

Progetto di Sistemi Distribuiti

Dott. Borsoi Diego
Dott. Callegari Filippo
DMIF, University of Udine, Italy
(also the authors are distributed)

Version $\frac{1}{4}$, 9 maggio 2021

Indice

Indice	a
1 Introduzione	1
1.1 Descrizione del problema	1
1.2 Struttura dell'implementazione	1
1.3 Trasparenze	2
1.4 Algoritmi	2
1.5 Testing	2
1.6 Piano di sviluppo	3
2 Analisi	5
2.1 Requisiti Funzionali	5
2.2 Requisiti Non Funzionali	6
3 Progetto	7
3.1 Architettura logica	7
3.2 Protocolli ed algoritmi	9
3.3 Architettura fisica e deployment	15
3.4 Piano di sviluppo	15
4 Implementazione	17
4.1 Software e hardware	17
4.2 Utilizzo dei "Supervisors"	17
4.3 Heartbeat (module: HB_in)	18
4.4 Memoria dei parametri (module: state_server)	20
4.5 Esecutore dei comandi (module: rules_worker)	21
4.6 Ambiente e la comunicazione con i nodi	22
A Supervisors	25
B Heartbeat	27
C Memoria dei parametri	31

Capitolo 1

Introduzione

Il progetto in questione riguarda la creazione di un sistema distribuito per la comunicazione fra dispositivi all'interno di una rete sparsa, tramite l'utilizzo di eventi.

1.1 Descrizione del problema

Ogni dispositivo corrisponde ad un nodo della rete ed è caratterizzato da:

- **Id** : numero univoco del nodo.
- **Stato** : tupla di variabili che rappresentano delle specifiche caratteristiche del nodo (es. temperatura di un sensore, stato di accensione di una termocoppia, ecc).
- **Regole** : insieme di regole del tipo ECA (Event Condition Action) che possono attivarsi a seguito di un evento inviato al nodo. Queste regole possono essere di due tipi: locali, l'azione modifica solamente lo stato del nodo in cui si attiva, oppure globale, l'azione viene inviata a tutti i nodi della rete perché venga letta ed, in caso la valutazione della guardia associata sia positiva, eseguita.

$$\{event; condition; action \mid \text{if } guard \text{ then } action\}$$

Lo stato della rete si evolverà ogni qual volta un evento verrà attivato, andando a sua volta ad innescare eventuali nuovi eventi e creando quindi una sequenza di azioni a cascata.

1.2 Struttura dell'implementazione

La rete in questione ha una struttura a mesh sparsa (cioè ogni nodo sarà al più connesso a un numero di nodi molto basso, rispetto alla totalità). I nodi sono idempotenti in modo tale da avere un sistema fortemente decentralizzato. Le

varie comunicazioni fra i nodi sono eseguite al di sopra di connessioni TCP, in tal modo possiamo garantire la consegna di ogni messaggio nell'ordine prestabilito. Per quanto concerne invece le comunicazioni riguardanti il sistema di *heartbeat* (3.2), queste vengono eseguite utilizzando connessioni UDP.

1.3 Trasparenze

Di seguito sono descritte le trasparenze che convergono e sono implementate dal sistema:

Trasparenza ai fallimenti: nel momento in cui un nodo fallisce/si disconnette, il resto della rete continua a funzionare normalmente.

Trasparenza alla scalabilità: la rete può espandersi in dimensione senza che il funzionamento dei nodi vari.

Trasparenza alla mobilità: un nodo può spostarsi all'interno della rete senza che il funzionamento suo e degli altri nodi vada a modificarsi.

1.4 Algoritmi

Il sistema implementa solamente due algoritmi:

- **Flooding Algorithm:** viene usato inizialmente per la comunicazione di un'azione a tutta la rete nel momento in cui in un nodo una regola globale o di transazione viene attivata.
- **Lamport clock (modificato):** viene utilizzata una versione modificata del Lamport clock per identificare i vari *flood* eseguiti; questo clock viene incrementato solamente dall'invio (o ricezione) di un messaggio, e non dalle azioni interne ad un nodo;
- **Distributed Spanning Tree:** viene usato per ridurre il traffico di rete quando la rete si è "stabilizzata". Questo algoritmo mira a creare un albero di copertura nella rete.

Il sistema non implementa particolari algoritmi essendo che si vuole realizzare una rete distribuita dove ogni nodo conosce esclusivamente i vicini ed evolve il suo stato solamente a causa di eventi ricevuti tramite dei messaggi.

1.5 Testing

Per testare il sistema verrà utilizzata un'entità *Ambiente*, la quale simulerà:

- la creazione iniziale della rete, caricando da dei file appositi la struttura degli stati, la lista di regole e la conformazione della rete
- la scoperta di nuove connessioni

- variazioni di variabili legate all'ambiente (es. temperatura registrata da un sensore)
- fallimenti di nodi
- ritardi nell'invio di messaggi fra nodi

Ogni modulo verrà testato singolarmente ed infine verranno eseguiti dei test completi del sistema.

1.6 Piano di sviluppo

Le future fasi di sviluppo seguiranno il seguente ordine:

1. Riunione con il committente per convalidare la risoluzione del problema
2. Implementazione ambiente virtuale per la gestione dei nodi
3. Implementazione della struttura del nodo
4. Implementazione sistema *heartbeat*
5. Implementazione del sistema algoritmico
6. Test totale
7. Validazione

Capitolo 2

Analisi

In questo capitolo vengono descritti nel dettaglio requisiti funzionali e non funzionali della soluzione proposta.

2.1 Requisiti Funzionali

I requisiti funzionali individuati sono:

- **Categorizzazione di un nodo:** ogni nodo ha un tipo il quale ne identifica lo stato e le sue regole;
- **Modifica delle regole di un nodo:** ogni tipo di nodo può avere le sue regole, codificabili attraverso la programmazione dello stesso;
- **Modifica dello stato di un nodo:** ogni evento permette di avere o degli *effetti locali*, degli *effetti globali* o degli *effetti transizionali*:
 - *effetti locali*: la regola va a modificare lo stato interno;
 - *effetti globali*: la regola può modificare lo stato delle variabili interne, e può generare un evento sugli altri nodi;
 - *effetti transizionali*: regola simile a quelle globali, ma con la differenza che l'esecuzione avverrà in tutti i nodi coinvolti “contemporaneamente”;
- **Aggiunta dinamica di un nodo:** un nodo può essere aggiunto alla rete in qualsiasi momento senza perturbarne la dinamicità, limitando l'aggiornamento ai vicini a cui si collega.
- **Esecuzione di un'azione ricevuta dai vicini:** nel momento in cui un nodo riceve un messaggio dai propri vicini esso va a verificare la soddisfacibilità della guardia (se presente) e nel caso di una valutazione positiva viene eseguita l'azione associata, andando quindi a modificare il proprio stato.

- **Attivazione di una regola:** ogni qual volta avviene un cambiamento nello stato di un nodo, viene eseguito un controllo delle regole, per vedere se gli eventi generati possano attivare una o più regole del nodo; nel caso in cui una regola venga attivata, in base al tipo (locale, globale o transazione) viene portata a termine l'azione corrispondente.

2.2 Requisiti Non Funzionali

I requisiti non funzionali individuati sono:

- **Decentralizzazione:** nessun nodo ha il controllo dell'ordine degli eventi, grazie al fatto che ogni nodo è idempotente;
- **Tolleranza ai guasti:** poiché tutti i nodi sono idempotenti, nel momento in cui un nodo si scollega dalla rete, la rete rimanente continua ad operare normalmente;
- **Etereogenità:** fintanto che i nodi aggiunti utilizzano il protocollo descritto, qualunque nodo di qualunque tipo (hardware o categoria) potrà essere aggiunto alla rete;
- **Scalabilità:** l'aggiunta dinamica dei nodi alla rete permette di scalare orizzontalmente con estrema facilità;
- **Trasparenze:** le trasparenze implementate sono quelle descritte al capitolo precedente (paragrafo 1.3).

Capitolo 3

Progetto

In questo capitolo vengono descritti in modo più approfondito l'architettura del progetto, i moduli, i protocolli e gli algoritmi utilizzati.

3.1 Architettura logica

Essendo tutti i nodi costruiti al medesimo modo, di seguito presentiamo la struttura di uno singolo di essi.

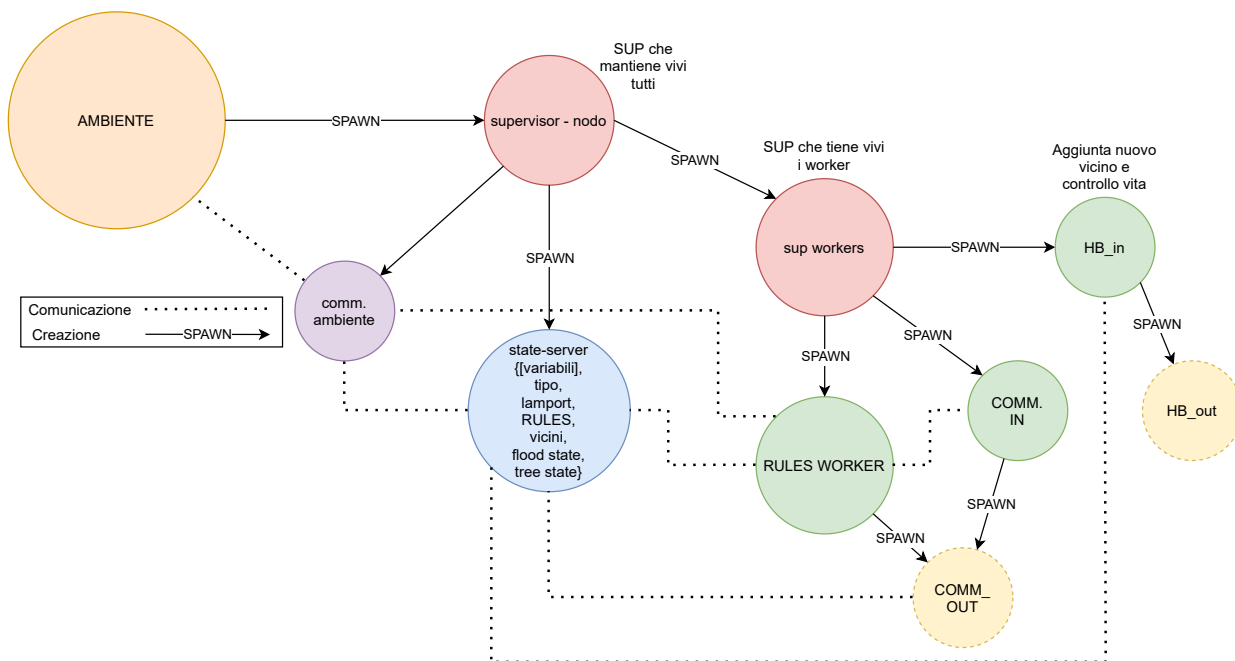


Figura 3.1: Struttura gerarchica dei moduli di un nodo e visualizzazione delle connessioni fra di essi.

Come si può vedere dalla figura 3.1, ogni nodo è formato dai seguenti moduli:

- **Supervisor nodo:** modulo che cerca di mantenere sempre operativi gli altri moduli interni. Se questo componente si riavvia equivarrebbe ad un riavvio del nodo, e quindi la conseguente perdita della modifica agli stati interni. I moduli da lui controllati sono:
 - **State server:** questo modulo manterrà tutte le variabili locali al nodo, che verranno modificate durante l’operatività dello stesso. Queste variabili sono:
 - * stato del nodo;
 - * lista delle regole associate al nodo;
 - * tipo del nodo;
 - * stato delle waves;
 - * informazioni sui nodi vicini;
 - * id univoco del nodo;
 - **Supervisor dei workers:** modulo che si occupa di gestire i moduli a lui dipendenti, riavviandoli in caso di “crash”. Questi sono:
 - * **heartbeat sense:** si occupa di controllare la “vitalità” dei vicini, di istanziare le connessioni con essi e di mantenere la “spanning tree” della rete. Ogni qual volta che deve comunicare con dei vicini affini, creerà un processo effimero (HB_out);
 - * **communication in:** si occupa di gestire tutti i messaggi in ingresso relativi al nodo in questione;
 - * **rules worker:** si occupa di applicare le azioni ricevute dai nodi vicini ed eventualmente eseguire una delle regole a lui locali al momento dell’attivazione; lui si occuperà anche di propagare le regole che generano degli eventi verso gli altri nodi, interpellando il modulo “*communication out*” (modulo apposito per l’invio dei messaggi ai nodi vicini).

Come è stato accennato in precedenza verrà sviluppato un ulteriore modulo chiamato “**Ambiente**”: questo modulo permette di simulare le interazioni che avverrebbero nel mondo reale. Nel dettaglio, le funzionalità del modulo sono:

- il “discovery” dei vicini;
- cambiamento delle variabili non dipendenti dalle regole (temperatura dell’ambiente/GPS/...);
- simulazioni di disservizi di rete;
- simulazione di guasti (temporanei o non) di un nodo;
- topologia della rete.

Dal punto di vista del nodo ci troviamo quindi costretti ad aggiungere un ulteriore modulo fittizio (*comunicazione ambiente*) per permettere la comunicazione con l'ambiente.

3.2 Protocolli ed algoritmi

Di seguito verranno descritti nel dettaglio i vari protocolli ed algoritmi utilizzati.

Controllo e attivazione delle regole

Nel momento in cui il sistema riceve un'azione da eseguire (dopo aver controllato la validità e che appartenga ad una wave non ancora ricevuta) si innesca la seguente serie di azioni:

1. il modulo *communication IN* invia l'azione da eseguire al modulo *rules worker*;
2. quest'ultimo utilizza una funzione dello *state server* per modificare lo stato del nodo in accordo all'azione ricevuta;
3. il *rules worker* esegue quindi un controllo sulle regole andando ad identificare quali possono essere attivate dalla modifica appena eseguita;
4. per ciascuna regola che viene attivata viene testata la condizione e in caso di risultato positivo:
 - se la regola è del tipo *locale*, viene eseguita l'azione associata
 - se invece la regola è del tipo *globale*, viene eseguita l'azione associata (in caso di guardia con valutazione positiva) e viene passata ad un processo di *communication OUT*, istanziato appositamente, che genera una nuova wave di messaggi inviando ai vicini la nuova azione.

Heartbeat

Il protocollo di "heartbeat" serve per mantenere consistente lo stato dei vicini di un nodo. Questo infatti controlla la loro vitalità e sarà componente chiave per l'aggiunta di un nuovo nodo.

L'algoritmo si suddivide quindi in tre componenti:

- *ECHO*: similmente al protocollo ICMP, si fa una richiesta di echo al vicino. Se questa "*ECHO_RQS*" andrà a buon fine, il primo nodo che istanzia una richiesta di "*ECHO*" riceverà un pacchetto di "*ECHO_RPL*". Definiamo come τ il tempo che intercorre tra un messaggio di *ECHO* ed un'altro. Se un nodo non rispondesse entro 2τ ad una *ECHO_RQS*, questo verrà considerato come non più collegato. Ogni *ECHO_RPL* conterrà il clock attuale del vicino.

- *ADD_NEW_ND*: similmente al protocollo DHCP, nella fase di aggiunta di un nodo alla rete, il nuovo nodo si annuncia al suo vicino “fisico”, chiedendo le informazioni essenziali per poter partecipare attivamente alla rete. Questo sarà spiegato più in dettaglio nella prossima sezione.
- *TREE_STATE*: questa componente si focalizza sulla gestione dell’albero di copertura della rete. Sarà analizzato in dettaglio nei paragrafi successivi.

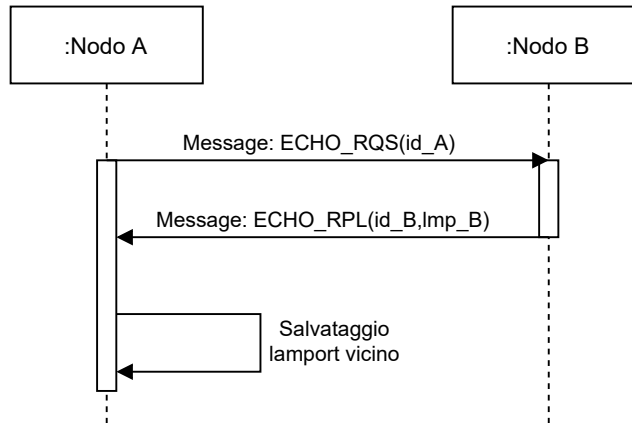


Figura 3.2: Sequence diagram dei messaggi usati per il sistema di *heartbeat*.

Aggiunta di un nuovo nodo

L’aggiunta di un nodo è una parte complicata del sistema: bisogna tener conto della possibilità che la rete si partizioni. Questo si può verificare nel caso in cui un nodo si riavvii. Il partizionamento della rete è visto come caso particolare di aggiunta di un nodo alla rete.

L’aggiunta di un nuovo nodo si compone dei seguenti passi:

1. *ADD_NEW_ND*: il nuovo nodo manda la richiesta di aggiunta alla rete a tutti i suoi vicini inviando il proprio id;
2. *ADD_NEW_NB*: il nodo che deve aggiungere il nuovo nodo risponde con un messaggio contenente (*clock*, *id*).

A questo punto, una volta che per ogni vicino ho i suoi parametri di clock, i casi possibili sono 2:

1. **tutti i nodi hanno medesimo clock**: imposto il mio clock al clock comune;
2. **esiste un clock massimo**: imposto il mio clock al valore massimo, e rispondo con *UPD_LMP* a tutti i miei vicini che non hanno il clock al

massimo. Questi a loro volta manderanno a tutti i loro vicini, con clock diverso da quello scelto, il nuovo valore.

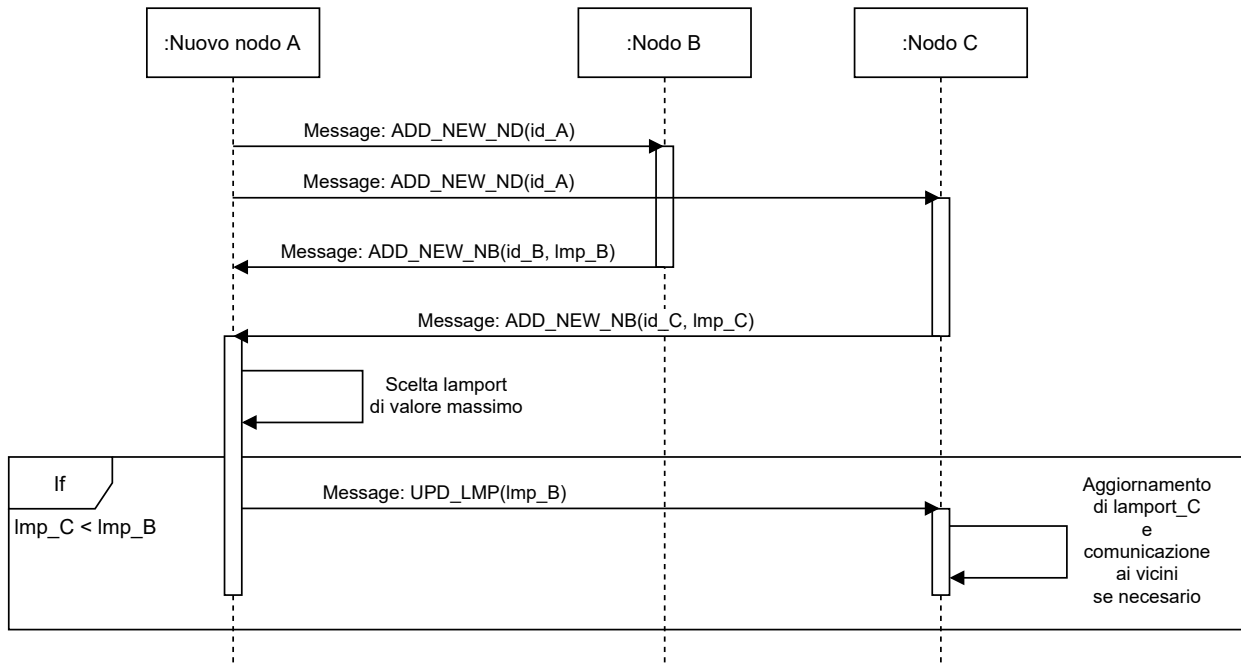


Figura 3.3: Sequence diagram dei messaggi usati per laggiunta di un nuovo nodo alla rete.

Two-phases Rules (Transazioni)

Questo è un protocollo definito come “two-phases commit”, ovvero permette l’esecuzione in contemporanea su più parti di almeno un’azione. Questo protocollo è stato modificato rispetto l’originale, in quanto, non avendo l’obbligo di avere sempre uno stato consistente tra i nodi, ci permette di semplificarlo. Le fasi sono quindi le seguenti:

1. *discovery*: troviamo tutti i nodi nella rete che sono interessati ad eseguire la transazione. Il nodo iniziatore manda nella rete un messaggio di richiesta di transazione, il quale contiene delle condizioni per partecipare. Se un nodo è interessato, risponderà all’iniziatore dicendogli di essere interessato.
2. *start*: il nodo iniziatore darà inizio alla transazione per i nodi interessati.

All’interno di questo protocollo, come possiamo vedere nell’immagine 3.4, vengono tenuti dei timeout di controllo, al fine di non rimanere bloccato nello

stato di transizione: se un nodo decide di partecipare alla transazione, allora non potrà eseguire altre regole.

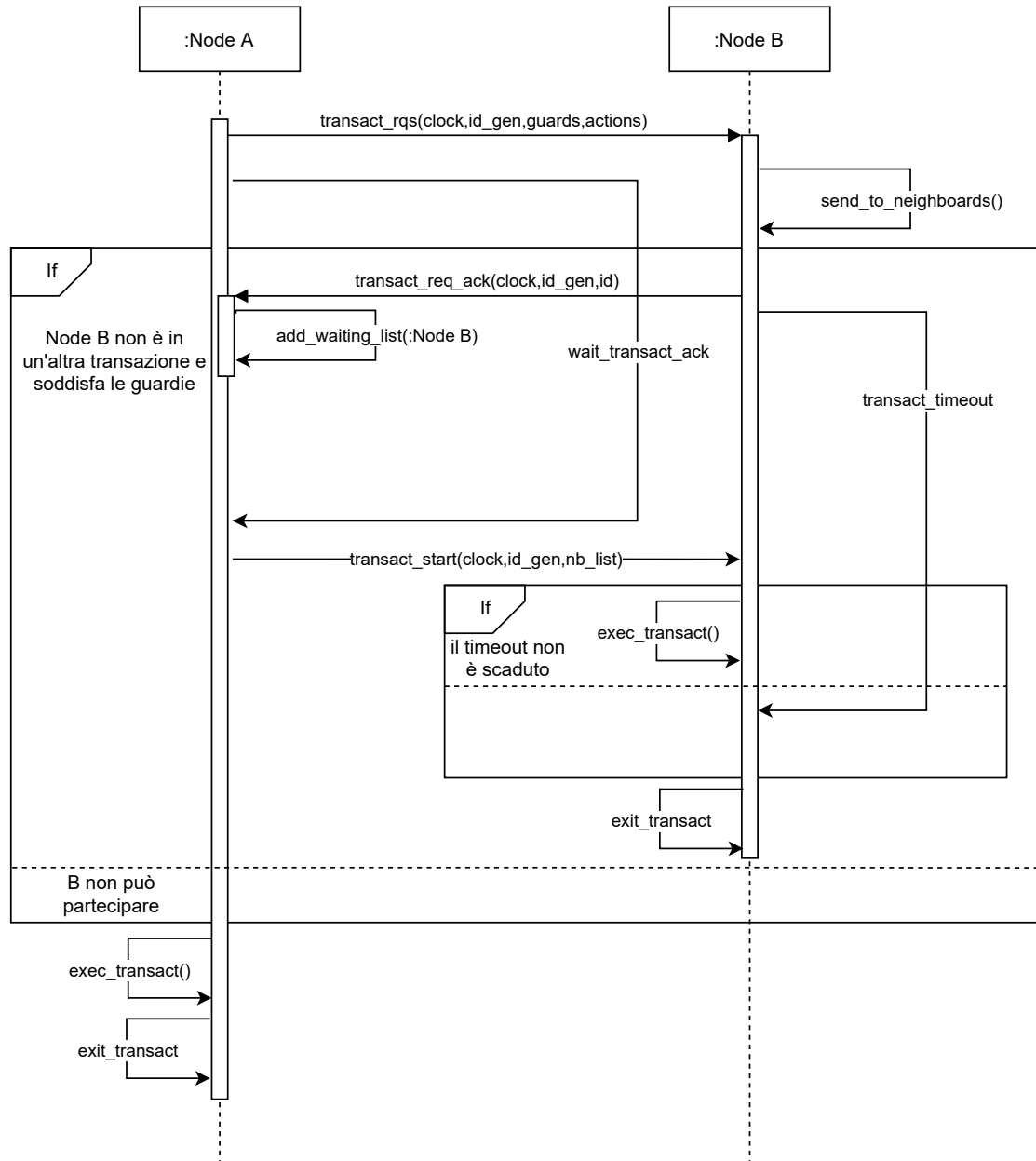


Figura 3.4: Sequence diagram del protocollo adottato per le regole di transaction.

Algoritmo: Flooding

Ogni qual volta si crea una regola che genera un evento “globale”, ogni nodo spedisce ai suoi vicini un messaggio contenente un’azione da eseguire. Questi, una volta ricevuto, aggiorneranno il loro “wave count” (tenuto dal clock, spiegato successivamente), e passerà all’esecuzione dell’istruzione condizionata contenuto nel messaggio. Al fine di evitar la presenza di messaggi vecchi nella rete, ogni messaggio conterrà la coppia (*clock, id_gen*): questo verrà salvato localmente nel nodo ricevente, e verrà mantenuto in memoria al fine di verificare se un messaggio ricevuto non è già stato eseguito. Ogni nodo quindi spedisce una copia del messaggio a tutti i suoi vicini (a patto non l’abbia già ricevuto in passato), meno a quello da cui l’ha ricevuto. Queste strategie applicate faranno sì che i messaggi circolanti nella rete siano il minor numero possibile, e nell’eventuale creazioni di cicli nella rete, dovuta alla topologia della stessa, siano soppressi alla prima occasione utile.

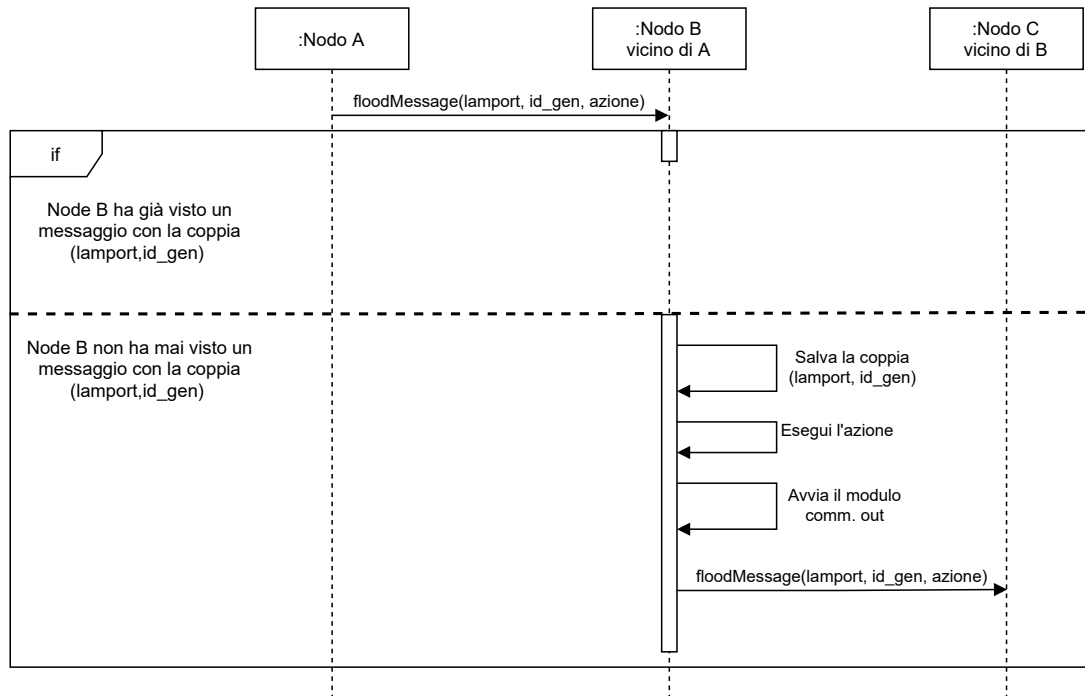


Figura 3.5: Sequence diagram dei messaggi usati per l’algoritmo di flooding.

Algoritmo: Lamport modificato

Il “Clock di Lamport” ci permette di dare una certa conseguenza temporale alle azioni. Questo sarà un numero intero crescente, ed identificherà ogni wave generata. Alla generazione di una wave ogni nodo userà il suo clock interno,

incrementato, per identificarla. Questo ci permette, come già abbiamo spiegato, di controllare la propagazione dei messaggi. Sia quindi Lmp_i il clock della nuova wave ricevuta e Lmp_n il clock del nodo corrente. Le azioni intraprendibili sono:

- $Lmp_i = Lmp_n + 1$: eseguo immediatamente l'azione in essa contenuta ed agguirno il valore di Lmp_n ;
- $Lmp_i > Lmp_n + 1$: aspetto un preimpostato timeout (es: 5τ) prima di eseguire l'azione di Lmp_i , in modo tale da poter ricevere le wave mancanti, e quindi mantenere l'ordine causale delle azioni;
- $Lmp_i < Lmp_n$: vado ad eseguire l'azione solo nel caso in cui tutte le variabili coinvolte non siano già state modificate da una wave ad essa successiva, cioè con $clock > Lmp_i$.

Algoritmo: Distributed Spanning Tree

L'albero di copertura (o spanning tree) è un protocollo distribuito tra i nodi ideato per limitare il numero di messaggi nella rete. Questo protocollo risulta essere una versione modificata del più noto *STB* (IEEE 802.1aq). Tutto questo protocollo viene eseguito sempre dal sistema di heartbeat, in quanto strettamente correlati.

Questo protocollo si divide in 3 fasi:

- **discovery**: in questa fase, il nodo che si sta connettendo esplora la rete, e dice a tutti di essere la radice. A seconda delle condizioni, esplicitate in dettagli nel capitolo successivo, potrà effettivamente esserlo o meno;
- **listening**: in questa fase, l'albero risulta stabile. Ogni τ la radice "eletta" manda dei pacchetti di keep alive ai figli. Nel momento in cui un la radice cesserà la sua permanenza nella rete, si passerà alla fase successiva;
- **transition**: in questa fase, per almeno uno componente della rete, la radice risulta non più raggiungibile. In questa fase i nodi coinvolti, non possono accettare la vecchia radice. Questa fase avviene quando un nodo perde la sua "radice", o quando lo spanning tree perde la radice che lo genera.

La scelta della radice avviene in maniera deterministica, sfruttando l'identificatore univoco di cui ogni nodo è dotato. Definiamo $\langle saved_root, saved_len, saved_nh \rangle$ la tupla che identifica i dati di radice attualmente salvata, la mia attuale distanza e attraverso quale vicino riesco a raggiungerla, mentre definiamo come $\langle root, len, nh \rangle$ la tupla proposta dai miei vicini. A questo punto la scelta per un nodo della sua radice avviene come segue:

- $saved_root > root$: la radice proposta è inferiore alla mia e sarà scelta come mia nuova radice;
- $saved_root = root \wedge saved_len + 1 < len$: il mio vicino ha una distanza peggiore della mia, lo avviso che il mio percorso potrebbe essere migliore;

- $saved_root = root \wedge saved_len > len + 1$: la mia distanza è peggiore di quella proposta, mi sposto sul mio vicino, in maniera di esser più vicino alla radice;
- $saved_root = root \wedge saved_len = len + 1 \wedge nh_saved > nh$: per questioni di determinismo, se la distanza e la radice proposta sono identiche, scelgo come “porta” verso la radice il mio vicino con id minore degli altri.

A questo punto possiamo scegliere anche a quali dei nostri possiamo inviare i messaggi: poichè ogni nodo è informato se è un nodo scelto o meno, se inoltriamo i messaggi solo ai vicini in cui sono stato scelto o al vicino che io ho scelto, posso raggiungere tutta la rete con un numero minimo di messaggi. A tal fine ogni vicino potrà assumere 3 vari stati:

- *root_port*: vicino scelto per raggiungere la radice;
- *active*: vicino che mi ha scelto per raggiungere la radice;
- *disabled*: vicino che non dipende da me.

3.3 Architettura fisica e deployment

Per quanto riguarda l’architettura fisica è necessario l’utilizzo di microcalcolatori, come dei “Raspberry” o degli “Arduino”. Ogni nodo corrisponderebbe fisicamente ad una di queste schede, avendo quindi la possibilità d’utilizzare più nodi per il medesimo “apparato”.

Non è strettamente necessario l’ausilio di microcalcolatori: potrei concentrare all’interno di un singolo calcolatore più nodi, a patto che siano gestiti in maniera consona.

Come descritto in precedenza, questi comunicherebbero con protocollo TCP ed UDP, non interessandoci quindi di tutta la rete sottostante.

Visto l’esempio a cui abbiamo pensato, si ritiene ideale l’utilizzo di connessioni wireless.

3.4 Piano di sviluppo

Le **funzionalità di base** che verranno implementate sono:

- comunicazione tra nodi;
- sistema di flooding;
- sistema di heartbeat;
- gestione dell’aggiunta di un nodo alla rete;
- gestione del riavvio di un nodo nella rete;
- impostazione iniziale dei nodi (ambiente);

- programmazione dei nodi;
- sistema basico d'esecuzione delle regole.

Sono state inoltre individuate le seguenti **funzionalità avanzate**:

- sistema avanzato d'esecuzione delle regole;
- salvataggio dello stato del nodo su files interni al controllore;
- riprogrammazione dinamica del nodo;
- implementazione di un *calculus* locale.

Capitolo 4

Implementazione

In questo capitolo tratteremo le scelte implementative avvenute nel corso del progetto. ci focalizzeremo man mano nei vari dettagli, cercando di spiegarne le motivazioni delle varianti adottate dai vari algoritmi o protocolli standard.

4.1 Software e hardware

Per l'implementazione non si è seguita una vera e propria scelta hardware, ma si è tenuto conto i vincoli, quali scarsa capacità computazionale e scarsa memoria.

A livello software la scelta è ricaduta su Erlang, in quanto linguaggio fortemente orientato ai threads, funzionale e con spiccata capacità alle connessioni. Altra peculiarità sono le librerie messe a disposizione dallo stesso tramite OTP (Open Telecom Platform). Di queste librerie sfrutteremo fortemente dei “behavior module”, quali *gen_server* e *supervisor*, i quali danno a disposizione l'astrazione del meccanismo del client/server e di supervisione dei threads del processo in vita. I moduli sviluppati più interessanti saranno quindi descritti in dettaglio nei paragrafi successivi.

4.2 Utilizzo dei “Supervisors”

Come già accennato, il meccanismo dei supervisor permette di gestire tutti i threads che verranno creati per i vari processi del nodo. Come si può notare in figura 3.1, abbiamo due supervisor: uno principale che rappresenta il vero e proprio nodo ed un secondo per mantenere la vitalità di tutti i processi interni legati alla comunicazione (dal modulo “rules_worker” al più semplice “HB_out”).

Il primo supervisor, generato dal modulo “*supervisor_nodo*”, provvede a sopprimere alla necessità di mantenere in vita tutti quei dati essenziali al funzionamento del nodo. Al suo interno infatti verranno istanziate le tabelle “ets” dove vengono salvate tutte le informazioni relative alle funzionalità del nodo (lo stato dell'albero, il clock, le regole,...), verrà istanziato il thread designato alla

comunicazione con l'ambiente, e per finire anche il thread che si occuperà di manipolare in maniera *sicura* i parametri del nodo. Ultimo, ma non per importanza, verrà generato come figlio del supervisor generale del nodo, il supervisor dei “workers”, che si occuperà di mantenere in vita i thread per l'heartbeat e per chi interpreterà le regole.

Codice d'esempio si trova nell'appendice A.

4.3 Heartbeat (module: HB_in)

Come descritto precedentemente, nel modulo *HB_in* possiamo trovare tutti i meccanismi legati alla gestione della rete e delle sue funzionalità come tale. Nello specifico, questo si occupa solamente di controllare la “vitalità” di un nodo e di creare un albero di copertura della rete stabile, con tutti i suoi protocolli annessi.

Heartbeat - consistenza della rete

La prima funzionalità (appunto, garantire la consistenza della rete), serve a garantire che se un nodo a me vicino si spenga o diventi non più raggiungibile, allora non devo più considerarlo come mio vicino. La consistenza della rete avviene tramite alcune variabili di stato interne al modulo, identificabili dalla tupla:

$$< neighb_clocks, neighb_state >$$

Il primo valore della tupla identifica una mappa di tipo $< id, clock >$, e mi serve nel momento in cui ho un aggiornamento del clock interno dovuto al ricongiungimento di una rete partizionata: per cercar di garantire la consistenza delle wave di regole che si propagheranno nella rete, dovrò cercare di allineare tutti i componenti della rete al medesimo clock. Il secondo valore della tupla identifica una mappa del tipo $< id, stato >$: per controllare la vitalità di un nodo devo infatti verificarne la sua presenza attraverso il meccanismo di “echo” descritto nel paragrafo 3.2. Alla creazione del nodo suppongo che tutti i miei vicini siano “vivi”, e sono tali fino a prova contraria (ovvero quando non rispondono più agli echo). Gli *stati* assumibili dai vicini in questa mappa quindi saranno:

$$stato(\eta : neighb) = \begin{cases} alive, & \text{se } \eta \text{ ha risposto all'ultimo echo} \\ maybe_death, & \text{se } \eta \text{ non ha risposto all'ultimo echo} \\ death, & \text{se } \eta \text{ non ha risposto agli ultimi due echo} \end{cases}$$

con η vicino di ogni nodo.

Si rimanda all'appendice B per ulteriori esemplificazioni.

Un vicino si dichiara *maybe_death* dopo $\tau = 5s$, e *death* dopo $2 * \tau$.

Heartbeat - spanning tree

Come abbiamo già accennato in precedenza, il protocollo di spanning tree prende spunto dal più famoso *spanning tree protocol* adottato dai bridge.

Questo protocollo basa il suo funzionamento sull'interscambio di tuple:

1. $\langle tree_state, \langle root, len, nh \rangle \rangle$: questa tupla indica lo stato dell'albero per il nodo che la propaga. Essa infatti contiene chi è la radice per quel nodo, quanto essa è distante e chi è che manda fuori il pacchetto;
2. $\langle tree_ack, id \rangle$: questo pacchetto avvisa il nodo vicino che verrà usato per "instradare" pacchetti nella rete;
3. $\langle tree_rm_rp, id \rangle$: avvisa la mia vecchia *root_port* che non la utilizzo più per instradare pacchetti.

Per ulteriori delucidazioni su come vengono intraprese le scelte si rimanda al paragrafo 3.2.

La necessità di avvisare i nodi vicini che li usiamo come root port o meno nasce dalla particolare situazione in cui ci troviamo: poiché non abbiamo un vero e proprio ambiente dove un nodo non designato al raggiungimento della radice comunica con tutti i nodi di quella zona, siamo costretti ad avvisare se un nodo è una *root_port* o meno. Questo passaggio complica la rete, in quanto crea un "overhead" di messaggi normalmente non contemplato dal protocollo originale.

Ulteriormente, poiché tutti i nodi non possono accorgersi in "tempo reale" della situazione delle loro "porte" (dato che sono tutte connessioni punto-punto), si richiede l'utilizzo di ulteriori clock per la correttezza dell'albero: uno per i timeout della vitalità della radice ed uno per evitare di utilizzare immediatamente una radice che si ipotizzava fosse morta. L'utilizzo di questi clock creano altri messaggi di saturazione della "rete emulata", desincronizzano i vari nodi, e peggiorandone le performance, imponendo necessariamente una potenza di calcolo più grande in quanto non sia una rete fisica.

Nel nostro specifico caso, viene fissato $\bar{\tau} = 2s$, come timeout $2 * \bar{\tau}$, come timeout per accettare di nuovo la vecchia radice $3 * \bar{\tau}$.

Heartbeat - comunicazioni in uscita

Per le comunicazioni in uscita, come possiamo notare da alcune linee di codice dell'appendice B, interpelliamo il modulo *HB_out*: questo modulo gestisce le comunicazioni in uscita, compresi gli errori di comunicazione. Parametri richiesti per la comunicazione verso l'esterno sono:

- *pacchetto*: il messaggio già confezionato da inviare
- *vicini*: una lista iterabile di vicini a cui mandarlo (per precisione, al loro heartbeat).

Questo thread viene direttamente creato dal modulo di *HB_in*, ed appena esaurisce la lista cessa la sua esistenza.

Si rimanda al file in `src/HB_out.erl` per vedere in dettaglio il suo funzionamento.

4.4 Memoria dei parametri (module: *state_server*)

Ogni nodo, come esplicitato nei requisiti funzionali (paragrafo 2.1), necessita d'essere programmabile e d'avere un insieme di variabili di stato proprie. Oltre a questo, abbiamo necessità di conoscere e di poter interpretare delle regole definite all'avvio del nodo. Inoltre, poiché abbiamo a che fare con un ambiente virtualizzato, necessitiamo di sapere quali sono i vicini con cui collabora.

Tutte queste informazioni, come già esplicitato nel capitolo precedente e nei paragrafi precedenti, sono salvate all'interno di uno "state server", il quale opera su delle tabelle persistenti (a meno di un riavvio totale di un nodo), create e fornite dal supervisor del nodo, al fine di non aver problematiche in caso di fallimenti.

Lo state server è incentrato in un'ottica client/server: per questo motivo il suo behavior module è *gen_server*.

Lo state server si occupa quindi di accedere in maniera sicura e non concorrentiale alle variabili a suo carico, al fine di non generare situazioni di deadlock o starvation.

Lo stato dello *state_server* sarà quindi composto dalla tupla

$$< vars_table, rules_table, neighb_table, node_params_table >$$

dove ogni campo equivale a:

- *vars_table*: tabella delle variabili;
- *rules_table*: tabella delle regole;
- *nb_table*: tabella dello stato dei vicini;
- *node_params_table*: tabella dei parametri interni al nodo.

I valori associati alla terza componente sono strettamente correlati a quanto affrontato nel paragrafo 4.3: al suo interno troveremo delle tuple con i valori legati allo spanning tree per cui sapremo se inoltrare o meno un messaggio a quel nodo. L'inoltro di un messaggio verso un nodo avviene solo nel caso in cui il vicino sia in stato *active* o *root_port*.

I valori associati al quarto membro della tupla sono invece legati ad informazioni chiave per identificare il nodo di appartenenza. Al suo interno possiamo trovare alcune voci, quali ad esempio: il suo identificativo univoco, che tipo di nodo è per la rete, quante "*rules wave*" ha visto, il suo clock interno e chi è la sua radice. Queste informazioni sono identificative ed univoche per un nodo, e sono il suo stato interno.

Il significato puntuale della prima e della seconda componente di questa tupla verranno affrontati in miglior dettaglio nel paragrafo 4.5.

Poiché lo state server è l'unico a poter manipolare le variabili del nodo, al suo interno avrà un piccolo "interprete" per manipolare e controllare le variabili. Questo avviene in maniera ricorsiva mediante l'utilizzo di una funzione ricorsiva chiamata "check_condition" (si rimanda all'appendice C per ulteriori delucidazioni).

4.5 Esecutore dei comandi (module: rules_worker)

blablablablablabla
step per il parsing delle regole blablablabla

propagazione delle regole

blablablabla
mi piace che vada fuori tutto subito e poi vedo se posso fare qualcosa
blablablablabla

Organizzazione delle regole

effetti locali/glob/trans
blablablablabl
struttura in appendice??

Regola locali

ci sono note particolari?

Regola globali

ci sono note particolari?

Regola transazioni

ci sono note particolari?

calcolabilità

code prioritarie
(AQ:fifo,PQ:lifo,OHT:che ne so io, forse fifo?)
blablablablabla
influenza del clock
blablablabla
calcolabilità ricorsiva
blablablabla

??

che altro devo mettere?

4.6 Ambiente e la comunicazione con i nodi

funzionalità dell'ambiente come:

```

1  %% client functions
2  -export([ignore_neighb/2, kill_node/1, ...(?)]).
3
4  -record(ambiente_state, {
5      graph,      %ets del grafo
6      id_spwn,
7      comm_spwn,
8      id_sup_node
9  }).

```

come comunico con il nodo?

comm_ambiente aiutaci tu!

```

1  %% client functions
2  -export([add_neighb/2, ignore_neighb/2]).
3
4  -record(comm_ambiente_state, {
5      name,
6      server_name,
7      rules_worker_name,
8      hb_name
9  }).
10
11  init([Name, Server_name, Rules_worker_name, HB_name, Id]) ->
12      ambiente ! {nodo_avviato, Name, {Id, HB_name}}, % IMPORTANTE: l'ambiente deve essere registrato sotto il
13      self() ! {start_discovery},
14      {ok, #comm_ambiente_state{name = Name, server_name = Server_name, rules_worker_name = Rules_worker_name,
15
16  handle_cast({add_neighb, Node = {_Node_ID, _Node_HB_name}}, State = #comm_ambiente_state{server_name = Ser
17      state_server:add_neighb(Server, Node),
18      {noreply, State};
19  handle_cast({ignore_neighb, Neighb}, State = #comm_ambiente_state{server_name = Server}) ->
20      state_server:ignore_neighb(Server, Neighb),
21      {noreply, State};
22  handle_cast(_Request, State = #comm_ambiente_state{}) ->
23      {noreply, State}.
24
25  handle_info({start_discovery}, State = #comm_ambiente_state{server_name = SN, hb_name = HBN}) ->

```



```
26  % richiedo all'ambiente di inviarmi la lista dei vicini
27  receive
28      {discover_neighbs, _Neighbs_list = []} ->
29      ok;
30      {discover_neighbs, Neighbs_list = [_ | _]} ->
31          io:format("Ricevuta lista di vicini: ~p.~n", [Neighbs_list]),
32          state_server:add_neighbs(SN, Neighbs_list)
33  end,
34  HBN ! {neighb_ready},
35  {noreply, State};
```

Appendice A

Supervisors

```
1 init({Id, Tipo}) ->
2   State_tables = create_table(Id, Tipo),
3   Server_name = list_to_atom(atom_to_list(Id) ++ "_server"),
4   Rules_worker_name = list_to_atom(atom_to_list(Id) ++ "_rules_worker"),
5   Comm_ambiente_name = list_to_atom(atom_to_list(Id) ++ "_comm_ambiente"),
6   HB_name = list_to_atom(atom_to_list(Id) ++ "_heartbeat_in"),
7   MaxRestart = 1,
8   MaxRestartPeriod = 5,
9   SupFlags = #{strategy => one_for_one,
10                intensity => MaxRestart, period => MaxRestartPeriod},
11   ChildSpecs = [
12     #{id => state_server,
13       start => {state_server, start_link, [Server_name, Id, State_tables]},
14       restart => permanent,
15       shutdown => infinity,
16       type => worker,
17       modules => [state_server]},
18     #{id => sup_workers,
19       start => {supervisor_workers, start_link,
20                [Id, Server_name, Rules_worker_name, HB_name]},
21       restart => permanent,
22       shutdown => infinity,
23       type => supervisor,
24       modules => [supervisor_workers]},
25     #{id => comm_ambiente,
26       start => {comm_ambiente, start_link,
27                [Comm_ambiente_name, Server_name, Rules_worker_name, HB_name, Id]},
28       restart => permanent,
29       shutdown => infinity,
30       type => worker,
31       modules => [comm_ambiente]}
32   ],
33   {ok, {SupFlags, ChildSpecs}}.
```

Appendice B

Heartbeat

```
1  -record(hb_state, {
2      id,                %il mio id
3      hb_name,           %il mio heartbeat
4      server_name,       %il mio state server
5      neighb_clocks,     %mappa dei vicini
6      neighb_state,      %mappa dello stato dei vicini
7      i_am_root,         %se sono la root true
8      is_root_alive,     %se la radice è viva è true
9      oldroot            %se la radice è morta per un po non posso riaccettarla
10 })
11
12 start_link(Id, Server_name, HB_name) ->
13     Pid = spawn_link(?MODULE, init, [Id, Server_name, HB_name]),
14     {ok, Pid}.
15
16 init(Id, Server_name, HB_name) ->
17     register(HB_name, self()),
18     State = #hb_state{id = Id, hb_name = HB_name, server_name = Server_name},
19     {ok, Clock} = state_server:get_clock(Server_name),
20     case Clock of
21         -1 ->
22             receive
23                 {neighb_ready} -> % aspetto che il mio comm_ambiente abbia aggiornato la lista dei vicini
24                     ok
25             end;
26         _ -> ok
27     end,
28     {ok, Neighbs} = state_server:get_neighb_hb(Server_name),
29     case {Clock, Neighbs} of
30         {-1, []} -> % siamo gli unici nella rete
31             state_server:update_clock(Server_name, 0),
32             Neighbs_clocks = maps:new();
33         {-1, _} -> %non ho ancora controllato i vicini
34             Neighbs_clocks = enter_network(Neighbs, State);
35         _ -> %il supervisor mi ha resuscitato
36             Neighbs_clocks = maps:from_list([Node, -1] || Node <- Neighbs)
37     end,
38     Neighbs_state = maps:from_list([Key, alive] || Key <- maps:keys(Neighbs_clocks)),
39     % fingo la fine di un timer per dare inizio al protocollo, essendo is_root_alive inizializzato a false
40     self() ! {start_tree},           %parto con l'albero
41     self() ! {start_echo},          %parto con gli echo
```

```

42     listen(State#hb_state{neighb_clocks = Neighbs_clocks,
43               neighb_state = Neighbs_state,
44               i_am_root = false,
45               is_root_alive = false,
46               oldroot = undefined}).
47
48 [...]
49
50 enter_network(Neighbs, _State = #hb_state{id = Id, hb_name = HB_name, server_name = Server_name}) ->
51   spawn(hb_OUT, init, [Server_name, {add_new_nd, Id, HB_name}, Neighbs]), % invia un messaggio ad ogni vic
52
53   Neighbs_clocks = maps:from_list([{Node, -1} || Node <- Neighbs]), % crea una mappa per il salvataggio de
54   New_neighbs_clocks = wait_for_all_neighbs(Neighbs_clocks, Server_name),
55   check_clock_values(New_neighbs_clocks, HB_name, Server_name),
56   New_neighbs_clocks.
57
58 [...]
59
60   {add_new_nd, Id_sender, Id_hb_sender} ->
61   %%      io:format("~p: Ricevuta richiesta connessione alla rete di ~p.~n", [Id, Id_sender]),
62   state_server:add_neighb(Server_name, {Id_sender, Id_hb_sender}),
63   {ok, Clock} = state_server:get_clock(Server_name),
64   try
65     Id_hb_sender ! {add_new_nb, HB_name, Clock}
66   catch
67     _:_ -> ok
68   end,
69   % aggiorno le due mappe usate presenti nello stato
70   New_NC = maps:put(Id_hb_sender, Clock, NC),
71   New_NS = maps:put(Id_hb_sender, alive, NS),
72   listen(State#hb_state{neighb_clocks = New_NC, neighb_state = New_NS});
73
74 [...]
75
76   {start_tree} ->
77   io:format("~p: Sono entrato in start_tree.~n", [Id]),
78   {ok, {Id_root, Dist, _ID_RP}} = state_server:get_tree_state(Server_name),
79   {ok, Neighbs_hb} = state_server:get_neighb_hb(Server_name),
80   spawn(hb_OUT, init, [Server_name, {tree_state, HB_name, {Id_root, Dist, Id}}, Neighbs_hb]),
81   if
82     Id_root == Id ->
83       New_Im_root = true,
84       self() ! {tree_root_keep_alive_timer_ended};
85     true ->
86       New_Im_root = false
87   end,
88   listen(State#hb_state{i_am_root = New_Im_root, is_root_alive = false});
89 [...]

```

```

90
91     {tree_state, HB_sender, {Id_root, Dist, Id_sender}} -> % messaggio di aggiornamento della radice dell'albero
92     io:format("~p: Ricevuto tree_state: ~p.~n", [Id, {Id_root, Dist, Id_sender}]),
93     {ok, {Saved_root, Saved_dist, Saved_RP}} = state_server:get_tree_state(Server_name),
94     {ok, Neighbs_map} = state_server:get_neighb_map(Server_name), % mappa con elementi {Id_nodo -> HB_nodo}
95     Neighbs_hb = maps:values(Neighbs_map),
96     if
97         Saved_root < Id_root ->
98         spawn(hb_OUT, init, [Server_name, {tree_state, HB_name, {Saved_root, Saved_dist, Id}}, [HB_sender]]),
99         New_Im_root = Im_root,
100         New_is_alive = Root_alive;
101     (Saved_root == Id_root) andalso (Saved_dist + 1 < Dist) ->
102     spawn(hb_OUT, init, [Server_name, {tree_state, HB_name, {Saved_root, Saved_dist, Id}}, [HB_sender]]),
103     New_Im_root = Im_root,
104     New_is_alive = Root_alive;
105     (Saved_root > Id_root) andalso (Oldroot /= Id_root) ->
106     % aggiorno lo stato dell'albero salvato
107     state_server:set_tree_state(Server_name, {Id_root, Dist + 1, Id_sender}),
108     % avviso la nuova route port che la uso come tale
109     spawn(hb_OUT, init, [Server_name, {tree_ack, Id}, [HB_sender]]),
110     % avviso la vecchia route port che non la uso più
111     if
112         Saved_RP == Id ->
113         ok;
114         true ->
115         spawn(hb_OUT, init, [Server_name, {tree_rm_rp, Id}, [maps:get(Saved_RP, Neighbs_map)]]),
116     end,
117     % avviso i vicini che ho cambiato porta
118     spawn(hb_OUT, init, [Server_name,
119         {tree_state, HB_name, {Id_root, Dist + 1, Id}},
120         Neighbs_hb -- [HB_sender]]),
121     New_Im_root = false,
122     New_is_alive = false,
123     erlang:send_after(4000, self(), {tree_keep_alive_timer_ended, Id_root});
124     ((Saved_root == Id_root) andalso (Saved_dist > Dist + 1))
125     or
126     ((Saved_root == Id_root) andalso (Saved_dist == Dist + 1) andalso (Id_sender < Saved_RP)) ->
127     % aggiorno lo stato dell'albero salvato
128     state_server:set_tree_state(Server_name, {Id_root, Dist + 1, Id_sender}),
129     % avviso la nuova route port che la uso come tale
130     spawn(hb_OUT, init, [Server_name, {tree_ack, Id}, [HB_sender]]),
131     % avviso la vecchia route port che non la uso più
132     spawn(hb_OUT, init, [Server_name, {tree_rm_rp, Id}, [maps:get(Saved_RP, Neighbs_map)]]),
133     % avviso i vicini che ho cambiato porta
134     spawn(hb_OUT, init, [Server_name,
135         {tree_state, HB_name, {Id_root, Dist + 1, Id}},
136         Neighbs_hb -- [HB_sender]]),
137     New_Im_root = false,

```

```
138         New_is_alive = Root_alive;
139     true ->
140         New_Im_root = Im_root,
141         New_is_alive = Root_alive
142     end,
143     listen(State#hb_state{i_am_root = New_Im_root, is_root_alive = New_is_alive});
```

Appendice C

Memoria dei parametri

```
1  %% API
2  -export([start_link/3]).
3
4  %% gen_server callbacks
5  -export([init/1, handle_call/3, handle_cast/2, handle_info/2, terminate/2, code_change/3]).
6
7  %% client functions
8  -export([exec_action/3, exec_action_from_local_rule/3, check_rule_cond/3, check_trans_guard/3]).
9  -export([update_clock/2, get_clock/1, get_rules/1]).
10 -export([get_neighb/1, get_neighb_hb/1, add_neighb/2, add_neighbs/2, rm_neighb/2, rm_neighb_with_hb/2, check_neighb/2,
11 -export([get_tree_state/1, reset_tree_state/1, set_tree_state/2, set_tree_active_port/2, rm_tree_active_port/2]).
12 -export([get_active_neighb/1, get_active_neighb_hb/1]).
13 -export([ignore_neighb/2, get_ignored_neighb_hb/1]).
14
15 -record(server_state, {
16     id,
17     vars_table,
18     rules_table,
19     neighb_table,
20     node_params_table,
21     lost_connections
22 }).
23
24 %%%=====
25 %%% API
26 %%%=====
27
28 start_link(Name, Id, State_tables) when is_atom(Name) ->
29     gen_server:start_link({local, Name}, ?MODULE, [Id, State_tables], []);
30
31 [...]
32
33 check_condition(Cond, VT, PT) ->
34     case Cond of          %operazioni binarie
35     {} ->
36         true;
37     {Op, Var1, Var2} -> % operazioni di arit  2: lt, lte, gt, gte, eq, neq
38         % nel caso in cui le variabili siano atomi (quindi vere e proprie variabili) devo ottenere il valore corrispondent
39         % altrimenti sono dei semplici numeri e quindi li uso cos  come sono
40         Real_var1 = if
41             is_atom(Var1) ->
```

```

42         [{Var1, Var1_value, _Var1_clock}] = ets:lookup(VT, Var1),
43         Var1_value;
44     true ->
45         Var1
46     end,
47     Real_var2 = if
48         is_atom(Var2) ->
49         [{Var2, Var2_value, _Var2_clock}] = ets:lookup(VT, Var2),
50         Var2_value;
51     true ->
52         Var2
53     end,
54     % ritorna il valore corrispondente all'operazione eseguita sulle variabili
55     case Op of
56     lt ->
57         Real_var1 < Real_var2;
58     lte ->
59         Real_var1 =< Real_var2;
60     gt ->
61         Real_var1 > Real_var2;
62     gte ->
63         Real_var1 >= Real_var2;
64     eq ->
65         Real_var1 == Real_var2;
66     neq ->
67         Real_var1 /= Real_var2;
68     %%      land ->
69     %%      check_condition(Real_var1,VT,PT) andalso check_condition(Real_var2,VT,PT);
70     %%      lor ->
71     %%      check_condition(Real_var1,VT,PT) orelse check_condition(Real_var2,VT,PT);
72     _ ->
73         io:format("State server - check_condition condizione non riconosciuta: ~p.~n", [Cond]),
74         false
75     end;
76     {OpU, Tpc} -> %operazioni unarie
77     case OpU of
78     tpe ->
79         [{tipo, Tpc}] == ets:lookup(PT, tipo);
80     ntp ->
81         [{tipo, Tpc}] /= ets:lookup(PT, tipo);
82     %%      lnot ->
83     %%      Real_var1 = if
84     %%          is_atom(Var1) ->
85     %%          [{Var1, Var1_value, _Var1_clock}] = ets:lookup(VT, Var1),
86     %%          Var1_value;
87     %%          true ->
88     %%          Var1
89     %%      end,

```

```
90  %%          check_condition(Real_var1, VT, PT);
91  _ ->
92      io:format("State server - check_condition condizione non riconosciuta: ~p.~n", [Cond]),
93      false
94  end;
95  _ ->
96      io:format("State server - check_condition condizione non riconosciuta: ~p.~n", [Cond]),
97      false
98  end.
```
