

ALMA MATER STUDIORUM · UNIVERSITÀ DI BOLOGNA

Scuola di Ingegneria e Architettura
Dipartimento di Informatica – Scienza e Ingegneria (DISI)
Corso di Laurea in Ingegneria Informatica

Tesi di Laurea in:
Amministrazione di Sistemi T

**SVILUPPO DI UN'INFRASTRUTTURA VIRTUALE
PER L'EROGAZIONE DI SERVIZI DI CALCOLO
CON SLURM**

Relatore:
Prof. Marco Prandini

Presentata da:
Massimo Valerio Zerbini

Correlatore:
Ing. Andrea Giovine

Appello V
Anno Accademico 2023/2024

*“It is only when they go wrong, that machines
remind you how powerful they are.”*

— Clive James

Abstract

Nel presente elaborato di tesi viene descritto e analizzato il processo di sviluppo di un'infrastruttura virtuale SLURM, per l'integrazione di specifici servizi di calcolo in ambito *High Performance Computing* (HPC). In particolare, si implementano la condivisione delle risorse di computazione, la definizione di una priorità di *scheduling* e la federazione di due *cluster* virtuali. La progettazione dell'ambiente avviene applicando il principio di *Infrastructure as Code* (IaC), tramite strumenti come Vagrant e Ansible. Durante il processo di sviluppo, si affrontano e analizzano le varie difficoltà legate alla corretta configurazione dell'intera infrastruttura distribuita. Le funzionalità integrate vengono infine convalidate tramite l'esecuzione di test mirati.

Indice

1	Introduzione	9
2	Stato dell'arte	11
2.1	Virtualizzazione	11
2.1.1	Vagrant	11
2.2	Provisioning	12
2.2.1	Ansible	13
2.3	DNS & DHCP	13
2.3.1	dnsmasq	14
2.4	Base di dati	14
2.4.1	MariaDB	15
2.5	Workload Management	15
2.5.1	SLURM	15
3	Analisi progettuale	19
4	Implementazione	23
4.1	Infrastruttura	23
4.1.1	Macchine virtuali	23
4.1.2	Configurazione di rete	25
4.1.3	DHCP & DNS	26
4.1.4	Accounting	27
4.1.5	Cluster SLURM	31
4.2	Condivisione di risorse	32
4.3	Priorità di scheduling	33
4.3.1	Impostazione del DB	34
4.4	Federazione di cluster	36
4.4.1	Impostazione del DB	37
5	Risultati	39
5.1	Test della condivisione di risorse	39
5.2	Test della priorità di scheduling	40
5.3	Test della federazione di cluster	41
6	Conclusioni	43
6.1	Sviluppi futuri	44

Bibliografia	47
Acronimi	49
Elenco delle figure	51
Ringraziamenti	53

1 | Introduzione

Il contesto dell'High Performance Computing (HPC) rappresenta l'avanguardia per la risoluzione di problemi complessi e computazionalmente avanzati. Esso consiste nel coordinamento di più unità di elaborazione (raggruppate in “*cluster*”) per fornire esecuzioni ad alte prestazioni, ricorrendo tipicamente al calcolo simultaneo [1]. La nascita del primo standard di programmazione parallela risale al 1990, ma in tempi recenti la potenza di calcolo ha raggiunto ordini di grandezza significativi per la simulazione di scenari reali.

L'applicazione di ambienti HPC viene principalmente associata all'ambito della ricerca scientifica, ma il suo utilizzo è esteso a qualunque disciplina che necessiti di calcolo a elevate prestazioni; alcuni esempi sono:

- la riproduzione e lo studio del clima globale (Climatologia);
- la risoluzione di equazioni fluidodinamiche (Fisica);
- la ricerca di metodi di conservazione per antichi testi e scritture (Archeologia);
- lo studio delle proteine per la cura di malattie degenerative (Medicina).

La progettazione e la manutenzione di un ambiente HPC coinvolgono tecnologie e aspetti tipici dell'ingegneria informatica, tra cui l'amministrazione di sistemi. È infatti di cruciale importanza la corretta impostazione della comunicazione tra i diversi nodi interessati.

In questo elaborato di tesi viene descritto e analizzato il procedimento di sviluppo di un'infrastruttura virtuale distribuita, dedicata all'implementazione di diverse funzionalità inerenti al contesto HPC. In particolare, verranno integrate:

- la **condivisione delle risorse di computazione** disponibili, al fine di raggiungere il pieno utilizzo della potenza di calcolo;
- la **priorità di *scheduling*** nei confronti di un determinato utente, ristretta a una specifica risorsa di computazione;
- la **federazione di *cluster***, ossia il coordinamento tra molteplici installazioni autonome per l'esecuzione di calcoli.

In primo luogo, si introdurranno le tecnologie e i particolari strumenti scelti per l'implementazione, discutendone il funzionamento e l'architettura.

In seguito, verrà progettato e definito l'ambiente virtuale IaC, affrontando problematiche come la configurazione di rete, la risoluzione DNS e la registrazione delle attività in un Database.

In ultimo, verranno condotti test mirati per la verifica di ciascuna delle funzionalità implementate.

2 | Stato dell'arte

2.1 Virtualizzazione

Una Virtual Machine (VM) (in italiano, *macchina virtuale*) rappresenta un'emulazione virtuale di un calcolatore, ottenuta tramite astrazione delle risorse hardware della macchina “*host*” (processo di virtualizzazione) [2]. Esistono diversi tipi di virtualizzazione:

- **Virtualizzazione nativa (o totale):** l'hardware viene completamente virtualizzato grazie all'uso di *hypervisor* (figura 2.1);
- **Virtualizzazione non nativa:** l'emulazione a livello software è basata su hardware diverso da quello fisicamente utilizzato, e necessita di codice aggiuntivo di basso livello per la compatibilità;
- **Paravirtualizzazione:** non viene simulato hardware in modo diretto, ma viene offerta una speciale Application Programming Interface (API) per l'interazione tra macchine virtuali e risorse fisiche;
- **Containerizzazione:** implementa la virtualizzazione a livello di *kernel* del sistema operativo, permettendo l'esecuzione contemporanea e isolata di istanze *user-space*, senza *overhead* dovuto all'emulazione hardware.

Ciascun tipo di virtualizzazione offre un diverso equilibrio tra flessibilità e prestazioni; in particolare, la virtualizzazione totale si è dimostrata specialmente adatta per:

- la condivisione di un sistema tra diversi utenti;
- l'isolamento dei programmi (e degli utenti) presenti sulle VM;
- l'emulazione di nuovo hardware per il raggiungimento di affidabilità, produttività e sicurezza.

2.1.1 Vagrant

Vagrant è un software *open source* per la gestione di VM, focalizzato sulla riproducibilità affidabile degli ambienti virtualizzati [3]. Nasce nel 2010 come progetto personale di Mitchell Hashimoto, ma cresce in popolarità dopo la fondazione di HashiCorp nel 2012.

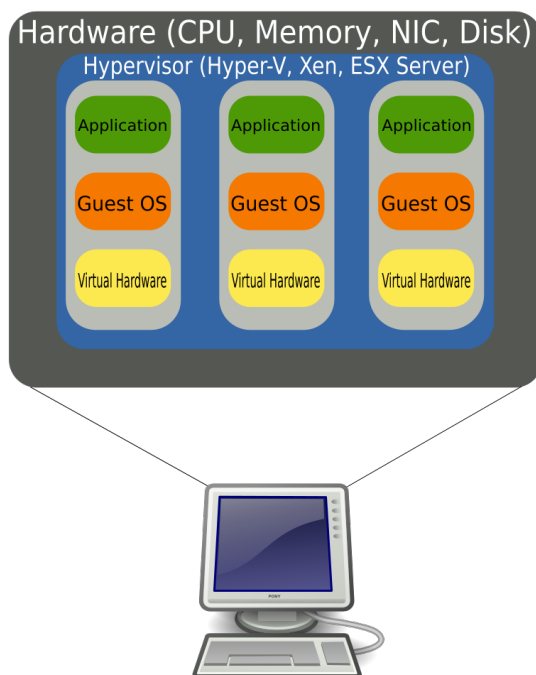


Figura 2.1: Virtualizzazione hardware totale [2]

Architettura

Vagrant dipende da software di virtualizzazione esterni (detti “*provider*” o “*hypervisor*”) come VirtualBox¹ e VMware², ma grazie al suo alto livello di astrazione, risulta portabile rispetto a essi.

Immagini di VM di base (denominate “*box*”) sono facilmente reperibili e condivisibili su depositi online: ciò rappresenta un punto di forza per Vagrant, che solleva gli sviluppatori dal laborioso processo di inizializzazione degli ambienti virtuali. L’uso della virtualizzazione hardware permette l’esecuzione di ambienti e sistemi operativi indipendenti dalla macchina *host*.

A tempo di esecuzione, Vagrant utilizza un file di configurazione (**Vagrantfile**) per ottenere i parametri definiti, quali immagine di base, specifiche hardware, cartelle condivise e indirizzi di rete locali.

2.2 Provisioning

Il *provisioning* rappresenta l’insieme delle azioni preparatorie necessarie a rendere un sistema o servizio pienamente operativo [4]. Nell’ambito dell’amministrazione di sistemi, esempi tipici sono installazione/aggiornamento di software, configurazione di rete e gestione degli utenti.

¹<https://www.virtualbox.org/>

²<https://www.vmware.com/>

2.2.1 Ansible

Ansible è un pacchetto di strumenti software *open source* per la configurazione e gestione automatizzata di ambienti informatici, tramite file di definizione (principio noto come Infrastructure as Code (IaC)) [5] [6]. Ansible nasce nel 2012 e viene acquisito dalla nota compagnia software RedHat nel 2015; è stato inizialmente sviluppato per sistemi Unix-like, ma ha esteso il proprio supporto anche a nodi Windows. Ansible è compatibile con computer fisici, macchine virtuali e interi ambienti *cloud*.

Architettura

Ansible si basa su interazione “*agentless*”, dove il controllore, mediante una semplice connessione Secure Shell (SSH), esegue le operazioni definite sui nodi interessati, senza quindi l’esigenza di particolare software sul sistema destinatario (escludendo il server SSH).

I nodi da gestire sono identificati tramite “*inventory*”, un file di testo che li elenca e li raggruppa. Per definire le operazioni da eseguire, Ansible utilizza i “*playbook*”, file descritti in linguaggio Yet Another Markup Language (YAML) di semplice interpretazione.

L’elenco dei moduli responsabili delle azioni sui nodi è costantemente aggiornato e consultabile nella documentazione ufficiale [7]. Essi sono sviluppati aspirando a un’alta astrazione dal sistema destinatario e all’idempotenza delle singole operazioni.

2.3 DNS & DHCP

Il protocollo Domain Name System (DNS) nasce in seguito all’esigenza di collegare indirizzi Internet Protocol (IP) a nomi di host, facilmente interpretabili e memorizzabili dagli esseri umani [8]. Per ciascuna zona di autorità è presente almeno un *name server*, dove vengono mantenute le corrispondenze tra *hostname* e indirizzi di competenza del dominio.

La risoluzione DNS è un procedimento trasparente alle applicazioni e agli utenti finali. La richiesta di traduzione del nome di host viene trasmessa al *resolver* locale, il quale ne controlla la presenza in memoria *cache*; in caso di non disponibilità del record, la richiesta viene inoltrata in autonomia ai *resolver* di livello superiore (figura 2.2).

Il protocollo Dynamic Host Configuration Protocol (DHCP) svolge invece la funzione di assegnamento automatico e dinamico di indirizzi IP, nell’ambito di una Local Area Network (LAN) [9].

DHCP lavora in architettura *client/server*: al momento della connessione alla rete locale, il cliente invia una richiesta in *broadcast*, e il server DHCP di competenza risponde assegnando un indirizzo IP locale all’interno di un intervallo configurabile.

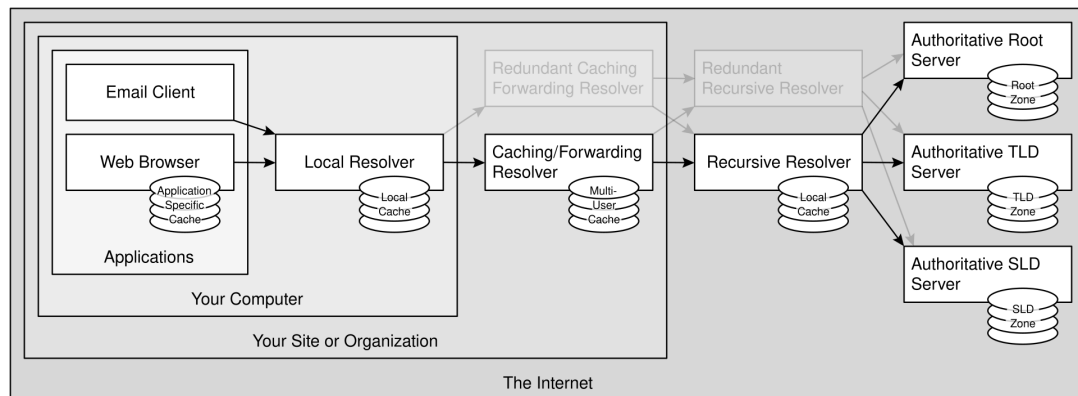


Figura 2.2: Sequenza di risoluzione DNS [8]

2.3.1 dnsmasq

dnsmasq è un software *open source* leggero e facilmente configurabile, in grado di fornire servizi DNS/DHCP su piccola scala [10].

In caso di richieste di risoluzione DNS, **dnsmasq** accede alla *cache* locale ed eventualmente inoltra l'operazione a un effettivo *resolver* ricorsivo. Per richieste DHCP, **dnsmasq** si affida invece alla configurazione locale per allocazioni statiche o dinamiche. Il servizio DHCP si integra in modo naturale con il servizio DNS, consentendo la registrazione automatica delle corrispondenze *hostname/address* in seguito al *leasing* degli indirizzi.

2.4 Base di dati

Un Database (DB) (in italiano, *base di dati*) è una collezione organizzata di dati accessibili tramite un Database Management System (DBMS), un software di interfacciamento per amministratori, applicazioni e utenti finali [11].

L'interazione con il DB è regolata da linguaggi specializzati:

- Data Control Language (DCL) per il controllo dell'accesso ai dati;
- Data Definition Language (DDL) per la definizione di dati e delle relazioni tra di essi;
- Data Manipulation Language (DML) per la manipolazione dei dati (inserimento, modifica, rimozione);
- Data Query Language (DQL) per la ricerca di informazioni estrapolate dai dati.

Structured Query Language (SQL) raggruppa le funzionalità di DDL, DML e DQL in un singolo linguaggio, rappresentando uno standard *de facto* nell'ambito delle basi di dati.

Le tipologie dominanti di DB sono:

- **Relazionale:** i dati sono strutturati in *tabelle* collegate univocamente tra loro tramite *relazioni* e *chiavi* identificative (figura 2.3); ciascuna tabella simboleggia un'entità: le colonne indicano gli *attributi*, mentre le righe rappresentano le sue *istanze*; grazie alle *transazioni*, i dati godono delle proprietà di Atomicity, Consistency, Isolation, Durability (ACID);
- **Not only SQL (NoSQL):** estensione dei DB relazionali, con supporto verso ulteriori strutture (chiave:valore, documenti, colonne, grafi, ...); è indicato per applicazioni a bassa latenza ed elevato *throughput*.

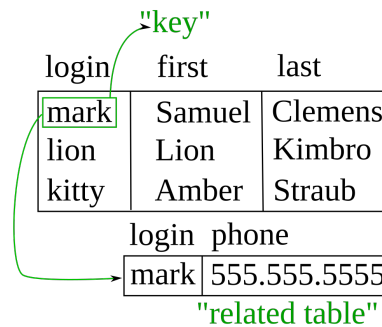


Figura 2.3: Modello relazionale [11]

2.4.1 MariaDB

MariaDB è un DBMS relazionale *open source*, basato sul diffuso MySQL, con il quale aspira a mantenere un'alta compatibilità nonostante l'introduzione di funzionalità avanzate, come nuovi *storage engine* [12].

2.5 Workload Management

Il concetto di *job scheduling* nasce con l'avvento dei primi sistemi informatici, in grado di eseguire una coda di schede perforate descriventi sequenze di istruzioni macchina [13]. I software moderni per la gestione di job sono tipicamente strutturati in architettura distribuita *master/agent*, con un singolo nodo di controllo e numerosi nodi di esecuzione.

2.5.1 SLURM

Simple Linux Utility for Resource Management (SLURM) è un software *open source* per la gestione e il coordinamento di nodi di calcolo, finalizzato alla pianificazione e l'esecuzione di programmi ("*jobs*") [14]. SLURM è altamente scalabile, rendendolo adatto anche in contesti High Performance Computing (HPC) a intenso utilizzo di risorse.

Vengono integrate tre funzioni principali:

- allocazione esclusiva (e non) delle risorse computazionali, per un periodo di tempo configurabile;

- inclusione di un *framework* per la sottomissione e il monitoraggio dell'esecuzione di programmi;
- gestione della contesa di risorse mediante code di esecuzione.

È possibile includere funzionalità aggiuntive, quali registrazione delle attività (“*accounting*”) e limitazione dell'uso di risorse, tramite il caricamento di *plugin* opzionali.

Architettura

Come illustrato in figura 2.4, un *cluster* SLURM è coordinato dal demone controllore `slurmctld`, il cui compito è gestire le richieste di esecuzione e allocare le risorse disponibili. Ciascun nodo di calcolo esegue `slurmd` per comunicare con il controllore ed effettuare i task richiesti. Il demone `slurmdbd`, connesso a un DBMS, diventa essenziale per la registrazione delle attività di molteplici cluster in un singolo DB; con esso vengono configurati gli *account* per il controllo degli utenti abilitati alla sottomissione di programmi. È inoltre possibile raggruppare i nodi di computazione in partizioni logiche assestanti e configurabili individualmente.

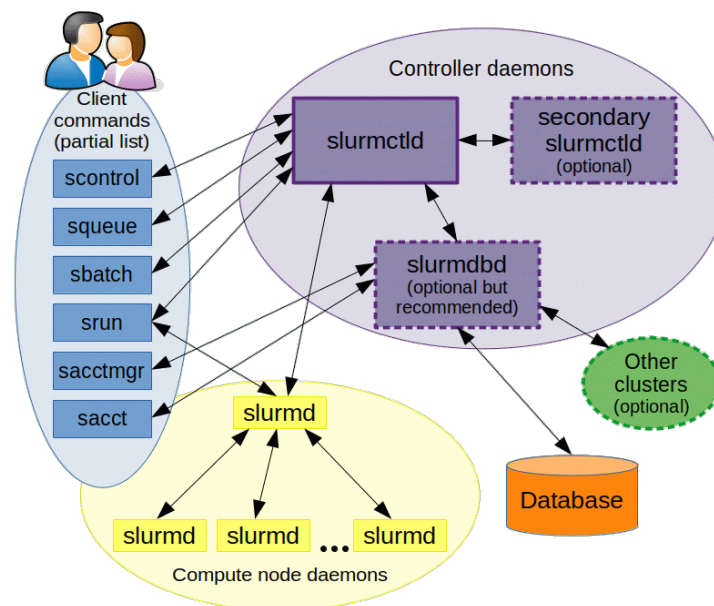


Figura 2.4: Componenti di un cluster SLURM [14]

Autenticazione

Le interazioni tra i vari demoni SLURM avvengono tramite Remote Procedure Call (RPC): è dunque necessario verificare l'autenticità della fonte del messaggio. Di default si utilizza MUNGE Uid 'N' Gid Emporium (MUNGE), un servizio per la creazione e validazione di credenziali utente, progettato appositamente per ambienti HPC [15].

La codifica/decodifica delle credenziali si basa su una chiave comune, condivisa ai soli nodi fidati. Un prerequisito fondamentale per il corretto funzionamento di MUNGE è la coerenza di utenti e gruppi in tutti i nodi interessati [16].

Condivisione di risorse

Di default, SLURM alloca in modo esclusivo interi nodi di computazione per l'esecuzione dei programmi: ciò significa che, anche in caso di basso utilizzo di risorse, nessun altro job verrà inoltrato allo stesso nodo. Per superare l'inefficienza di questo scenario, è possibile attivare la condivisione di risorse, quali Central Processing Unit (CPU) e Random Access Memory (RAM) [17]; configurazioni più granulari sono applicabili per le componenti della CPU, come *socket* e *core* (figura 2.5).

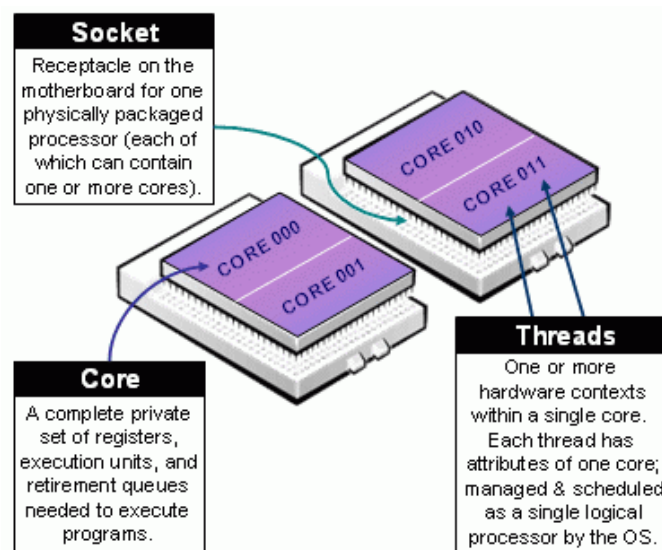


Figura 2.5: Definizioni di *socket*, *core* e *thread* [18]

Ulteriori tipologie di risorse, come le Graphics Processing Unit (GPU), vengono supportate sotto la classificazione di Generic Resource (GRES): per esse è necessaria la precisa definizione tramite file di configurazione. È inoltre possibile la loro condivisione tra i diversi job in coda [19].

Priorità di scheduling

SLURM offre numerosi metodi per impostare la priorità di *scheduling* tra i job; i principali sono:

- configurazione del **peso di molteplici fattori** (“*Multifactor Priority*”) per la selezione del job da eseguire, come ad esempio la sua dimensione, la quantità di risorse richieste e il tempo trascorso in coda [20];
- definizione di **Quality of Service (QoS)** per partizioni o singoli nodi [21], con la possibilità di limitare l'uso di risorse specifiche [22].

Federazione

SLURM permette la definizione di federazioni di *cluster*, dove la sottomissione dei job è inoltrata a tutti i componenti [23]. Ciascuno di essi tenta dunque di portare a termine la richiesta di esecuzione, in modo indipendente. Per il corretto coordinamento all'interno della federazione, è necessario che tutte le istanze di `slurmctld` (una per cluster) possano comunicare tra loro, oltre che con la singola istanza di `slurmdbd`, dedicata alla registrazione delle attività (figura 2.6).

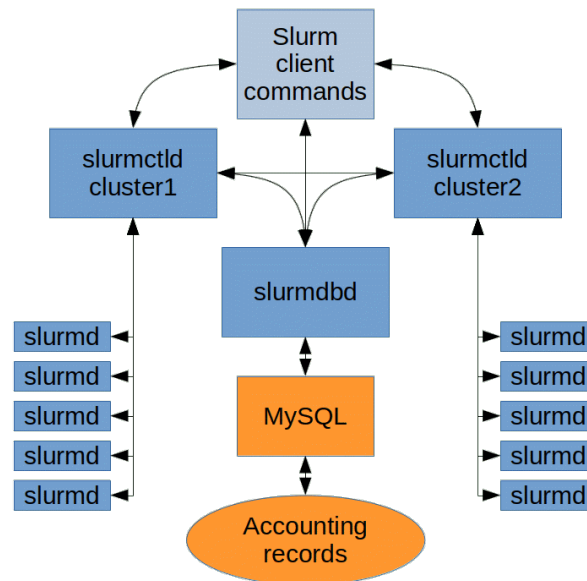


Figura 2.6: Comunicazione in una federazione SLURM [24]

3 | Analisi progettuale

Il progetto si pone di sviluppare, in stile Infrastructure as Code (IaC), un'infrastruttura di VM dedicata all'integrazione di particolari funzionalità SLURM. Di seguito sono analizzati i passi necessari.

Macchine virtuali

L'implementazione richiede anzitutto la scelta di un *hypervisor* (in grado di fornire virtualizzazione totale) e di una *box* Vagrant di riferimento; si è dunque optato per VirtualBox¹, nativamente supportato da Vagrant, e Debian Bookworm², un sistema operativo basato su Linux e noto per la sua stabilità ed efficienza.

Una funzionalità interessante di VirtualBox sono le *linked clone* [25], macchine virtuali basate sulle differenze con il disco della *box* di base, al fine di evitarne la copia totale (figura 3.1). In questo modo, l'*overhead* dovuto alla creazione di molteplici VM viene ridotto al minimo.

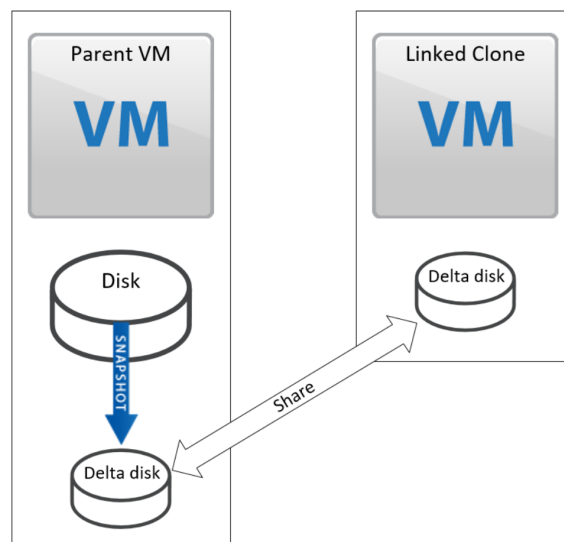


Figura 3.1: Relazione tra VM di base e *linked clone* [26]

Il *provisioning* delle macchine virtuali viene svolto da diversi *playbook* Ansible.

¹<https://www.virtualbox.org/>

²<https://www.debian.org/releases/bookworm/>

DHCP & DNS

SLURM basa le proprie comunicazioni su nomi di host, astraendo da indirizzi IP; è quindi necessaria la corretta risoluzione DNS all'interno della LAN. La scelta di `dnsmasq` offre inoltre funzionalità DHCP per più reti.

Accounting

Essendo prevista la registrazione delle attività (“*accounting*”) in ambienti multi-cluster SLURM, è necessaria la configurazione di un DB e relativo DBMS; è stato scelto MariaDB per familiarità con il contesto SQL. L'interazione tra il DBMS e il demone `slurmdbd`, responsabile dell'*accounting*, è configurabile seguendo la documentazione appropriata [27].

Topologia di rete

Una federazione SLURM, per definizione, è composta da più *cluster* in coordinamento tra loro; in contesti concreti, è verosimile che ciascuno di essi sia confinato nella propria LAN, dove la trasmissione verso l'esterno viene mediata da un *router* (o *gateway*). La figura 3.2 visualizza la topologia di rete per due *cluster* in federazione.

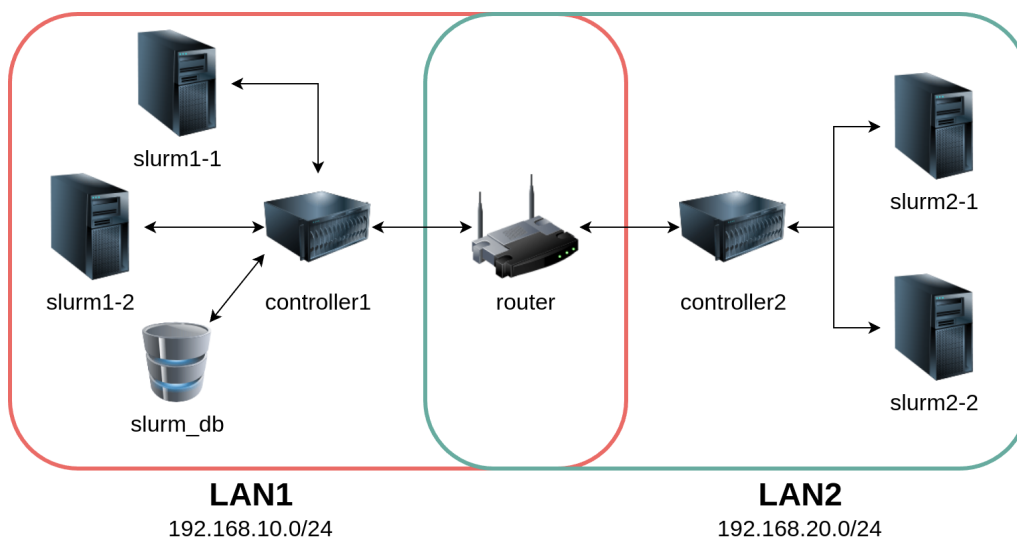


Figura 3.2: Topologia virtuale desiderata

In particolare, la configurazione della macchina virtuale `router` deve permettere la corretta comunicazione tra le due reti, inoltrando i pacchetti interessati.

Condivisione di risorse

La condivisione delle risorse disponibili su un nodo è configurabile grazie al *plugin* `select/cons_tres` [17]. Il progetto si pone, nello specifico, la condivisione di CPU, RAM e GPU; essendo quest'ultima catalogata come Generic Resource (GRES), è necessaria una configurazione aggiuntiva per il corretto interfacciamento [19].

Priorità di scheduling

Una volta attiva la condivisione di risorse, si può impostare la priorità di *scheduling*, ottenuta grazie alla definizione di una partizione SLURM separata e all'utilizzo del *plugin priority/multifactor* [20].

Sulla partizione viene applicata una Quality of Service (QoS) prioritaria [21], accessibile dall'utente desiderato e ristretta a una singola GPU del nodo di computazione [22].

Federazione di cluster

Ciascun *cluster* SLURM è descritto interamente da un proprio file di configurazione, che deve essere coerente su tutti i suoi nodi. Per la generazione di due installazioni, è dunque necessario definire due file separati, da inoltrare correttamente alle rispettive macchine virtuali.

Una prerogativa essenziale per l'impostazione di una federazione SLURM è la registrazione, sullo stesso DB, di tutti i cluster coinvolti, tramite l'uso di una singola istanza di `slurmdbd` [23].

4 | Implementazione

L'intero codice sviluppato per l'implementazione è consultabile pubblicamente su un *repository* Git¹.

4.1 Infrastruttura

4.1.1 Macchine virtuali

Il primo passo è la definizione delle macchine virtuali, seguendo il grafico e la nomenclatura presentati in figura 3.2; è necessario prestare attenzione ai seguenti aspetti:

- creazione delle VM come *linked clone* della *box* di base;
- corretto assegnamento delle reti interne VirtualBox (LAN1 e LAN2);
- delega del *provisioning* a un *playbook* Ansible.

Il file di definizione Vagrant risulta così composto:

```
Vagrantfile
1 Vagrant.configure("2") do |config|
2   #
3   # VM BOX CONFIGURATION
4   config.vm.box = "debian/bookworm64"
5   config.vm.provider "virtualbox" do |vb|
6     vb.linked_clone = true
7     vb.memory = "1024"
8     vb.cpus = "4"
9   end
10  #
11  # VIRTUAL MACHINES
12  config.vm.define "router" do |machine|
13    machine.vm.hostname = "router"
14    machine.vm.network "private_network", virtualbox__intnet: "LAN1",
15      ↪ auto_config: false
16    machine.vm.network "private_network", virtualbox__intnet: "LAN2",
17      ↪ auto_config: false
18  end
19 end
```

¹<https://github.com/maxzrbn/Tirocinio-SLURM>

```

16 end
17 (1..2).each do |i|
18   config.vm.define "controller#{i}" do |machine|
19     machine.vm.hostname = "controller#{i}"
20     machine.vm.network "private_network", virtualbox__intnet: "LAN#{i}",
21       ↪ auto_config: false
22   end
23   (1..2).each do |j|
24     config.vm.define "slurm#{i}-#{j}" do |machine|
25       machine.vm.hostname = "slurm#{i}-#{j}"
26       machine.vm.network "private_network", virtualbox__intnet: "LAN#{i}",
27         ↪ auto_config: false
28     end
29   end
30 end
31 #
32 # PROVISIONING
33 config.vm.provision "ansible" do |ansible|
34   ansible.playbook = "site.yml"
35 end

```

Ciascun nodo è quindi virtualmente dotato di 4 CPU e 1 Gigabyte (GB) di RAM. Da notare l'uso di iterazioni per la definizione compatta di molteplici VM.

Il *provisioning* Ansible viene definito in `site.yml`, attribuendo “ruoli” alle varie macchine virtuali; ciascuno di essi esegue una particolare sequenza di istruzioni.

```

----- site.yml -----
1 ---
2 - hosts: router
3   become: true
4   roles:
5     - common
6     - router
7
8 - hosts: controller1
9   become: true
10  roles:
11    - common
12    - mariadb
13    - slurmcommon
14    - cluster1
15    - slurmctlr
16
17 - hosts: controller2
18   become: true

```



```
19  roles:
20      - common
21      - slurmcommon
22      - cluster2
23      - slurmctlr
24
25  - hosts: slurm1-1:slurm1-2
26      become: true
27      roles:
28          - common
29          - slurmcommon
30          - cluster1
31          - slurmwkr
32
33  - hosts: slurm2-1:slurm2-2
34      become: true
35      roles:
36          - common
37          - slurmcommon
38          - cluster2
39          - slurmwkr
```

4.1.2 Configurazione di rete

Svolgendo la funzione di DHCP/DNS server, la macchina virtuale **router** dispone di due interfacce di rete con indirizzi IP locali statici; è possibile configurarle grazie al seguente file, opportunamente inoltrato alla VM:

```
roles/router/files/eth_cfg
1  auto eth1
2  iface eth1 inet static
3      netmask 255.255.255.0
4      address 192.168.10.254
5
6  auto eth2
7  iface eth2 inet static
8      netmask 255.255.255.0
9      address 192.168.20.254
```

Affinché si possano effettivamente instradare i pacchetti tra le due LAN, è necessario attivare, su **router**, il *forwarding* a livello di *kernel*, modificando un particolare file di configurazione:

```
roles/router/tasks/main.yml
1  ...
2  - name: Allow IPv4 forwarding
3    ansible.builtin.lineinfile:
4      path: '/etc/sysctl.conf'
5      line: 'net.ipv4.ip_forward=1'
6      state: present
7      notify: Reload kernel config
8  ...
```

Le interfacce delle restanti VM vengono automaticamente configurate grazie al protocollo DHCP.

4.1.3 DHCP & DNS

Il servizio DHCP/DNS è svolto, per entrambe le reti, da `dnsmasq`, in esecuzione su `router`. La sua configurazione è la seguente:

```
roles/router/files/dnsmasq.conf
1  interface=eth1
2  interface=eth2
3  no-hosts
4  no-resolv
5
6  # NAMESERVER USED FOR THE INTERNET (the host running the VMs)
7  server=10.0.2.3
8
9  # ADDRESS RANGES FOR LEASING
10 dhcp-range=interface:eth1,192.168.10.1,192.168.10.253,12h
11 dhcp-range=interface:eth2,192.168.20.1,192.168.20.253,12h
12
13 # SPECIFY SELF AS NAMESERVER (for both LANs)
14 dhcp-option-force=interface:eth1,option:dns-server,192.168.10.254
15 dhcp-option-force=interface:eth2,option:dns-server,192.168.20.254
16
17 # AVOID ROUTING ALL TRAFFIC THROUGH THIS VM
18 dhcp-option=3
19
20 # STATIC ROUTE FOR CROSS-LAN COMMUNICATION
21 dhcp-option=121,192.168.20.0/24,192.168.10.254,192.168.10.0/24,192.168.20.254
```

L'opzione evidenziata specifica il *gateway* tra le reti, ossia l'instradamento dei pacchetti provenienti da una LAN e destinati all'altra.

4.1.4 Accounting

La registrazione delle attività (“*accounting*”) viene svolta dal demone `slurmdbd`, in esecuzione su `controller1`.

Database

Per integrare l’inserimento di password nel *provisioning*, si definiscono apposite variabili di esecuzione, eventualmente crittografabili con l’ausilio di Ansible Vault [7]:

```
----- secrets.yml -----
1 mariadb:
2   root_password: THIS_IS_THE_ROOT_DB_PASSWORD
3   slurm_password: THIS_IS_THE_SLURM_DB_PASSWORD
```

Al fine di isolare l’installazione del DB, si può utilizzare il software di containerizzazione `podman`², in grado di recuperare ed eseguire l’immagine di MariaDB opportunamente configurata:

```
----- roles/mariadb/tasks/main.yml -----
1 ...
2 - name: Pull a MariaDB container
3   become: true
4   become_user: 'mariadb'
5   containers.podman.podman_image:
6     name: 'mariadb:latest'
7
8 - name: Run MariaDB container
9   become: true
10  become_user: 'mariadb'
11  containers.podman.podman_container:
12    name: 'mariadb'
13    image: 'mariadb:latest'
14    state: started
15    ports: "127.0.0.1:3306:3306"
16    env:
17      MARIADB_ROOT_PASSWORD: "{{ mariadb.root_password }}"
18  ...
```

In questo modo MariaDB risulta in ascolto sulla porta Transmission Control Protocol (TCP) 3306 di `controller1`.

Seguendo poi la documentazione SLURM per il corretto interfacciamento tra il DBMS e `slurmdbd` [27], si generano DB e utente dedicati:

²<https://podman.io/>

```

roles/mariadb/tasks/main.yml
1  ...
2  - name: Ensure 'slurm_db' database is present
3    community.mysql.mysql_db:
4      name: slurm_db
5      state: present
6      login_host: 127.0.0.1
7      login_port: 3306
8      login_user: root
9      login_password: "{{ mariadb.root_password }}"
10   retries: 6
11   delay: 10
12
13  - name: Ensure 'slurm' MySQL user is present with password
14    community.mysql.mysql_user:
15      name: slurm
16      password: "{{ mariadb.slurm_password }}"
17      host: '%'
18      login_host: 127.0.0.1
19      login_user: root
20      login_password: "{{ mariadb.root_password }}"
21      state: present
22  ...

```

Si impostano poi i privilegi appropriati, assicurandosi che vengano salvati e applicati (operazione di “*flushing*”):

```

roles/mariadb/tasks/main.yml
1  ...
2  - name: Grant usage privilege to the user
3    community.mysql.mysql_user:
4      name: slurm
5      priv: '*.*:USAGE'
6      append_privs: yes
7      host: '%'
8      login_host: 127.0.0.1
9      login_user: root
10     login_password: "{{ mariadb.root_password }}"
11     state: present
12
13  - name: Grant all privileges on the database to the user
14    community.mysql.mysql_user:
15      name: slurm
16      priv: 'slurm_db.*:ALL'
17      append_privs: yes
18      host: '%'
19      login_host: 127.0.0.1

```

```

20     login_user: root
21     login_password: "{{ mariadb.root_password }}"
22     state: present
23
24 - name: Flush privileges (update & reload grant tables)
25   community.mysql.mysql_user:
26     login_user: root
27     login_password: "{{ mariadb.root_password }}"
28     check_implicit_admin: yes
29     append_privs: yes
30     login_host: 127.0.0.1
31     state: present
32     name: slurm
33     host: '%'
34     ...

```

A questo punto è possibile verificare manualmente l'accesso al DB, presentandosi come utente `slurm` (utilizzato poi da `slurmdbd`):

```

vagrant@controller1:~$ mariadb -u slurm --port=3306
↪ --password=THIS_IS_THE_SLURM_DB_PASSWORD slurm_db
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MariaDB connection id is 12
Server version: 11.5.2-MariaDB-ubu2404 mariadb.org binary distribution

Copyright (c) 2000, 2018, Oracle, MariaDB Corporation Ab and others.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

MariaDB [slurm_db]>_

```

slurmdbd

Si prosegue con l'installazione del demone `slurmdbd`:

```

_____ roles/mariadb/tasks/main.yml _____
1     ...
2 - name: Install 'slurmdb' package
3   ansible.builtin.apt:
4     name: 'slurmdbd'
5     update_cache: true
6     ...

```

Il suo file di configurazione è così composto:

```
----- slurmdbd.conf -----
1 AuthType=auth/munge
2 DbdHost=controller1
3 DbdPort=6819
4 DebugLevel=verbose
5 StorageHost=localhost
6 StorageLoc=slurm_db
7 StoragePass="{ { mariadb.slurm_password } }"
8 StoragePort=3306
9 StorageType=accounting_storage/mysql
10 StorageUser=slurm
11 LogFile=/var/log/slurm/slurmdbd.log
12 PidFile=/run/slurmdbd.pid
13 SlurmUser=slurm
-----
```

Da notare le opzioni:

- `AuthType=auth/munge` : attivazione del *plugin* MUNGE per l'autenticazione delle trasmissioni RPC;
- `DbdPort=6819` : porta TCP di ascolto per le comunicazioni con gli altri componenti SLURM.

Nel tentativo di avviare il demone si sono però presentati i seguenti errori:

```
vagrant@controller1:~$ sudo slurmdbd -D -vvv
slurmdbd: debug2: accounting_storage/as_mysql: init: mysql_connect() called for
↳ db slurm_db
slurmdbd: debug2: Attempting to connect to localhost:3306
slurmdbd: error: mysql_real_connect failed: 2002 Can't connect to local server
↳ through socket '/run/mysqld/mysqld.sock' (2)
slurmdbd: error: The database must be up when starting the MYSQL plugin.
↳ Trying again in 5 seconds.
```

Indagando sul messaggio evidenziato si è appreso che `slurmdbd`, specificata l'opzione `StorageHost=localhost`, tenta la connessione al DB tramite *socket* locale di default (`/run/mysqld/mysqld.sock`) [27]; eseguendo però il DBMS all'interno di un *container*, la posizione della *socket* diventa relativa al *filesystem* astratto:

```
vagrant@controller1:~$ sudo find / -type s | grep 'mysqld.sock'
/home/mariadb/.local/share/containers/storage/overlay/e7d63d4dddde910afaf1d257
↳ 4c3345fc74c42372bffc95f84528ead29c2f52d7e/diff/run/mysqld/mysqld.sock
```

È dunque necessario aggiornare il file di configurazione con `StorageHost=127.0.0.1`, in modo da forzare il demone all'utilizzo di porte TCP/IP.

4.1.5 Cluster SLURM

Il file di configurazione `slurm.conf` viene utilizzato da entrambi `slurmctld` e `slurmd`; in esso è descritto il generico comportamento del *cluster*:

```
roles/cluster1/templates/slurm.conf
1  # CLUSTER 1 CONFIG
2  #
3  ClusterName=cluster1
4  SlurmctldHost=controller1
5  MpiDefault=none
6  ProctrackType=proctrack/cgroup
7  ReturnToService=1
8  SlurmctldPidFile=/var/run/slurmctld.pid
9  SlurmctldPort=6817
10 SlurmdPidFile=/var/run/slurmd.pid
11 SlurmdPort=6818
12 SlurmdSpoolDir=/var/slurm/slurmd
13 SlurmUser=slurm
14 SlurmdUser=root
15 StateSaveLocation=/var/slurm/slurmctld
16 SwitchType=switch/none
17 TaskPlugin=task/affinity
18 #
19 # SCHEDULING
20 SchedulerType=sched/backfill
21 SelectType=select/linear
22 #
23 # LOGGING AND ACCOUNTING
24 AccountingStorageType=accounting_storage/slurmdbd
25 AccountingStorageUser=slurm
26 JobAcctGatherType=jobacct_gather/cgroup
27 SlurmctldLogFile=/var/log/slurm/slurmctld.log
28 SlurmdLogFile=/var/log/slurm/slurmd.log
29 #
30 # COMPUTE NODES
31 NodeName=slurm1-[1-2] CPUs=4 State=UNKNOWN
32 PartitionName=debug Nodes=ALL Default=YES MaxTime=INFINITE State=UP
```

Con l'opzione evidenziata, `slurmctld` attiva e delega la registrazione delle attività a `slurmdbd`, in esecuzione sul medesimo nodo (`controller1`).

La configurazione del secondo *cluster* presenta invece alcune differenze:

```
roles/cluster2/templates/slurm.conf
1  # CLUSTER 2 CONFIG
2  #
3  ClusterName=cluster2
```

```

4 SlurmctldHost=controller2
5 ...
6 AccountingStorageHost=controller1
7 ...
8 NodeName=slurm2-[1-2] ...
9 ...

```

L'opzione in evidenza, specificata solo nel *cluster* senza DB, indica il nodo da contattare per l'*accounting* delle attività, dove viene eseguito `slurmdbd`.

4.2 Condivisione di risorse

Per la condivisione delle risorse è necessario modificare opportunamente i file di configurazione; di seguito vengono indicate le variazioni relative al primo *cluster*, ma le differenze con il secondo interessano esclusivamente la nomenclatura dei nodi:

```

roles/cluster1/templates/slurm.conf
1 ...
2 SelectType=select/cons_tres
3 SelectTypeParameters=CR_CPU_Memory
4 ...
5 AccountingStorageTRES=gres/gpu
6 ...
7 GresTypes=gpu
8 NodeName=slurm1-[1-2] CPUs=4 RealMemory=1024 Gres=gpu:rtx2080ti:2 State=UNKNOWN
9 ...

```

In particolare, le opzioni indicano [17]:

- `SelectType=select/cons_tres`: *plugin* per la gestione delle risorse (“*consumable trackable resources*”), permettendo l'esecuzione di più job sullo stesso nodo;
- `SelectTypeParameters=CR_CPU_Memory`: risorse ordinarie da condividere (nel caso in questione, CPU e RAM);
- `AccountingStorageTRES=gres/gpu`: attivazione dell'*accounting* relativo alle GPU;
- `GresTypes=gpu`: tipologie di Generic Resource (GRES) da condividere;
- `RealMemory=1024`: memoria totale in Megabyte (MB) allocabile per nodo di computazione (opzione obbligatoria per la condivisione della memoria; deve essere minore o uguale alla RAM riconosciuta dal sistema operativo);
- `Gres=gpu:rtx2080ti:2`: specifiche GRES presenti nel nodo.

È inoltre fondamentale introdurre nei nodi di computazione due ulteriori file di configurazione, per il corretto interfacciamento con le GRES [19]:


```

roles/slurmwrk/templates/cgroup.conf
1 CgroupPlugin=autodetect
2 CgroupAutomount=yes
3 ConstrainCores=no
4 ConstrainRAMSpace=no
5 ConstrainDevices=yes

```

In `gres.conf` vanno specificati:

- il nome del nodo di computazione;
- le precise caratteristiche delle risorse presenti.

```

gres.conf
1 NodeName=slurm1-1 Name=gpu Type=rtx2080ti Count=2 File=/dev/nvidia[0-1]

```

Una volta configurate, è possibile verificare la disponibilità delle risorse tramite il comando `sinfo`:

```

vagrant@controller1:~$ sinfo -o "%15N %10c %10m %20G"
NODELIST          CPUS      MEMORY      GRES
slurm1-[1-2]       4         1024      gpu:rtx2080ti:2

```

4.3 Priorità di scheduling

L'applicazione di una priorità di *scheduling* necessita anzitutto di un contesto verosimile di utilizzo; a tal proposito vengono definiti utenti/gruppi di prova, tra cui quelli interessati dalla priorità (in questo caso, `user1`).

Affinché venga mantenuta la coerenza di utenti/gruppi all'interno del *cluster*, la loro definizione avviene con task comuni a tutti i nodi:

```

roles/slurmcommon/tasks/main.yml
1 ...
2 - name: Ensure test groups exist
3   ansible.builtin.group:
4     name: "{{ item.name }}"
5     gid: "{{ item.id }}"
6     state: present
7   loop:
8     - { name: 'user1', id: '2001' }
9     - { name: 'user2', id: '2002' }
10    - { name: 'user3', id: '2003' }
11    - { name: 'user4', id: '2004' }
12
13 - name: Ensure test users exist

```

```

14  ansible.builtin.user:
15      name: "{{ item.name }}"
16      uid: "{{ item.id }}"
17      group: "{{ item.name }}"
18      state: present
19  loop:
20      - { name: 'user1', id: '2001' }
21      - { name: 'user2', id: '2002' }
22      - { name: 'user3', id: '2003' }
23      - { name: 'user4', id: '2004' }
24  ...

```

Il file di configurazione SLURM è stato modificato nel seguente modo:

```

_____ roles/cluster1/templates/slurm.conf _____
1  ...
2  PriorityType=priority/multifactor
3  PriorityWeightQOS=10000
4  PriorityWeightTRES=GRES/gpu=100
5  ...
6  AccountingStorageEnforce=limits
7  JobAcctGatherType=jobacct_gather/cgroup
8  ...
9  PartitionName=gupart Nodes=slurm1-[1-2] QOS=guprio AllowAccounts=user1acct
   ↪ State=UP

```

In particolare, le opzioni specificano:

- `PriorityType=priority/multifactor`: *plugin* per la priorità a molteplici fattori [20]; per il caso in questione si sono privilegiate, tramite le opzioni `PriorityWeight-`:
 - l'appartenenza a Quality of Service (QoS) prioritarie;
 - la richiesta di GPU per l'esecuzione;
- `AccountingStorageEnforce=limits`: imposizione dei limiti sulle risorse specificati nel DB; implica la presenza di un'associazione (utente, account, partizione, ...) valida per tutte le richieste di esecuzione [22];
- `PartitionName=gupart`: partizione aggiuntiva ristretta all'account `user1acct`, con QoS `guprio` (definiti in seguito).

4.3.1 Impostazione del DB

È fondamentale la registrazione nel DB (tramite il comando `sacctmgr`) delle informazioni per il comportamento desiderato, ovvero:

- la QoS gpuprio [21];
- gli account user1acct (prioritario) e noprio [27].

```
vagrant@controller1:~$ sudo sacctmgr add qos gpuprio Priority=10
↪ GrpTRES=gres/gpu=1
Adding QOS(s)
  gpuprio
Settings
  Description      = gpuprio
  GrpTRES          = gres/gpu=1
  Priority         = 10
Would you like to commit changes? (You have 30 seconds to decide)
(N/y): y
vagrant@controller1:~$ sudo sacctmgr add account user1acct QOS=gpuprio
Adding Account(s)
  user1acct
Settings
  Description      = Account Name
  Organization     = Parent/Account Name
Associations =
  C = cluster1    A = user1acct
Settings
  QOS              = gpuprio
Would you like to commit changes? (You have 30 seconds to decide)
(N/y): y
vagrant@controller1:~$ sudo sacctmgr add account noprio
Adding Account(s)
  noprio
Settings
  Description      = Account Name
  Organization     = Parent/Account Name
Associations =
  C = cluster1    A = noprio
Would you like to commit changes? (You have 30 seconds to decide)
(N/y): y
```

L'account **noprio** risulta essenziale affinché gli utenti non prioritari possano richiedere l'esecuzione di job. Da notare il confinamento della priorità a una singola GPU, tramite l'attributo **GrpTRES** [22].

Gli utenti, durante la registrazione nel DB, vengono collegati a un account di default, pertinente al comportamento desiderato:

```
vagrant@controller1:~$ sudo sacctmgr add user user1 DefaultAccount=user1acct
Adding User(s)
  user1
```

```

Settings
  Default Account = user1acct
Associations =
  C = cluster1    A = user1acct          U = user1
Would you like to commit changes? (You have 30 seconds to decide)
(N/y): y
vagrant@controller1:~$ sudo sacctmgr add user user2,user3,user4
↪  DefaultAccount=noprio
Adding User(s)
  user2
  user3
  user4
Settings
  Default Account = noprio
Associations =
  C = cluster1    A = noprio             U = user2
  C = cluster1    A = noprio             U = user3
  C = cluster1    A = noprio             U = user4
Would you like to commit changes? (You have 30 seconds to decide)
(N/y): y

```

Il comando `sacctmgr` permette infine di verificare le informazioni registrate:

```

vagrant@controller1:~$ sacctmgr show qos format=name,priority,GrpTRES
      Name      Priority      GrpTRES
-----
normal          0
gpuprio         10    gres/gpu=1
vagrant@controller1:~$ sacctmgr show association user=user1,user2,user3,user4
↪  format=user,account,qos
      User      Account      QOS
-----
user2      noprio      normal
user3      noprio      normal
user4      noprio      normal
user1  user1acct      gpuprio

```

4.4 Federazione di cluster

L'implementazione di una federazione SLURM si basa sulla corretta impostazione dei singoli *cluster*, svolta in sezione 4.1.5.

La seguente opzione, non essenziale per il funzionamento in sé, imposta la visualizzazione federata per i comandi `sacctmgr` e `scontrol` [23]:

```

1  ...
2  FederationParameters=fed_display
3  ...

```

4.4.1 Impostazione del DB

La registrazione di entrambi i *cluster* sullo stesso DB è un requisito fondamentale per l'impostazione della federazione [23]:

```

vagrant@controller1:~$ sudo sacctmgr add federation testfederation
↪ clusters=cluster1,cluster2
Adding Federation(s)
  testfederation
Settings
  Cluster      = cluster1
  Cluster      = cluster2
Would you like to commit changes? (You have 30 seconds to decide)
(N/y): y

```

Utilizzando `sacctmgr` e `scontrol`, si è dunque confermata la corretta comunicazione tra di essi:

```

vagrant@controller1:~$ sacctmgr show federation
Federation    Cluster ID      Features      FedState
-----
testfeder+    cluster1  1              ACTIVE
testfeder+    cluster2  2              ACTIVE
vagrant@controller1:~$ scontrol show federation
Federation: testfederation
Self:        cluster1:127.0.0.1:6817 ID:1 FedState:ACTIVE Features:
Sibling:      cluster2:192.168.20.194:6817 ID:2 FedState:ACTIVE Features:
↪ PersistConnSend/Recv:Yes/Yes Synced:Yes

```


5 | Risultati

L'accertamento della funzionalità dei servizi SLURM proposti richiede, in primo luogo, l'avvio dell'intera infrastruttura. Tramite il comando “**vagrant up**”, la creazione e il *provisioning* delle macchine virtuali vengono automatizzati in stile Infrastructure as Code (IaC).

Il comando **sinfo** convalida la disponibilità dei *cluster* aderenti alla federazione:

```
vagrant@controller1:~$ sinfo
PARTITION CLUSTER AVAIL TIMELIMIT NODES STATE NODELIST
debug*    cluster1  up    infinite     1  idle slurm1-[1-2]
debug*    cluster2  up    infinite     1  idle slurm2-[1-2]
```

5.1 Test della condivisione di risorse

Per verificare la condivisione delle risorse presenti su un nodo, vengono allocati quattro job da 1 CPU (-c1) e 256 MB di RAM (--mem=256), specificandone lo svolgimento su un singolo nodo (-N1) da 4 CPU e 1 GB di RAM:

```
vagrant@controller1:~$ srun -N1 -c1 --mem=256 sleep 60 &
[1] 18254
vagrant@controller1:~$ srun -N1 -c1 --mem=256 sleep 60 &
[2] 18264
vagrant@controller1:~$ srun -N1 -c1 --mem=256 sleep 60 &
[3] 18274
vagrant@controller1:~$ srun -N1 -c1 --mem=256 sleep 60 &
[4] 18284
```

Grazie al comando **squeue** è possibile confermare l'esecuzione contemporanea dei job (tutti in stato *running*), a dimostrazione della condivisione di processori e memoria:

```
vagrant@controller1:~$ squeue
JOBID PARTITION NAME USER ST TIME NODES NODELIST(REASON)
1 debug sleep vagrant R 0:15 1 slurm1-1
```

2	debug	sleep	vagrant	R	0:14	1 slurm1-1
3	debug	sleep	vagrant	R	0:13	1 slurm1-1
4	debug	sleep	vagrant	R	0:12	1 slurm1-1

Allo stesso modo, vengono allocati due job richiedenti 1 GPU ciascuno (`--gpus=1`) verificando la loro esecuzione parallela e quindi la condivisione delle GPU:

```
vagrant@controller1:~$ srun -N1 --gpus=1 --mem=256 sleep 60 &
[1] 18294
vagrant@controller1:~$ srun -N1 --gpus=1 --mem=256 sleep 60 &
[2] 18304
vagrant@controller1:~$ squeue
JOBID PARTITION    NAME      USER  ST        TIME  NODES NODELIST(REASON)
    5      debug    sleep    vagrant  R         0:08      1 slurm1-1
    6      debug    sleep    vagrant  R         0:07      1 slurm1-1
```

5.2 Test della priorità di scheduling

Per confermare la priorità di *scheduling* su una GPU, da parte di un determinato utente, vengono allocati tre job non prioritari e un job lanciato da `user1` (sulla partizione prioritaria `gpupart`):

```
vagrant@controller1:~$ sudo -u user2 srun -N1 --gpus=1 --mem=256 sleep 60 &
[1] 18314
vagrant@controller1:~$ sudo -u user3 srun -N1 --gpus=1 --mem=256 sleep 60 &
[2] 18324
vagrant@controller1:~$ sudo -u user4 srun -N1 --gpus=1 --mem=256 sleep 60 &
[3] 18334
vagrant@controller1:~$ sudo -u user1 srun -N1 --partition=gpupart --gpus=1
↪ --mem=256 sleep 60 &
[4] 18344
vagrant@controller1:~$ squeue
JOBID PARTITION    NAME      USER  ST        TIME  NODES NODELIST(REASON)
    9      debug    sleep    user4  PD         0:00      1 (Resources)
    8      debug    sleep    user3  R         0:03      1 slurm1-1
    7      debug    sleep    user2  R         0:17      1 slurm1-1
   10  gpupart    sleep    user1  PD         0:00      1 (QOSGrpGRES)
```

Il job 10 risulta quindi in concorrenza di risorse con il job 9 (entrambi in stato *pending*). Attendendo la liberazione di una GPU (ossia il termine del job 7), la partizione prioritaria ha la precedenza di *scheduling*:


```
vagrant@controller1:~$ squeue
```

JOBID	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST(REASON)
9	debug	sleep	user4	PD	0:00	1	(Resources)
8	debug	sleep	user3	R	0:50	1	slurm1-1
10	gpupart	sleep	user1	R	0:04	1	slurm1-1

5.3 Test della federazione di cluster

Grazie all'opzione `--clusters` in fase di sottomissione, è possibile specificare una lista di *cluster* sulla quale tentare di eseguire il job richiesto (ad esempio, il programma `hostname`):

```
vagrant@controller1:~$ srun --clusters=cluster2 -N2 hostname
slurm2-1
slurm2-2

vagrant@controller2:~$ srun --clusters=cluster1 -N2 hostname
slurm1-1
slurm1-2
```

È importante notare come, per entrambe le sottomissioni di esempio:

- il controllore di partenza non appartenga al *cluster* specificato per l'esecuzione;
- venga richiesta l'esecuzione del job su due nodi distinti (opzione `-N2`).

Il risultato di `hostname` conferma l'inoltro dei job tra i due *cluster* e quindi il corretto comportamento della federazione.

6 | Conclusioni

Il progetto di tesi si è concentrato sullo sviluppo di un'infrastruttura virtuale, al fine di implementare particolari servizi di calcolo SLURM. Si sono scelti strumenti di gestione e configurazione automatica di macchine virtuali (nello specifico, Vagrant & Ansible), ottenendo una struttura Infrastructure as Code (IaC). Grazie all'uso di Git¹, un software *open source* per il controllo delle versioni, il processo di sviluppo è stato supportato da:

- cronologia delle modifiche;
- progettazione non lineare delle funzionalità.

In primo luogo, si sono dovute integrare le caratteristiche di base per l'operatività dell'infrastruttura, come assegnamento DHCP e risoluzione DNS (per entrambe le reti virtuali); a tal fine si è scelto lo strumento **dnsmasq**, di semplice impostazione.

Seguendo la documentazione appropriata, si sono poi configurati i vari componenti di un'installazione SLURM distribuita, prestando particolare attenzione alla configurazione di **slurmdbd** e DBMS associato.

Si è dunque passati all'implementazione dei particolari servizi posti come obiettivo del progetto.

La condivisione delle risorse di computazione permette il loro assegnamento dinamico in base alle necessità dei job da eseguire; si è definita la condivisione di:

- Central Processing Unit (CPU);
- Random Access Memory (RAM);
- Graphics Processing Unit (GPU).

In particolare, il corretto interfacciamento con la GPU, catalogata da SLURM come Generic Resource (GRES), ha richiesto ulteriore configurazione.

Per implementare la priorità di *scheduling*, riservata a un determinato utente e ristretta a una specifica risorsa, si è aggiunta una partizione SLURM dedicata; essa risulta:

- prioritaria su una singola GPU, grazie alla definizione di Quality of Service (QoS);

¹<https://git-scm.com/>

- accessibile esclusivamente dall'utente interessato, in accordo con l'*accounting* impostato nel DB.

L'attivazione della federazione di *cluster* SLURM riguarda essenzialmente la loro registrazione sul DB in comune; una volta inseriti, le richieste di esecuzione vengono inoltrate tra di essi, a seconda della necessità.

Si è infine svolta la fase di verifica delle funzionalità implementate, tramite la sottomissione mirata di programmi e l'accertamento della loro esecuzione.

6.1 Sviluppi futuri

Il modello di sviluppo IaC include numerosi vantaggi a beneficio dell'integrazione di un servizio o un'infrastruttura; tra i più rilevanti si ricordano:

- rapidità ed efficienza di sviluppo;
- riproduzione coerente degli ambienti;
- riduzione degli errori di configurazione dovuti alla gestione manuale.

Sono dunque possibili diversi sviluppi per l'infrastruttura progettata in questo elaborato, ad esempio:

- l'ampliamento delle funzionalità di calcolo disponibili, come *preemption* (sospensione di un job a bassa priorità per l'esecuzione di uno a priorità maggiore [28]) e *gang scheduling* (sospensione alternata dell'esecuzione di più job su una stessa risorsa [29]);
- l'implementazione di prassi *DevOps*, come test automatizzati, provocati dalla modifica dell'infrastruttura (principio noto come Continuous Integration & Continuous Delivery (CI/CD));
- l'applicazione in contesti reali, per la concreta erogazione di calcolo ad alte prestazioni.

Bibliografia

- [1] *High Performance Computing (Wikipedia)*. URL: https://it.wikipedia.org/wiki/High_performance_computing (visitato il giorno 20/02/2025).
- [2] *Virtualization (Wikipedia)*. URL: <https://en.wikipedia.org/wiki/Virtualization> (visitato il giorno 21/02/2025).
- [3] *Vagrant (Wikipedia)*. URL: [https://it.wikipedia.org/wiki/Vagrant_\(software\)](https://it.wikipedia.org/wiki/Vagrant_(software)) (visitato il giorno 21/02/2025).
- [4] *Provisioning (Wikipedia)*. URL: [https://en.wikipedia.org/wiki/Provisioning_\(technology\)](https://en.wikipedia.org/wiki/Provisioning_(technology)) (visitato il giorno 22/02/2025).
- [5] *Ansible (Wikipedia)*. URL: [https://en.wikipedia.org/wiki/Ansible_\(software\)](https://en.wikipedia.org/wiki/Ansible_(software)) (visitato il giorno 22/02/2025).
- [6] *Infrastructure as Code (Wikipedia)*. URL: https://en.wikipedia.org/wiki/Infrastructure_as_code (visitato il giorno 23/02/2025).
- [7] *Ansible Documentation*. URL: <https://docs.ansible.com/> (visitato il giorno 22/02/2025).
- [8] *Domain Name System (Wikipedia)*. URL: https://en.wikipedia.org/wiki/Domain_Name_System (visitato il giorno 23/02/2025).
- [9] *Dynamic Host Configuration Protocol (Wikipedia)*. URL: https://en.wikipedia.org/wiki/Dynamic_Host_Configuration_Protocol (visitato il giorno 23/02/2025).
- [10] *Dnsmasq (Wikipedia)*. URL: <https://en.wikipedia.org/wiki/Dnsmasq> (visitato il giorno 24/02/2025).
- [11] *Database (Wikipedia)*. URL: <https://en.wikipedia.org/wiki/Database> (visitato il giorno 24/02/2025).
- [12] *MariaDB (Wikipedia)*. URL: <https://en.wikipedia.org/wiki/MariaDB> (visitato il giorno 24/02/2025).
- [13] *Job Scheduler (Wikipedia)*. URL: https://en.wikipedia.org/wiki/Job_scheduler (visitato il giorno 25/02/2025).

- [14] *SLURM Overview*. URL: <https://slurm.schedmd.com/overview.html> (visitato il giorno 25/02/2025).
- [15] *SLURM Authentication Guide*. URL: <https://slurm.schedmd.com/authentication.html> (visitato il giorno 26/02/2025).
- [16] *MUNGE Documentation*. URL: <https://github.com/dun/munge/wiki> (visitato il giorno 25/02/2025).
- [17] *SLURM Consumable Resources Guide*. URL: https://slurm.schedmd.com/cons_tres.html (visitato il giorno 26/02/2025).
- [18] *SLURM Multi-core/Multi-thread Support*. URL: https://slurm.schedmd.com/mc_support.html (visitato il giorno 26/02/2025).
- [19] *SLURM Generic Resources Support*. URL: <https://slurm.schedmd.com/gres.html> (visitato il giorno 26/02/2025).
- [20] *SLURM Priority Multifactor Guide*. URL: https://slurm.schedmd.com/priority_multifactor.html (visitato il giorno 27/02/2025).
- [21] *SLURM Quality of Service Guide*. URL: <https://slurm.schedmd.com/qos.html> (visitato il giorno 27/02/2025).
- [22] *SLURM Resource Limits Guide*. URL: https://slurm.schedmd.com/resource_limits.html (visitato il giorno 27/02/2025).
- [23] *SLURM Federation Guide*. URL: <https://slurm.schedmd.com/federation.html> (visitato il giorno 28/02/2025).
- [24] *SLURM Network Configuration Guide*. URL: <https://slurm.schedmd.com/network.html> (visitato il giorno 28/02/2025).
- [25] *Vagrant Linked Clones*. URL: <https://developer.hashicorp.com/vagrant/docs/providers/virtualbox/configuration#linked-clones> (visitato il giorno 01/03/2025).
- [26] *Cloning VMs*. URL: <https://itgaiden.com/2019-02-25-cloning-vms-in-vsphere-part-1/> (visitato il giorno 01/03/2025).
- [27] *SLURM Accounting Guide*. URL: <https://slurm.schedmd.com/accounting.html> (visitato il giorno 02/03/2025).
- [28] *SLURM Preemption Guide*. URL: <https://slurm.schedmd.com/preempt.html> (visitato il giorno 06/03/2025).
- [29] *SLURM Gang Scheduling Guide*. URL: https://slurm.schedmd.com/gang_scheduling.html (visitato il giorno 06/03/2025).

Acronimi

ACID Atomicity, Consistency, Isolation, Durability	15
API Application Programming Interface	11
CI/CD Continuous Integration & Continuous Delivery	44
CPU Central Processing Unit	17
DB Database	14
DBMS Database Management System	14
DCL Data Control Language	14
DDL Data Definition Language	14
DHCP Dynamic Host Configuration Protocol	13
DML Data Manipulation Language	14
DNS Domain Name System	13
DQL Data Query Language	14
GB Gigabyte	24
GPU Graphics Processing Unit	17

GRES Generic Resource	17
HPC High Performance Computing	5
IaC Infrastructure as Code	5
IP Internet Protocol	13
LAN Local Area Network	13
MB Megabyte	32
MUNGE MUNGE Uid 'N' Gid Emporium	16
NoSQL Not only SQL	15
QoS Quality of Service	17
RAM Random Access Memory	17
RPC Remote Procedure Call	16
SLURM Simple Linux Utility for Resource Management	15
SQL Structured Query Language	14
SSH Secure Shell	13
TCP Transmission Control Protocol	27
VM Virtual Machine	11
YAML Yet Another Markup Language	13

Elenco delle figure

2.1	Virtualizzazione hardware totale [2]	12
2.2	Sequenza di risoluzione DNS [8]	14
2.3	Modello relazionale [11]	15
2.4	Componenti di un cluster SLURM [14]	16
2.5	Definizioni di <i>socket</i> , <i>core</i> e <i>thread</i> [18]	17
2.6	Comunicazione in una federazione SLURM [24]	18
3.1	Relazione tra VM di base e <i>linked clone</i> [26]	19
3.2	Topologia virtuale desiderata	20

Ringraziamenti

*Ringrazio i miei relatori di tesi, il Prof. Prandini e l'Ing. Giovine,
per il supporto e l'opportunità di approfondire l'argomento.*

*Ringrazio la mia famiglia, per il loro continuo sostegno e
incoraggiamento verso la realizzazione dei miei obiettivi.*

*Ringrazio infine gli amici di studio, di liceo, di tennistavolo,
per avermi accompagnato durante questo percorso.*