

Infrastrutture di calcolo su GRID in Italia

IFAE 2007
Napoli, 12 aprile

Enzo Miccio
INFN/CNAF - CERN

Sommario

- ▶ In prospettiva
- ▶ Progetti attuali e futuri
- ▶ Lo stato attuale
 - ▶ GRID in produzione, ora
 - ▶ GRID in preparazione, per LHC
- ▶ Conclusioni

Infrastruttura GRID

La disponibilità di servizi di rete e **protocolli standard** su uno strato fisico costituito da **risorse** di calcolo, spazio di storage e supercomputers, reso accessibile e condivisibile da un'adeguata organizzazione

Un po' di storia

- ▶ All'inizio del 2000 viene approvato il progetto **INFN Grid** (*grid.infn.it*)
 - ▶ 20 sedi, un centinaio di persone coinvolte
 - ▶ collaborazione fra fisici, ingegneri, computer scientist...
 - ▶ la motivazione principe era rispondere alla sfida del *computing problem* per LHC, ma sin dall'inizio il progetto è stato aperto alle esigenze provenienti da altri campi di ricerca come quello biomedico o di osservazione terrestre, come pure ad applicazioni industriali
 - ▶ è stato il punto di partenza per lo sviluppo dell'infrastruttura GRID in Italia

Un po' di storia

- ▶ Nel 2001, grazie all'esperienza di INFN Grid, nasce il progetto **Grid.it**
 - ▶ vengono stanziati fondi dal ministero e vengono coinvolte altre istituzioni oltre all'INFN (CNR, ASI...)
 - ▶ l'obiettivo è quello di fornire le basi per un'infrastruttura comune a tutta l'area di ricerca italiana (IRA)

Un po' di storia

- ▶ Nel 2001, nell'ambito del V Programma Quadro europeo e in collaborazione con il CERN, diversi paesi europei e alcune industrie, INFN Grid lancia **DataGrid** (cern.ch/eu-datagrid), pietra miliare verso la costruzione di un'infrastruttura a livello dell'Area di Ricerca Europea (ERA).
- ▶ Il progetto sfocia in **EGEE** (**E**nabling **G**rids for **E**-scienc**E**, eu-egee.org), un progetto quadriennale (2004-08) finanziato nell'ambito del VI Programma Quadro europeo e orma a metà della sua seconda fase.

La strategia

- ▶ Sviluppo di middleware e infrastruttura all'interno di progetti europei ed internazionali (accedendo a fondi extra-INFN: EU, MIUR...), passando spesso per il coordinamento del CERN (DataGrid, EGEE, EGEE II, WLCG...), spesso promossi dallo stesso INFN
- ▶ Collaborazione internazionale (Open Science Grid, Open Grid Forum...) improntata a garantire l'interoperatività globale dei servizi sviluppati e l'adozione di standard internazionali
- ▶ Sviluppi nazionali del middleware nelle aree non coperte da progetti UE

Oggi

► L'infrastruttura di produzione GRID italiana

- più di 40 centri di ricerca coinvolti
- le risorse sono raggiungibili attraverso servizi specifici per ciascuna VO
- la maggior parte di essi (~30) sono coinvolti anche a livello internazionale (EGEE/LCG)
- gli altri sono accessibili attraverso servizi di grid su scala nazionale



Oggi

► Portale operativo:



<http://grid-it.cnaf.infn.it/>

Documentazione per l'utente

Documentazione per gli amministratori

Repository per il software

Monitoring

Sistema a ticket per la notifica di problemi

FAQ e supporto

- News
- Organisation
 - ▶ People & tasks
 - ▶ Deployment
 - ▶ Meetings
 - ▶ Internal doc
- Access to the grid
 - ▶ Enable your UI [upd]
 - ▶ Get your certificate
 - ▶ Register to a VO
 - ▶ Use the grid
 - ▶ Applications
- Manage your site
 - ▶ Installation
 - ▶ Upgrade
 - ▶ Releases
 - ▶ Quattor
 - ▶ Quattor test [upd!]
 - ▶ More doc
 - ▶ Test & cert
- Grid status
 - ▶ Job monitor
 - ▶ GSTAT (Cnaf)
 - ▶ Grid services
 - ▶ Calendar
 - ▶ Downtime Advices
- Support
 - ▶ Ticketing System
 - ▶ Knowledge base
 - ▶ Broadcast Advice
- EGEE SA1
 - ▶ EGEE SA1 Italy
 - ▶ INFN SC Report [upd]
 - ▶ EGEE SA1 Europe
 - ▶ SAM Admin (Cnaf) [new]
 - ▶ Gridops.org [new]
- Search
- Links

Welcome to the INFN and Grid.it Production



This is the official web site of

INFN-GRID is a research project implementing and widespread international research projects

The national Grid.it project is implemented in various environments. Important installations include Astronomy, Biology, Geology

We're coordinating our efforts to build a Grid.

Our efforts are evaluated in terms of the best standards in building blocks over which

Read how real users are using

Latest news

Domani (FP6)

- ▶ Garantire l'evoluzione del Middleware Grid Open Source verso standards internazionali (OMII-Europe)
- ▶ Contribuire alle attività informatiche di Ricerca e Sviluppo
- ▶ Coordinare l'espansione di EGEE nel mondo (EUMedGrid, Eu-IndiaGrid, EUChinaGrid...)
- ▶ Sostenere l'allargamento di EGEE a nuove comunità scientifiche
 - ▶ GRIDCC (Applicazioni real time e controllo apparati)
 - ▶ BionfoGrid (Bionformatici; Coordinato dal CNR)
 - ▶ LIBI (MIUR; Bionfomatici in Italia)
 - ▶ Cyclops (Protezione Civile)

Oggi

- ▶ La fase di R&S preliminare è ormai ampiamente superata
- ▶ L'infrastruttura di GRID è ormai funzionante a livello di produzione per oltre 20 VO e migliaia di job al giorno
- ▶ In particolare gli esperimenti di LHC (CMS, ATLAS) non potrebbero più fare a meno, *oggi*, della GRID
- ▶ Il successo scientifico stesso di LHC è strettamente vincolato al successo della GRID
- ▶ Il successo di GRID è strettamente vincolato alla sua capacità di soddisfare le richieste di LHC

LHC e GRID

► Motivazioni

- il CERN da solo può fornire solo una frazione delle risorse necessarie
- decine di istituti possono contribuire con risorse
- necessità di integrare tali risorse

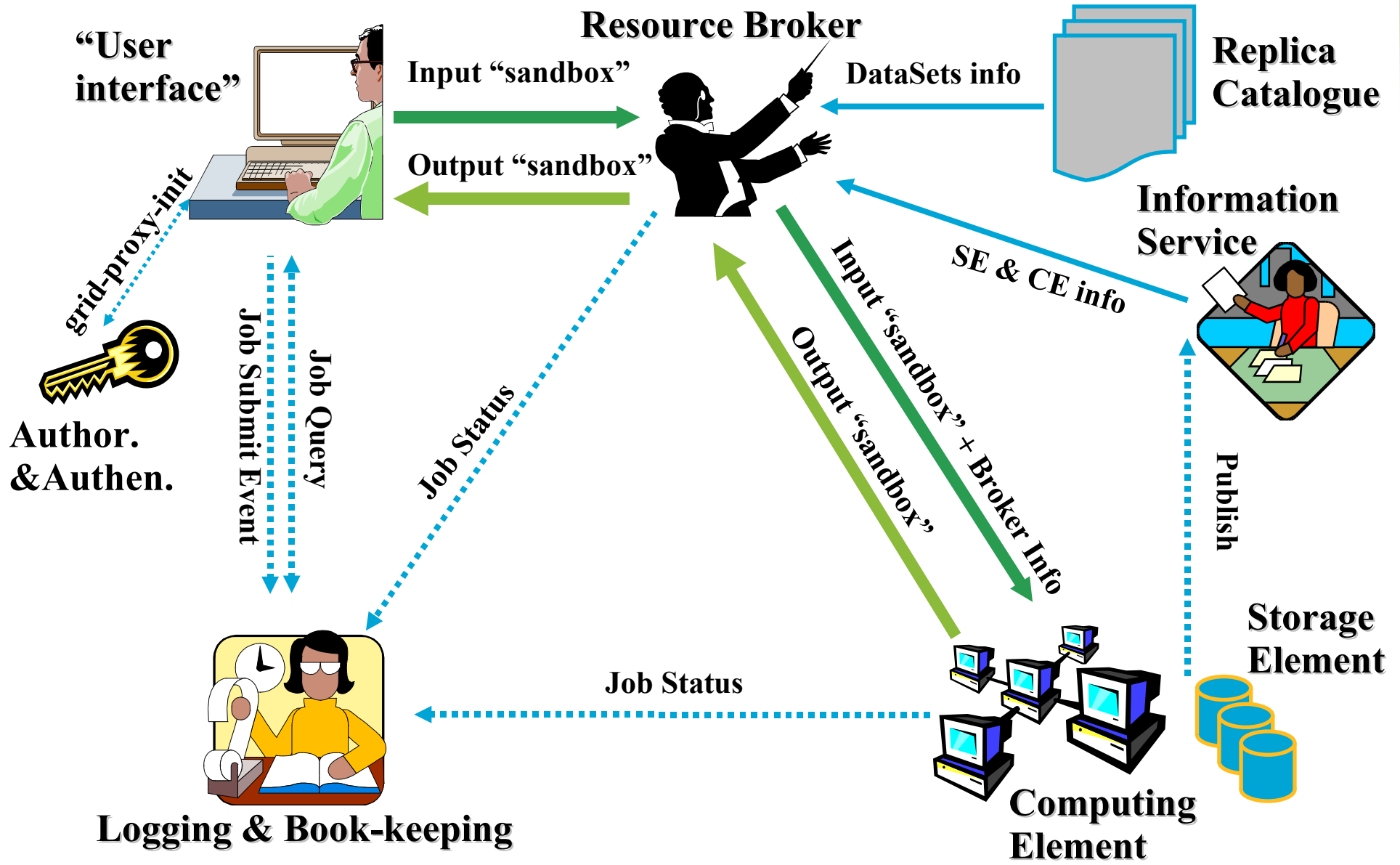
► Requisiti

- accesso uniforme e “user-friendly”
 - interfacce standard a risorse eterogenee
 - protocolli standard per l’accesso ai dati
- ottimizzazione dell’accesso ai dati
 - distribuzione intelligente e il più possibile automatica dei dati
 - assegnazione intelligente delle risorse di calcolo richieste

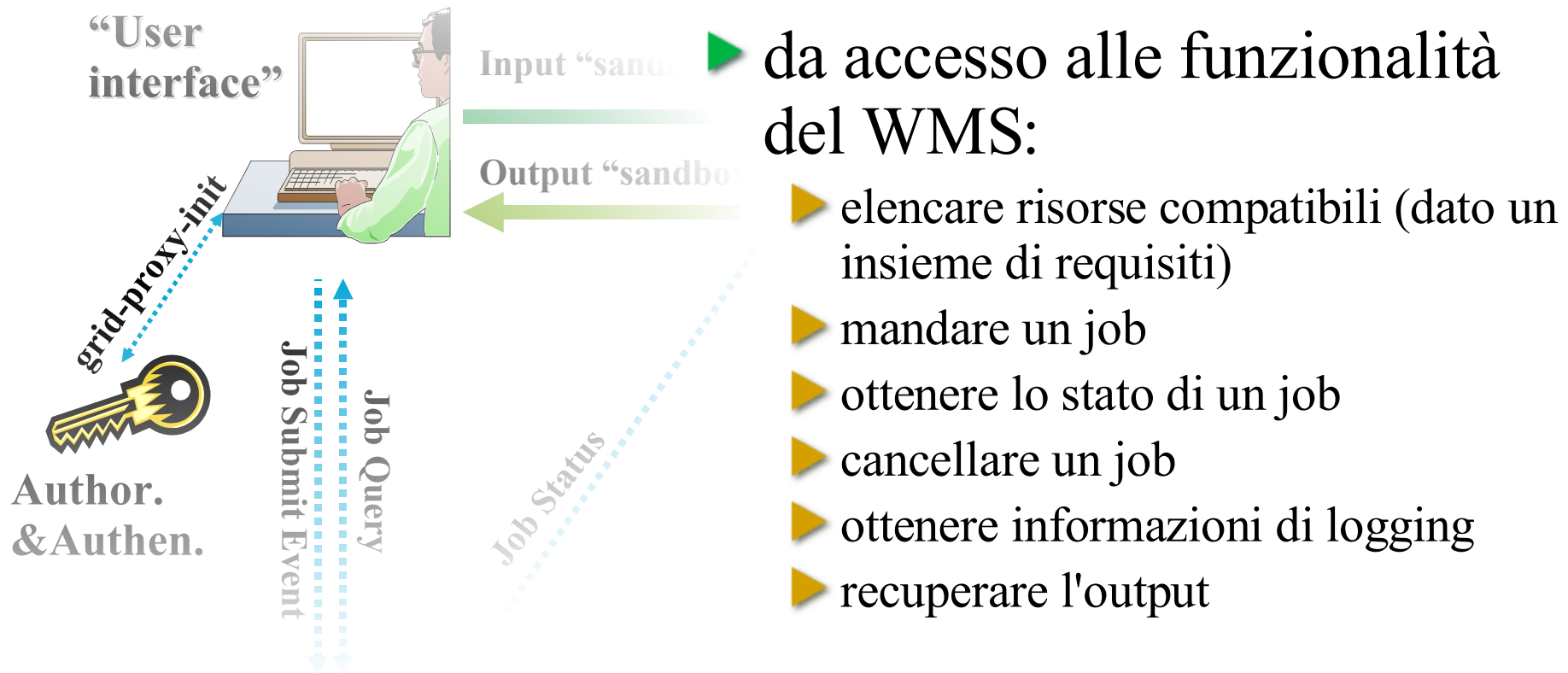
LHC e GRID

- ▶ Componenti principali:
 - ▶ Servizi di autenticazione e autorizzazione
 - ▶ Workload Management
 - ▶ Information System
 - ▶ Data Management
 - ▶ Computing Element
 - ▶ Monitoraggio

La vita di un job



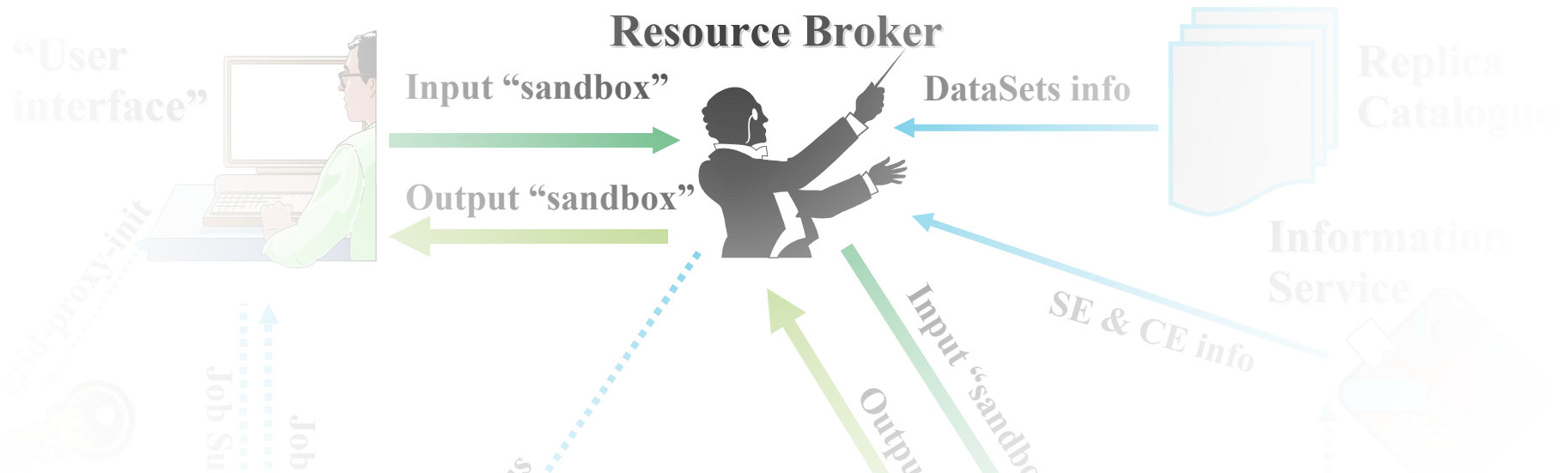
User Interface



▶ Autenticazione e autorizzazione

- ▶ Gestione delle VO e dei privilegi degli utenti
- ▶ corrispondenza tra utenti (identificati via certificato X.509) e account locali

Workload Management



- ▶ accetta e gestisce i job degli utenti
- ▶ seleziona le risorse più appropriate
- ▶ tiene traccia di quello che succede ai job
- ▶ restituisce l’output all’utente

Information System

- ▶ Fornisce in tempo reale lo stato della Grid (servizi e risorse)
- ▶ usato dal WMS
 - ▶ per sapere quali risorse di calcolo possono soddisfare le richieste degli utenti



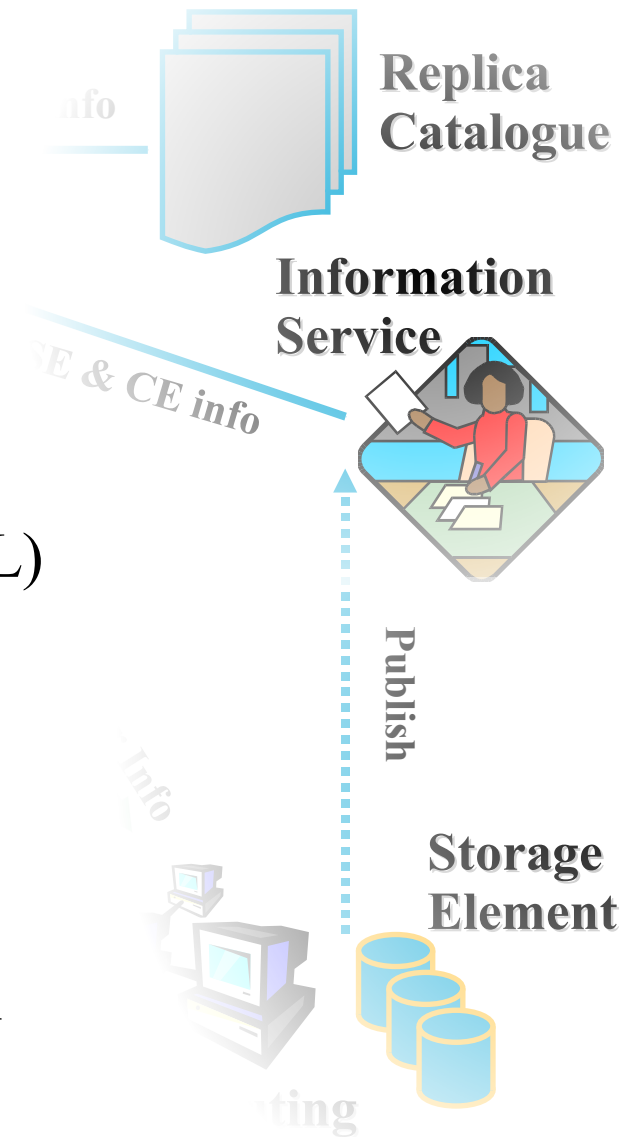
Data Management

▶ Catalogo dei file

- ▶ contiene le mappe tra
 - ▶ *Logical File Names (LFN)*
 - ▶ *Globally Unique Identifiers (GUID)*
 - ▶ *Physical File Names (PFN)*
- ▶ è centralizzato (backend Oracle o MySQL)

▶ Storage Element

- ▶ Essenzialmente un disk server (eventualmente front-end a un sistema di mass storage)



Monitoraggio

- ▶ visualizza lo stato presente e passato della Grid
- ▶ consente di diagnosticare i problemi
- ▶ modello gerarchico (come per l'IS)
 - ▶ l'informazione viene generata sulle singole macchine da appositi sensori
 - ▶ poi viene raccolta a livello di sito
 - ▶ infine viene spedita e immagazzinata in un database centrale



Logging & Book-keeping

- ▶ diversi sistemi attualmente in funzione, con finalità spesso diverse
 - ▶ R-GMA
 - ▶ GridICE
 - ▶ Site Functional Tests

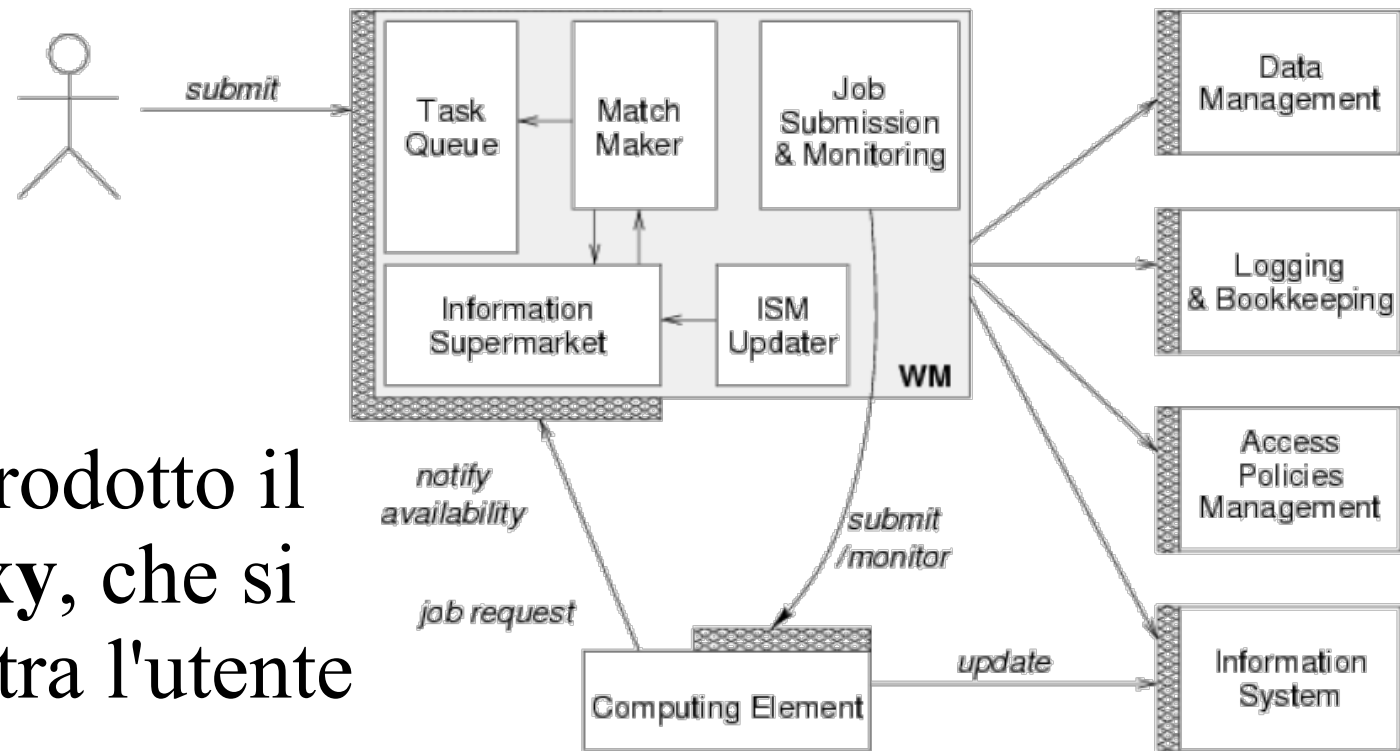
GridICE >> Site::ALL

Site ▼	Computing Resources												Storage Resources		
	General											Network			
	GK#	Q#	RunJob	WaitJob	SlotLoad	MH#	Power	WN#	CPU#	CPUload	Available	Total	%		
CNR-ILC-PISA	1	6	4	5	100%	4	18K	2	4	100%	380.9 GB	381.6 GB	0%		
ENEA-INFO	1	3	0	0	43%	18	72K	15	15	2%	13.8 GB	18.3 GB	24%		
ESA-ESRIN	1	3	0	0	0%	8	12K	6	6	0%	12.3 GB	16.9 GB	27%		
INAF-Trieste	1	4	0	0	0%	9	35K	8	16	0%	463.2 GB	522.2 GB	11%		
INFN-BARI	1	4	50	1	37%	43	378K	36	72	34%	311 GB	698.5 GB	55%		
INFN-BOLOGNA	1	6	0	0	-	7	102K	5	18	5%	33.6 GB	71.3 GB	50%		
INFN-BOLOGNA-CMS	1	2	0	0	0%	8	54K	6	12	0%	1.2 TB	1.7 TB	30%		
INFN-CAGLIARI	1	6	7	0	25%	14	92K	11	22	5%	910.5 GB	938.1 GB	3%		
INFN-CATANIA	1	5	0	0	-	95	1M	76	224	4%	2.3 TB	3.5 TB	55%		
INFN-CNAF	1	33	7	80	75%	27	38K	4	8	58%	960.2 GB	1.7 TB	44%		
INFN-FERRARA	1	4	1	0	6%	13	82K	9	18	11%	10 GB	32.9 GB	70%		
INFN-FIRENZE	1	6	22	4	39%	18	379K	16	64	36%	511.3 GB	722.1 GB	29%		
INFN-FRASCATI	1	3	5	0	83%	5	34K	3	6	68%	1.1 TB	1.2 TB	14%		
INFN-GENOVA	1	3	0	0	0%	6	42K	4	8	0%	337.9 GB	341.7 GB	1%		
INFN-LECCE	1	3	0	0	0%	4	11K	2	4	1%	15.2 GB	17.9 GB	15%		
INFN-LNL-2	1	5	8	0	7%	77	776K	70	140	6%	752.8 GB	2 TB	63%		
INFN-LNS	1	5	0	0	0%	6	90K	4	16	0%	475.8 GB	476.2 GB	0%		
INFN-MILANO	1	13	2	8	96%	32	274K	29	58	80%	1.7 GB	169.1 GB	99%		
INFN-NAPOLI	1	6	35	0	97%	27	251K	24	48	81%	115.2 GB	870.3 GB	87%		
INFN-NAPOLI-ATLAS	1	6	11	3	92%	-	-	-	-	-	-	-	-		
INFN-NAPOLI-CMS	1	6	3	0	25%	8	59K	6	12	25%	169 GB	417.8 GB	60%		
INFN-NAPOLI-VIRGO	1	7	0	2	0%	3	4K	1	2	0%	3.4 GB	7 GB	52%		
INFN-PADOVA	1	19	96	82	100%	60	531K	50	100	91%	9.2 TB	10.8 TB	14%		
INFN-PERUGIA	1	4	14	0	61%	24	200K	22	44	66%	201.7 GB	216.5 GB	7%		
INFN-PISA	1	3	0	3	0%	5	12K	2	4	0%	8.2 GB	9.3 GB	12%		
INFN-PISA2	1	6	15	0	89%	16	75K	13	23	54%	854.6 GB	4 TB	79%		
INFN-ROMA1	1	21	12	0	29%	26	280K	23	50	11%	1.6 TB	3 TB	47%		
INFN-ROMA1-CMS	1	2	1	6	10%	10	54K	5	12	0%	1.7 TB	1.9 TB	6%		
INFN-ROMA1-VIRGO	1	2	4	0	40%	8	35K	6	12	42%	194.6 GB	238.1 GB	18%		
INFN-ROMA2	-	1	3	1	0%	6	38K	4	8	0%	-	-	-		
INFN-ROMA3	1	3	0	0	0%	6	38K	4	8	0%	955.2 GB	956.7 GB	0%		
INFN-T1	2	13	306	0	86%	8	-	-	-	-	373 TB	559.3 TB	33%		
INFN-TORINO	1	9	2	0	4%	27	248K	24	48	0%	479.6 GB	2 TB	77%		
INFN-TRIESTE	1	6	2	11	100%	3	4K	1	2	50%	15 GB	26.1 GB	42%		
ITB-BARI	1	6	0	0	21%	14	294K	12	48	16%	51.9 GB	58.7 GB	12%		
SNS-PISA	1	6	0	0	0%	5	24K	3	6	0%	182 GB	203.2 GB	10%		
SPACI-LECCE-IA64	1	6	5	0	32%	9	15K	7	14	41%	18.7 GB	29.4 GB	36%		
SPACI-NAPOLI	1	6	0	0	0%	4	8K	3	3	0%	-	-	-		
SPACI-NAPOLI-IA64	1	6	1	79	26%	62	251K	60	120	11%	47.1 GB	63.4 GB	26%		
UNI-PERUGIA	1	5	0	0	0%	15	52K	13	26	0%	-	-	-		
TOTAL	# 40	41	265	614	284	35%	740	6M	589	1301	23%	398.4 TB	598.4 TB	33%	

Dal middleware LCG...

- ▶ Il Resource Broker di LCG si è dimostrato robusto, ma
 - ▶ il codice è ormai congelato
 - ▶ niente nuove feature
 - ▶ difficoltà di bug fix
 - ▶ la sottomissione via Network Server è troppo lenta
 - ▶ può richiedere decine di secondi per job, se il RB è carico
 - ▶ il rate massimo di sottomissione è limitato
 - ▶ l'esperienza mostra che non possono essere gestiti più di ~7000 jobs/day
 - ▶ non supporta il rinnovo dei VOMS proxy
 - ▶ i VOMS proxy sono ormai diventati uno standard *sine qua non* poiché permettono un'autorizzazione *fine-grained* (data access, job priorities)

...al middleware gLite



- ▶ Viene introdotto il **WMPProxy**, che si inserisce tra l'utente e il WMS vero e proprio, ottimizzando la gestione dei job

gLite WMS: vantaggi

- ▶ Task queue interna:
 - ▶ Se non ci sono risorse disponibili che corrispondono alle richieste di un job, questo può essere mantenuto in coda per un tentativo in un secondo momento
 - ▶ “Shallow resubmission”: il job viene ri-sottomesso se il fallimento è avvenuto prima di raggiungere il Worker Node
- ▶ *Information Supermarket*
 - ▶ Può sottomettere jobs basandosi su informazioni raccolte da più parti e raccolte in una cache locale

gLite WMS: vantaggi

- ▶ Bulk submission
 - ▶ Collections: insieme di job indipendenti
- ▶ Job sandbox
 - ▶ Condivisione dell'input sandbox per le collection
 - ▶ Download/upload delle sandbox via GridFTP, https, http
- ▶ Autenticazione più rapida via WMPProxy
- ▶ Match-making più rapido
 - ▶ *Bulk* match-making

gLite WMS: vantaggi

- ▶ Tempi di risposta più rapidi per l'utente
- ▶ Job throughput più alti

gLite WMS sotto test

- ▶ Servizio sperimentale
 - ▶ urgenza di avere un servizio funzionante e prestante
 - ▶ necessità di poterlo già usare in produzione (ATLAS)
- ▶ Sotto intensa attività di testing dall'estate 2006
 - ▶ Test congiunti ATLAS & CMS
 - ▶ Sono state usate come WMS poche macchine costantemente sotto controllo (CERN, Milano, CNAF) e aggiornate tutte con la stessa configurazione
 - ▶ Ciclo di test-patch-deploy molto a stretto giro

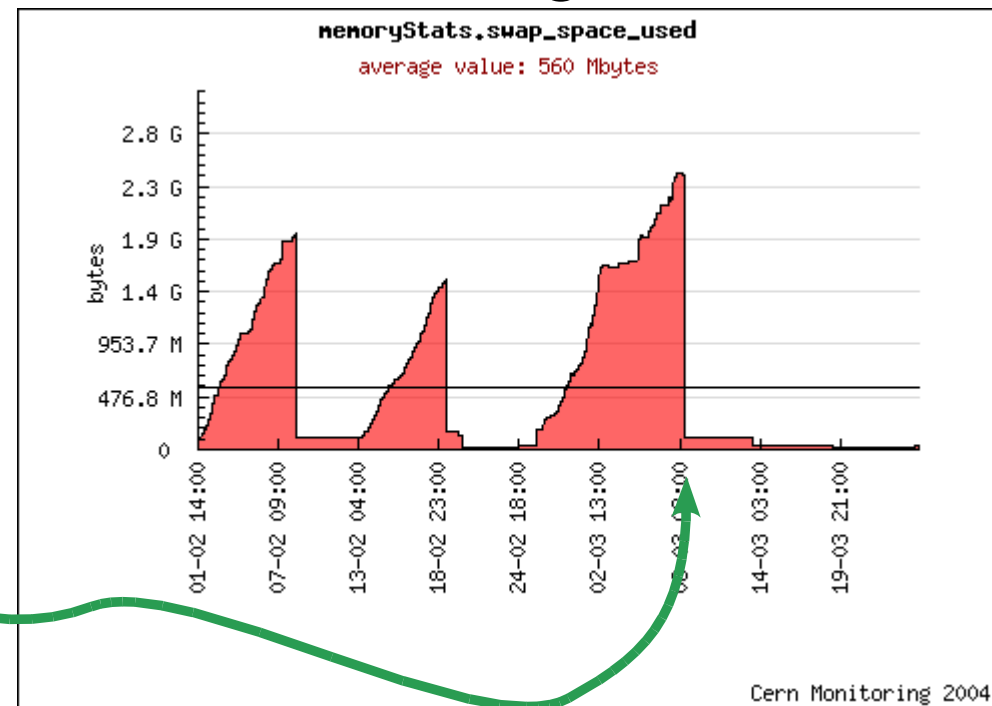
gLite WMS: problemi

▶ *memory leak*

▶ L'uso della memoria cresce linearmente in gLite WMS 3.0

▶ mantenerlo in funzione diventa difficile (restarts, reboots...)

▶ problema risolto in gLite WMS 3.1



gLite WMS: problemi

- ▶ Job che vanno in stallo
 - ▶ Problema serio in gLite WMS 3.0
 - ▶ 15% di jobs su 15K jobs/day
 - ▶ Peggiora considerevolmente sotto carico
 - ▶ Situazione migliore in gLite WMS 3.1
 - ▶ ~5% di jobs su 15K jobs/day
 - ▶ Il problema è stato individuato nel Condor DAGMAN
 - ▶ meccanismo di gestione delle collezioni
 - ▶ Problema risolto
 - ▶ come? ...rimuovendo Condor DAGMAN!
 - ▶ le collezioni sono gestite con un meccanismo più semplice
 - ▶ gli ultimi test non mostrano più alcun job in stallo

gLite WMS sotto test

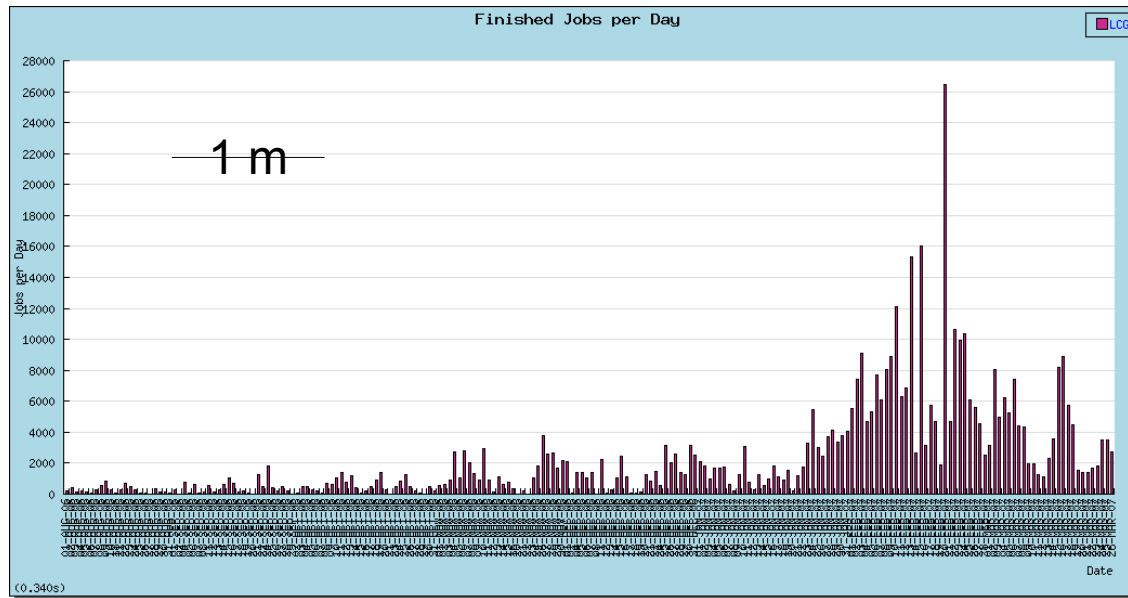
- ▶ bulk submission
 - ▶ target di 15k jobs/day per una settimana
 - ▶ Effettivi requirement di ATLAS
 - ▶ Semplici “HelloWorld” job, bulk submission
 - ▶ Limite di *Shallow Retry* pari a 5
- ▶ ultimi risultati
 - ▶ 10K job/day
 - ▶ meno dell'1% di job persi (in stallo)
 - ▶ ritmo sostenuto stabilmente per una settimana

Site	Submit	Wait	Ready	Sched	Run	Don (S)	Don (E)	Don (F)	Abort	Clear	Canc
ce05-lcg.cr.cnaf.infn.it	0	0	0	0	0	14188	0	0	12	0	0
ce06-lcg.cr.cnaf.infn.it	0	0	0	0	0	14160	0	0	40	0	0
cclcgceli02.in2p3.fr	0	0	0	0	0	13781	0	0	19	0	0
ce04.pic.es	0	0	0	0	0	14340	0	1	259	0	0
ce-fzk.gridka.de	0	1	0	2	0	10954	0	0	2043	0	0
lcgce01.gridpp.rl.ac.uk	0	0	0	0	0	12946	0	0	54	0	0
lcgce01.triumf.ca	0	2	0	0	0	12384	0	0	14	0	0
ce113.cern.ch	0	0	0	0	3	12809	0	0	388	0	0
ce114.cern.ch	0	0	0	0	0	13519	0	0	281	0	0
ce115.cern.ch	0	0	0	0	0	13919	0	1	80	0	0

gLite WMS sotto test

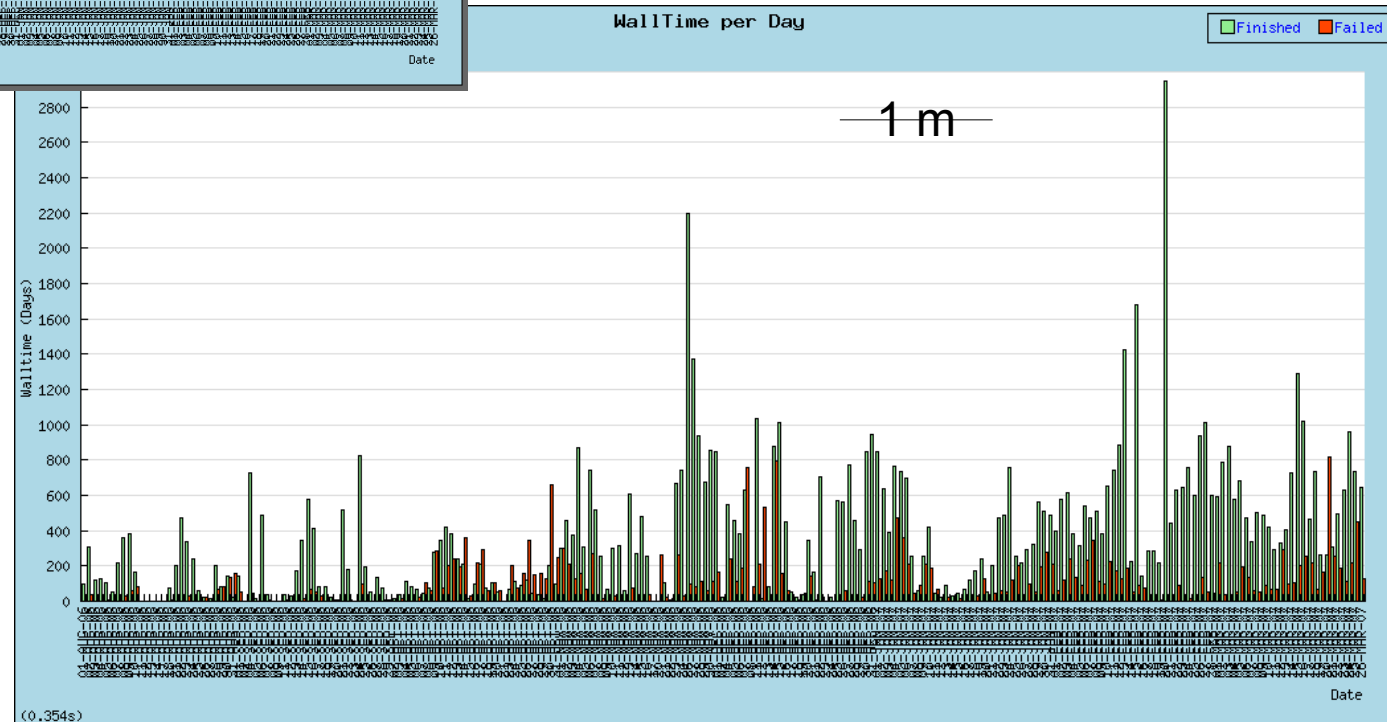
- ▶ non-bulk submission
 - ▶ viene sottomessa vera applicazione CMSSW
 - ▶ Limite di *Shallow Retry* pari a 3
 - ▶ Interrogazione frequente dello status dei job con *retrieving* automatico dell'output dalla UI
- ▶ ultimi risultati
 - ▶ raggiunto limite del ritmo di sottomissione *consecutivo* (~6k job/day)
 - ▶ limite superabile via sottomissione *parallela*, anche da una stessa UI
 - ▶ >10 jobs/day

...e in produzione



- ▶ job terminati vs giorni (ATLAS)
 - ▶ Significativo incremento nell'ultimo mese
 - ▶ raggiunti i 20k job/day

- ▶ Wall-Clock-Time
 - ▶ il tempo perso in job falliti (rosso) è tipicamente basso
 - ▶ occasionali aumenti dovuti a validazioni e sporadici incidenti



gLite WMS: stato attuale

- ▶ I problemi più grossi sono stati risolti e testati
- ▶ Gli altri problemi minori sono stati compresi e risolti
- ▶ Resta da effettuare i test con queste ultime patch

Conclusioni

- ▶ Grid è già funzionante
- ▶ Grid è già *determinante*
- ▶ Grid ha ancora sfide da affrontare