

Continuità operativa dei servizi informatici presso i LNGS

LNGS – Servizio Calcolo e Reti



I servizi informatici dei Laboratori Nazionali del Gran Sasso, con esperimenti che acquisiscono dati 24 ore al giorno ed utenti da tutto il mondo devono essere sempre usufruibili. Senza sistemi di supervisione automatica e di ridondanza i disservizi vengono, di solito, individuati abbastanza velocemente se si presentano durante le ore lavorative. Se però hanno luogo durante la notte o durante il fine settimana possono passare ore o giorni prima che vengano rilevati, affrontati e risolti.

Progetto

- Sistema di monitoring centralizzato
- Alta affidabilità dei servizi



Per questo è nata l'esigenza di realizzare un sistema di monitoring centralizzato, con particolare riguardo alla supervisione dei parametri ambientali del centro di calcolo, e ridondare parte dei servizi in maniera da "nascondere" all'utente malfunzionamenti e crash di sistema.

Un lavoro di questo tipo richiede un grande dispendio di tempo ed energia, almeno in una fase iniziale.

L'occasione per realizzare questo progetto ci è stata da una borsa di studio per diplomati della durata di un anno legata al progetto POR Abruzzo e finanziata dal Fondo Sociale Europeo che è stata bandita, dietro nostra proposta, su questo argomento.

Monitoring - Requirements

- Monitoring del sistema di alimentazione elettrica
- Monitoring del sistema di raffreddamento
- Monitoring dei sistemi di storage, di calcolo, di rete
- Monitoring dei servizi
- Allarmistica via e-mail e SMS
- Azioni in caso di problemi



La realizzazione di un sistema di supervisione centralizzato e` stato il primo compito che ci siamo dati.

Il sistema di supervisione deve tenere sotto controllo elementi infrastrutturali, apparati di rete, sistemi di storage, server, servizi.

La notifica degli allarmi deve essere immediata e diretta alle persone che hanno in carico il sistema malfunzionante: e-mail e SMS ci sono sembrati i veicoli piu` adatti.

Il sistema di monitoring deve anche essere in grado di compiere azioni per tentare di risolvere o porre rimedio ad alcuni tipi di problemi.

Monitoring - Modello

- Un nodo raccoglie informazioni, le elabora e si occupa della notifica via e-mail e SMS e di eventuali azioni sugli apparati supervisionati
- Comunicazione su rete ethernet
- Software open source standard
- GUI per la visualizzazione dello stato del sistema e possibilmente per la configurazione



L'architettura di massima del sistema di monitoring prevede un unico nodo, eventualmente ridondato, che si occupi della raccolta dei dati, dell'elaborazione dell'informazione, della notifica degli allarmi e possa agire sui sistemi monitorati (ad esempio spegnendoli).

L'host che si occupa del monitoring dovrebbe poter raggiungere sia la LAN LNGS che la LAN privata dedicata ai servizi interni.

Tutta la comunicazione deve avvenire su rete ethernet usando i protocolli della suite TCP/IP. Il software dovrebbe essere basato sul sistema operativo linux, dovrebbe essere scritto usando un linguaggio interpretato di largo uso e presente comunemente su tutte le distribuzioni linux (perl, python..) e dovrebbe appoggiarsi solo su programmi e librerie universalmente presenti sui sistemi linux, qualunque sia la distribuzione usata. Una eventuale interfaccia grafica che permetta di compilare un file di configurazione (eventualmente in xml) sarebbe desiderabile.

Monitoring - Scelte

- Nagios (<http://www.nagios.org/>)
- Cacti (<http://www.cacti.net/>)
- Shutdown (fatto in casa)



Queste sono state le nostre scelte:

Nagios serve per la rilevazione dei problemi, per l'allarmistica via e-mail e SMS, per la visualizzazione delle mappe.

Cacti ci permette la visualizzazione dello storico di partametri importanti (carico delle ups, temperature, carico dei server)

Shutdown, un software scritto da Mario Cimini, si occupa di spegnere ordinatamente le macchine interne al centro di calcolo in caso di sovratemperatura o di mancata fornitura ENEL per piu` di 5 minuti.

Monitoring – sala calcolo

- Controllo stato degli UPS su rete nascosta
- Controllo sensori di temperatura
- Spegnimento automatico dei server in situazioni critiche
- Notifica anche al di fuori del servizio
- Inclusione delle farm di esperimento nel sistema di monitoring e di shutdown

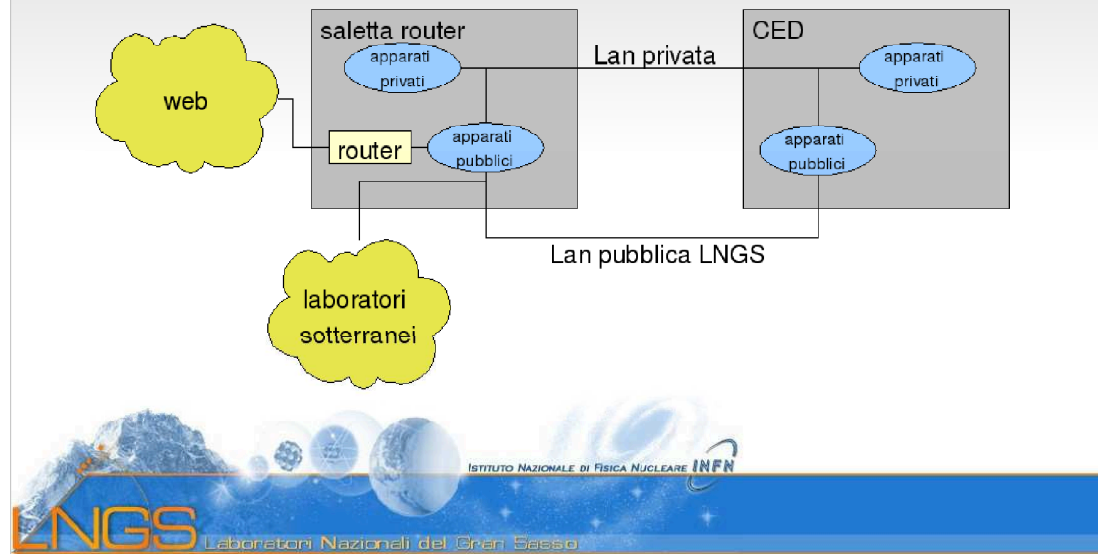


Per quanto riguarda il centro di calcolo ci siamo occupati di acquistare delle schede che permettono il controllo via rete (snmp) dei tre UPS che, in parallelo, alimentano l'edificio principale dei laboratori esterni, incluso lo stesso centro di calcolo. Nel centro di calcolo abbiamo anche installato due sensori di temperatura interrogabili via rete. Le informazioni ottenute dagli apparati descritti qui sopra sono rese disponibili anche a chi si occupa degli impianti elettrici e di refrigerazione e sono elaborate dal sistema di shutdown che, in casi critici, avvia la procedura di spegnimento delle macchine presenti all'interno del centro di calcolo.

L'inclusione degli host di esperimento nel sistema di monitoring centrale e di spegnimento automatico sta riscuotendo un buon interesse e sarà realizzata al più presto .

Alta affidabilità infrastruttura

- Click to add an outline

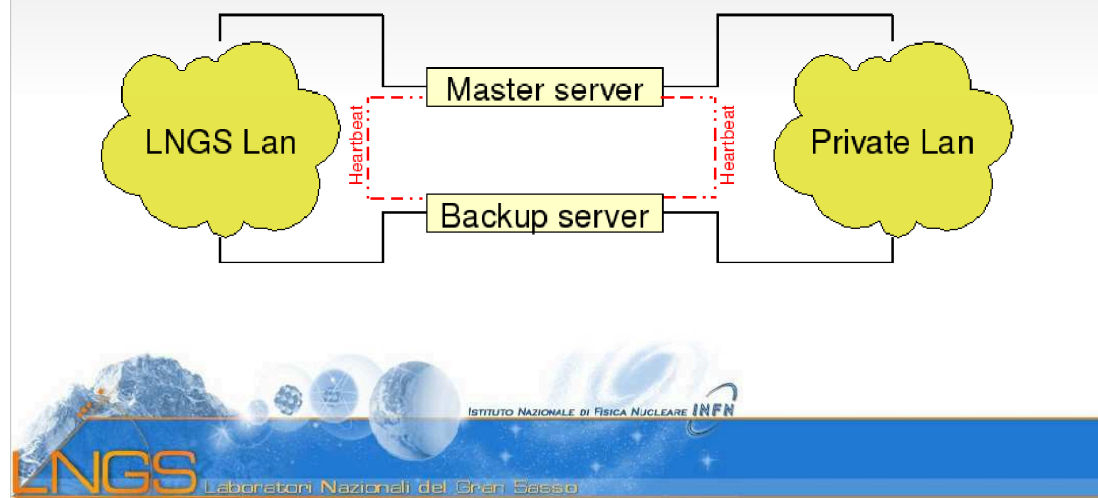


L'infrastruttura di rete dei LNGS, per quel che ci riguarda, comprende due locali dedicati ad ospitare apparati di rete e sistemi di calcolo e di storage. I due locali sono situati in edifici diversi e serviti da UPS, gruppi elettrogeni e sistemi di raffreddamento separati e sono connessi tra loro da due reti fisiche distinte.

Questa architettura è adatta a realizzare la ridondanza dei servizi perché ci permette di distribuire gli elementi dei cluster ad alta affidabilità in luoghi fisicamente lontani pur mantenendo la doppia connessione di rete necessaria per l'heartbeat tra gli host.

Alta Affidabilita` Heartbeat v.1

Click to add an outline



Heartbeat (<http://www.linux-ha.org/Heartbeat>) e` un software open source per l'implementazione di sistemi in regime di alta disponibilita`.

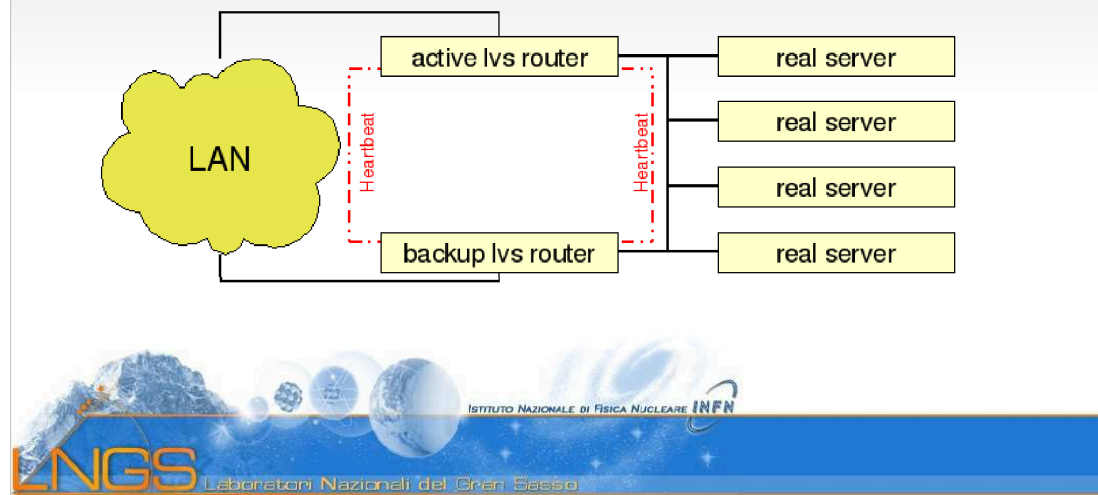
Ogni host che fornisce un servizio ha un "alter ego" col quale comunica attraverso due canali separati (normalmente di rete, ma anche seriali) e che e` pronto a sostituire l'host primario in caso di fallimento di quest'ultimo acquisendone gli indirizzi di rete ed i servizi offerti.

Heartbeat non e` in grado di reagire in caso ci siano problemi sul servizio offerto dal server principale ma il server stesso rimanga attivo e presente sulla rete.

Lo utilizziamo per ridondare i server web.

Alta Affidabilita` LVS

- Click to add an outline



LVS (<http://www.linuxvirtualserver.org/>) e` un software open source che permette di offrire servizi in regime di alta disponibilita` e con bilanciamento del carico.

Le richieste ai server reali passano attraverso un host chiamato "lvs director" che puo` essere ridondato con un sistema simile a quello di Heartbeat.

Tipicamente "dietro" al director vi sono coppie di server, ognuna delle quali offre un servizio specifico.

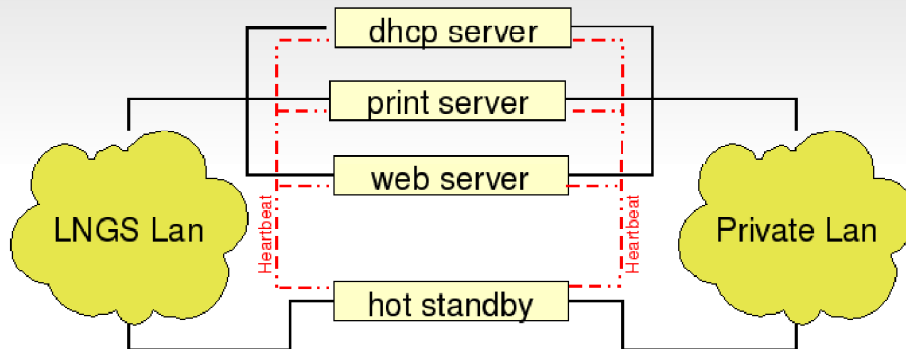
Il director distribuisce le richieste ai server reali usando algoritmi di bilanciamento del carico. Il director smette di dirigere le richieste verso gli host che non sono piu` capaci di fornire il servizio richiesto, rendendo cosi` invisibili agli utenti le failure di singoli host.

Lo utilizziamo per: DNS, proxy server, web server.

Internet Explorer ha problemi a connettersi con cluster LVS di web server se li contatta a livello 2.

Alta Affidabilita` Heartbeat v.2

- Click to add an outline



Con la versione 2 Heartbeat fornisce nuove feature:

- la possibilita` di creare dei cluster di piu` di due macchine
- un sistema (primitivo) di monitoring dei servizi
- un interfaccia grafico per la configurazione del cluster

L'interfaccia grafico pero` non consente un controllo completo del sistema e non e` semplice modificare a mano il file di configurazione (xml).

Alta Affidabilita` altre tecniche

- Alta affidabilita` integrata in alcuni servizi di rete di base:

nis, kerberos, mysql, smtp, DNS



Nel corso delle indagini che sono state necessarie per capire quali servizi possano essere resi piu` sicuri con sistemi di alta affidabilita` abbiamo anche indagato sui sistemi di alta affidabilita` inclusi nel codice del software di alcuni servizi: kerberos, dns, mysql, smtp, nis.

kerberos

Il sistema di ridondanza dei servizi di autenticazione di kerberos prevede l'esistenza di almeno due server di autenticazione: un master ed uno o piu` slave che periodicamente scaricano dal master il database aggiornato.

Il sistema funziona correttamente e in maniera trasparente all'utente a condizione che il file di configurazione di kerberos sui client indichi il nome di almeno due server kerberos del dominio o che i client siano configurati per chiedere il nome dei server kerberos ai name server locali.

nis

Il sistema di ridondanza dei servizi di autenticazione ed autorizzazione di NIS prevede l'esistenza di almeno due server di autenticazione: un master ed uno o piu` slave che periodicamente scaricano dal master il database aggiornato.

Il sistema funziona correttamente e in maniera trasparente all'utente a condizione che il file di configurazione di nis sui client indichi il nome di almeno due server nis del dominio o che i server nis possano essere trovati con richieste broadcast.

mysql

abbiamo implementato un sistema di replica in configurazione master-master sui server web (quello principale e quello di "scorta") in modo da avere i database sempre aggiornati e da poterli usare anche in caso di fallimento del web server principale ed entrata in attivita` di quello di "scorta".

dns

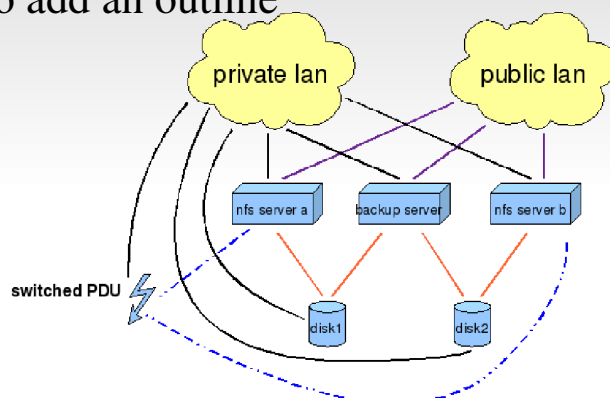
il sistema di ridondanza dei name server per un dominio funziona se il file di configurazione del dns del client elenca almeno due server dns. In caso di caduta del server che appare per primo nell'elenco, ogni richiesta di risoluzione di nomi viene passata al server successivo solo dopo essere andata in time out sul primo. Questo provoca gravi rallentamenti di tutte le applicazioni di rete, rendendo questo sistema di ridondanza praticamente inutilizzabile, anche con la migliore configurazione possibile dei client. Si e` deciso pertanto di usare LVS per creare un server DNS "virtuale". Le richieste dirette a questo server vengono smistate dal director LVS (come descritto nella sezione relativa a LVS stesso) verso i server DNS reali, rendendo trasparente all'utente il malfunzionamento di uno di essi.

smtp

il sistema di ridondanza dei server di posta elettronica basato sui record MX dei name server funziona bene ed e` probabilmente il modo migliore di offrire un servizio smtp in regime di alta disponibilita`.

Alta Affidabilita` - storage

Click to add an outline



Per quanto riguarda il sistema di condivisione dello storage in rete (NFS), che spesso ci aveva dato problemi, oltre ad aver cercato un setup che fosse il piu` stabile possibile, abbiamo provato ad utilizzare Heartbeat v2 per creare un cluster come quello in figura.

Abbiamo visto pero` che la configurazione di Heartbeat v2 e` molto complicata e non permette un controllo completo del sistema risultante. Abbiamo cosi` scelto di scrivere noi stessi un sistema di monitoring dei due (o piu`) server nfs attivi che gira su un server di backup. In caso di problemi il server di backup acquisisce le risorse del server malfunzionante dopo averlo fermato e spento attraverso un PDU (ciabatta) controllabile da remoto.

Alta Affidabilita` rete locale e accesso alla WAN

- Ridondanza fisica dei rami di backbone della LAN
- RSTP
- Linea di emergenza dai laboratori sotterranei verso Bologna



Il lavoro di realizzazione della ridondanza sulla LAN e di messa in funzione della linea di collegamento di emergenza alla rete GARR ha interessato e sta interessando un altro gruppo di persone all'interno Servizio Calcolo e Reti dei LNGS.

La linea di emergenza verso il GARR parte dai laboratori sotterranei e va ad attestarsi a Bologna passando dal versante teramano del Gran Sasso.

In questo modo il percorso fisico e` completamente separato da quello della linea di rete principale (dai laboratori esterni a Roma Tizii) e da quello della linea di backup (dai laboratori esterni al POP dell'Aquila).

Documentazione

- Content Management System
- Documentazione su carta
- Mappe fisiche di rete, apparati, server
- Mappe logiche dei servizi



L'implementazione di un sistema di monitoring capillare su rete ethernet e la realizzazione di cluster per l'alta affidabilità fanno aumentare il numero degli host e degli apparati presenti in rete e ancora di più gli indirizzi IP presenti sulla rete stessa.

Aumentano anche i percorsi di rete di interconnessione tra le macchine.

Il fatto che alcuni indirizzi di rete non risiedano sempre sullo stesso host e che la virtualizzazione sia usata in maniera massiccia crea molto confusione agli amministratori di sistema che non sanno più dove siano fisicamente i loro server.

Un sistema di documentazione sempre aggiornato è indispensabile!

Risorse Umane

- Condivisione dell'informazione
- Lavoro in gruppo
- Manutenzione della documentazione



Molto spesso la mancata conoscenza di come sono realizzati i servizi informatici e del loro stato all'interno dello stesso gruppo di persone che li gestisce e` causa del prolungarsi delle interruzioni del servizio.

La comunicazione tra persone, il lavoro di gruppo e la manutenzione della documentazione sono altrettanto importanti dei sistemi automatici di monitoring e dei cluster ad alta affidabilita`.

Conclusioni

- Miglioramento della qualità del servizio
 - Monitoring
 - Cluster HA
 - Sistemi di ridondanza interni ai servizi di rete
 - Verifica di scelte passate
 - Documentazione
- Aumenta la complessità del sistema



In generale, alla fine di questo anno di lavoro, possiamo dire che la qualità dei servizi informatici offerti è migliorata notevolmente.

Forse il fattore che pesa di più nel miglioramento della qualità del servizio è il sistema di monitoring che permette la notifica dei problemi in tempo reale. Pesano anche, oltre ai cluster HA, la documentazione, una corretta configurazione dei sistemi di ridondanza interni ai servizi di rete (ove esistenti) e la verifica che si è fatta su alcune scelte passate che non risultavano soddisfacenti.

A controbilanciare tutti questi fattori è la aumentata complessità del sistema, in particolare per quanto riguarda il numero di host sulla rete ed i percorsi di rete di interconnessione.

Persone

- Mario Cimini – borsista POR

Suo è il merito di gran parte del lavoro qui descritto. Sta scrivendo la sua tesi di laurea di primo livello in informatica su questo argomento.

- Fabio di Bernardini – contratto a progetto

Pur essendo impegnato su altri temi ha dato un contributo molto valido proponendo soluzioni che spesso sono state adottate e partecipando alle discussioni.

- Stefano Stalio - dipendente

ha coordinato il gruppo e ha lavorato al progetto quando non impegnato nella normale attività lavorativa.

