

RightSizing-SpikeServer

Contents

1	Caso di Studio	2
2	Obiettivi dello studio	2
3	Modello concettuale	2
3.1	Spike server	3
3.2	Web Server	3
4	Modello delle specifiche	3
4.1	Componenti fisici del sistema:	3
4.2	Specifiche del carico di lavoro:	3
4.3	Logica di controllo:	4
4.4	Metrica di valutazione delle prestazioni:	4
4.5	Gestione della fine della simulazione:	4
5	Modello Computazionale	4
5.1	Stato del sistema	4
5.2	Descrizione del job	5
5.3	Identificazione degli eventi	5
5.4	Logica di controllo e routing	5
5.5	Configurazione dei parametri delle distribuzioni	5
5.6	Struttura del simulatore	6
5.6.1	Parallelismo	7
5.7	Raccolta delle metriche	7
5.8	Scheduler	9
6	Verifica	9
6.1	Verifica del routing basato su soglia SI_{max}	9
6.1.1	Verifica con $SI_{max} = 0$	9
6.1.2	Verifica con $SI_{max} = \infty$	10
6.2	Verifica delle metriche prestazionali	10
6.2.1	Verifica dell'utilizzazione al variare di SI_{max}	10
6.2.2	Verifica del throughput	11
7	Validazione	13
8	Risultati	13
8.1	Obiettivo 1	14
8.2	Obiettivo 2	15
8.3	Obiettivo 3	17

1 Caso di Studio

Il sistema oggetto di studio è un'architettura di data center per un Internet Service Provider, progettata per gestire dinamicamente le fluttuazioni di carico e garantire la Quality of Service (QoS) ottimizzando al contempo l'uso delle risorse.

Il problema principale affrontato è il "right-sizing", ovvero come evitare sia il sovradimensionamento (spreco di risorse) sia il sottodimensionamento (violazione degli SLA e degrado delle prestazioni), specialmente in presenza di fluttuazioni di carico a breve e lungo termine.

L'architettura proposta, come descritto nel caso di studio 6.2 del libro di testo "Performance Engineer", si basa su un livello di scaling verticale che gestisce i picchi di carico improvvisi e di breve durata. Questo livello introduce uno Spike Server dedicato. Un Load Controller monitora un indicatore di picco (Spike Indicator, SI), definito come il numero di richieste concorrenti totali presenti su un Web Server (tutte le richieste vengono gestite con un scheduler processor sharing).

Il comportamento del sistema seguirebbe quanto descritto:

- Quando l'indicatore SI supera una soglia di allarme SI_{max} , le nuove richieste in arrivo non vengono più inviate al Web Server congestionato, ma vengono reindirizzate allo Spike Server.
- Quando il carico sul Web Server diminuisce e SI scende al di sotto della soglia, il routing delle richieste torna alla normalità.

2 Obiettivi dello studio

Lo studio si pone l'obiettivo di analizzare e validare l'efficacia del modello di autoscaling basato su spike server attraverso la simulazione. Gli obiettivi specifici sono:

- Determinare il valore di SI_{max} più alto possibile (provando vari valori di SI) che mantenga comunque il tempo di risposta medio $E[R] \leq 8$ secondi, con un tasso di arrivo iperesponenziale con media di 6.66 req/s (400 req/min).
- Analizzare come varia il tempo di risposta medio al variare del coefficiente di variazione degli arrivi iperesponenziali (fluttuazioni a breve termine) per carichi di lavoro crescenti (da 1 req/s a 12 req/s).
- Verificare come cambia il contesto (sempre sotto carichi crescenti e al variare di SI_{max}) se lo spike server ha un tasso di servizio doppio rispetto al webserver principale, invece che uguale come nel caso base.

3 Modello concettuale

Il modello descritto può essere schematizzato come in figura 1.

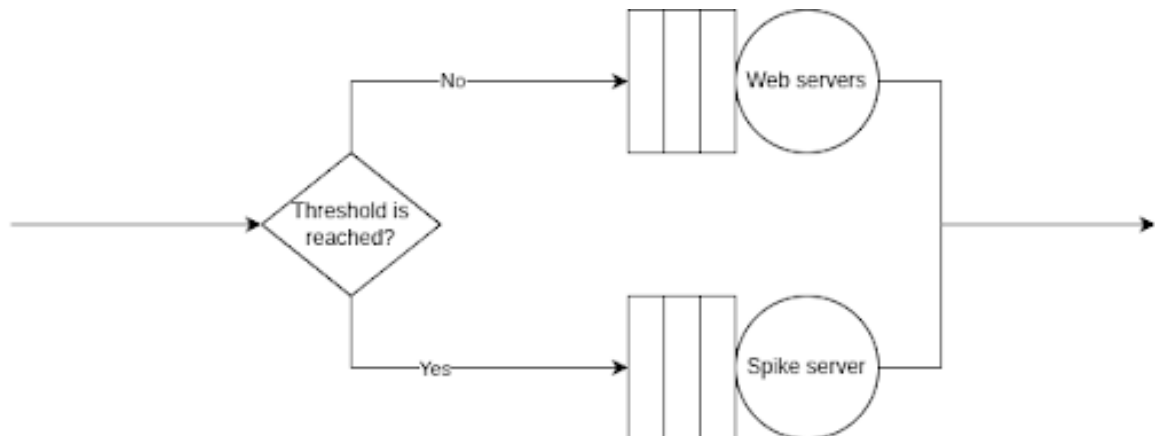


Figure 1: Modello concettuale del sistema di autoscaling gerarchico con Spike Server

I job arrivano e a seconda del livello di pienezza dei webservers. Nonostante il sistema possa sembrare troppo semplice per una analisi simulativa, il fatto che il routing non sia probalistico lo rende molto complesso da analizzare in modo analitico. il routing dei job, dipende strettamente dallo stato del webserver nel momento del routing.

Spiegazione tempi di servizio esponenziale: Nel caso di studio affrontato nel libro viene utilizzata una distribuzione iperesponenziale nei tassi di servizio per modellare il fatto che ad un server arrivano job di dimensione molto variabile. Tale comportamento potrebbe essere modellato anche utilizzando delle size dei job differenti anziché agire sul tasso di servizio. Tuttavia ho deciso di attenermi al testo originale e utilizzare anche io dei tassi di servizio iperesponenziali per modellare questo comportamento per poter confrontare meglio i risultati ottenuti.

3.1 Spike server

Lo Spike Server avrà la stessa potenza rispetto al webserver nei primi 2 obiettivi. Nell'obiettivo 3, invece, avrà un tasso di servizio doppio rispetto a quello del webserver principale, in modo da poter smaltire rapidamente i picchi di carico. Per il calcolo di SI_{max} ottimale è importante sottolineare che lo spike server è bene usarlo il meno possibile, in quanto è una risorsa più costosa rispetto al webserver. Quindi anziché cercare di minimizzare il tempo di risposta medio del sistema, si cerca di rispettare lo SLA (8 secondi) utilizzando il meno possibile lo spike server.

3.2 Web Server

Il web server sarà solamente uno, per semplicità di modellazione. Per questo motivo non vengono gestite le fluttuazioni a lungo termine, in quanto queste richiederebbero un sistema con più webserver per utilizzare una scaling orizzontale.

4 Modello delle specifiche

Queste variabili matematiche, sono la rappresentazione dello stato del sistema:

- $SI(t) \in [0, SI_{max}]$: Spike indicator, indica il numero di job presenti nel web server al tempo t .
- $n_{spike}(t) \in [0, \infty)$: numero di job presenti nello spike server al tempo t .

per cui lo stato del sistema è definito come la tupla $S(t) = (SI(t), n_{spike}(t))$.

4.1 Componenti fisici del sistema:

- Web server principale: server principale che gestisce le richieste in arrivo.
- Spike server: server secondario che viene attivato quando il carico sul web server principale supera una certa soglia (SI_{max}).

Componenti logici del sistema:

- Load Controller: componente che monitora l'indicatore di picco (Spike Indicator, SI) e decide il routing delle richieste in base al valore di SI e alla soglia SI_{max} .

4.2 Specifiche del carico di lavoro:

- Arrivi:
 - Valore base: Processo di arrivo iperesponenziale con $cv = 4$ e media 0.15 secondi (6.66 req/s)
 - Valore stress: Processo di arrivo iperesponenziale con carico variabile da 1 req/s a 12 req/s
- Servizio: Distribuzione iperesponenziale con $cv = 4$ e media 0.16 secondi (web server e spike server), media 0.08 secondi (spike server nell'obiettivo 3)

- Scheduling: Processor Sharing

Vale la pena notare che, nel caso base del webserver, siccome $\rho > 1$, il sistema è instabile e senza spike server il tempo di risposta divergerebbe a infinito.

4.3 Logica di controllo:

- Arrivi nuovi job:
 - Se $SI < SI_{max}$: invia al webserver principale
 - Se $SI \geq SI_{max}$: invia allo spike server
- Completamento job:
 - Decrementa il contatore SI se il job era nel webserver principale
 - Decrementa il contatore n_spike se il job era nello spike server

Nella logica originale del caso di studio, il valore di SI parte da SI_{max} e viene decrementato ad ogni completamento di un job. La logica da me adottata, nonostante sia inversa, è perfettamente equivalente.

4.4 Metrica di valutazione delle prestazioni:

Lo SLA da rispettare è un tempo di risposta medio $E[R] \leq 8$ secondi.

4.5 Gestione della fine della simulazione:

Dato che il sistema viene studiato a regime, esso fa riferimento ad un orizzonte temporale prefissato che va dai 120 secondi (per eliminare il bias della fase transitoria in cui le code sono ancora vuote) ai 1200 secondi di simulazione. Questo viene fatto nel report di riferimento e quindi per confrontare al meglio i risultati si è deciso di adottare lo stesso approccio. Inoltre, non verrà fatta allo scadere del tempo di simulazione una pulizia delle code, ma si considereranno solo i job completati entro il tempo di simulazione, per il calcolo delle metriche medie aggregate (es. $R0$), coerentemente con l'approccio adottato nel caso di studio di riferimento.

5 Modello Computazionale

Il sistema di simulazione è implementato in Python e segue un modello ad eventi discreti. Sono state utilizzate alcune file della libreria di Steve Park & Dave Geyer, tradotti in python da Philip Steele. In particolare sono state utilizzate le librerie per la generazione di numeri casuali per il multistream (rngs.py) e per le distribuzioni (rvgs.py).

5.1 Stato del sistema

Le variabili di programmazione che rappresentano lo stato del sistema sono:

- web_jobs : list: è una lista che contiene i job in fase di processamento nel web server, la sua lunghezza rappresenta l'indicatore SI
- spike_jobs : list: è una lista che contiene i job in fase di processamento nello spike server.

5.2 Descrizione del job

Ogni job è rappresentato da una classe Job:

```
class Job:
    def __init__(self,
                  arrival_time,
                  service_demand,
                  is_spike=False):
        self.arrival_time = arrival_time
        self.service_demand = service_demand
        self.remaining_work = service_demand
        self.is_spike = is_spike
```

Dove:

- arrival_time: tempo di arrivo del job
- service_demand: tempo di servizio totale richiesto dal job
- remaining_work: lavoro rimanente da completare (utile per lo scheduling processor sharing)
- is_spike: booleano che indica se il job è stato assegnato allo spike server

5.3 Identificazione degli eventi

- arrivo di un nuovo job
- completamento di un job nel web server
- completamento di un job nello spike server

5.4 Logica di controllo e routing

Ogni volta che arriva un nuovo job:

- Se $\text{len}(\text{web_jobs}) \leq SI_{max}$, il job viene aggiunto a web_jobs.
- Se $\text{len}(\text{web_jobs}) > SI_{max}$, il job viene aggiunto a spike_jobs.

Ogni volta che job completa rimuove dalla lista corrispondente il job (web_jobs o spike_jobs).

5.5 Configurazione dei parametri delle distribuzioni

- Stream 0: Inter-arrivi (Iperesponenziale, media 0.15s, cv=4).
- Stream 1: Servizio Web Server (Iperesponenziale, media 0.16s, cv=4).
- Stream 2: Servizio Spike Server (Iperesponenziale, media 0.08s, cv=4).

La distribuzione iperesponenziale è stata implementata utilizzando l'esponenziale e l'uniforme della libreria rvgs.py. Con l'uniforme calcoliamo la probabilità di scegliere una delle due esponenziali, e con l'esponenziale calcoliamo il servizio vero e proprio (Figura 2).

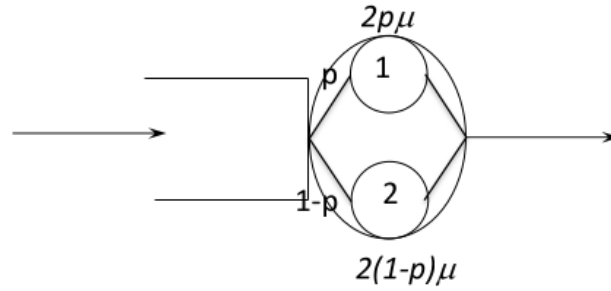


Figure 2: Generazione di una variabile casuale iperesponenziale

In particolare date in input la media e il coefficiente di variazione cv . Sapendo che $g(p) = cv^2$, possiamo calcolare p_1 , p_2 , m_1 e m_2 come:

$$c = \frac{cv^2 - 1}{cv^2 + 1}$$

$$p_1 = \frac{1 + c}{2}$$

$$p_2 = 1 - p_1$$

$$m_1 = \frac{mean}{2 \cdot p_1}$$

$$m_2 = \frac{mean}{2 \cdot p_2}$$

e quindi utilizzare un semplice if per scegliere quale esponenziale utilizzare in base ad un campione uniforme:

```
u = Uniform(0,1)
if u < p1:
    return Exponential(m1)
else:
    return Exponential(m2)
```

Da notare che per ogni campionamento dell'esponenziale si utilizzano, quindi, due estrazioni dallo stream del generatore di numeri casuali.

5.6 Struttura del simulatore

Il simulatore è strutturato come una classe Simulator che gestisce l'inizializzazione, l'esecuzione e la raccolta delle metriche della simulazione. La classe ha una serie di costanti di configurazione:

```
START      = 0.0
BIAS_PHASE = 120.0
STOP       = 1200.0
INFINITY   = 1e15
SEED       = 8
REPLICAS   = 100
N_PROCESSES = 12
```

Dove:

- START: tempo di inizio della simulazione
- BIAS_PHASE: durata della fase transitoria (120 secondi)
- STOP: tempo di fine della simulazione (1200 secondi)

- INFINITY: un valore molto grande usato per inizializzare i tempi di evento (deve essere molto più grande di STOP)
- SEED: seme per il generatore di numeri casuali, che è stato utilizzato insieme alla libreria rngs.py per tutti i generatori di numeri casuali
- REPLICAS: numero di repliche della simulazione (100 repliche)
- N_PROCESSES: numero di processi paralleli per l'esecuzione delle repliche (12 processi)

5.6.1 Parallelismo

Per velocizzare l'esecuzione delle repliche della simulazione, è stata implementata una logica di multiprocessing per sfruttare tutti i core della mia CPU disponibili. Tuttavia per farlo ho dovuto prendere alcune accortezze che permettessero di non minare né l'indipendenza delle repliche né la riproducibilità dei risultati. In particolare:

- ad ogni processo vengono assegnati 3 stream dei 256 disponibili, in modo che ogni processo abbia il proprio spazio di numeri casuali indipendenti dagli altri processi, ma in maniera tale che potessero procedere in parallelo senza che la non deterministicità del parallelismo influisse sull'ordine di estrazione dei numeri casuali.
- Per il numero di repliche introdotte, ogni processo non si avvicinava nemmeno lontanamente al limite dello stream assegnato.

$$MODULUS/STREAM \approx 8388607$$

Ho comunque implementato un meccanismo in cui una volta sfiorato il limite dello stream, gli venivano assegnati 3 nuovi stream. Ed essendo il numero di processi 12, e ogni processo usa 3 stream, questo è tranquillamente possibile senza sfiorare il limite dei 256 stream e senza ulteriori meccanismi di controllo.

- Sempre per garantire la riproducibilità dei risultati, ad ogni processo vengono assegnate sempre lo stesso numero di repliche da eseguire, in modo tale che le ultime repliche non vengano eseguite da processi diversi (e quindi stream diversi) in base all'assegnazione dinamica di essi.

5.7 Raccolta delle metriche

Per raccogliere le metriche è stata utilizzata una classe Track:

```
class Track:
    def __init__(self):
        self.area_node_web = 0.0
        self.area_node_spike = 0.0

        self.area_busy_web = 0.0
        self.area_busy_spike = 0.0

        self.completed_web = 0
        self.completed_spike = 0

        self.scaling_actions = 0
```

Dove:

- area_node_web: area sotto la curva del numero di job nel web server

$$\int_{BIAS_PHASE}^{STOP} N_{web}(t) dt$$

- **area_node_spike**: area sotto la curva del numero di job nello spike server

$$\int_{BIAS_PHASE}^{STOP} N_{spike}(t)dt$$

- **area_busy_web**: area sotto la curva del tempo di utilizzo del web server

$$\int_{BIAS_PHASE}^{STOP} U_{web}(t)dt$$

- **area_busy_spike**: area sotto la curva del tempo di utilizzo dello spike server

$$\int_{BIAS_PHASE}^{STOP} U_{spike}(t)dt$$

- **completed_web**: numero di job completati nel web server
- **completed_spike**: numero di job completati nello spike server
- **scaling_actions**: numero di azioni di scaling (numero di volte in cui lo spike server è stato attivato)

Gli integrali sono calcolati ad ogni evento in maniera incrementale durante la simulazione.

Le metriche finali raccolte sono state:

- **web_response_time**: tempo di risposta medio nel web server calcolato come l'integrale del numero di job completati nel web server diviso per il numero di job completati nel web server

$$\frac{area_node_web}{completed_web}$$

- **spike_response_time**: tempo di risposta medio nello spike server

$$\frac{area_node_spike}{completed_spike}$$

- **total_response_time**: tempo di risposta medio totale

$$\frac{area_node_web + area_node_spike}{completed_web + completed_spike}$$

- **utilization_web**: utilizzo medio del web server

$$\frac{area_busy_web}{STOP - BIAS_PHASE}$$

- **utilization_spike**: utilizzo medio dello spike server

$$\frac{area_busy_spike}{STOP - BIAS_PHASE}$$

- **throughput_web**: throughput medio del web server

$$\frac{completed_web}{STOP - BIAS_PHASE}$$

- **throughput_spike**: throughput medio dello spike server

$$\frac{completed_spike}{STOP - BIAS_PHASE}$$

- **throughput_total**: throughput medio totale

$$\frac{completed_web + completed_spike}{STOP - BIAS_PHASE}$$

- **scaling_actions**: numero di azioni di scaling

Per quanto riguarda le metriche finali su tutte le repliche, è stato utilizzato l'algoritmo di Welford per il calcolo delle medie e varianze e poi degli intervalli di confidenza al 95%.

5.8 Scheduler

Siccome lo scheduler è di tipo processor sharing, il tempo prima del prossimo evento da schedulare è calcolato come:

$$\min\{next_arrival, time_to_complete_web_job, time_to_complete_spike_job\}$$

Dove:

- *next_arrival*: tempo rimanente al prossimo arrivo
- *time_to_complete_web_job*: tempo rimanente al completamento del job con il minor lavoro rimanente nel web server
- *time_to_complete_spike_job*: tempo rimanente al completamento del job con il minor lavoro rimanente nello spike server

Per calcolare i tempi rimanenti al completamento dei job, non basta trovare il job con il minor lavoro rimanente, perché essendo in processor sharing il tempo di completamento dipende anche dal numero di job in servizio. Quindi il tempo di completamento del job con il minor lavoro rimanente è calcolato come:

$$time_to_complete_web_job = min_remaining_work \times len(web_jobs)$$

$$time_to_complete_spike_job = min_remaining_work \times len(spike_jobs)$$

Una volta trovato il tempo del prossimo evento, si avanza il tempo di simulazione di tale intervallo e si aggiorna il lavoro rimanente di tutti i job in servizio sottraendo il lavoro svolto in tale intervallo:

$$j.remaining_work = j.remaining_work - \frac{time_advance}{len(web_jobs)} \forall j \in web_jobs$$

$$j.remaining_work = j.remaining_work - \frac{time_advance}{len(spike_jobs)} \forall j \in spike_jobs$$

6 Verifica

Ci sono varie verifiche che si possono fare andando a controllare se tutta una serie di caratteristiche volute si riflettano veramente nel modello simulativo costruito e quindi nei risultati ottenuti. In particolare, andrò ad effettuare le seguenti verifiche:

- Verifica del corretto funzionamento di routing basato su soglia si SI_{max} .
 - Controllo che impostando $SI_{max} = 0$ venga utilizzato solo lo spike server.
 - Controllo che impostando $SI_{max} = \infty$ non venga mai utilizzato lo spike server.
- Verifica delle metriche prestazionali:
 - Verifico che aumentando SI_{max} l'utilizzazione del web server aumenta
 - Verifico che il throughput sia pari al tasso di arrivo quando il sistema è stabile e uguale al tasso di servizio quando il sistema è sovraccarico.

6.1 Verifica del routing basato su soglia SI_{max}

6.1.1 Verifica con $SI_{max} = 0$

Impostando $SI_{max} = 0$ si può verificare che tutti i job vengano instradati allo spike server. Infatti, come mostrato in figura 3, l'utilizzazione del web server è nulla, mentre tutto il carico viene gestito dallo spike server.

```

Simulazione completata con i seguenti parametri:
  SI_max: 0
  Arrival Mean: 0.15
  Web Mean: 0.16
  Spike Mean: 0.16
  Coefficiente di Variazione: 4.0
Statistiche raccolte con intervallo di confidenza al 95%:
spike_response_time      : 12.9032 +/- 1.1099 (Var: 32.069075)
total_response_time      : 12.9032 +/- 1.1099 (Var: 32.069075)
utilization_web          : 0.0000 +/- 0.0000 (Var: 0.000000)
utilization_spike        : 0.9910 +/- 0.0028 (Var: 0.000205)
throughput_web           : 0.0000 +/- 0.0000 (Var: 0.000000)
throughput_spike         : 6.5423 +/- 0.0339 (Var: 0.029877)
throughput_total         : 6.5423 +/- 0.0339 (Var: 0.029877)
scaling_actions          : 44.9200 +/- 7.3205 (Var: 1394.983434)

```

Figure 3: Verifica routing con $SI_{max} = 0$

6.1.2 Verifica con $SI_{max} = \infty$

Impostando $SI_{max} = \infty$ si può verificare che nessun job venga instradato allo spike server. Infatti, come mostrato in figura 4, l'utilizzazione dello spike server è nulla, mentre tutto il carico viene gestito dal web server.

```

Simulazione completata con i seguenti parametri:
  SI_max: 1000000000000000.0
  Arrival Mean: 0.15
  Web Mean: 0.16
  Spike Mean: 0.16
  Coefficiente di Variazione: 4.0
Statistiche raccolte con intervallo di confidenza al 95%:
web_response_time        : 13.1808 +/- 1.0418 (Var: 28.251530)
total_response_time      : 13.1808 +/- 1.0418 (Var: 28.251530)
utilization_web          : 0.9914 +/- 0.0028 (Var: 0.000202)
utilization_spike        : 0.0000 +/- 0.0000 (Var: 0.000000)
throughput_web           : 6.5389 +/- 0.0351 (Var: 0.032065)
throughput_spike         : 0.0000 +/- 0.0000 (Var: 0.000000)
throughput_total         : 6.5389 +/- 0.0351 (Var: 0.032065)
scaling_actions          : 0.0000 +/- 0.0000 (Var: 0.000000)

```

Figure 4: Verifica routing con $SI_{max} = \infty$

6.2 Verifica delle metriche prestazionali

6.2.1 Verifica dell'utilizzazione al variare di SI_{max}

Impostando un carico di lavoro fisso e variando il valore di SI_{max} , si può verificare che l'utilizzazione del web server aumenti al crescere di SI_{max} , come mostrato in figura 5.

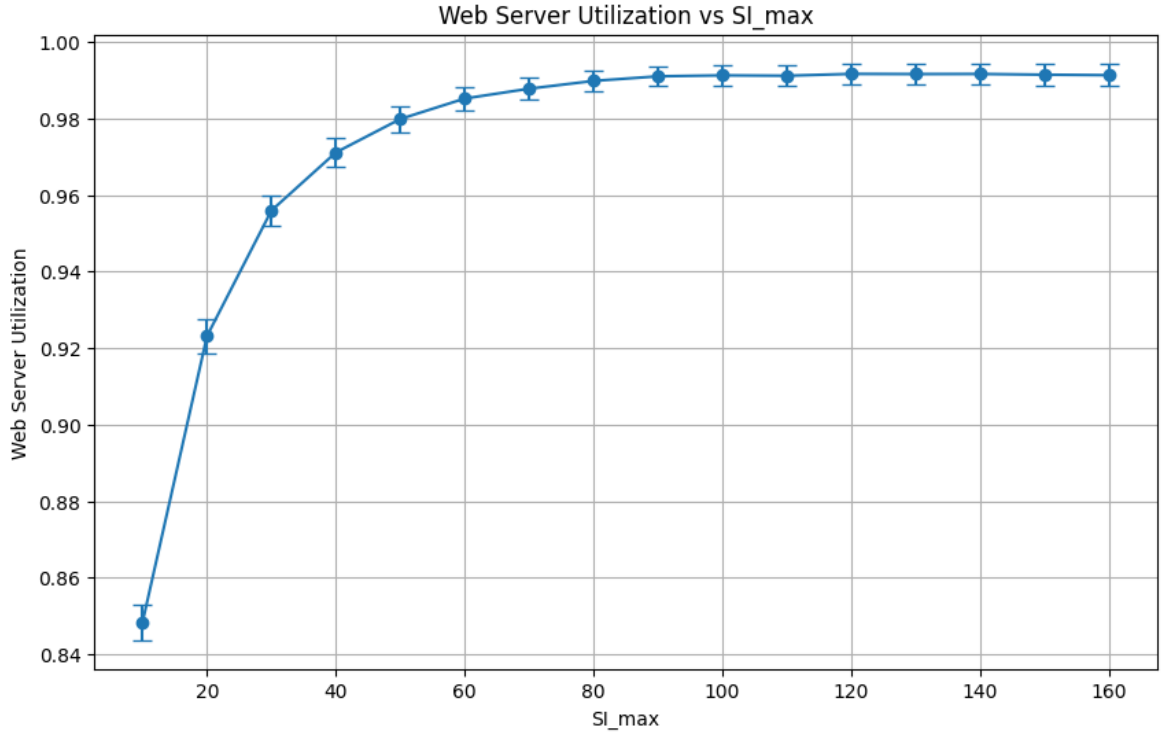


Figure 5: Verifica dell'utilizzazione del web server al variare di SI_{max}

Table 1: Verifica dell'utilizzazione del web server al variare di SI_{max}

SI_{max}	U_{web}	U_{web} CI95
10	0.8484	0.0047
20	0.9231	0.0045
30	0.9559	0.0040
40	0.9711	0.0037
50	0.9799	0.0034
60	0.9852	0.0031
70	0.9878	0.0030
80	0.9898	0.0028
90	0.9910	0.0025
100	0.9913	0.0026
110	0.9912	0.0026
120	0.9917	0.0026
130	0.9916	0.0027
140	0.9916	0.0028
150	0.9914	0.0028
160	0.9913	0.0028

6.2.2 Verifica del throughput

Variando il carico di lavoro da 1 req/s a 12 req/s e impostando $SI_{max} = \infty$, in modo da misurare solo il web server, si può verificare che il throughput del sistema non superi mai il suo limite teorico, che nel caso considerato è uguale a $\min(\mu_{web}, \lambda_{web})$. Tuttavia, come evidenziato in Figura 6, il throughput misurato $X_{measured} = N_{completed}/T_{sim}$ supera il limite teorico in condizioni di sovraccarico ($\rho > 1$). L'analisi di questa anomalia ha identificato la causa nella gestione della fine della simulazione (stopping

condition).

Il protocollo di simulazione adottato, per mantenere stretta aderenza al caso di studio di riferimento, interrompe la simulazione a $T_{stop} = 1200s$ senza attendere lo svuotamento delle code. Questa scelta metodologica introduce un bias sistematico quando il sistema è instabile:

- Il calcolo del throughput considera solo i job completati entro T_{stop} .
- A causa della distribuzione iperesponenziale dei servizi ($cv=4$), i job con richieste di servizio brevi tendono a completare rapidamente, mentre i job con richieste lunghe ("heavy tail") tendono ad accumularsi in coda e rimanere incompleti allo scadere del tempo.

Di conseguenza, il rapporto tra job completati e tempo trascorso risulta sovrastimato rispetto alla capacità reale del server di smaltire lavoro medio.

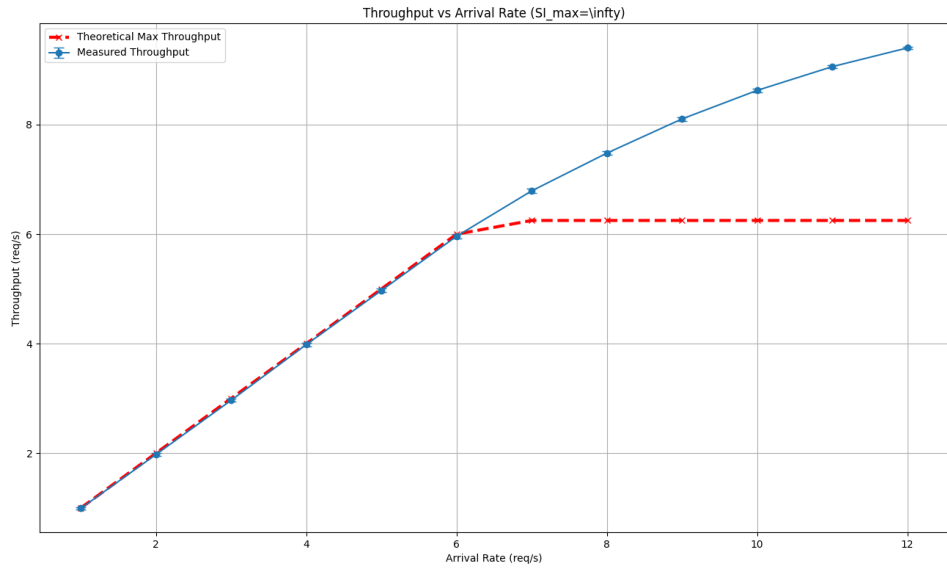


Figure 6: Throughput misurato che supera il valore teorico

Sebbene si sia scelto di mantenere questo approccio per garantire la comparabilità dei risultati con il report originale, è fondamentale notare che il simulatore si comporta correttamente se si considera la relazione fondamentale:

$$X_{calculated} = U_{web} \cdot \mu_{web}$$

Come mostrato in Figura 7 e nella Tabella 2, il throughput calcolato a partire dall'utilizzazione ($X_{calculated}$) rispetta perfettamente il limite teorico, confermando che la logica interna del simulatore è corretta e che la discrepanza in Figura 6 è un artefatto della metrica di misurazione scelta dallo studio di riferimento, non un errore implementativo. È bene notare che questa discrepanza emerge solo in condizioni di sovraccarico estremo come questi ($SI_{max} = \infty$ e $\lambda > \mu$) e non influenza i risultati principali dello studio, che si concentrano su scenari stabili o moderatamente carichi.

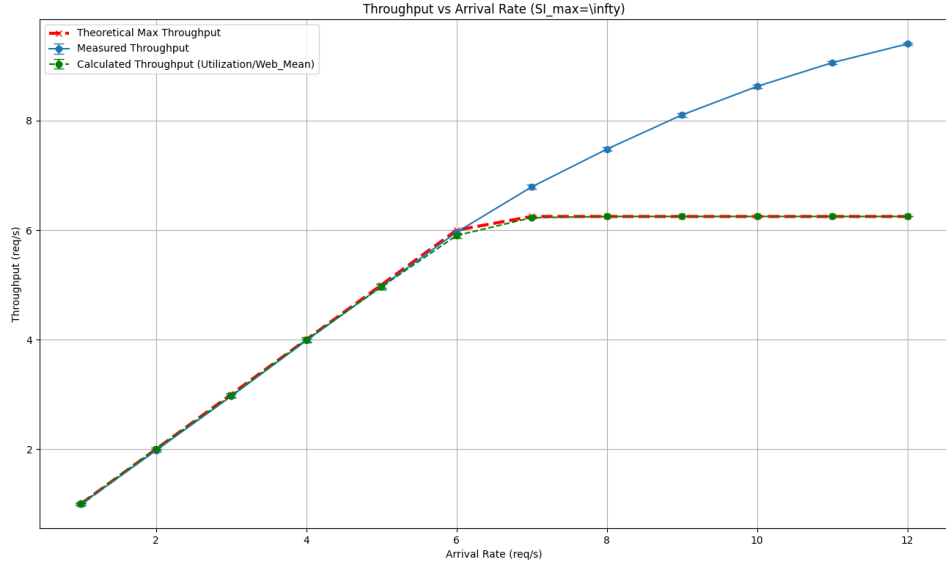


Figure 7: Throughput calcolato a partire dall'utilizzazione

Table 2: Verifica del Throughput

λ	$X_{measured}$	$X_{measured}$ CI95	$X_{theoretical}$	$X_{calculated}$	$X_{calculated}$ CI95
1	0.9923	0.0158	1.0000	1.0056	0.0223
2	1.9767	0.0226	2.0000	2.0027	0.0297
3	2.9658	0.0263	3.0000	2.9825	0.0388
4	3.9840	0.0311	4.0000	3.9999	0.0428
5	4.9719	0.0349	5.0000	4.9650	0.0511
6	5.9621	0.0396	6.0000	5.9011	0.0453
7	6.7908	0.0381	6.2500	6.2283	0.0123
8	7.4797	0.0361	6.2500	6.2491	0.0008
9	8.1046	0.0348	6.2500	6.2498	0.0002
10	8.6257	0.0311	6.2500	6.2501	0.0001
11	9.0594	0.0300	6.2500	6.2499	0.0001
12	9.4018	0.0267	6.2500	6.2499	0.0001

7 Validazione

Dato che in questo studio non sono disponibili dati reali con cui confrontare i risultati ottenuti dalla simulazione, la validazione del modello simulativo confronterà i risultati ottenuti con quelli ottenuti dallo studio di riferimento nella prossima sezione del report.

8 Risultati

Gli esperimenti condotti mirano a raggiungere gli obiettivi delineati nella sezione 2.

8.1 Obiettivo 1

Come si può evincere dalla figura 8, il tempo di risposta medio $E[R]$ rimane al di sotto della soglia di 8 secondi per valori di SI_{max} fino a circa 90. Esattamente come suggerito nello studio di riferimento, il valore ottimale di SI_{max} che massimizza l'utilizzo del web server mantenendo il rispetto dello SLA si attesta intorno a 80-90. Considerando che con l'intervallo di confidenza con 90 si sfora, il valore più prudente da adottare per SI_{max} tra i due risulta essere 80.

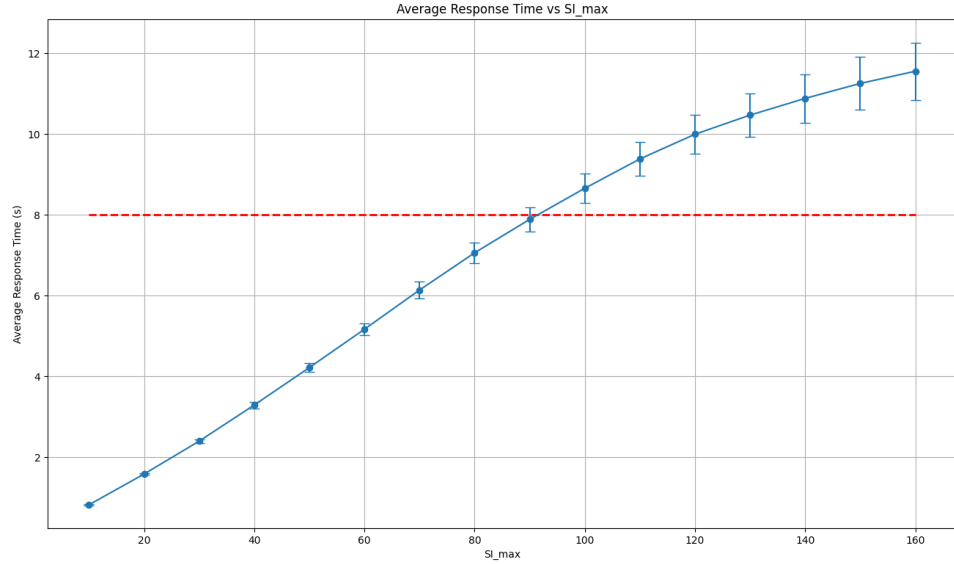


Figure 8: Tempo di risposta medio in funzione di SI_{max}

Table 3: Tempo di risposta medio in funzione di SI_{max}

SI_{max}	$E[T_s]$	$E[T_s]$ CI95
10	0.8174	0.0092
20	1.5783	0.0225
30	2.3941	0.0464
40	3.2887	0.0794
50	4.2172	0.1102
60	5.1665	0.1543
70	6.1383	0.2045
80	7.0609	0.2506
90	7.8901	0.2986
100	8.6595	0.3719
110	9.3868	0.4216
120	9.9987	0.4868
130	10.4731	0.5419
140	10.8868	0.5989
150	11.2576	0.6543
160	11.5601	0.7126

8.2 Obiettivo 2

Come mostrato in figura 9, accade lo stesso fenomeno descritto nello studio di riferimento: al crescere del tasso di arrivo, il tempo di risposta medio aumenta all'aumentare degli arrivi, fino a che non si arriva a superare la soglia SI_{max} . A questo punto, il tempo di risposta comincia a calare fino a che anche lo spike server si satura e il tempo di risposta ricomincia a salire per valori superiori a 11.

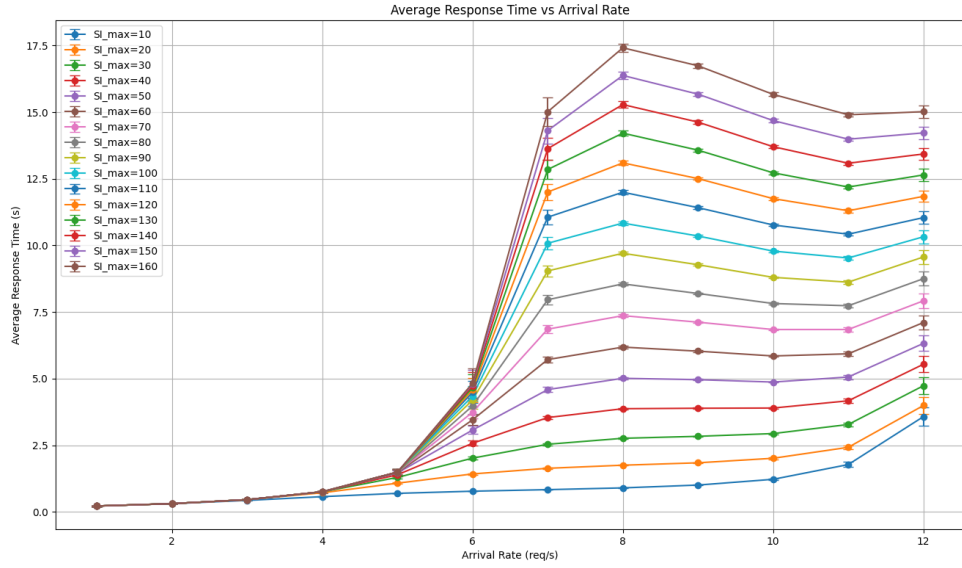


Figure 9: Tempo di risposta medio in funzione del tasso di arrivo e di SI_{max}

Table 4: Tempo di risposta medio in funzione del tasso di arrivo e di SI_{max}

SI max	Arrival Rate	Total Response Time Mean	Total Response Time CI95
10	1.00	0.2210	0.0054
10	2.00	0.3092	0.0074
10	3.00	0.4312	0.0086
10	4.00	0.5726	0.0094
10	5.00	0.6964	0.0093
10	6.00	0.7766	0.0099
10	7.00	0.8354	0.0088
10	8.00	0.9009	0.0101
10	9.00	1.0068	0.0184
10	10.00	1.2228	0.0363
10	11.00	1.7764	0.0994
10	12.00	3.5733	0.3537
20	1.00	0.2210	0.0054
20	2.00	0.3103	0.0073
20	3.00	0.4599	0.0122
20	4.00	0.7168	0.0187
20	5.00	1.0781	0.0269
20	6.00	1.4247	0.0264
20	7.00	1.6363	0.0212
20	8.00	1.7521	0.0200
20	9.00	1.8421	0.0176
20	10.00	2.0123	0.0333
20	11.00	2.4231	0.0726
20	12.00	3.9912	0.3099
30	1.00	0.2210	0.0054
30	2.00	0.3103	0.0073
30	3.00	0.4592	0.0123
30	4.00	0.7511	0.0263
30	5.00	1.2941	0.0478
30	6.00	2.0178	0.0532
30	7.00	2.5375	0.0383
30	8.00	2.7620	0.0280
30	9.00	2.8358	0.0211
30	10.00	2.9386	0.0323
30	11.00	3.2799	0.0724
30	12.00	4.7331	0.3056
40	1.00	0.2210	0.0054
40	2.00	0.3103	0.0073
40	3.00	0.4592	0.0123
40	4.00	0.7562	0.0288
40	5.00	1.3933	0.0676
40	6.00	2.5757	0.0961
40	7.00	3.5394	0.0618
40	8.00	3.8720	0.0359
40	9.00	3.8894	0.0256
40	10.00	3.8973	0.0250
40	11.00	4.1682	0.0757
40	12.00	5.5398	0.2982
50	1.00	0.2210	0.0054
50	2.00	0.3103	0.0073
50	3.00	0.4592	0.0123
50	4.00	0.7568	0.0302
50	5.00	1.4611	0.0859
50	6.00	3.0724	0.1533
50	7.00	4.5896	0.0939
50	8.00	5.0136	0.0403
50	9.00	4.9609	0.0283
50	10.00	4.8720	0.0257
50	11.00	5.0615	0.0810
50	12.00	6.3221	0.2870

8.3 Obiettivo 3

Nell'obiettivo 3 si va ad duplicare la potenza dello spike server, per verificare se questo evita il saturamento dello spike server e quindi il successivo aumento del tempo di risposta medio. Come mostrato in figura 10, lo spike server non arriva più a saturarsi e il tempo di risposta continua a diminuire all'aumentare del tasso di arrivo.

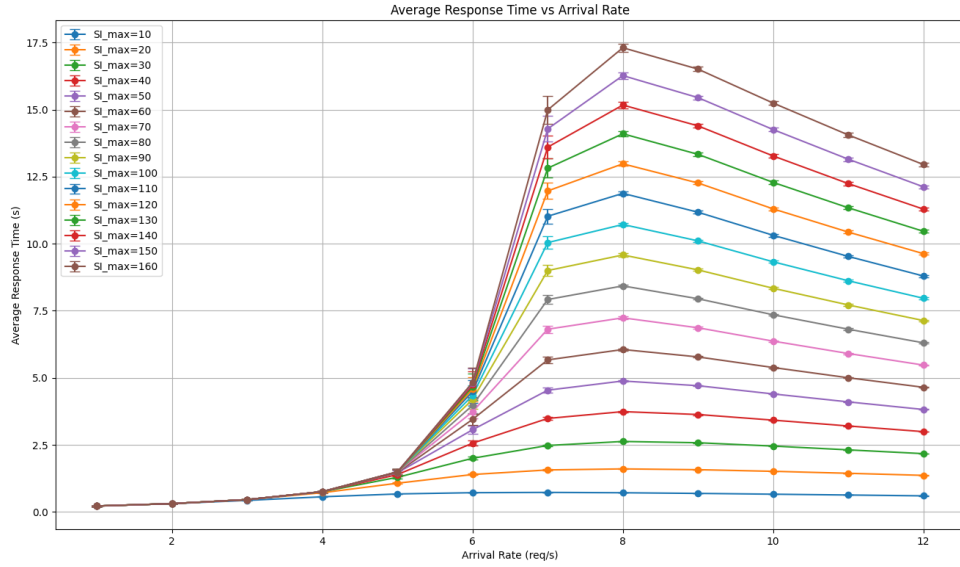


Figure 10: Tempo di risposta medio in funzione del tasso di arrivo e di SI_{max} con spike server potenziato

Table 5: Tempo di risposta medio in funzione del tasso di arrivo e di SI_{max} con spike server potenziato

SI max	Arrival Rate	Total Response Time Mean	Total Response Time CI95
10	1.00	0.2210	0.0054
10	2.00	0.3091	0.0074
10	3.00	0.4301	0.0085
10	4.00	0.5646	0.0090
10	5.00	0.6704	0.0084
10	6.00	0.7181	0.0077
10	7.00	0.7281	0.0057
10	8.00	0.7158	0.0046
10	9.00	0.6905	0.0041
10	10.00	0.6618	0.0032
10	11.00	0.6308	0.0027
10	12.00	0.5999	0.0026
20	1.00	0.2210	0.0054
20	2.00	0.3103	0.0073
20	3.00	0.4598	0.0122
20	4.00	0.7161	0.0186
20	5.00	1.0709	0.0264
20	6.00	1.3961	0.0250
20	7.00	1.5663	0.0187
20	8.00	1.6048	0.0141
20	9.00	1.5752	0.0093
20	10.00	1.5147	0.0072
20	11.00	1.4411	0.0055
20	12.00	1.3640	0.0053
30	1.00	0.2210	0.0054
30	2.00	0.3103	0.0073
30	3.00	0.4592	0.0123
30	4.00	0.7510	0.0263
30	5.00	1.2917	0.0473
30	6.00	2.0012	0.0522
30	7.00	2.4779	0.0359
30	8.00	2.6277	0.0240
30	9.00	2.5801	0.0159
30	10.00	2.4555	0.0101
30	11.00	2.3147	0.0091
30	12.00	2.1737	0.0088
40	1.00	0.2210	0.0054
40	2.00	0.3103	0.0073
40	3.00	0.4592	0.0123
40	4.00	0.7562	0.0288
40	5.00	1.3925	0.0674
40	6.00	2.5644	0.0950
40	7.00	3.4861	0.0594
40	8.00	3.7402	0.0323
40	9.00	3.6328	0.0223
40	10.00	3.4226	0.0131
40	11.00	3.2056	0.0130
40	12.00	2.9940	0.0126
50	1.00	0.2210	0.0054
50	2.00	0.3103	0.0073
50	3.00	0.4592	0.0123
50	4.00	0.7568	0.0302
50	5.00	1.4608	0.0858
50	6.00	3.0650	0.1522
50	7.00	4.5398	0.0914
50	8.00	4.8830	0.0371
50	9.00	4.7069	0.0263
50	10.00	4.3996	0.0175
50	11.00	4.1035	0.0174
50	12.00	3.8203	0.0164