

mésocentre de calcul de franche-comté



UNIVERSITÉ DE FRANCHE-COMTÉ

mésocentre de calcul de franche-comté





Expérience d'un jeune mésocentre

L. Philippe
10 juin 2010





Plan

- Situation de départ / contraintes
- Solution installée
- Utilisation actuelle
- Bilan et évolutions



Projet

- Financements acquis
- Pré-projet
- Pas de personnel
- Peu d'expertise
- Locaux insuffisamment climatisés



Franche-Comté
Conseil régional



Utilisateurs

- Etablissements franc-comtois + industriels
- Quelques utilisateurs avertis :
 - Utilisation de clusters locaux
 - Centres de calcul
- Beaucoup d'utilisateurs non avertis
 - Septiques
 - Ne savent pas ce que peuvent en faire
 - Windows
- Participation à l'élaboration du CCTP

Choix d'une solution clef en main

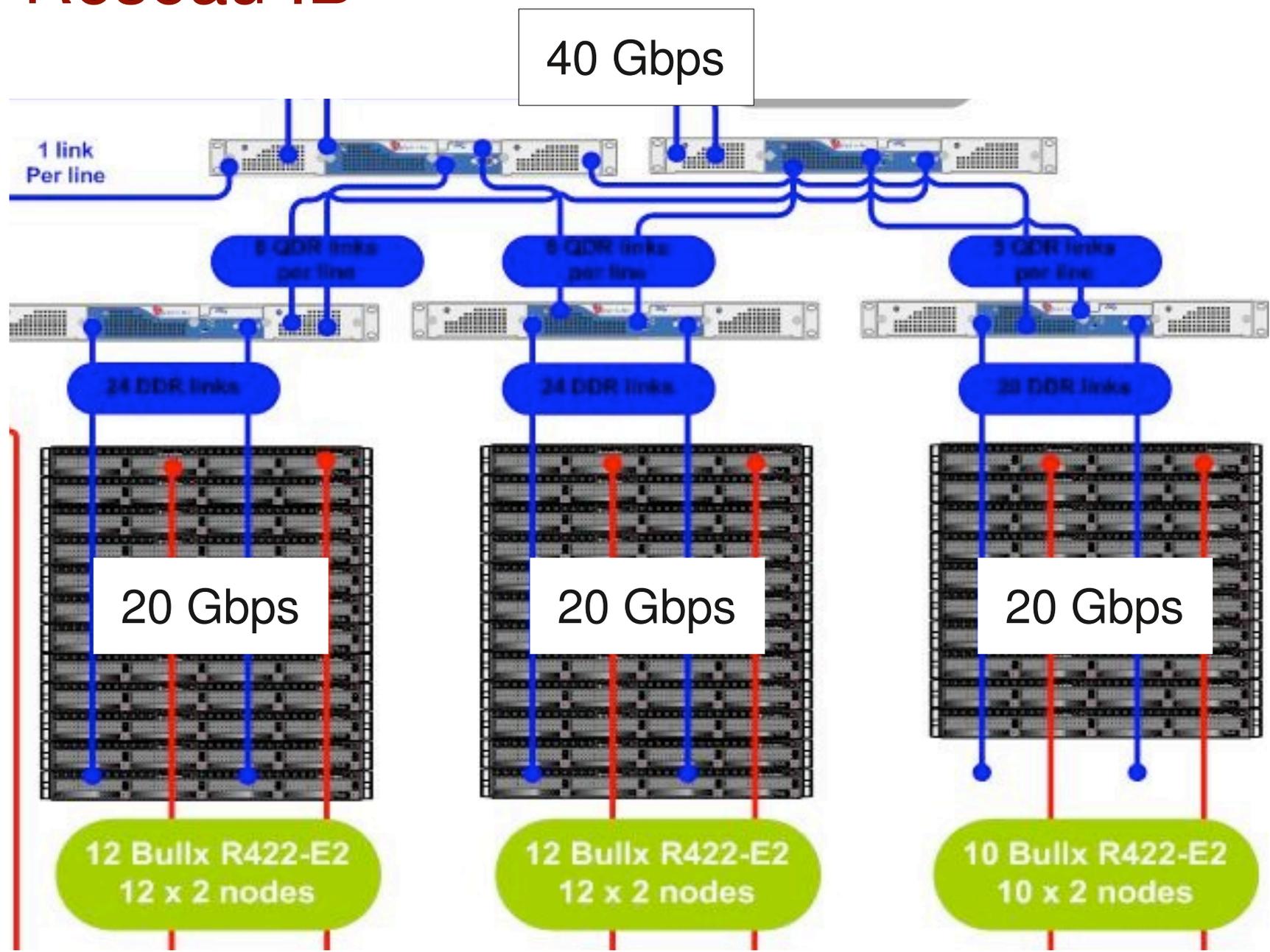


Cluster - Calcul

- Machine Bull
- 68 noeuds de calcul : 5,5 Tflops
 - R422-E2 : bi-proc, Xeon X5560, 2,66 Ghz
 - 544 cœurs
 - 32 avec 24 Go, 1 avec 96 Go, 35 avec 12 Go
- 1 (2) noeud GPGPU Telsa s1070 : 4 Tflops
- Réseau InfiniBand DDR / QDR :
 - Non bloquant DDR
- Non ondulés



Réseau IB





Cluster - Admin

- 5 noeuds hors files d'attentes:
 - 1 noeud de connexion
 - 2 noeuds interactifs
 - 1 noeud serveur NFS
 - 1 admin
 - Ondulés
- Baies de stockage :
 - Système de fichiers partagés NFS 20TB
 - Système de backup 10TB
- OS Bullx Cluster Suite (Base RHEL 5)
- Système de batch : SGE



Portes froides

- Installation
- Deux armoires climatisées
- Faux plancher
- Commandes :
 - Nagios
 - Redémarrage
- Efficace
- Evolution ?

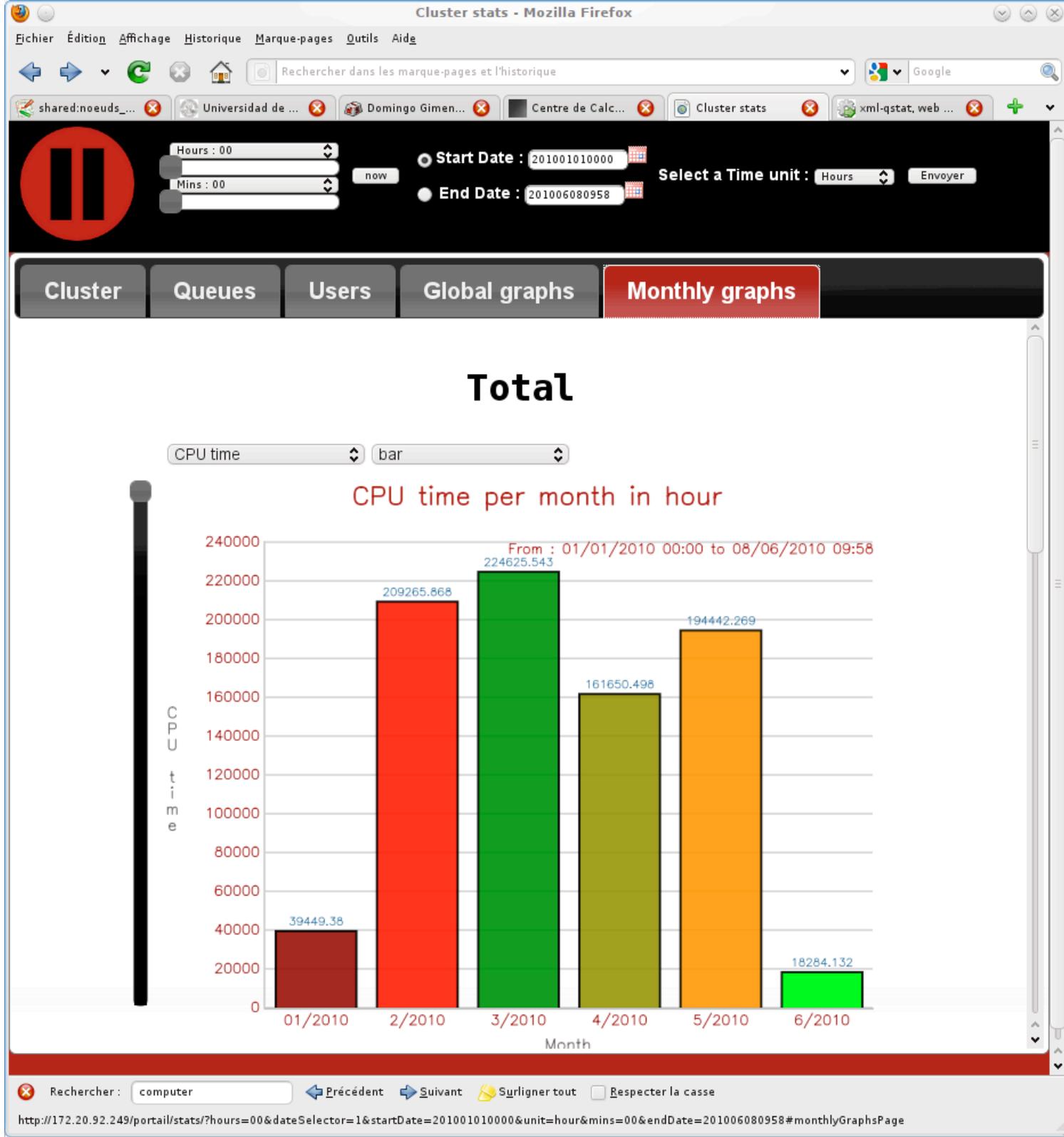




Logiciels administration

- Bull cluster suite :
 - Redhat Enterprise
 - Cluster DB : référence mais statique
 - BSM : Bull System Manager
= Nagios
- Ganglia
- Xmlqstat : détails sur les jobs
- Outil maison exploite SGE







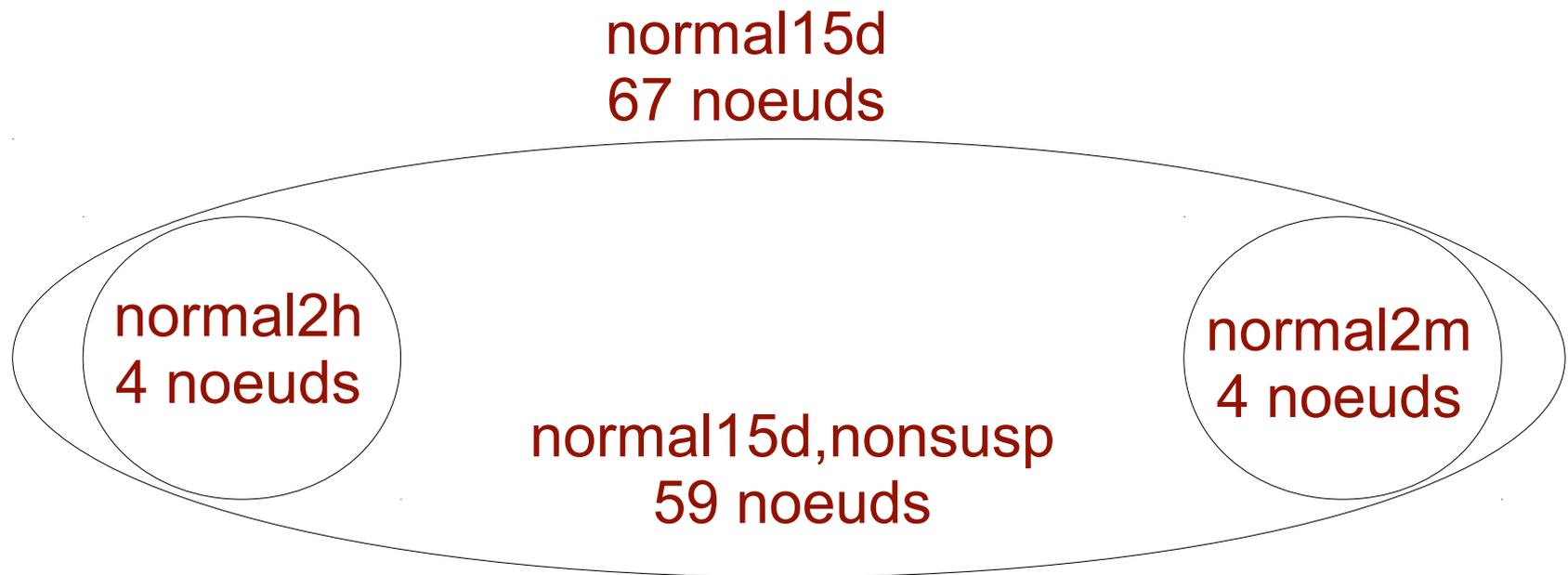
Files d'attente

- Gestionnaire SGE
- Contraintes :
 - Permettre l'exécution de jobs courts
 - Demande pour une exécution longue : 2 mois
 - Noeuds spécifiques :
 - Tesla
 - 96 Go de mémoire
 - Jobs parallèles
 - Soumissions interactives



Files d'attente

- Files spécifiques : tesla, 96 Go, inter
- Files suspensives : pb jobs parallèles

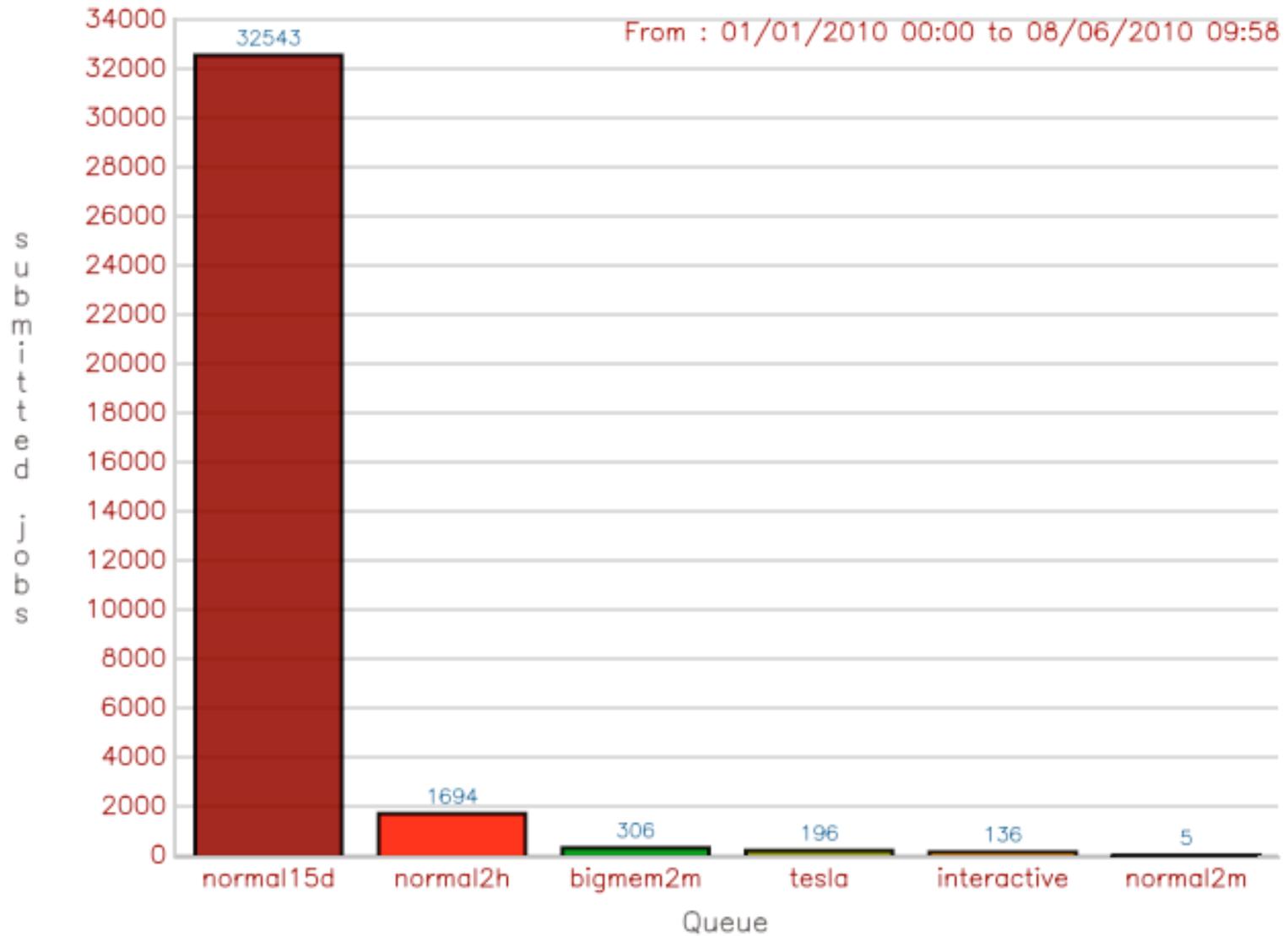


- Environnements parallèles : openmp, mpi



Files d'attente

submitted jobs per queue





Accès

- Réseau privé université
 - SSH
 - VPN
- Annuaire LDAP
- Modules pour paramétrer environnement:
 - Logiciels
 - MPI
 - Langages



Logiciels installés

- Libres :
 - Distribution, bibliothèques (Blas, Lapack, Petsc, ...)
 - Gamess, Ab-Initio, Namd, Meep, Espresso, FFTW, R, Grace, ...
- Commerciaux :
 - Intel Cluster toolkit
 - Molpro, Gaussian, VASP, MAGMA
 - Matlab :
 - 7 toolboxes + compilateur
 - Comsol multiphysics :
 - Structural mechanics



Portail - Mésocentre de calcul de Franche-Comté - Mozilla Firefox

Echier Édition Affichage Historique Marque-pages Outils Aide

http://172.20.92.249/portail/src/php/portail.php

shared:noeuds_du_cl... Universidad de Murcia Domingo Gimenez C... Centre de Calcul de ... Cluster stats Portail - Mésocentre ...



mésocentre de calcul de Franche-Comté

DÉCONNEXION 

PORTAIL - MÉSOCENTRE DE CALCUL DE FRANCHE-COMTÉ

Explorateur Jobs Script Informations

Parent Rafraîchir Vignettes Transférer Nv. Rép. Nv. Fich. Télécharger Éditer Lancer Revenir Copier Déplacer Supprimer

/SGE/scripts Default Files

Répertoires	Nom du fichier	Taille	Type	Modifié le
Default Files	myJob.87109.out	27.54 Ko	Document texte	19/05/2010 10:43
	myJob_201005191043.sge	210 o	Document texte	19/05/2010 10:42
	myJob_201005201205.sge	268 o	Document texte	20/05/2010 12:06
	myJob_201006011356.sge	224 o	Document texte	01/06/2010 13:56
	Test_epis.87102.out	1.77 Ko	Document texte	19/05/2010 10:25
	Test_epis_201005191026.sge	230 o	Document texte	19/05/2010 10:25

Détails Recherche



Nom : Test_epis_201005191026.sge
Taille : 230 o
Type : Document

http://172.20.92.249/SGE/scripts



Portail - Mésocentre de calcul de Franche-Comté - Mozilla Firefox

Fichier Édition Affichage Historique Marque-pages Outils Aide

http://172.20.92.249/portail/src/php/portail.php

shared: noeuds_du_cl... Universidad de Murcia Domingo Gimenez C... Centre de Calcul de ... Cluster stats Portail - Mésocentre ...

mésocentre de calcul de franche-comté

DÉCONNEXION

PORTAIL - MÉSOCENTRE DE CALCUL DE FRANCHE-COMTÉ

Explorateur Jobs **Script** Informations

Le but de cette page est de vous permettre de générer un script SGE.
 Vous pourrez retrouver vos script grâce À l'explorateur dans le répertoire SGE/scripts.
 La sortie et les erreurs se trouvent quant à eux dans le répertoire SGE/out.

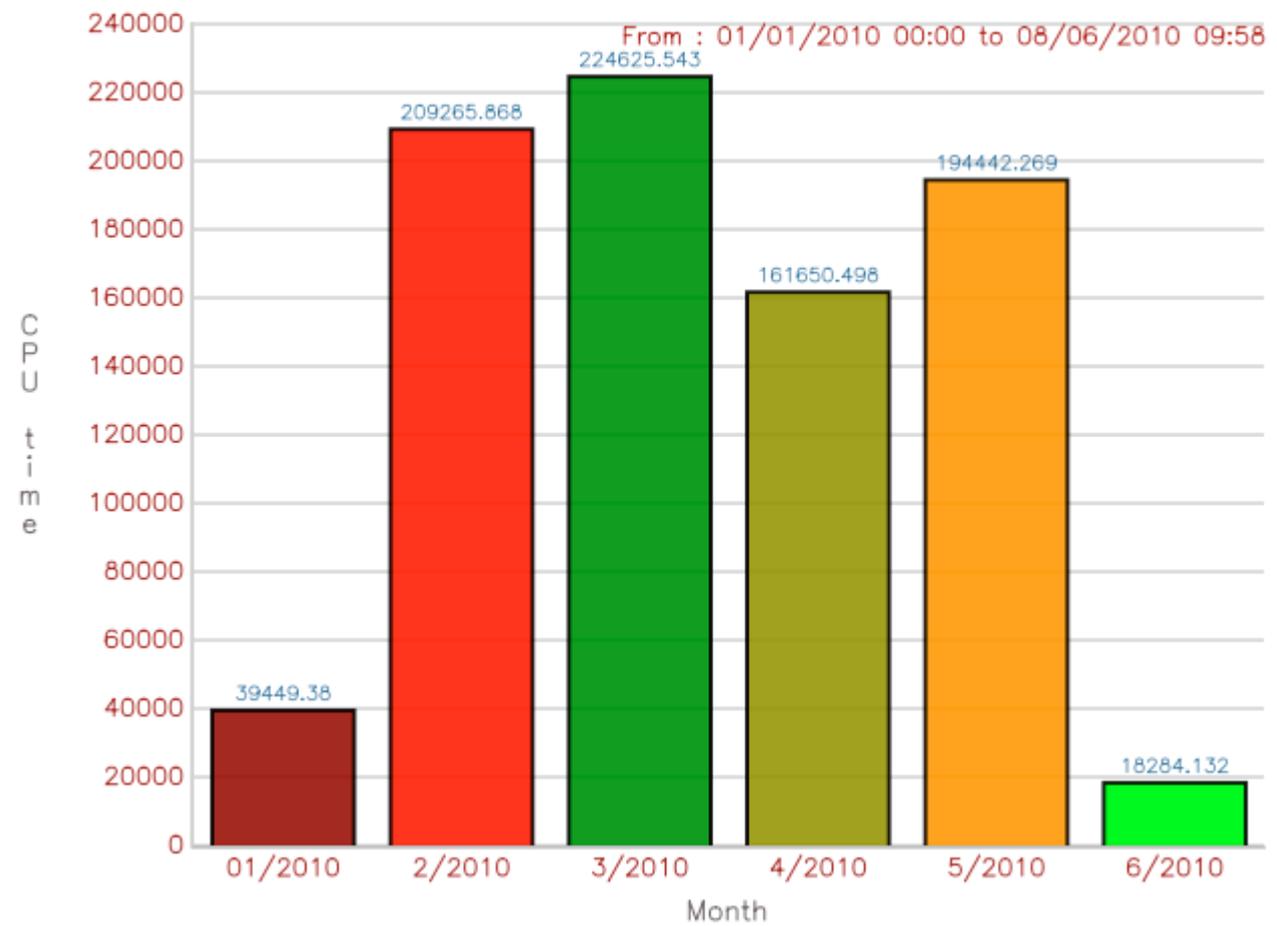
Ressource	Valeur
Type de Job <ul style="list-style-type: none"> séquentiel : un seul slot parallèle : plusieurs slots 	<input type="radio"/> séquentiel <input checked="" type="radio"/> parallèle
Environnement parallèle	impi_tight Intel MPI
Nombre de slots	<input type="text"/>
Queue	normal2h
Mémoire utilisée	<input type="range"/> 9279 MO
Temps maximum d'exécution <ul style="list-style-type: none"> plus ce temps est court, plus la priorité du job sera importante le job sera tué s'il dépasse ce temps 	<input type="text" value="2:00:00"/> h (format H:MM:SS)
Nom du Job	<input type="text" value="myjob"/>
Nom du fichier ou stocker les erreurs	<input type="text" value="\$JOB_NAME.\$JOB_ID.err"/>
Nom du fichier ou stocker la sortie	<input type="text" value="\$JOB_NAME.\$JOB_ID.out"/>
Chemin complet de l'application	<input type="text"/>
Paramètres de l'application	<input type="text"/>
Commandes supplémentaires (commandes système):	<input type="text"/>

Terminé



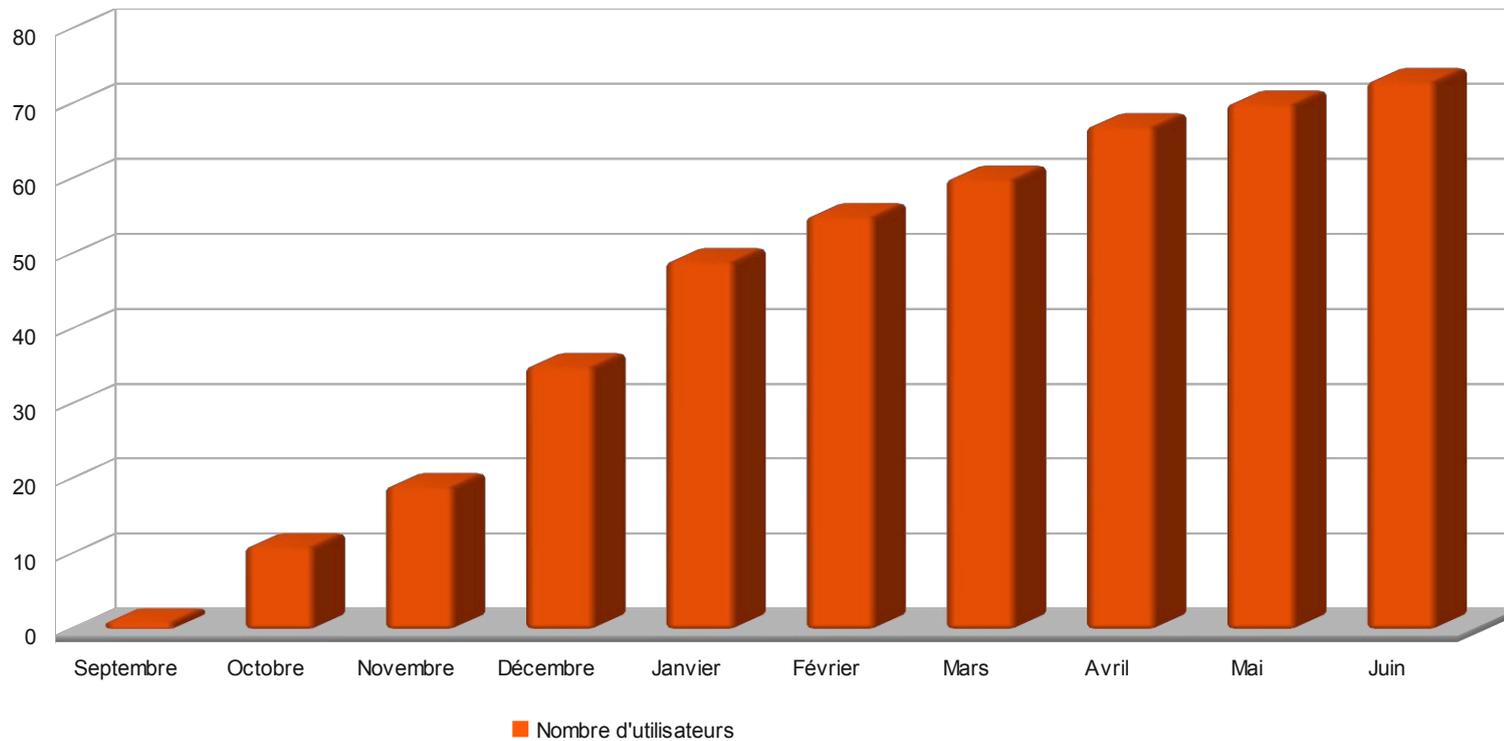
Utilisation du cluster

CPU time per month in hour





Evolution du nombre d'utilisateurs



● Environ 30 réels



Problèmes

- Matériels :
 - Panne courant
 - Composants: disques, mémoires, carte mère
 - Mise à jour SGE
- MPI :
 - 3 versions : Bull MPI, Intel MPI, OpenMPI
 - Problèmes :
 - Orphan daemon : fin d'exécution (Namd, VASP)
 - Recompilation des codes avec OpenMPI
- Fichiers Windows = dos2unix



Bilan

- Fonctionnel rapidement
- Nombre d'heures
- Nombre d'utilisateurs :
 - Croissance rapide
 - Utilisation des anciens clusters
 - Contacts industriels
- Formation :
 - 5 jours : Optimisation, OpenMP, MPI
 - 1 à 2 jours : sujet spécifique
- Financements ?



Evolutions

- Très vite confronté aux problèmes d'évolution
- Réseau :
 - Ajouts limités : 72 noeuds
 - Baisser la bande passante:
 - Facteur bloquage $\frac{1}{2}$ = 90 noeuds
 - Sous-clusters = hétérogénéité
 - Adaptation des files ?
 - Achat nouveau switch :
 - 3 x 36 = 180 noeuds
 - 96 ports = 216 noeuds



Evolutions 2

- Matériel :
 - Noeud interactif 32 coeurs
 - Ajout de noeuds de calcul
 - Espace de stockage 40 To
- Logiciels :
 - Système : changer de distribution
 - Tests :
 - Machines virtuelles
 - Wine