

Un nouvel outil de monitoring du point de vue du site

By Julia Andreeva, Max Boehm, Adrian Casajus, Benjamin Gaidioz, Costin Grigoras, Lukasz Kokoszkiewicz, **Elisa Lanciotti**, Ricardo Rocha, Pablo Saiz, Irina Sidorova, Andrei Tsaregorotsev.
CERN, IT Dept.

Réunion des sites LCG-France
Grenoble 27-28 Novembre 2008

Plan

- Introduction
 - Motivation et objectif de l'activité
- Structure du nouvel outil de monitoring
 - Flux de l'information de la source jusqu'à l'affichage final
- Objet du monitoring:
 - Les activités principales des VOs aux sites
- L'affichage: Gridmap
 - Structure du Gridmap
 - Apparence du premier prototype
- État courant de l'activité et questions ouvertes

Monitoring des activités pendant le CCRC08

- Après l'expérience du CCRC08, impressions des administrateurs des sites:
 - Les principaux outils de monitoring utilisés pendant cet expérience étaient des outils spécifiques pour les VOs. Dans la plupart des cas, ils ont bien fonctionné, mais ils ne sont pas très faciles à utiliser pour quelqu'un qui est externe à la VO
 - Les sites qui servent plusieurs Vos ont du apprendre à utiliser des outils très différents entre eux: temps d'apprentissage nécessaire et peu d'efficacité
 - Comme résultat, souvent les administrateurs des sites ne savaient pas bien si son site contribuait à l'activité de la VO comme on s'y attendait
- Les administrateurs des sites voudraient avoir la possibilité de comparer la performance du site de point de vue de la VO avec la vision qu'ils obtiennent de ses propres systèmes de monitoring

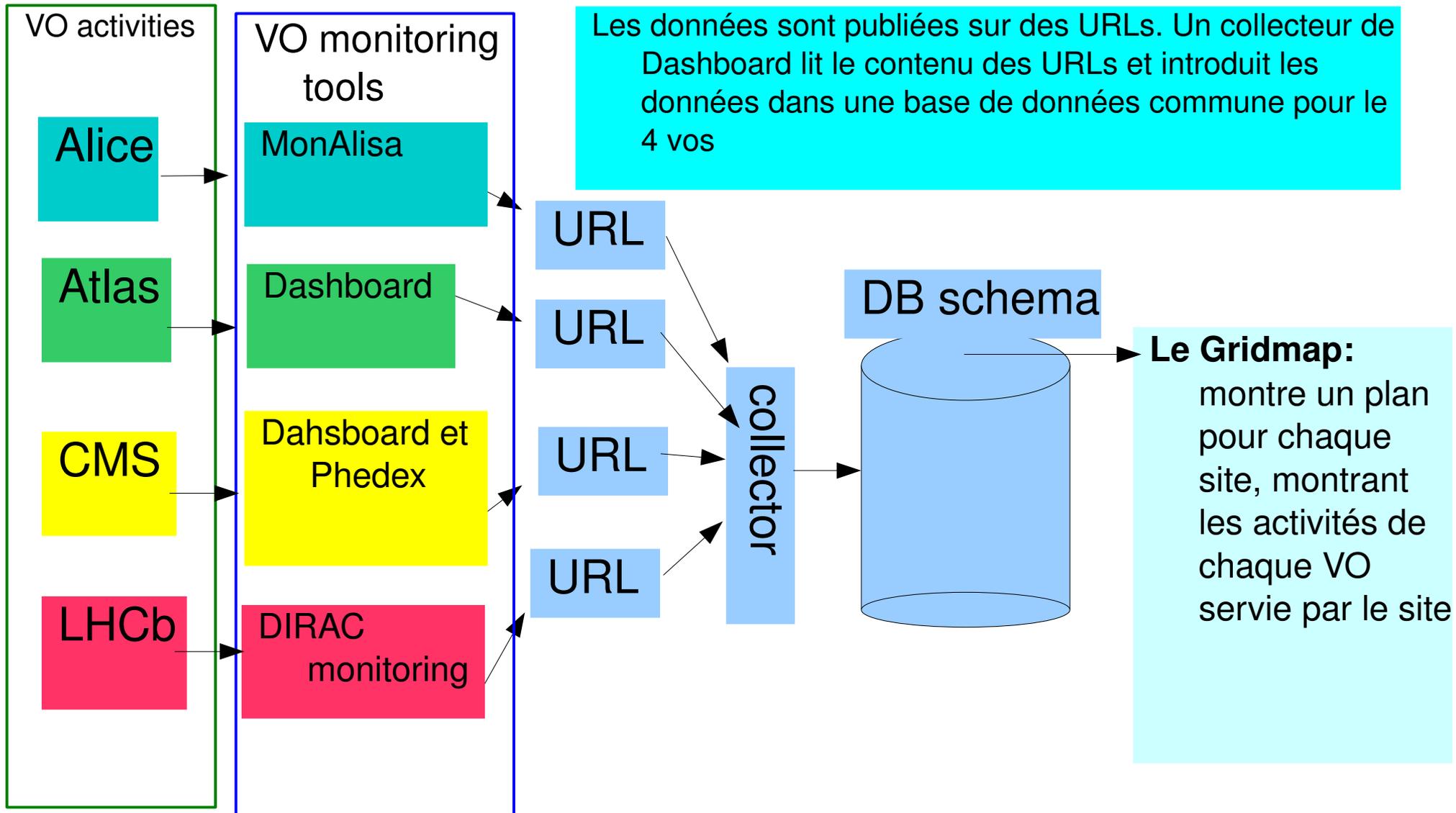
Objectif de cette activité

Fournir un nouvel outil de monitoring qui devrait:

- À partir de la même vue, donner une vue générale de l'état des services dans le site. Il devrait être un outil simple à utiliser, même pour quelqu'un externe à la VO, et qui n'a pas de connaissances spécifiques de l'expérience --> **un aperçu général des activités du site à partir de une unique vue**
- Cet outil obtiendra l'information des outils de monitoring particuliers de chaque VO (Dashboard, MonALISA, Dirac, Phedex) et montrera le résultats essentiels de une façon simplifiée, en fournissant des liens vers les sources de l'information --> **un outil de haut niveau**
- Il montrera l'information en utilisant la technologie de Gridmap --> visualisation **immédiate des problèmes grâce aux plans avec couleurs**

Une requête ultérieure des sites c'est aussi de obtenir une définition claire des objectifs des VOs relativement aux activités qui se déroulent au site: cette information devrait être montrée par cet outil à coté des valeurs observées pour permettre une évaluation des performances

Flux de l'information



Objets du monitoring: fournis par les outils de monitoring particuliers de chaque VO

- L'information sur les activités des VO est obtenue des différents outils de monitoring spécifiques de la VO :
 - Dashboard, MonALISA, DIRAC, Phedex
- Pour construire une base de données commune aux 4 VO il est nécessaire:
 - Identifier des grandeurs qui sont communes à toutes les VO
 - Définir clairement le signifié de chaque grandeur
 - Vérifier que ces grandeurs sont utiles pour le site et que elles sont disponibles dans les outils de monitoring

 **Effort de coordination:** Définir un group de grandeurs qui ont le même signifié pour toutes le VO

Construire un modèle de données et dessiner un schéma pour la base de données commune -> une seule et unique base de données pour le grandeurs des différentes expériences

Activités et grandeurs objet du monitoring

ACTIVITÉS

état général du site

job processing

transfert de données

...

GRANDEURS.

Nombre de jobs, nombre de jobs achevés (pendant la dernière heure)

wall time et temps CPU (si possible en KSI2K)

...

taux de réussite (moyenne de la dernière heure)

taux de transfert (moyen dans la dernière heure)

...

Il y a 2 activités principales: job processing et transfert de données Leurs grandeurs sont rassemblés pour les 4 expériences LHC

Pour chacune des activités principales il y a un groupe de activités secondaires: elles sont optionnelles et dépendent de la VO et de son modèle de calcul

Job Processing: MC production, data reconstruction, user analysis, SAM tests...

Data transfer: production transfer, transfer t0-t1, transfer t1-t1...

Le collecteur

- Un nouveau collecteur a été développé dans le cadre du Dashboard, en analogie avec le collecteur du module Site Status Board (développé par Pablo Saiz pour CMS)
- Le collecteur se exécute périodiquement (toutes les heures), lit les grandeurs publiées dans les URLs par les outils de monitoring des Vos, et les introduit dans la base de données

http://.....

site,activity,metric,actual,value, pledged_value, status, time, URL
CERN,mc prod, number of jobs..
...

http://.....

site,activity,metric,actual,value, pledged_value,
status, time, URL
CERN,mc prod, number of jobs..
...

collector

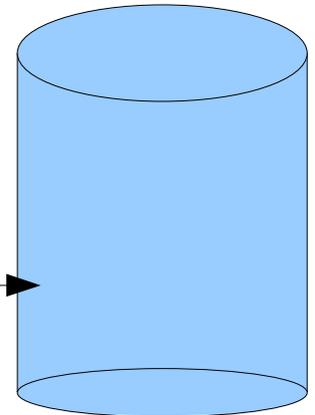
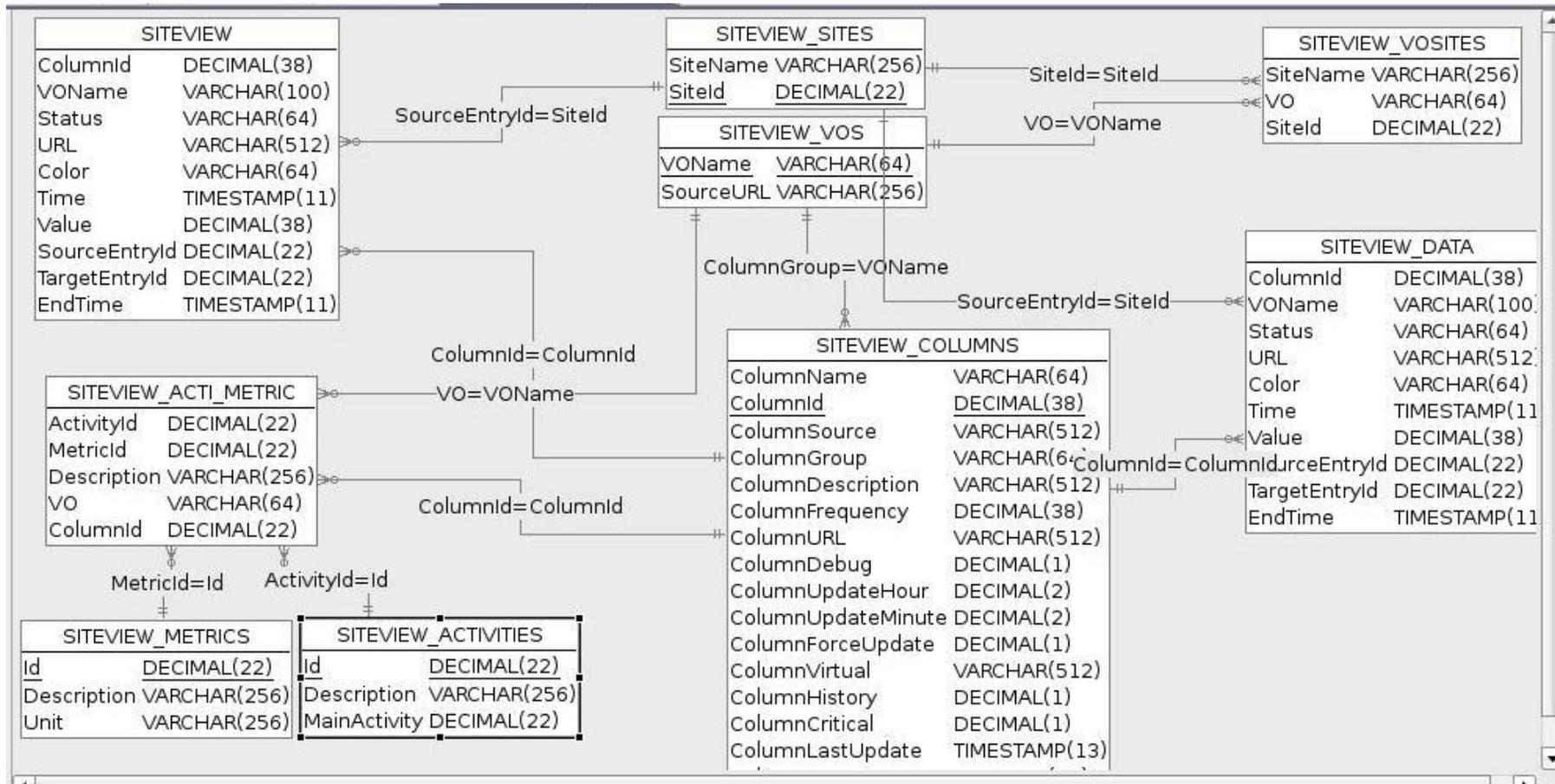


Schéma de la base des données

Implémentation sur un serveur de test de Oracle



Très flexible: d'autres activités et grandeurs peuvent être ajoutées ensuite, selon le besoin de la VO

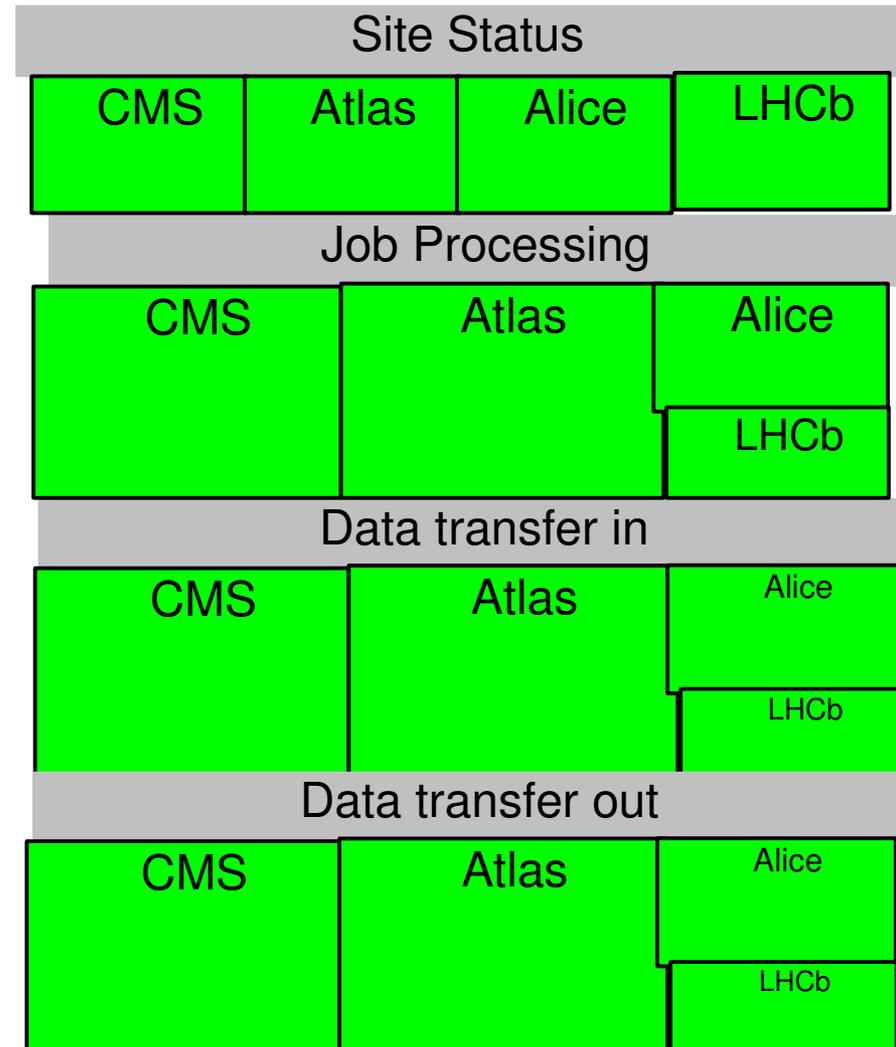
Le Gridmap

- Un Gridmap pour chaque site
- Dans le Gridmap du site, la hiérarchie est la suivante:
 - premier niveau: Activités
 - second niveau: VO

Le Gridmap principal ne montre que les activités principales, pour toute VO servie par le site

Les dimensions du rectangle qui représente la VO est proportionnel à une grandeur qui caractérise son activité (i.e. le nombre de jobs parallèles). Un offset constant est ajouté, à fin de pouvoir toujours visualiser toutes le Vos, même si elles ont une très faible activité en ce moment.

L'état (=COLOR) devrait être fourni para la VO. Chaque VO peut avoir une définition de état indépendante



L'état du site

- Il s'agit d'une évaluation générale de l'état du site, du point de vue de la VO. Si l'état est 'ok' (vert), alors le site est en bon état pour la VO, même s'il n'y a aucun job en exécution
- Pour LHCb et Atlas l'état du site est calculé à partir des tests SAM spécifiques pour la VO
- Pour CMS non seulement les tests SAM, mais aussi le taux de réussite de l'activité de production, en ne prenant en compte que les échecs dus à un problème du site (vu le exit code). Et aussi si le site est visible sur le BDII
- Pour Alice il est calculé dans le cadre de Monalisa
- **Important:** si l'état du site apparaît rouge ça ne signifie pas nécessairement que le site est le responsable. Même si la définition de état est le plus possible relié au site, il peut arriver certaines fois que le site apparaît rouge pour des problèmes de la VO. Même en ce cas là, il est utile que le site soit au courant du problème

Un aperçu du prototype pour FZK

- Le Gridmap principal ne montre que les activités principales, pour toutes les Vos servies par le site

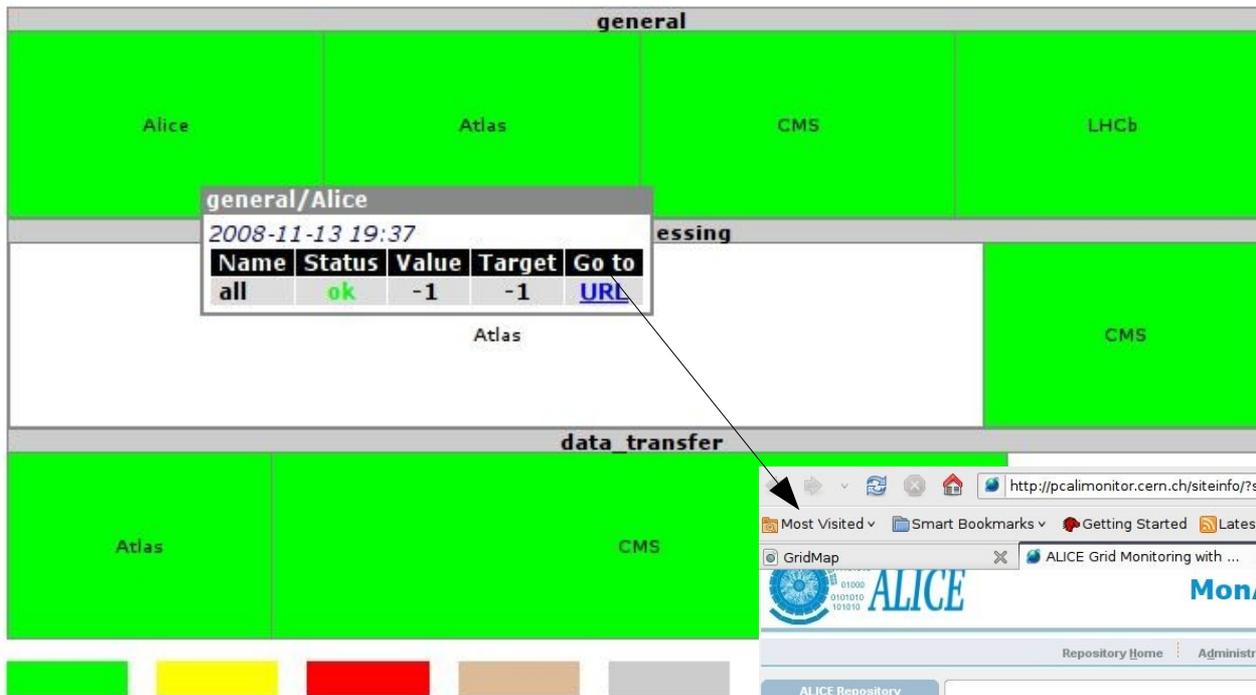
Siteview GridMap Test Page



- En pointant la souris sur la barre d'en-tête il apparaît une fenêtre avec les noms de Vos, l'état et la date de la dernière mise a jour

Lien a l'état du site pour Alice

Siteview GridMap Test Page



Les infobulles fournissent un lien vers la source de l'information

MonALISA information Version: 1.8.8 (JDK 1.6.0_06)
Running on: alice-fzk.gridka.de
Administrator: Kilian Schwarz <K.Schwarz@gsi.de>

Service health NTP: SYNC, offset: 0s

Name	Status	Size	Used	Free	Usage	No of files	Type	ADD test
ALICE::FZK::dCache	OK	423.8 TB	1.026%	419.4 TB	4.346 TB	0 B	SRM	OK
ALICE::FZK::dCache sink	OK	2.728 PB	7.618%	2.521 PB	212.8 TB	0 B	SRM	OK

Lien pour l'état du site pour CMS: Site Status Board

Siteview GridMap Test Page

general/CMS
2008-11-13 19:37

Name	Status	Value	Target	Go to
all	ok	-1	-1	URL

Site Status for the CMS sites

OK Degraded Warning Down Maintenance Inactive

A partir d'ici on a tous les liens nécessaires pour continuer la recherche

Put the mouse over any column header to get the description of the column
 Clicking on a column header will display the evolution of that column over the last 24 hours
 : information is more than 24h old
[Back to the index](#)

Site Name	Visible	JobRobot	SAM TESTS		Production	Analysis	Site usage		Phedex		CMSSW	Open issues	Maintenance	Maintenance	
			CE	SRM			Running	Pending	# Links	In rate					Out rate
T1_DE_FZK	OK	n/a	OK	OK	100%(1839)	77%(2544)	296	1744	ok	56	43	OK	info	GOCDDB-info	n/a

Lien pour l'état du site pour les VOs Atlas et LHCb: test SAM spécifiques

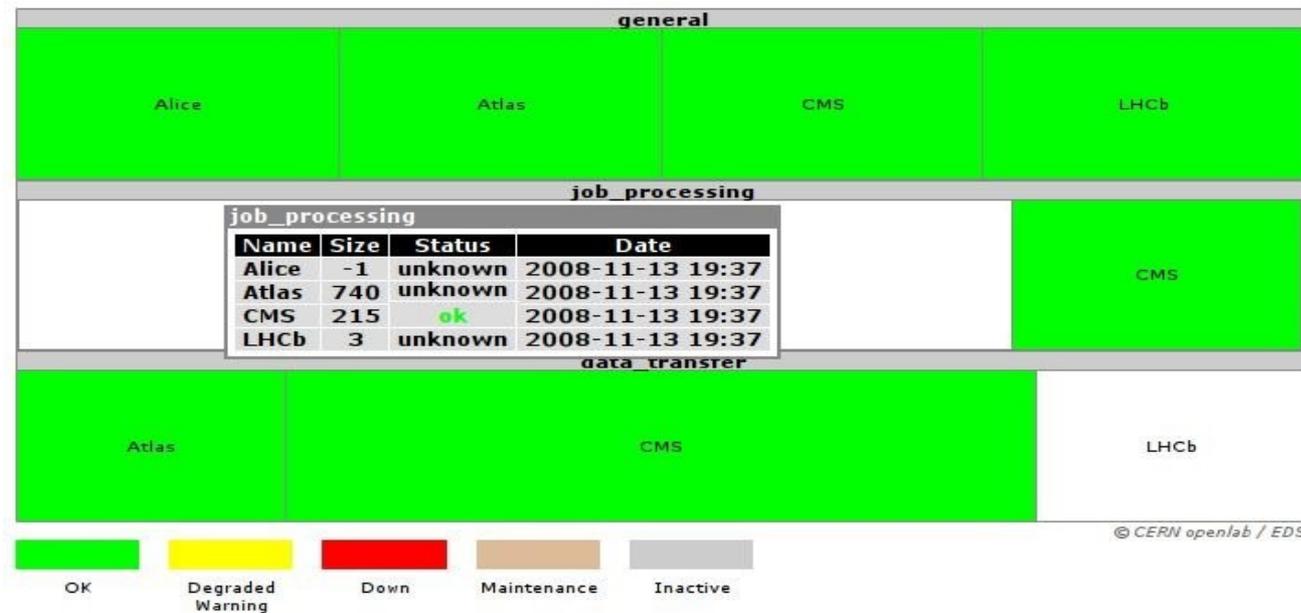
Siteview GridMap Test Page

The screenshot displays the SAM visualization interface for LHCb. At the top, a 'general' grid map shows Alice, Atlas, CMS, and LHCb in green. Below it, a 'job_processing' section shows Atlas and CMS. A browser window shows the URL: <http://dashb-lhcb-sam.cern.ch/dashboard/request.py/historicalserviceavailability?mode=serviceavl&algoid=5&sit>. The main interface includes a legend for status (OK, Degraded Warning, Down, Maintenance) and a control panel for 'VO view' and 'Site View'. The 'Site View' panel shows filters for 'Service Availability', 'Algorithm', 'Time Range' (Last 24 Hours), 'Sites' (Tier0 + Tier1s, CERN-PROD, FZK-LCG2, IN2P3-CC, INFN-T1, NIKHEF-ELPROD, RAL-LCG2), and 'Service Types' (All service types, Select all, CE, FTS, LFC, LFC_C, LFC_L, RB). A 'Show Results' button is present. Below the control panel, a 'Service Availability' chart shows a green grid for '24 Hours from 2008-11-12 19:00 to 2008-11-13 19:00 UTC' for various sites. The chart lists sites like FZK-LCG2 - CE - ce-1-fzk.gridka.de through FZK-LCG2 - FTS - fts2-fzk.gridka.de. At the bottom, there is a section titled 'Algorithm for calculating the Site and Service Availability'.

Activité de job processing

- En pointant la souris sur l'en-tête: information sur les Vos qui sont servies par le site.
- On montre le nombre de jobs en exécution (c'est le paramètre qui détermine les dimensions du rectangle), l'état (qui détermine la couleur), et la date de la dernière mise à jour

Siteview GridMap Test Page

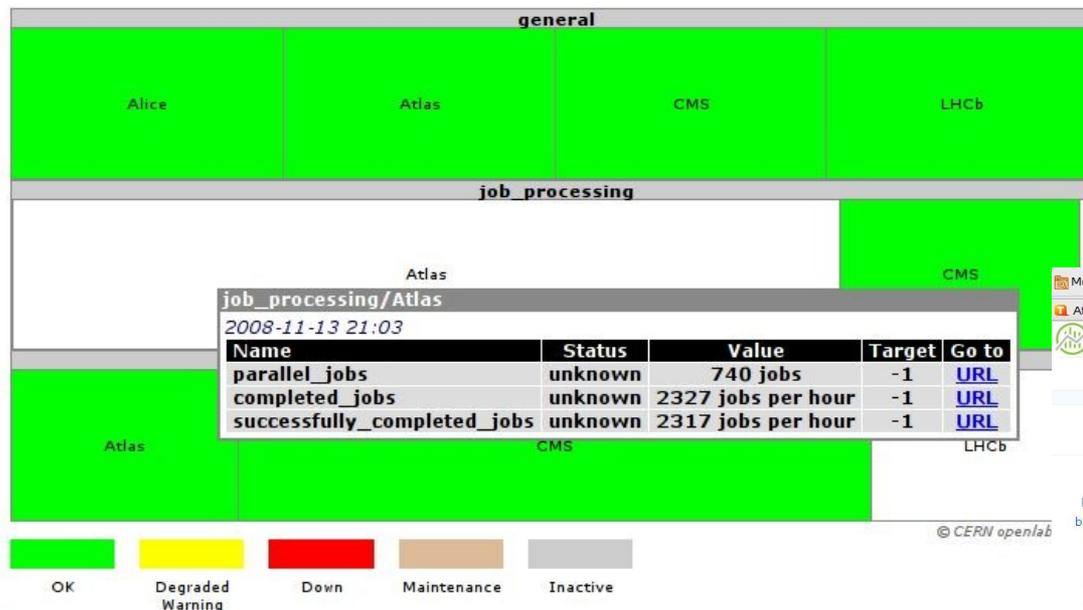


- **Important:** si l'état est rouge, il n'indique pas un problème du site. Par contre, il indique un problème relié à l'activité de la VO sur ce site

Job processing: Atlas

- Toutes les grandeurs relatives à l'activité de job processing (jusqu'à maintenant que MC production, mais bientôt aussi user analysis et autres types de taches) sont montrées
- Le lien affiché renvoie vers la page de Dashboard relative à la production MC au site

Siteview GridMap Test Page



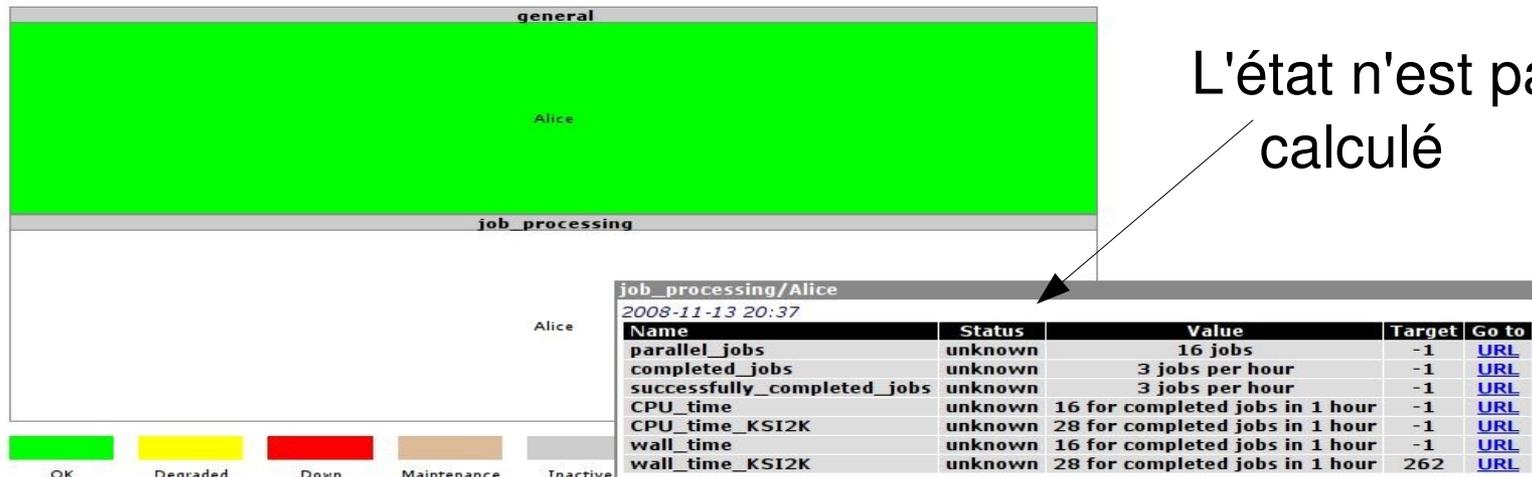
Topology:



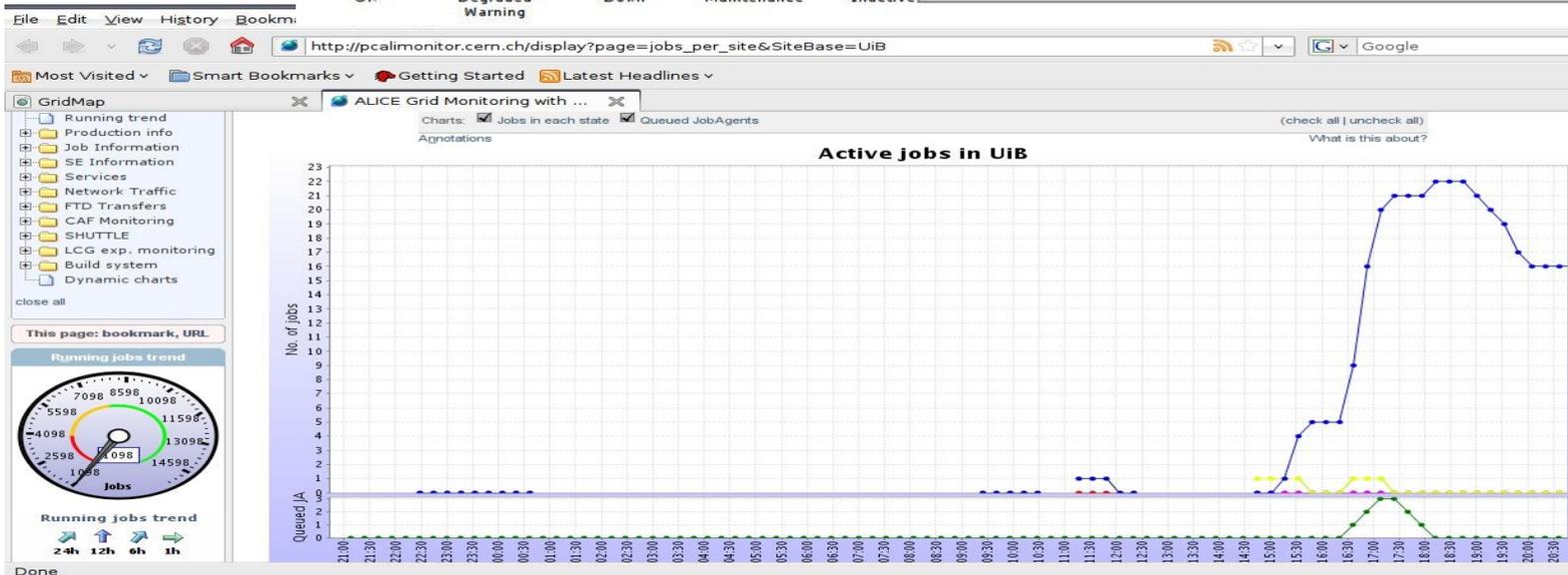
Job Processing: Alice

- Toutes les grandeurs relatives au job processing sont affichées
- Le wall time est affiché, **normalisé en ksi2k** avec aussi la valeur attendue

Siteview GridMap Test Page



L'état n'est pas calculé



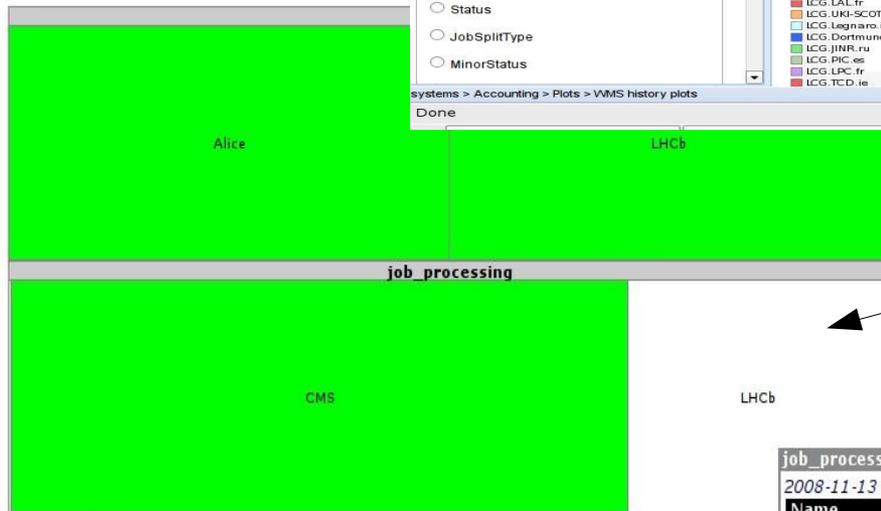
Le lien dans le infobulle renvoie a la URL de Monalisa pour ce site et cet activité

Job processing LHCb

- Job processing LHCb. Lien vers le history plot du Dirac WMS



Siteview GridMap Test Page



L'etat bientôt sera calculé sur la base du taux de réussite

Le temps est affiché en seconds (il n'est pas normalisé)

Pas de valeur attendue

Name	Status	Value	Target	Go to
parallel_jobs	unknown	107 jobs	-1	URL
completed_jobs	unknown	21 jobs per hour	-1	URL
successfully_completed_jobs	unknown	21 jobs per hour	-1	URL
CPU_time	unknown	1821425 for completed jobs in 1 hour	-1	URL
wall_time	unknown	1828682 for completed jobs in 1 hour	-1	URL

Topology:
 {
 "job_processing": ["job_processing/Alice", "job_processing/CMS", "job_processing/LHCb"],

Le second niveau du Gridmap

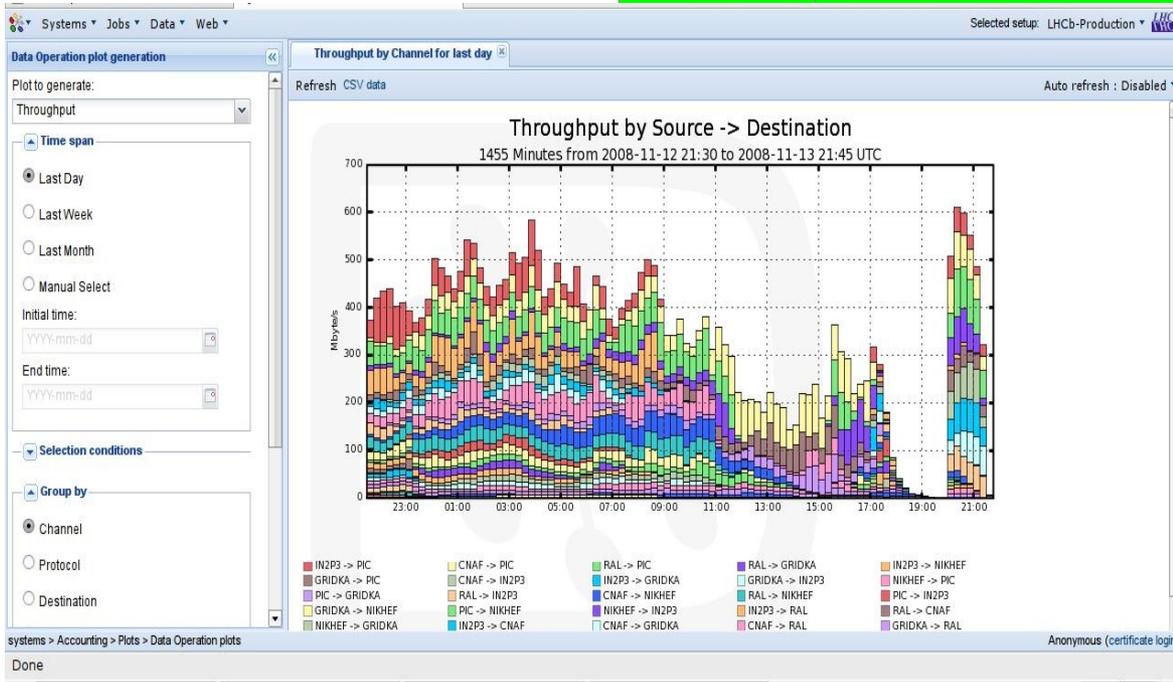
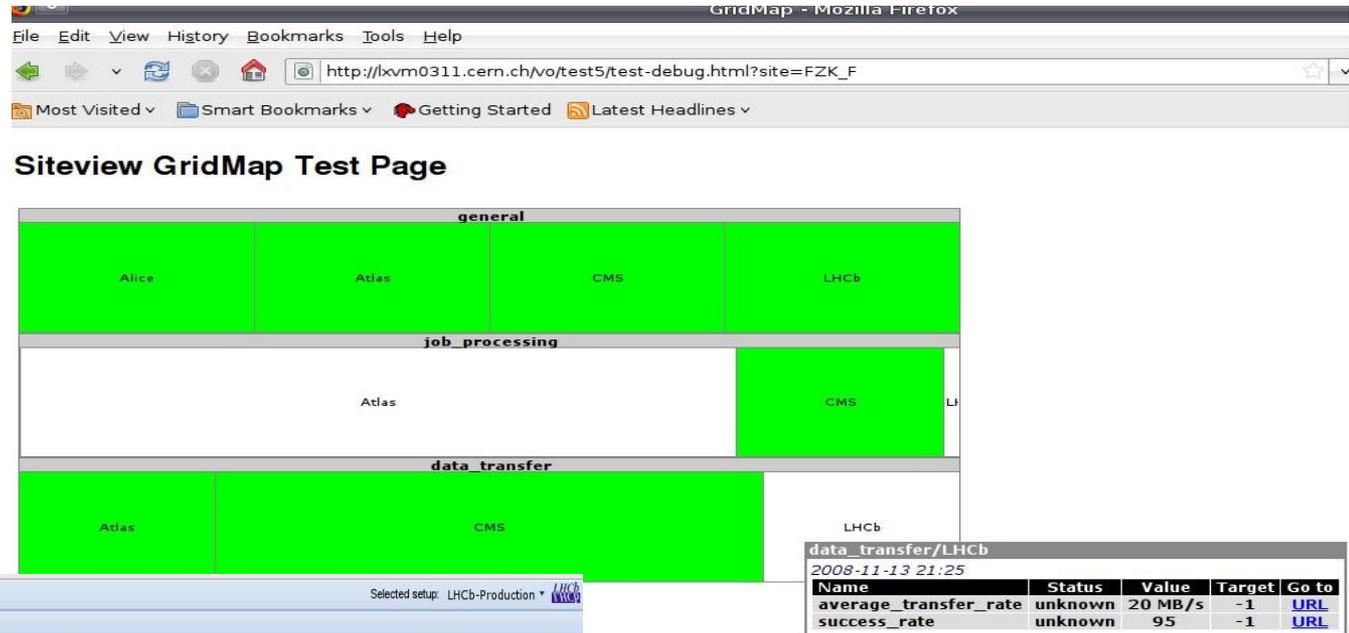
Siteview GridMap Test Page



- Par ex. activité de job processing de la VO CMS.
- En cliquant sur le plan, un second plan apparaît avec les activités secondaires. Par ex. pour CMS: reprocessing, MC production, user analysis, ...
- L'infobulle montre le nombre de jobs en exécution pour chaque activité secondaire

Transfert des données

- Dans la dernière version on montre le flux de données entrantes et sortantes séparément
- Le Gridmap de second niveau montre les différents types de transfert



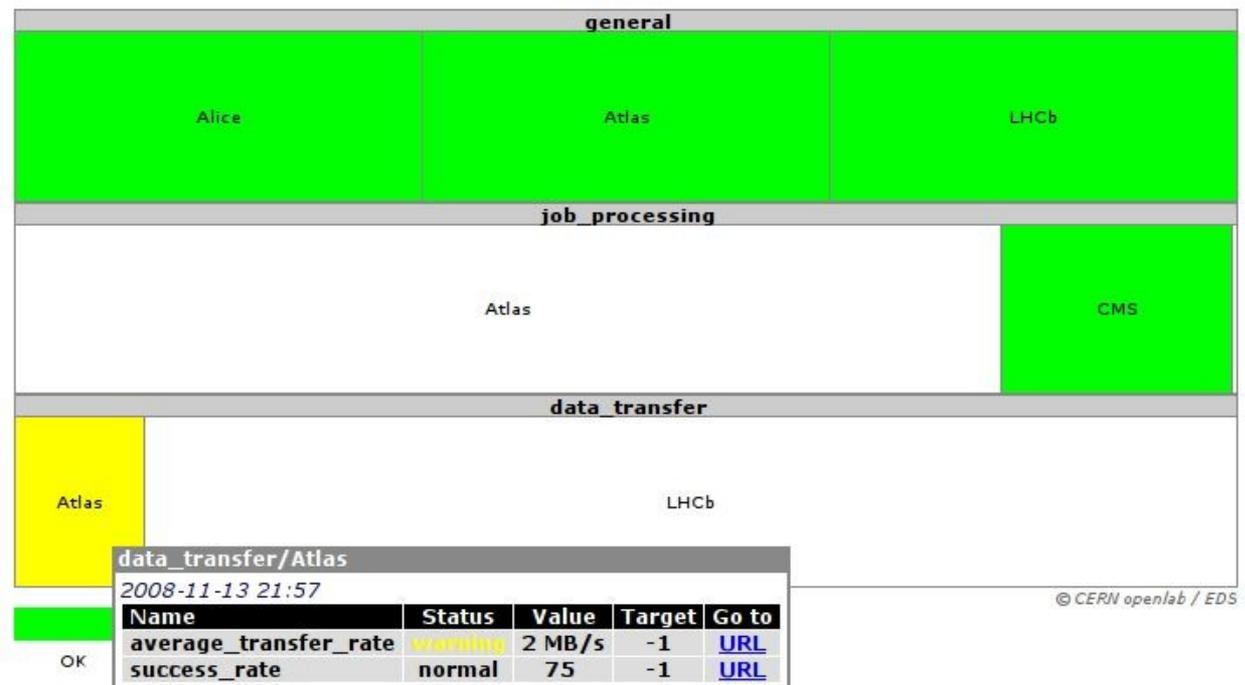
- en développement pendant cette semaine. Une vue plus récente sera montrée pendant la demo

Transfert des données pour Atlas

- L'état et la valeur attendue sont affichés
- Pour le transferts de T0 aux T1, l'état est calculé pour le taux de transfert moyen en comparant avec la valeur attendue: il peut être ok ou warning
- Pour le taux de réussite l'état est calculé suivant cette règle:

- $SR < 20\%$: critical
- $20 \leq SR < 50\%$: warning
- $50 \leq SR < 80\%$: normal
- $SR > 80\%$ good

Siteview GridMap Test Page



État actuel et plans pour l'avenir

- On a défini un groupe de activités et de grandeurs
- On a dessiné un schéma pour la base de données commune et on a vérifié sa validité pour nos objectifs
- On a développé un collecteur qui lit l'information publiée sur des URL et introduit les données dans le schéma commun
- Le Vos ont déjà fourni la plupart des grandeurs, mais quelques grandeurs sont encore manquantes
- Souvent les valeurs attendues et l'état ne sont pas fournis par les VOs

Partie facile
du travail:
l'outil a
une
structure
très
simple

pas très
facile...

Difficile, mais
réalisable

Plans: une version de test du Gridmap complet
avant la fin de l'année

Questions ouvertes

- Souvent c'est difficile pour les VO de afficher les valeurs attendues pour les grandeurs et l'état pour les activités
- Souvent la valeur attendue pour une grandeur n'est pas disponible, ou bien elle est disponible mais elle se trouve dans un autre système ou base de données qui n'a aucune connexion avec l'outil de monitoring, donc la valeur observée ne peut pas être comparée avec la valeur attendue
- En général, le Vos sont très réticentes à déclarer un état pour son activité au site

Quelques observations à ce propos:

Si l'activité d'une VO sur un site apparaît rouge, cela ne signifie pas que le site est responsable, mais c'est plutôt la VO responsable du problème

L'état du site est montré dans le rectangle supérieure 'Site Status' : si cela est vert, alors le site est en bon état, même s'il n'y a aucun job en exécution

Cet outil n'est que un moyen informatif, qui a comme but de simplifier l'accès a l'information donnée par les outil de monitoring des Vos et d'aider les administrateurs des sites à dépister les problèmes de une façon simple et rapide

Remerciements

- Ce projet ne peut avancer que grâce a la collaboration des différentes Vos qui publient les données nécessaires pour peupler la base de données commune
- Merci a: P.Saiz (CMS), C.Grigoras (Alice), R.Rocha and B. Gaidioz (Atlas), A. Casajus and A.Tsaregorotsev (LHCb), W.Ollivier (LHCb and Atlas) pour sa contribution fondamentale

Feedback to: elisa.lanciotti@cern.ch