

LCG-France Tier-1 @ Analysis Facility Réunion Mensuelle de Coordination

Relevé des Conclusions

18 décembre 2008 – 13h30

Présents:

- Suzanne Poulat [SP]
- Lionel Schwarz [LS]
- Ghita Rahal [GR]
- Farida Fassi [FF]
- Nelli Pukhaeva [NP]
- Catherine Biscarat [CB]
- Julien Devemy [JD]
- Philippe Olivero [PHO]
- Rolf Rumler [RR]
- Marc Hausard [MH]
- Pierre-Emmanuel Brinette [PEB]
- Pierre Girard [PG]
- David Bouvet [DAB]
- Yvan Calas [YC]
- Hélène Cordier [HC]
- Fabio Hernandez [FH]

Président : FH **Secrétaire** : FH

Agenda: http://indico.in2p3.fr/conferenceDisplay.py?confId=1094

Début de réunion: 13h35

1. Introduction

Cette réunion a pour objet de faire le point sur l'état d'avancement de la plate-forme de surveillance des services basée sur NAGIOS, d'une part, et de présenter une récapitulatif des problèmes récents rencontrés lors des différents exercices de transfert et de traitement des données des expériences LHC.

Le format de cette réunion est légèrement différent à celui utilisé habituellement. FH a attaché une présentation à l'agenda de la réunion mais ces slides n'ont pas été présentés afin de consacrer le temps de la réunion au traitement des problèmes les plus prioritaires. FH annonce qu'à partir de la réunion mensuelle du mois de janvier 2009, HC présentera régulièrement une synthèse des problèmes soulevés par les responsables des services lors des réunions hebdomadaires de l'activité transverse Grid.

2. Etat d'avancement de la plate-forme de surveillance des services

MH et PEB présentent en 2 parties (voir support attaché à l'agenda de la réunion) l'état de la plate-forme NAGIOS et l'état des sondes de surveillance. Les points essentiels de cette présentation sont ci-dessous :

CC-IN2P3 1/5



- Une plate-forme NAGIOS de production est en place, composée de 2 machines en configuration redondante. Les différentes versions des fichiers de configuration et des sondes elles mêmes sont déposées dans un référentiel SVN.
- La plate-forme de développement et de validation de sondes NAGIOS est composée d'une machine. Elle est configurée pour supporter plusieurs utilisateurs travaillant simultanément sur des sondes ou sur la configuration du serveur. Les modifications validées et intégrées dans le référentiel SVN sont automatiquement propagées vers la plate-forme NAGIOS de production.
- Les sondes NAGIOS des services grid produites et packagées par WLCG sont exécutées sur une machine dédiée séparée de la plate-forme de production. Ceci permet de découpler l'évolution des sondes WLCG et celle des sondes locales.
- La documentation relative à l'installation et à l'utilisation de NAGIOS, se trouve dans le wiki des opérations, à l'adresse https://cctools2.in2p3.fr/operations/wiki/doku.php?id=docservices:expert:nagios:start. On y trouve un guide d'écriture des sondes à l'usage des responsables des différents services ainsi que la liste des sondes actuellement en production. Un espace pour collecter le retour des utilisateurs des sondes est aussi prévu.
- L'état instantané des services surveillés par NAGIOS peut être visualisé à l'adresse http://ccnagios.in2p3.fr/nagios (l'utilisateur et le mot de passe sont à demander à nagiosmaster@cc.in2p3.fr)
- Le calendrier pour les étapes à venir est comme suit :
 - O Phase de pré-production (22/12/2008): l'ensemble des sondes existantes seront activées avec comme destinataires des notifications les responsables de NAGIOS et les membres du groupe Exploitation. Dans cette phase, les sondes NAGIOS seront exécutées en parallèle avec les sondes NGOP afin de valider que la couverture des services surveillés par ces 2 outils est identique.
 - Phase de production (vers le 01/02/2009): dans cette phase, les notifications NAGIOS seront envoyées, en plus des responsables de l'exploitation, aux responsables des différents services surveillés, et un message RLS sera aussi émis.
- MH est chargé de contacter les responsables des différents services afin de mettre en place avec eux les sondes pertinentes. Il s'agit d'un processus itératif qui impliquera non seulement les responsables des services, mais aussi les responsables de l'exploitation et les ingénieurs d'astreinte.

3. Synthèse des problèmes récents rencontrés par les expériences LHC au CC-IN2P3

GR présente une synthèse des principaux problèmes rencontrés par les expériences LHC pendant les derniers mois. La présentation, préparée collectivement par les expertes calcul des 4 expériences, est attachée à l'agenda de la réunion.

Un résumé des problèmes et les actions identifiées suit :

• Forte charge sur la zone AFS utilisée pour l'installation du logiciel des expériences : problème observé par ATLAS et LHCb. Une solution est en phase finale de validation entre ATLAS et Xavier Canehan. Dans cette solution, l'installation d'une nouvelle version du logiciel de l'expérience suppose une phase de pré- et post-processing spécifiques au site qui ont pour objectif la création et la réplication des volumes AFS correspondants. D'après les experts AFS du site, cette solution est satisfaisante et permet un passage à l'échelle. Une fois la phase de validation avec ATLAS sera terminée, le même mécanisme sera utilisé pour les 4 expériences LHC.

CC-IN2P3 2/5



D'autre part, une plus forte réactivité des experts AFS du site est nécessaire pour réduire les latences d'intervention observées dans ce service. L'augmentation du nombre de personnes susceptibles d'intervenir sur ce sujet est souhaitée.

ACTION [GR, LA, FF, NP]: demander aux experts AFS du site la généralisation de la solution de réplication de volumes AFS et coordonner la mise en place du mécanisme, en interaction avec les responsables de l'installation du logiciel de chacune des 4 expériences.

 Débit insuffisant pour la copie des données HPSS (bande) dans les disques gérés par dCache :

Des tests de staging (organisé ou non) des données entre HPSS et dCache ont montré des limitations dans le débit atteignable avec la configuration actuelle de ces outils. Cette copie de données est indispensable pour le traitement et retraitement des données stockées sur cartouches, notamment pour les phases de reprocessing. Des exercices réalisés par ATLAS et CMS ont mis en évidence ce problème qui nécessite une amélioration de l'interaction entre dCache et HPSS. [Une analyse des résultats observés par CMS en février 2008 est disponible à l'adresse

http://indico.in2p3.fr/conferenceDisplay.py?confId=800].

Afin d'optimiser la lecture, des stratégies différentes d'organisation des données sur bande sont envisagées par chaque expérience [voir par exemple l'orientation de CMS à l'adresse https://twiki.cern.ch/twiki/bin/view/CMS/DMWMPG Namespace]. Aussi, des stratégies différentes pour déclencher le rapatriement des données de la bande (HPSS) vers le disque (dCache) sont évaluées (staging organisé ou pré-staging, staging à la demande, etc.). Il paraît indispensable de tirer partie de l'organisation des données ainsi que du mécanisme de déclenchement de la copie vers la bande, dans l'interaction de dCache et HPSS.

Par ailleurs, PEB informe que Jonathan étudie actuellement un outil de ONL pour l'ordonnancement des copies des fichiers entre le disque HPSS et la bande HPSS, utilisé aussi à BNL.

ACTION [FH]: convoquer à une réunion les experts des expériences LHC, les experts HPSS et les experts dCache afin de comprendre précisément le mode de fonctionnement envisagé par chaque expérience en ce qui concerne l'organisation des données et le modèle de déclenchement de la copie des données de la bande HPSS vers le disque dCache. Ces informations permettront d'identifier les possibilités techniques offertes par le couple dCache-HPSS pour satisfaire ces demandes et établir un plan d'action.

Incohérence du catalogue dCache avec le contenu réel des serveurs de fichiers : Le problème est mis en évidence par des jobs voulant accéder à des fichiers gérés et catalogués par dCache mais qui n'existent pas réellement, ce qui suppose une intervention manuelle très consommatrice en temps aussi bien de côté support que de côté experts dCache. Deux causes ont été identifiées qui expliqueraient cette incohérence, mais la liste des fichiers impactés n'est Plusieurs pistes ont été évoquées : audit complet du catalogue de dCache suivi d'audits incrémentaux, améliorer la séparation entre les composants de dCache utilisés par chaque expérience afin que la charge provoquée par une expérience n'impacte pas négativement les autres.

LS a expliqué quelques possibilités supplémentaires : la configuration en base de données séparées pour le catalogue des fichiers de chaque expérience sera possible avec Chimera, qui sera introduit dans une version ultérieure de dCache. Cette possibilité sera étudiée.

ACTION [LS]: identifier les moyens nécessaires pour établir un audit complet du catalogue dCache et des audits réguliers qui nous permettent d'assurer la cohérence entre le catalogue central et le contenu des serveurs des fichiers dCache. Un audit complet sera programmé rapidement même si cette opération implique un arrêt

CC-IN2P3 3/5



prolongé du service dCache. L'objectif est d'assainir la situation du catalogue dCache. Le début de l'année semble une période propice pour une telle opération.

- Correspondance entre les limitations réelles des services et le paramétrage de BQS: Les ressources BQS sont utilisées pour la régulation de l'exécution des jobs des expériences LHC. Des ressources ont été définies pour plusieurs services (dCache, HPSS, SPS, AFS, ...). Les experts LHC du site ne sont systématiquement pas informés des changements des valeurs limites de ces ressources. PHO informe qu'au sein de l'équipe Opérations un travail est en cours pour améliorer la traçabilité et la lisibilité des modifications faites de ces valeurs limites. L'outil elog servira à stocker ces informations qui seront à disposition des ingénieurs du Centre. D'autre part, le suivi de l'utilisation des ressources BQS par expérience est disponible à l'adresse http://cctools2.in2p3.fr/mrtguser/info_manips.php.
- Jobs ATLAS restés en attente lors de la soumission au Computing Element : La cause de ce problème n'a pas été identifiée. Des modifications dans la plage des ports utilisée par le CE ainsi qu'une mise à jour du middleware ont été implémentées sans réel succès pour ce problème particulier. Une mise à jour matérielle et logicielle de la VO Box ATLAS est prévue, ce qui implique une mise à jour de PANDA et tous ses composants (y compris CONDOR) utilisés pour la soumission des jobs.

 ACTION [CB, GR]: mettre à jour la VO Box de ATLAS en SL4 et réaliser des tests de soumission des jobs PANDA aux CEs du site.
- Problèmes d'échange d'information entre les experts responsables des services et les experts locaux des expériences Des amélioration sont nécessaires dans ce domaine. Il a été convenu par l'ensemble des responsables des services présents à la réunion de l'utilisation de l'outil elog pour diminuer l'échange d'informations par courriel et améliorer la visibilité et traçabilité des informations relatives aux événements opérationnels de chacun des services. Les services concernés sont : dCache, HPSS, LFC, FTS, CEs, système d'information, NAGIOS. ACTION [JD]: modifier la configuration de l'instance elog de sorte que les responsables de chacun des services ci-dessus puisse y déposer l'information relative à son service. Comme convenu, elog sera le référentiel d'informations opérationnelles de chaque service. Les échanges de courriels avec ces informations devraient être limités. Le point informations est à l'adresse: http://cctools2.in2p3.fr/elog/ d'entrée de ces **ACTION** [LS, PEB, DB, PG, MH]: s'assurer que l'information des événements opérationnels quotidiens est disponible sur elog.
- Manque d'outils de monitoring à plusieurs niveaux de détail, permettant de suivre l'activité d'une expérience. Exemples des informations qui seraient utiles: évolution dans le temps du nombre de requêtes sur dCache de chaque VO, évolution du nombre d'accès aux pools dCache (par VO et par type de protocole utilisé, dcap, srm, ...), nombre de transferts gridFTP, etc.
- Manque de privilèges suffisants pour les experts LHC du site pour accéder efficacement et à distance aux logs des jobs en exécution, stockés sur les worker nodes.

CC-IN2P3 4/5



ACTION [GR] : demander l'activation de la consultation à distance des logs des jobs en exécution pour les experts LHC du site.

4. Prochaine réunion

Jeudi 15 janvier 2009, 13h30-15h30, Amphi

Fin de réunion : 16h35

CC-IN2P3 5/5