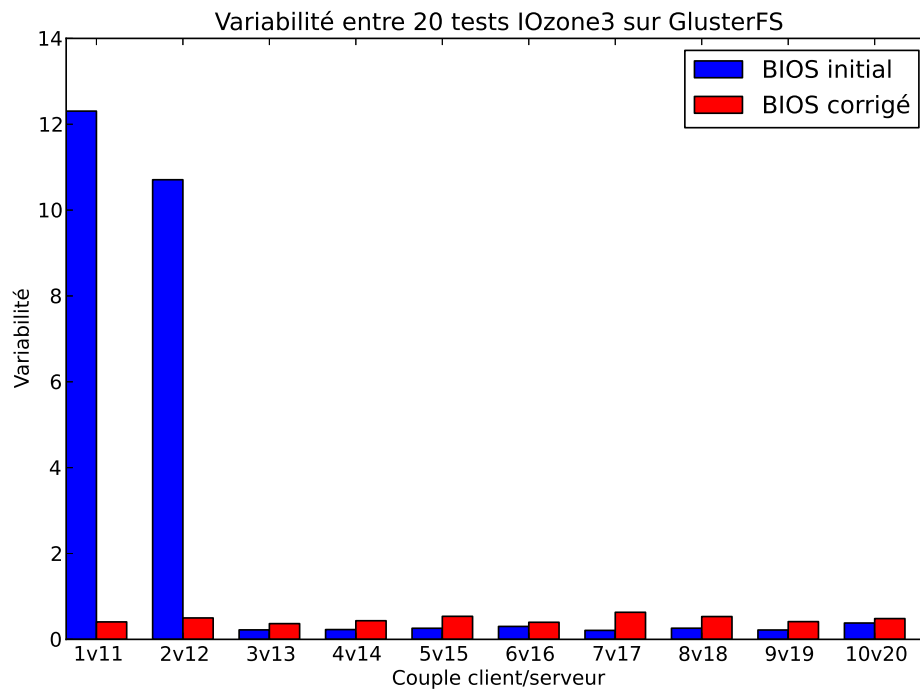


FIGURE 2 – Série de 20 tests IOzone3 entre un client et un serveur, pour deux couples différents de machines.

4.3. La variabilité : de la jeunesse des BIOS à la vieillesse des composants

Des variabilités temporelles n'affectent pas uniquement les équipements récents. Le CBP fonctionne essentiellement sur des équipements d'âge «canonique», déclassés des centres de calcul. Leur exploitation quotidienne, sur des expériences, menées sur ces trois dernières années, suggère de multiples origines à cette variabilité sur les temps de calcul : cela va de l'état des ventilateurs, du vieillissement des alimentations, du taux de remplissage de messages de la pile IPMI, de la tenue dans le temps de la pâte thermique entre processeurs et radiateur, de la position dans la baie, de la «maladie des condensateurs», ... : ces différences influençaient les environnements électrique et thermique, mais nous ne nous doutions pas de l'importance de telles influences, essentiellement lorsque les réseaux ou de grosses quantités de mémoire sont sollicités.

4.4. Conclusion SIDUS

Outil d'administration ou de recherche ? Avec SIDUS, il n'y a plus de système sur les disques durs embarqués (et donc des Watts et des BTU à récupérer, de la maintenance en moins), plus de procédures lourdes d'installation ou d'administration : c'est donc indubitablement un outil d'administration de parc, quelle que soit la machine ou sa destination. Sa seule exigence : un démarrage en réseau. Pour quel bénéfice ? Du temps d'ingénieur à reventiler sur des tâches plus en relation avec le calcul scientifique.

De plus, face aux évolutions matérielles permanentes, c'est un outil permettant une comparaison rapide entre des machines différentes, des BIOS dissemblables, des noyaux distincts, voire

des environnements climatiques variés. C'est donc, par conséquent, un outil d'investigation qui élimine toute variabilité sur le socle logiciel.

Ensuite, par sa simplicité, SIDUS facilite l'accès aux ressources de calcul scientifique. Ce n'est pas dans l'appropriation de la puissance de calcul que réside la difficulté d'apprentissage : c'est dans la maîtrise des outils. Démarrer un système complet graphique, pleinement opérationnel, en quelques secondes, comparable à un nœud de calcul, sans toucher quoi que ce soit à la machine hôte, favorise cette transition. Elle permet en outre l'exploitation de ressources locales (CPU & RAM) souvent largement surdimensionnées pour leur usage courant. Cette approche COMOD, certainement l'avenir d'un *sparse computing*, reste à généraliser pour rationaliser au mieux l'usage d'un parc informatique surtout en cette période de restriction budgétaire. Le portage de son poste de travail vers les équipements dédiés est instantanée : SIDUS accélère donc appropriation et intégration pour les acteurs de la recherche exploitant le calcul scientifique.

Quelles perspectives pour SIDUS ? Le mettre à disposition sur tout le site de l'ENS-Lyon, poursuivre dans la lignée de 2013[6][8][9] sa promotion dans la communauté HPC et sur la place lyonnaise, permettre son accès via un accès VPN, l'adapter sur des structures de grilles ou expérimentales (sur Grid'5000) : autant de projets à l'étude. Un autre axe de travail vise aussi le remplacement de la glu AUFS par OverlayFS, nouvellement intégrée aux noyaux Linux les plus récents.

Bibliographie

1. Corden (D. M. J.) et Kreitzer (D.). – Consistency of floating-point results using the intel(r) compiler or why doesn't my application always give the same answer? – <http://software.intel.com/en-us/articles/consistency-of-floating-point-results-using-the-intel-compiler>, août 2012.
2. Kreitzer (D.). – Home - abinit. – <http://www.abinit.org>, janvier 2014.
3. Norcott (W.). – Iozone filesystem benchmark. – <http://www.iozone.org/>, janvier 2014.
4. Quemener (E.). – Site institutionnel de sidus. – <http://www.cbp.ens-lyon.fr/sidus/>, janvier 2014.
5. Quemener (E.) et Corvellec (M.). – Extreme os deduplication using sidus. *Linux Journal*, no235, nov 2013, pp. 100–111.
6. Quemener (E.) et Corvellec (M.). – Os deduplication with sidus (single-instance distributing universal system). – http://conference.scipy.org/scipy2013/presentation_detail.php?id=199, juin 2013.
7. Quemener (E.) et Taulelle (L.). – Le projet distonet : le stockage distribué du cluster au poste de travail. In : *Journées Réseaux 2011*. – Paris, France, novembre 2011.
8. Quemener (E.) et Taulelle (L.). – Déduplication extrême avec SIDUS : un premier pas vers la reproductibilité ? In : *Journées SUCCES 2013*. – Paris, France, novembre 2013.
9. Quemener (E.) et Taulelle (L.). – Déduplication extrême d'os avec sidus. In : *Journées Réseaux 2013*. – Montpellier, France, décembre 2013.
10. Quemener (E. Y.). – Exemple de rootaufs. – <http://www.cbp.ens-lyon.fr/sidus/rootaufs>, janvier 2014.
11. Schembri (N. A.). – Google code de rootaufs. – <http://code.google.com/p/>

rootaufs/, septembre 2010.