

U-LITE: una proposta per il futuro calcolo scientifico ai LNGS

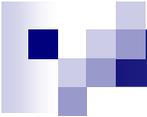
S. Parlati

LNFS 28 novembre 2011



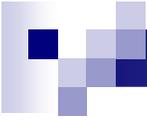
25 anni di calcolo scientifico ai LNGS

- Il Servizio di Calcolo e Reti
 - Nazzareno Taborogna (da 25 anni nel servizio)
 - Sandra Parlati (da 15 anni)
 - Stefano Stalio (da 12 anni)
 - Roberto Giuliani (da 12 anni)
 - Piero Spinnato (da 5 anni)
- Supporto di G. Di Carlo
- Andrea Donati nel servizio fino al 2005
- Borsisti e laureandi



25 anni di calcolo scientifico ai LNGS

- Anni '80 e '90: sistema **fortemente centralizzato** basato su cluster VMS e rete DECNET; Macro, EAS-Top, Gallex, LVD
- Dal 2000: diversificazione dei sistemi operativi (Digital Unix, Linux..); ambienti di calcolo **eterogenei** tra esperimenti: tendenza degli esperimenti a far ricorso a risorse di calcolo indipendenti (es. Borexino, Icarus, Opera).
- Il Servizio ha continuato a offrire un ambiente di calcolo completo (interattivo, batch, storage e backup) agli esperimenti (es LVD, Luna..). **Storage e backup** centralizzato utilizzato anche da esperimenti con farm di calcolo propria.
- 2010-2011: modello condiviso, gestito centralmente; nuove tecniche (virtualizzazione) permettono di avere ambienti eterogenei sullo stesso hardware -> **U-LITE** modello di calcolo per i futuri esperimenti.



Stato del calcolo scientifico ai LNGS

- Il Servizio di calcolo gestisce:
 - Sistemi di **storage** multipli per un totale di 150TB
 - Due librerie di nastri (120 slot ognuna) per il **backup e l'archiviazione** situate in due diversi edifici per aumentare l'affidabilità
 - La cella **AFS** /afs/lngs.infn.it
 - Servizio di **public login**: circa 600 utenti
 - Cluster **LSF**
 - **Kerberos e LDAP** per utenti LNGS e per gli utenti delle farm di esperimento (xenon, warp, darkside, lucifer, gerda, lvd, teo): circa 1000 utenti registrati.
 - Librerie, software scientifico, etc...

2 sale calcolo distinte, site in edifici differenti, per aumentare l'affidabilità dei sistemi e minimizzare il rischio di perdita di dati

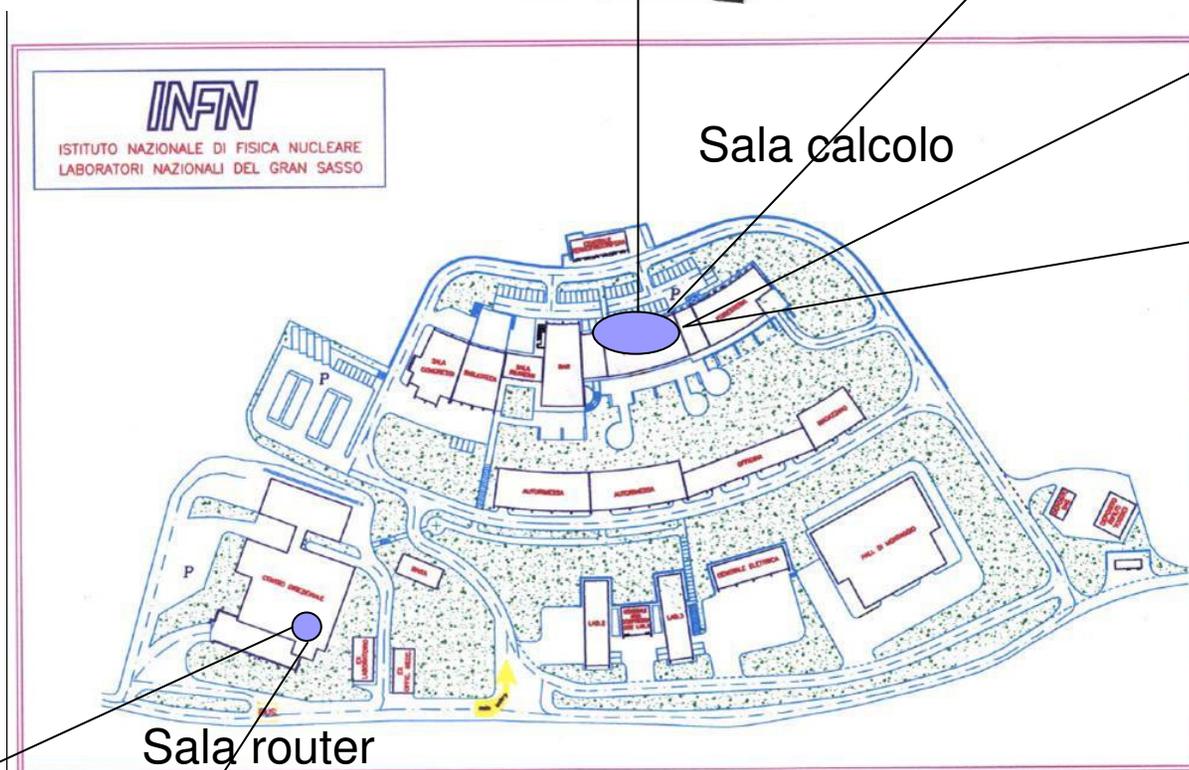
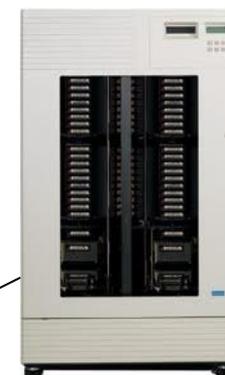
Storage



Cpu



Backup



Sala calcolo

Sala router

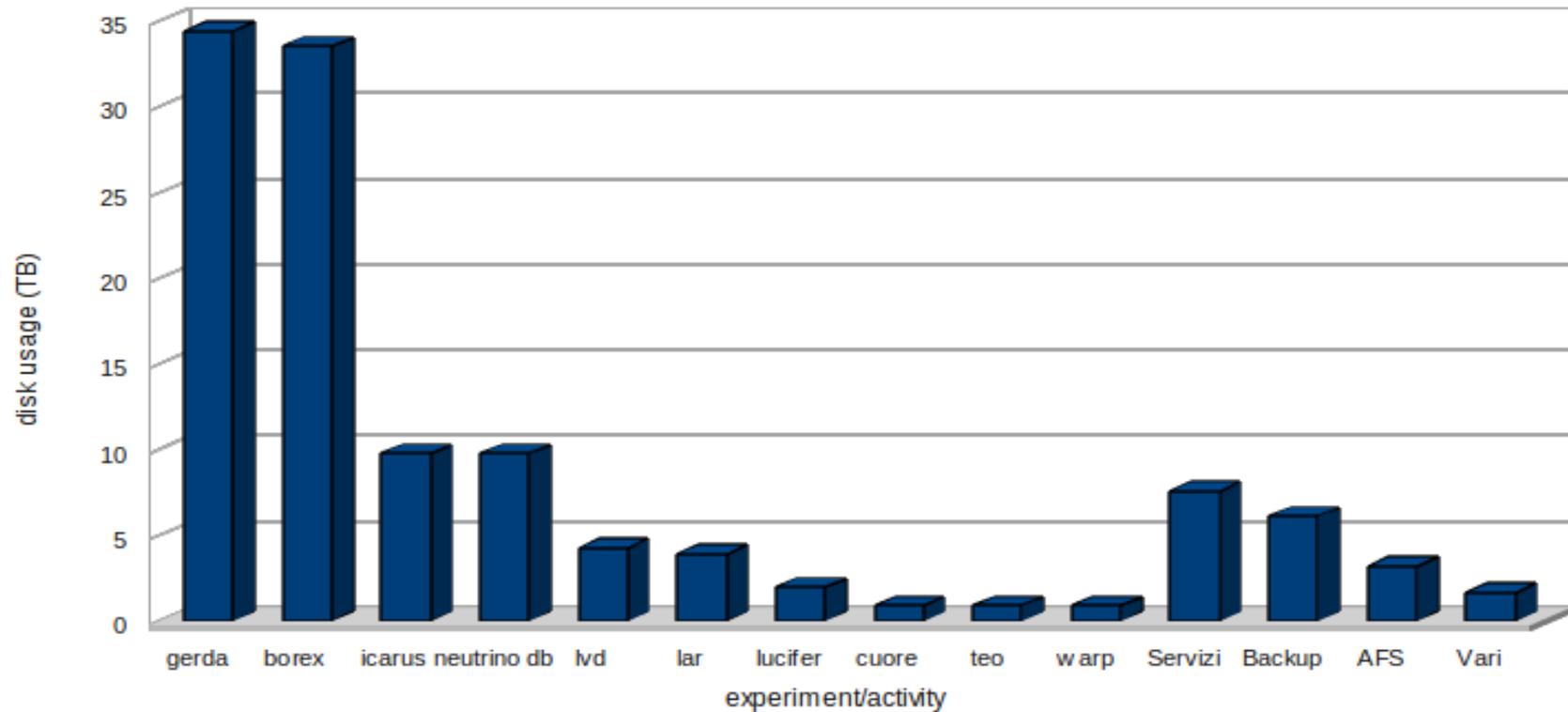


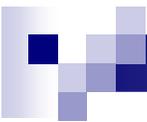
Servizi



Entrambi i locali sono dotati di condizionamento e sono sotto gruppo di continuità e generatore diesel in caso di back out.

Es: storage gestito centralmente



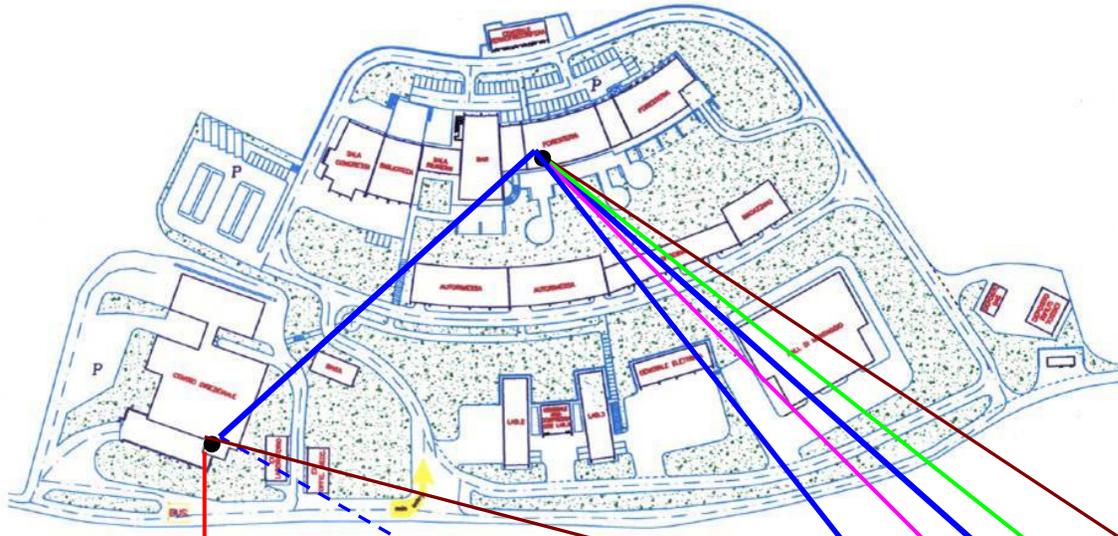


Stato del calcolo scientifico ai LNGS (2)

- L'affidabilità della rete e' fondamentale quando i dati sono in un tunnel sotto la montagna e distanti 10Km!
- Il Servizio di calcolo gestisce la rete locale in ogni suo aspetto, dando grande importanza all'affidabilità dei collegamenti con il sito sperimentale:
 - Percorsi fisici ridondati e con cammini differenti
 - Link aggregation per aumentare la capacità e l'affidabilità dei link
 - Spanning tree: percorsi di backup pronti ad attivarsi nel caso di failure dei link principali
 - Centro-stella ridondati ai lab esterni e sotterranei e VRRP per ridondare il routing interno
- In aggiunta alla LAN condivisa, alcuni esperimenti utilizzano link punto-punto dedicati per il trasferimento dati dal DAQ allo storage.
- Il Servizio gestisce il link geografico con il GARR

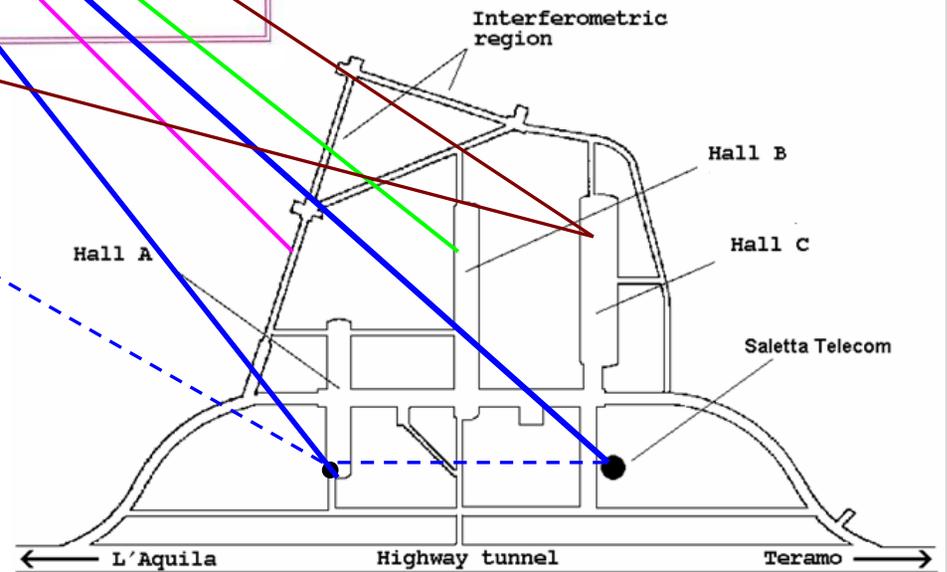
INFN

ISTITUTO NAZIONALE DI FISICA NUCLEARE
LABORATORI NAZIONALI DEL GRAN SASSO



- LAN LNGS 12Gb/s
percorsi principali +
percorsi STP 1 o 2 Gb
- Borexino 2Gb/s
- Icarus 3Gb/s
- Xenon 1Gb/s

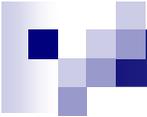
2 Gb/s GARR-X
(da agosto 2011)





Stato del calcolo scientifico ai LNGS (3)

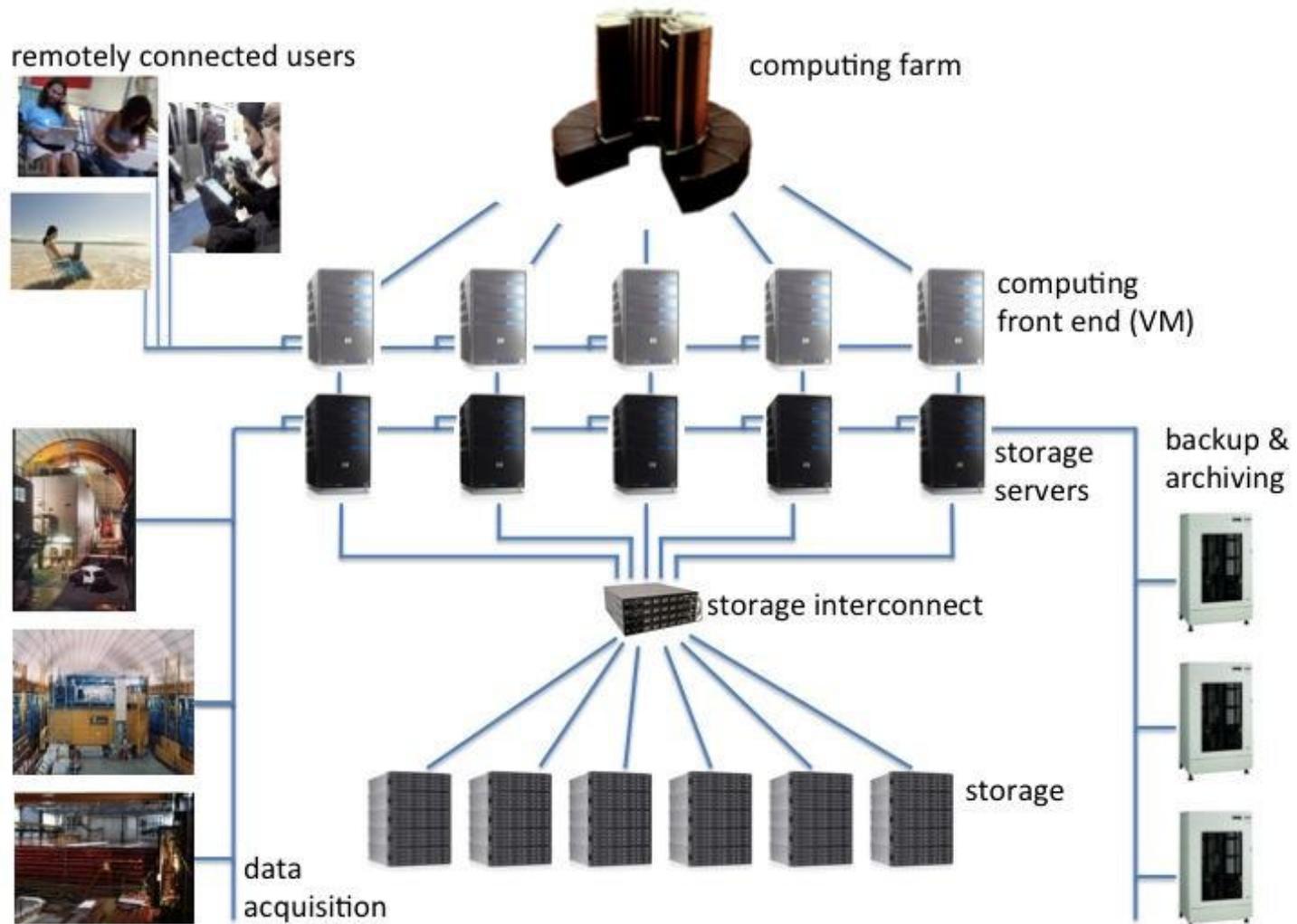
- Tutti i servizi di rete (NFS, AFS, DNS, DHCP, Radius, etc..) lavorano in regime di alta affidabilità.
- Tutti i sistemi (calcolo, rete, servizi) sono controllati da un sistema di monitoring centralizzato (Nagios) che invia alert via mail o via sms in caso di guasti o malfunzionamenti.

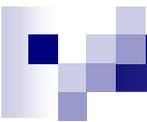


Motivazioni per una nuova proposta

- Offrire un **miglior servizio agli esperimenti**: creare un ambiente di calcolo completo dove e' più necessario ossia lì **dove il dato è acquisito**.
- Risorse umane: gestione a carico del servizio di calcolo con **personale esperto on-site** permette agli esperimenti di risparmiare risorse umane ed previene problemi causati dalla gestione di risorse di calcolo da parte di personale non qualificato o non motivato.
- Aspetto **economico**: la condivisione delle risorse riduce le spese infrastrutturali e ne ottimizza l'uso.
- La **virtualizzazione** rende possibile la condivisione dello stesso HW da parte di ambienti sw diversi. L'esperienza acquisita con la virtualizzazione dei servizi ci e' tornata utile nella progettazione di un nuovo modello per il calcolo scientifico.

U-Lite: Unified LNGS IT Infrastructure



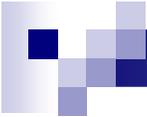


U-Lite: gli ingredienti (1)

- **Storage servers:**

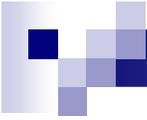
- 1 o 2 per ciascun esperimento. Copiano i dati dal DAQ sul sistema di storage, eseguono la preanalisi o la riduzione dei dati, salvano i dati sul sistema di backup e presentano i dati alle macchine di front-end e alla computing farm.
- Per il ruolo delicato che rivestono, a metà tra il DAQ e l'offline, la gestione degli storage server è a carico degli esperimenti, con possibilità per il Servizio di calcolo di intervenire con diritti amministrativi in caso di necessità.

- **Computing front-end:** 1 o più per esperimento. Macchine per il login interattivo, per l'analisi dati interattiva, per sviluppo del software e per la sottomissione dei job al computing cluster.



U-Lite: gli ingredienti (2)

- **Computing farm:**
 - più host virtuali (KVM) che girano su un cluster di server fisici multicore.
 - Ogni gruppo/esperimento ha i propri nodi di calcolo basati su un **template di VM, sviluppato dalle collaborazioni** sulla base delle proprie esigenze.
 - I job di simulazione o analisi vengono sottomessi ad un sistema di code (Torque/Maui) che li esegue sui nodi di calcolo virtuali.
- **Storage:** sistemi RAID ridondati (doppio controller, doppio alimentatore).
- **Backup e archiviazione:** librerie di nastri (attualmente LTO4); software di backup open source (bacula); formato di scrittura dei dati aperto.
- **Autenticazione e Autorizzazione:** utenti gestiti centralmente attraverso servizi Kerberos e LDAP del laboratorio.



Stato attuale di U-Lite

- Tutte le componenti HW e SW di U-Lite sono testate e pronte all'uso: un primo nucleo di U-Lite e' già operativo da aprile 2011 (6 host fisici, 2x4 core ciascuno, 10.25 HepSpec per core, costituiscono il cluster di calcolo)
- Storage e backup sono **consolidati**; il cluster di calcolo e' operativo, il sw di gestione dei job e' sotto sviluppo.
- Il sistema e' scalabile
- Il pool HW crescerà con l'arrivo di nuovi esperimenti in U-LITE.
- Sono disponibili strumenti per il **monitoring on-line** dei job e **l'accounting**.
- Alcuni gruppi cominciano a testare/usare U-Lite (es:Gerda)

On-line job monitoring di U-Lite

File Edit View History Bookmarks Tools Help

http://qmaster.lnfs.infn.it/cgi-bin/status

mysql arithmetic in query

La Republic... Nagios U-Lite comp... job statistics... INFN Web SquirrelMail ... Proxmox Virt... template.php

CRM Monitor

User	Running jobs	Max slots
DEFAULT	0	100
gbruno	1	100
gdicarlo	0	100
giordano	0	100
pandola	33	100
spinnato	0	100
stalio	0	100

Jobs run on Fri 25 November 2011

Queued	62
Run	44
Completed	23
Deleted	0

Jobs run on Thu 24 November 2011

Queued	73
Run	72
Completed	60
Deleted	2

Jobs run on Tue 22 November 2011

Queued	17
Run	17
Completed	17
Deleted	0

Jobs run on Mon 07 November 2011

Queued	1
Run	1
Completed	1
Deleted	0

Jobs run on Fri 04 November 2011

Queued	9
Run	8
Completed	8
Deleted	1

Server running on host qmaster
Start time Thu Nov 24 15:08:54 2011
Local time Fri Nov 25 12:54:06 2011

Accounting records
Last Week
Last Month
Last Quarter

Torque is running	11/25/2011 12:53:31:0100:PBS_Server:job:1159.qmaster.lnfs.infn.it:dequeuing from teo-long. state COMPLETE
Maiui is running	11/25 12:53:33 INFO: scheduling complete. sleeping 30 seconds
CRM is running	Fri Nov 25 12:53:55 2011 - jobs (queued/running): 18/33. nodes (down/up): 4/13. free resources (slots/cores/RAM): 17/-2/14325
CRM last log entry	Fri Nov 25 12:53:34 2011 - jobs (queued/running): 18/33. nodes (down/up): 4/13. free resources (slots/cores/RAM): 17/-2/14325

Node	RAM	CPUS	Server	Last Op	Idle	State	Properties	Job	Run Time	Queue	Owner	Group
ge-login	4096	4	hnode00	none	0	job-exclusive	gerdanode,ge-login,always_on	1122.1123.1124.1125	02:42:45	gerda-long	pandola	gerda
vnnode001	4096	4	hnode00	none	50	free	teonode,vnnode001,always_on					
vnnode002	4096	4	hnode05	none	10	free	teonode,vnnode002,always_on					
vnnode003	4096	4	hnode02	start	80	free	teonode,vnnode003					
vnnode004	4096	4	hnode02	start	80	free	teonode,vnnode004					
vnnode005	4096	4	hnode04	none	0	offline.job-exclusive	gerdanode,vnnode005	1098.1119.1120.1121	02:45:33	gerda-long	pandola	gerda
vnnode006	4096	4	hnode02	none	0	offline.job-exclusive	gerdanode,vnnode006	1115.1116.1117.1118	02:43:03	gerda-long	pandola	gerda
vnnode007	2048	2	hnode05	stop	0	down	gsnode,vnnode007					
vnnode008	2048	2	hnode06	stop	0	down	gsnode,vnnode008					
vnnode009	4096	4	hnode05	stop	0	down	lvdnode,vnnode009					
vnnode010	4096	4	hnode06	none	0	job-exclusive	lvdnode,vnnode010,always_on	1148	01:19:52	lvd-long	gbruno	lvd
vnnode011	4096	4	hnode02	stop	0	down	gerdanode,vnnode011					
vnnode012	4096	4	hnode01	none	0	offline.job-exclusive	gerdanode,vnnode012	1111.1112.1113.1114	02:42:16	gerda-long	pandola	gerda
vnnode013	4096	4	hnode05	start	0	offline.job-exclusive	gerdanode,vnnode013	1126.1127.1128.1129	02:41:33	gerda-long	pandola	gerda
vnnode014	4096	4	hnode03	none	0	offline.job-exclusive	vnnode014,gerdanode	1107.1108.1109.1110	02:42:49	gerda-long	pandola	gerda
vnnode015	4096	4	hnode01	none	0	offline.job-exclusive	vnnode015,gerdanode	1103.1104.1105.1106	02:43:24	gerda-long	pandola	gerda
vnnode016	4096	4	hnode03	none	0	offline.job-exclusive	gerdanode,vnnode016	1099.1100.1101.1102	02:43:23	gerda-long	pandola	gerda

Server	Slots (used/free/total)	CPU cores (used/free/total)	RAM (used/free/total)	CPU speed
hnode00	2/2/4	8/0/8	8192-13/8179	2332.000
hnode01	2/2/4	8/0/8	8192-13/8179	2332.000
hnode02	4/4/8	13-5/8	13312/3119/16431	1595.781
hnode03	2/2/4	8/0/8	8192-13/8179	2332.000
hnode04	1/3/4	4/0/4	4096/8212/12308	2327.501
hnode05	2/2/4	8/0/8	8192-13/8179	2332.000
hnode06	2/2/4	5/3/8	5120/3046/8166	2327.598

Queue	Priority	Max Time	Group	Nodes	Jobs (running/queued/total)
gerda-long	200	9999:00:00	gerda	gerdanode	32/18/50
gerda-short	100	02:00:00	gerda	gerdanode	0/0/0
gerda-xpress	300	00:20:00	gerda	gerdanode	0/0/0
gs-long	200	9999:00:00	gs	gsnode	0/0/0
gs-short	100	02:00:00	gs	gsnode	0/0/0
gs-xpress	300	00:20:00	gs	gsnode	0/0/0
lvd-long	200	9999:00:00	lvd	lvdnode	1/0/1

Done

4 down 4 unknown 12 warnings 1 critical

Accounting

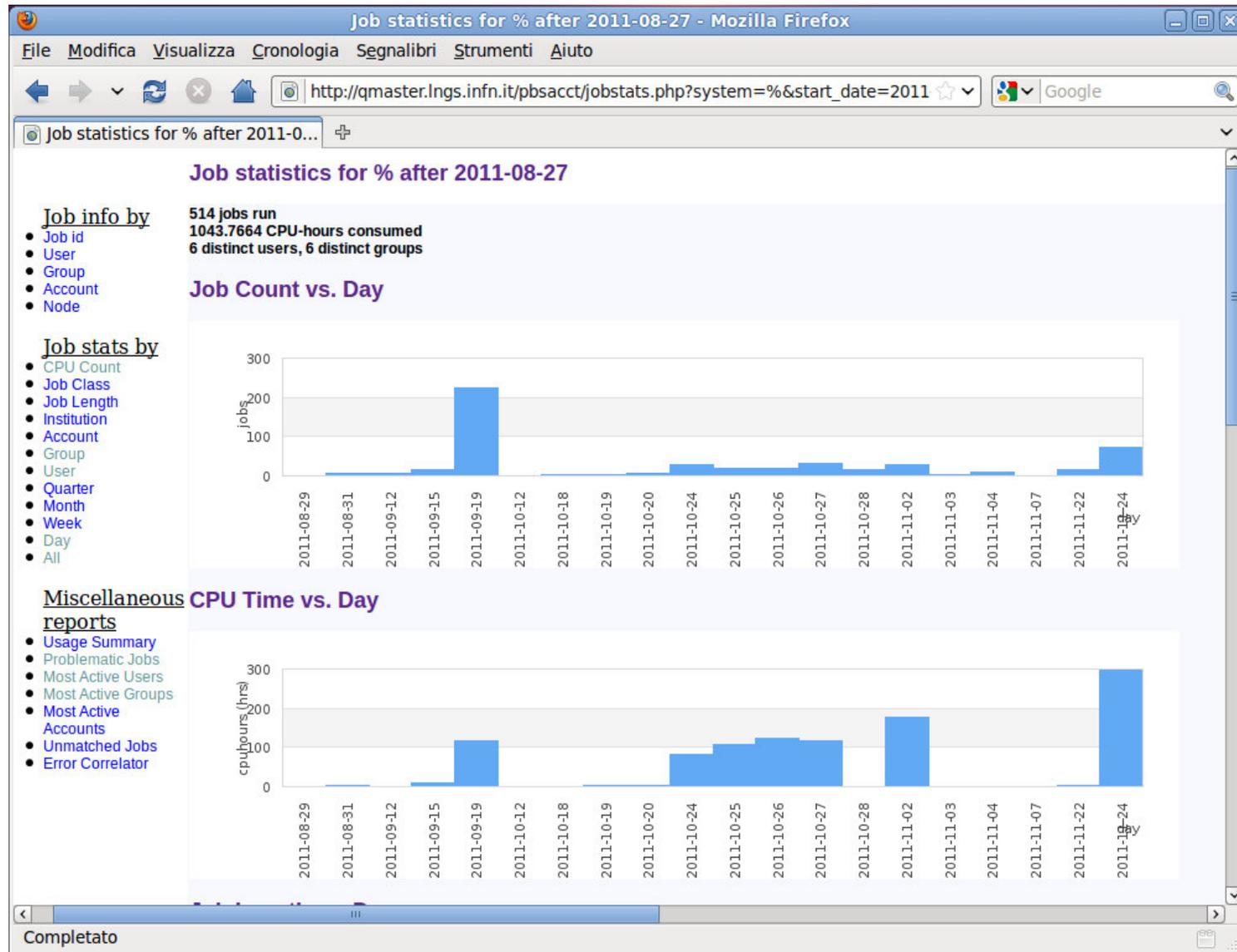
Most active groups on % between 2011-09-01 and 2011-11-25

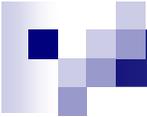
group	users	job count	CPU-hours
gerda	1	166	898.0006
macro	1	294	125.4928
200	1	17	7.9650
lngs	2	16	6.2875
lvd	1	11	4.2639
teo	1	2	0.0006

Bookmarkable URL for this report:
http://qmaster.lngs.infn.it/pbsacct/active-groups.php?system=%&start_date=20

Completo

Accounting





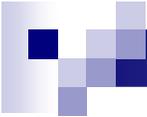
Chi e' coinvolto in U-Lite?

■ Responsabilità

- Gestione e coordinazione: S. Parlati
- Responsabilità tecniche: P. Spinnato and S. Stalio
- Supervisione e pianificazione: G. Di Carlo

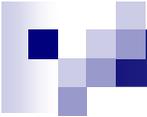
■ Staff tecnico

- Servizi di base: tutto il Servizio di Calcolo (5 persone)
- Dedicato staff: P. Spinnato e S. Stalio
- Necessità stimata di FTE:
 - **3 FTE**, distribuiti tra personale del servizio e personale in formazione (2 borsisti già richiesti per il 2012), per il setup di U-LITE
 - **1.5 FTE** distribuiti tra due o più membri del Servizio quando U-lite sarà pienamente operativo.



Modello di condivisione dei costi

- Il **Laboratorio**, attraverso il Servizio di Calcolo, ha finanziato l'attuale **infrastruttura** di U-Lite (rack, rete, SAN, librerie di tape) più un nucleo di macchine di calcolo, con il contributo della **CCR**.
- Costi vivi di **cpu, dischi e tape** sono a carico degli **esperimenti**;
 - dischi e tape sono ad accesso riservato
 - le macchine di calcolo si inseriscono in un ambiente condiviso; ogni collaborazione finanzia solo la cpu necessaria in media mentre i picchi di utilizzo sono coperti dalle risorse di altre collaborazioni o finanziate dal Servizio di Calcolo.
 - Il sistema batch di U-Lite prevede algoritmi di allocazione delle risorse per garantire sempre alle collaborazioni l'accesso alle risorse direttamente finanziate.
- U-Lite e' scalabile e stimiamo che l'infrastruttura attuale rimanga sostanzialmente invariata per una crescita di un fattore 10.
- Gli interventi fisiologici di mantenimento dell'infrastruttura sono a carico del laboratorio.



Vantaggi di U-Lite

- **Presenza stabile di personale esperto on-site per la gestione dell'intera infrastruttura.**
- Grande attenzione del personale all'**affidabilità** dei sistemi e **consapevolezza** dell'importanza dei dati (i LNGS sono il Tier-0 degli esperimenti!).
- **Libertà** per gli esperimenti di sviluppare l'ambiente di simulazione e analisi su qualunque piattaforma linux.
- Nessuna necessità di adeguarsi all'ambiente GRID
- Risparmio garantito grazie alla condivisione delle risorse



Links

- <http://u-lite.Ings.infn.it>
Sito principale di U-LITE
- <http://qmaster.Ings.infn.it>
Job monitoring e accounting
- <http://computing.Ings.infn.it>
Sito del Servizio di Calcolo e Reti LNGS