

Throughput and Heavy-traffic Optimality of General Load Balancing Algorithm in Cloud Networks

Xingyu Zhou

The Ohio State University

zhou.2055@osu.edu

October 24, 2016

1 Introduction and Motivation

- Why do we need an effective and fast cloud?
- Load Balancing and Previous Works

2 Throughput and Heavy-traffic optimality of General Load Balancing

- Model, Challenges, Contributions and Key insights
- Methodology and Sufficient Conditions
- Homogeneous Servers
- Heterogeneous Servers
- Conclusions

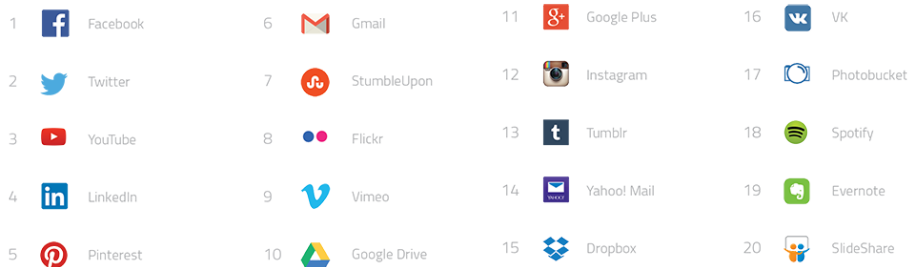
1 Introduction and Motivation

- Why do we need an effective and fast cloud?
- Load Balancing and Previous Works

2 Throughput and Heavy-traffic optimality of General Load Balancing

- Model, Challenges, Contributions and Key insights
- Methodology and Sufficient Conditions
- Homogeneous Servers
- Heterogeneous Servers
- Conclusions

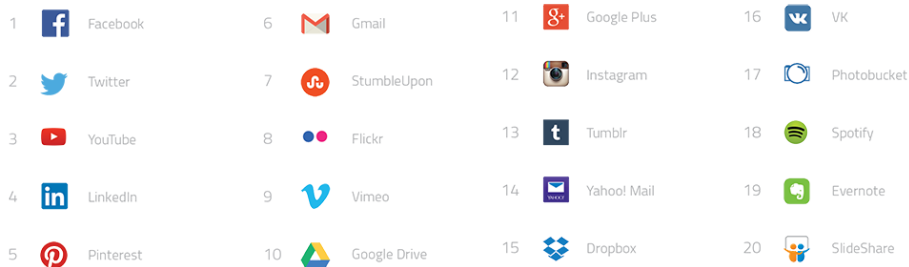
Cloud is everywhere in life



Source: <https://www.skyhighnetworks.com/>

Figure: Top 20 consumer cloud services in 2016

Cloud is everywhere in life



Source: <https://www.skyhighnetworks.com/>

Figure: Top 20 consumer cloud services in 2016

Question: What makes a good cloud service?

An Effective and Fast Cloud

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.
 - ▶ It determines how many units of information an application can process in a period of time.

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.
 - ▶ It determines how many units of information an application can process in a period of time.
 - ▶ High throughput means the application can handle more number of concurrent users.

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.
 - ▶ It determines how many units of information an application can process in a period of time.
 - ▶ High throughput means the application can handle more number of concurrent users.
- **Fast:** Low delay and response time

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.
 - ▶ It determines how many units of information an application can process in a period of time.
 - ▶ High throughput means the application can handle more number of concurrent users.
- **Fast:** Low delay and response time
 - ▶ It determines how much time an application takes to return a request to users.

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.
 - ▶ It determines how many units of information an application can process in a period of time.
 - ▶ High throughput means the application can handle more number of concurrent users.
- **Fast:** Low delay and response time
 - ▶ It determines how much time an application takes to return a request to users.
 - ▶ Low response time means a good user experience and also great profit for companies.

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.
 - ▶ It determines how many units of information an application can process in a period of time.
 - ▶ High throughput means the application can handle more number of concurrent users.
- **Fast:** Low delay and response time
 - ▶ It determines how much time an application takes to return a request to users.
 - ▶ Low response time means a good user experience and also great profit for companies.
 - ▶ Examples¹:
 - ★ Google: an half second increase in loading time drops traffic by 20%
 - ★ Amazon: every 100ms of latency costs them 1% in sales

¹Credit: Prof. Ness Shroff Mobihoc 2015 Talk

An Effective and Fast Cloud

- **Effective:** High utilization and throughput.
 - ▶ It determines how many units of information an application can process in a period of time.
 - ▶ High throughput means the application can handle more number of concurrent users.
- **Fast:** Low delay and response time
 - ▶ It determines how much time an application takes to return a request to users.
 - ▶ Low response time means a good user experience and also great profit for companies.
 - ▶ Examples¹:
 - ★ Google: an half second increase in loading time drops traffic by 20%
 - ★ Amazon: every 100ms of latency costs them 1% in sales

Question: How do we design an effective and fast cloud system?

¹Credit: Prof. Ness Shroff Mobihoc 2015 Talk

1 Introduction and Motivation

- Why do we need an effective and fast cloud?
- Load Balancing and Previous Works

2 Throughput and Heavy-traffic optimality of General Load Balancing

- Model, Challenges, Contributions and Key insights
- Methodology and Sufficient Conditions
- Homogeneous Servers
- Heterogeneous Servers
- Conclusions

Load Balancing in Cloud

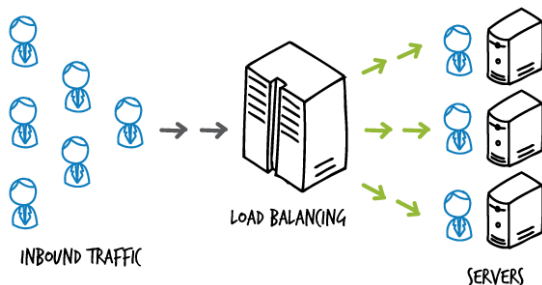


Figure: A typical model in cloud system

- **Load balancing:** Choose the right server(s) when requests coming.

Load Balancing in Cloud

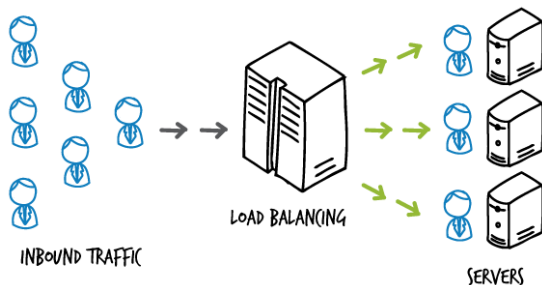


Figure: A typical model in cloud system

- **Load balancing:** Choose the right server(s) when requests coming.
 - ▶ It is the key to optimize resource use, maximize throughput, reduce response time in cloud system.

Load Balancing in Cloud

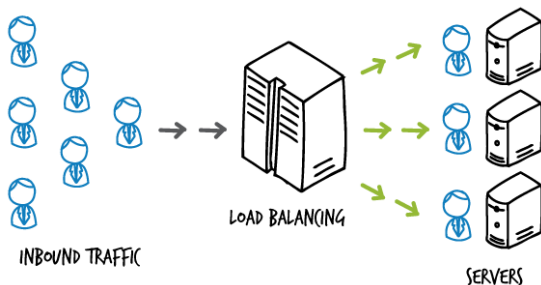


Figure: A typical model in cloud system

- **Load balancing:** Choose the right server(s) when requests coming.
 - ▶ It is the key to optimize resource use, maximize throughput, reduce response time in cloud system.
 - ▶ It becomes more and more critical due to explosive increase in the number of **servers** and **traffic** in cloud system.

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.
 - ▶ **Task arrival:** average-delay of Join-the-shortest-queue (JSQ), Power-of-d and Join-the-idle-queue (JIQ)[Mitzenmacher'96,'16, Lu'11]

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.
 - ▶ **Task arrival:** average-delay of Join-the-shortest-queue (JSQ), Power-of-d and Join-the-idle-queue (JIQ)[Mitzenmacher'96,'16, Lu'11]
 - ▶ **Job arrival:** Batch sampling is introduced to reduce sampling overhead. [Srikant'15, Ousterhout'13]

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.
 - ▶ **Task arrival:** average-delay of Join-the-shortest-queue (JSQ), Power-of-d and Join-the-idle-queue (JIQ)[Mitzenmacher'96,'16, Lu'11]
 - ▶ **Job arrival:** Batch sampling is introduced to reduce sampling overhead. [Srikant'15, Ousterhout'13]
- **Heavy traffic limit:** the utilization of server ρ goes to 1.

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.
 - ▶ **Task arrival:** average-delay of Join-the-shortest-queue (JSQ), Power-of-d and Join-the-idle-queue (JIQ)[Mitzenmacher'96,'16, Lu'11]
 - ▶ **Job arrival:** Batch sampling is introduced to reduce sampling overhead. Srikant'15, Ousterhout'13]
- **Heavy traffic limit:** the utilization of server ρ goes to 1.
 - ▶ **Brownian motion condition:** Sample-path heavy-traffic optimality in scaled time over a finite time interval of JSQ, power-of-d[Foschini et al,'78, Whiting'16]

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.
 - ▶ **Task arrival:** average-delay of Join-the-shortest-queue (JSQ), Power-of-d and Join-the-idle-queue (JIQ)[Mitzenmacher'96,'16, Lu'11]
 - ▶ **Job arrival:** Batch sampling is introduced to reduce sampling overhead. [Srikant'15, Ousterhout'13]
- **Heavy traffic limit:** the utilization of server ρ goes to 1.
 - ▶ **Brownian motion condition:** Sample-path heavy-traffic optimality in scaled time over a finite time interval of JSQ, power-of-d[Foschini et al,'78, Whiting'16]
 - ▶ **Lyapunov Drift condition:** First moment heavy-traffic optimality of JSQ, power-of-d. [Erylimaz'13, Maguluri'12]

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.
 - ▶ **Task arrival:** average-delay of Join-the-shortest-queue (JSQ), Power-of-d and Join-the-idle-queue (JIQ)[Mitzenmacher'96,'16, Lu'11]
 - ▶ **Job arrival:** Batch sampling is introduced to reduce sampling overhead. Srikant'15, Ousterhout'13]
- **Heavy traffic limit:** the utilization of server ρ goes to 1.
 - ▶ **Brownian motion condition:** Sample-path heavy-traffic optimality in scaled time over a finite time interval of JSQ, power-of-d[Foschini et al,'78, Whiting'16]
 - ▶ **Lyapunov Drift condition:** First moment heavy-traffic optimality of JSQ, power-of-d. [Erylimaz'13, Maguluri'12]

Questions: Can we generalize existing load balancing algorithm?

Previous Works

Throughput analysis is relatively easy, however, delay is really difficult for exact analysis. Two alternatives are asymptotic analysis:

- **Large system limit:** the number of servers N goes to infinity.
 - ▶ **Task arrival:** average-delay of Join-the-shortest-queue (JSQ), Power-of-d and Join-the-idle-queue (JIQ)[Mitzenmacher'96,'16, Lu'11]
 - ▶ **Job arrival:** Batch sampling is introduced to reduce sampling overhead. [Srikant'15, Ousterhout'13]
- **Heavy traffic limit:** the utilization of server ρ goes to 1.
 - ▶ **Brownian motion condition:** Sample-path heavy-traffic optimality in scaled time over a finite time interval of JSQ, power-of-d[Foschini et al,'78, Whiting'16]
 - ▶ **Lyapunov Drift condition:** First moment heavy-traffic optimality of JSQ, power-of-d. [Erylimaz'13, Maguluri'12]

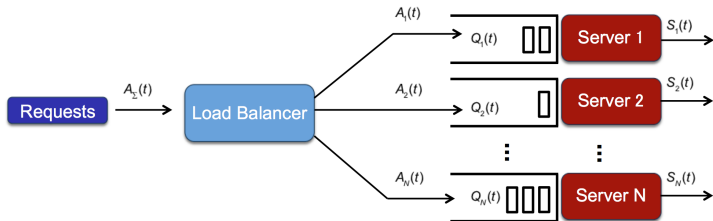
Questions: Can we generalize existing load balancing algorithm?

Question: Do we really need sampling for each arrival?

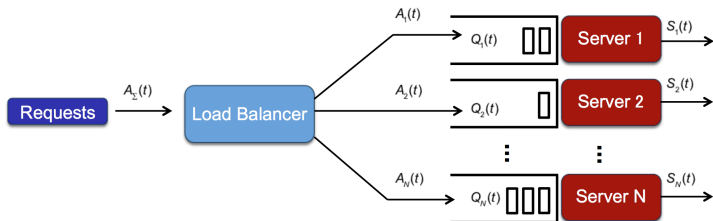
- 1 Introduction and Motivation
 - Why do we need an effective and fast cloud?
 - Load Balancing and Previous Works
- 2 Throughput and Heavy-traffic optimality of General Load Balancing
 - Model, Challenges, Contributions and Key insights
 - Methodology and Sufficient Conditions
 - Homogeneous Servers
 - Heterogeneous Servers
 - Conclusions

- 1 Introduction and Motivation
 - Why do we need an effective and fast cloud?
 - Load Balancing and Previous Works
- 2 Throughput and Heavy-traffic optimality of General Load Balancing
 - Model, Challenges, Contributions and Key insights
 - Methodology and Sufficient Conditions
 - Homogeneous Servers
 - Heterogeneous Servers
 - Conclusions

System Model and Assumptions

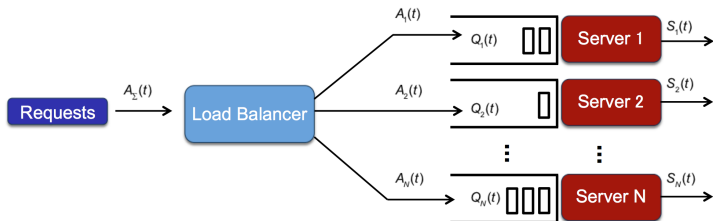


System Model and Assumptions



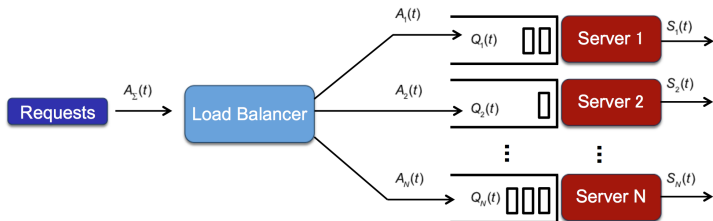
- The exogenous arrival process $A_{\Sigma}(t)$ with mean rate λ_{Σ} is i.i.d and also independent of service, queue length, and routing decision. $A_{\Sigma}(t) \leq A_{\max} < \infty$ for all $t \geq 0$

System Model and Assumptions



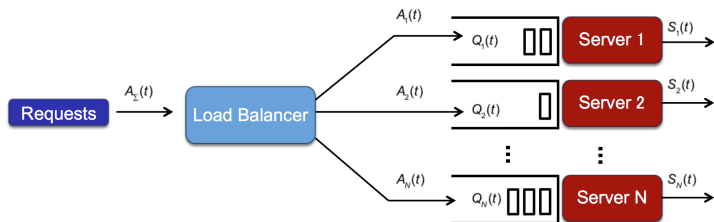
- The exogenous arrival process $A_\Sigma(t)$ with mean rate λ_Σ is i.i.d and also independent of service, queue length, and routing decision. $A_\Sigma(t) \leq A_{\max} < \infty$ for all $t \geq 0$
- Each service process $S_n(t)$ with mean rate μ_n is i.i.d and also independent of other service, arrival, queue length, and routing decision. $S_n(t) \leq S_{\max} < \infty$ for all $t \geq 0$ and all $n \in \{1, 2, \dots, N\}$, and $S_\Sigma(t) = \sum S_n(t)$ is the hypothetical total service with mean $\mu_\Sigma = \sum \mu_n$. $\epsilon = \mu_\Sigma - \lambda_\Sigma$ indicates how close the arrival rate to the boundary.

System Model and Assumptions



- The exogenous arrival process $A_\Sigma(t)$ with mean rate λ_Σ is i.i.d and also independent of service, queue length, and routing decision. $A_\Sigma(t) \leq A_{\max} < \infty$ for all $t \geq 0$
- Each service process $S_n(t)$ with mean rate μ_n is i.i.d and also independent of other service, arrival, queue length, and routing decision. $S_n(t) \leq S_{\max} < \infty$ for all $t \geq 0$ and all $n \in \{1, 2, \dots, N\}$, and $S_\Sigma(t) = \sum S_n(t)$ is the hypothetical total service with mean $\mu_\Sigma = \sum \mu_n$. $\epsilon = \mu_\Sigma - \lambda_\Sigma$ indicates how close the arrival rate to the boundary.
- Requests cannot be removed from server to server.

System Model and Assumptions



- The exogenous arrival process $A_\Sigma(t)$ with mean rate λ_Σ is i.i.d and also independent of service, queue length, and routing decision. $A_\Sigma(t) \leq A_{\max} < \infty$ for all $t \geq 0$
- Each service process $S_n(t)$ with mean rate μ_n is i.i.d and also independent of other service, arrival, queue length, and routing decision. $S_n(t) \leq S_{\max} < \infty$ for all $t \geq 0$ and all $n \in \{1, 2, \dots, N\}$, and $S_\Sigma(t) = \sum S_n(t)$ is the hypothetical total service with mean $\mu_\Sigma = \sum \mu_n$. $\epsilon = \mu_\Sigma - \lambda_\Sigma$ indicates how close the arrival rate to the boundary.
- Requests cannot be removed from server to server.
- Queue length dynamic is

$$\begin{aligned} Q_n(t+1) &= [Q_n(t) + A_n(t) - S_n(t)]^+ \\ &= Q_n(t) + A_n(t) - S_n(t) + U_n(t). \end{aligned} \tag{1}$$

Definitions of Throughput and Heavy-traffic Optimal

Definitions of Throughput and Heavy-traffic Optimal

Definition (Throughput Optimality)

In general, a load balancing policy is said to be throughput optimal if it can stabilize any set of arrival rates which can be stabilized by another policy. In our model, a load balancing policy is throughput optimal if it can stabilize any exogenous rate $\lambda_{\Sigma} < \mu_{\Sigma}$.

Definitions of Throughput and Heavy-traffic Optimal

Definition (Throughput Optimality)

In general, a load balancing policy is said to be throughput optimal if it can stabilize any set of arrival rates which can be stabilized by another policy. In our model, a load balancing policy is throughput optimal if it can stabilize any exogenous rate $\lambda_{\Sigma} < \mu_{\Sigma}$.

Definition (Equivalent Single Queue)

A single queue process $q(t)$ is said to be the equivalent single queue of the load balancing system if $a(t) = A_{\Sigma}(t)$ and $s(t) = S_{\Sigma}(t)$. Clearly, $q(t) \leq_{st} \sum Q_n(t)$

Definitions of Throughput and Heavy-traffic Optimal

Definition (Throughput Optimality)

In general, a load balancing policy is said to be throughput optimal if it can stabilize any set of arrival rates which can be stabilized by another policy. In our model, a load balancing policy is throughput optimal if it can stabilize any exogenous rate $\lambda_{\Sigma} < \mu_{\Sigma}$.

Definition (Equivalent Single Queue)

A single queue process $q(t)$ is said to be the equivalent single queue of the load balancing system if $a(t) = A_{\Sigma}(t)$ and $s(t) = S_{\Sigma}(t)$. Clearly, $q(t) \leq_{st} \sum Q_n(t)$

Definition (Heavy-traffic Optimality)

A load balancing policy is said to be heavy-traffic optimal if the expected steady-state sum queue length is asymptotically the same as the equivalent single queue when the arrival rate approaches to the capacity boundary.

Main Contributions

Main Contributions

- We give sufficient conditions for a policy to be throughput and heavy-traffic optimal under certain assumptions.

Main Contributions

- We give sufficient conditions for a policy to be throughput and heavy-traffic optimal under certain assumptions.
- We introduce a new load balancing policy, called **Join-Below-Average (JBA)**, which is both throughput and heavy-traffic optimal. Instead of joining the shortest queue, it needs only to randomly selects a queue with queue length below average to join.

Main Contributions

- We give sufficient conditions for a policy to be throughput and heavy-traffic optimal under certain assumptions.
- We introduce a new load balancing policy, called **Join-Below-Average (JBA)**, which is both throughput and heavy-traffic optimal. Instead of joining the shortest queue, it needs only to randomly select a queue with queue length below average to join.
 - ▶ it is often easy to obtain the sum queue length of a system, and hence the average queue length.

Main Contributions

- We give sufficient conditions for a policy to be throughput and heavy-traffic optimal under certain assumptions.
- We introduce a new load balancing policy, called **Join-Below-Average (JBA)**, which is both throughput and heavy-traffic optimal. Instead of joining the shortest queue, it needs only to randomly select a queue with queue length below average to join.
 - ▶ it is often easy to obtain the sum queue length of a system, and hence the average queue length.
 - ▶ sometimes, it is prohibited for a request to join the shortest queue due to data locality.

Main Contributions

- We give sufficient conditions for a policy to be throughput and heavy-traffic optimal under certain assumptions.
- We introduce a new load balancing policy, called [Join-Below-Average \(JBA\)](#), which is both throughput and heavy-traffic optimal. Instead of joining the shortest queue, it needs only to randomly select a queue with queue length below average to join.
 - ▶ it is often easy to obtain the sum queue length of a system, and hence the average queue length.
 - ▶ sometimes, it is prohibited for a request to join the shortest queue due to data locality.
- We show that a class of load balancing policies is throughput and heavy-traffic optimal.

Main Contributions

- We give sufficient conditions for a policy to be throughput and heavy-traffic optimal under certain assumptions.
- We introduce a new load balancing policy, called **Join-Below-Average (JBA)**, which is both throughput and heavy-traffic optimal. Instead of joining the shortest queue, it needs only to randomly select a queue with queue length below average to join.
 - ▶ it is often easy to obtain the sum queue length of a system, and hence the average queue length.
 - ▶ sometimes, it is prohibited for a request to join the shortest queue due to data locality.
- We show that a class of load balancing policies is throughput and heavy-traffic optimal.
 - ▶ we need only use any 'good' load balancing policy once every T time slots, for any finite T , and random routing in other time slots to achieve throughput and heavy-traffic optimality

Main Contributions

- We give sufficient conditions for a policy to be throughput and heavy-traffic optimal under certain assumptions.
- We introduce a new load balancing policy, called **Join-Below-Average (JBA)**, which is both throughput and heavy-traffic optimal. Instead of joining the shortest queue, it needs only to randomly selects a queue with queue length below average to join.
 - ▶ it is often easy to obtain the sum queue length of a system, and hence the average queue length.
 - ▶ sometimes, it is prohibited for a request to join the shortest queue due to data locality.
- We show that a class of load balancing policies is throughput and heavy-traffic optimal.
 - ▶ we need only use any 'good' load balancing policy once every T time slots, for any finite T , and random routing in other time slots to achieve throughput and heavy-traffic optimality
 - ▶ The 'good' policy can be any one time-slot sampling heavy traffic optimal policy, such as JBA, Power-of-d, and JSQ depending on different situations.

Key Insights

Key Insights

- We know that purely random routing is not heavy-traffic optimal, but, it is not so bad as we might think.

Key Insights

- We know that purely random routing is not heavy-traffic optimal, but, it is not so bad as we might think.
 - ▶ For homogeneous servers, it actually does no harm (of course no good) to heavy-traffic optimality.

Key Insights

- We know that purely random routing is not heavy-traffic optimal, but, it is not so bad as we might think.
 - ▶ For homogeneous servers, it actually does no harm (of course no good) to heavy-traffic optimality.
 - ▶ For heterogeneous servers, the harm it does will tend to zero as traffic becomes more heavy.

Key Insights

- We know that purely random routing is not heavy-traffic optimal, but, it is not so bad as we might think.
 - ▶ For homogeneous servers, it actually does no harm (of course no good) to heavy-traffic optimality.
 - ▶ For heterogeneous servers, the harm it does will tend to zero as traffic becomes more heavy.
- As the traffic gets heavier, it actually becomes more and more easy to achieve heavy-traffic optimality.

Key Insights

- We know that purely random routing is not heavy-traffic optimal, but, it is not so bad as we might think.
 - ▶ For homogeneous servers, it actually does no harm (of course no good) to heavy-traffic optimality.
 - ▶ For heterogeneous servers, the harm it does will tend to zero as traffic becomes more heavy.
- As the traffic gets heavier, it actually becomes more and more easy to achieve heavy-traffic optimality.
 - ▶ The upper bound on the sampling interval T gets larger when the traffic becomes heavier.

Key Insights

- We know that purely random routing is not heavy-traffic optimal, but, it is not so bad as we might think.
 - ▶ For homogeneous servers, it actually does no harm (of course no good) to heavy-traffic optimality.
 - ▶ For heterogeneous servers, the harm it does will tend to zero as traffic becomes more heavy.
- As the traffic gets heavier, it actually becomes more and more easy to achieve heavy-traffic optimality.
 - ▶ The upper bound on the sampling interval T gets larger when the traffic becomes heavier.
 - ▶ In some sense, load balancing becomes easier (kind of counterintuitive, but can be explained) when traffic becomes heavier.

1 Introduction and Motivation

- Why do we need an effective and fast cloud?
- Load Balancing and Previous Works

2 Throughput and Heavy-traffic optimality of General Load Balancing

- Model, Challenges, Contributions and Key insights
- **Methodology and Sufficient Conditions**
- Homogeneous Servers
- Heterogeneous Servers
- Conclusions

Main Ideas

- Choose the Lyapunov function $Z(\mathbf{Q}) \triangleq \|\mathbf{Q}\|_1^2 = \left(\sum_{n=1}^N Q_n\right)^2$, and set the mean drift to zero in steady-state to get bounds.

Main Ideas

- Choose the Lyapunov function $Z(\mathbf{Q}) \triangleq \|\mathbf{Q}\|_1^2 = \left(\sum_{n=1}^N Q_n\right)^2$, and set the mean drift to zero in steady-state to get bounds.
 - ▶ (R1) It requires the existence of a stationary distribution $\bar{\mathbf{Q}}$.

Main Ideas

- Choose the Lyapunov function $Z(\mathbf{Q}) \triangleq \|\mathbf{Q}\|_1^2 = \left(\sum_{n=1}^N Q_n\right)^2$, and set the mean drift to zero in steady-state to get bounds.
 - ▶ (R1) It requires the existence of a stationary distribution $\bar{\mathbf{Q}}$.
 - ▶ (R2) It requires bounded moments of the stationary distribution to avoid the situation $\infty - \infty$, i.e., $\mathbb{E} \left[\|\bar{\mathbf{Q}}\|^2 \right] \leq M$.

Main Ideas

- Choose the Lyapunov function $Z(\mathbf{Q}) \triangleq \|\mathbf{Q}\|_1^2 = \left(\sum_{n=1}^N Q_n\right)^2$, and set the mean drift to zero in steady-state to get bounds.
 - ▶ (R1) It requires the existence of a stationary distribution $\bar{\mathbf{Q}}$.
 - ▶ (R2) It requires bounded moments of the stationary distribution to avoid the situation $\infty - \infty$, i.e., $\mathbb{E} \left[\|\bar{\mathbf{Q}}\|^2 \right] \leq M$.
- Assume (R1) and (R2) satisfy, we obtain the expected sum queue length in steady-state

$$\epsilon \mathbb{E} \left[\sum_{n=1}^N \bar{Q}_n^{(\epsilon)} \right] = \frac{\zeta^{(\epsilon)}}{2} + \mathbb{E} \left[\|\bar{\mathbf{Q}}(t+1)\|_1 \|\bar{\mathbf{U}}(t)\|_1 \right] - \mathbb{E} \left[\|\bar{\mathbf{U}}(t)\|_1^2 \right] \quad (2)$$

Main Ideas

- Choose the Lyapunov function $Z(\mathbf{Q}) \triangleq \|\mathbf{Q}\|_1^2 = \left(\sum_{n=1}^N Q_n\right)^2$, and set the mean drift to zero in steady-state to get bounds.
 - ▶ (R1) It requires the existence of a stationary distribution $\bar{\mathbf{Q}}$.
 - ▶ (R2) It requires bounded moments of the stationary distribution to avoid the situation $\infty - \infty$, i.e., $\mathbb{E} \left[\|\bar{\mathbf{Q}}\|^2 \right] \leq M$.
- Assume (R1) and (R2) satisfy, we obtain the expected sum queue length in steady-state

$$\epsilon \mathbb{E} \left[\sum_{n=1}^N \bar{Q}_n^{(\epsilon)} \right] = \frac{\zeta^{(\epsilon)}}{2} + \mathbb{E} \left[\|\bar{\mathbf{Q}}(t+1)\|_1 \|\bar{\mathbf{U}}(t)\|_1 \right] - \mathbb{E} \left[\|\bar{\mathbf{U}}(t)\|_1^2 \right] \quad (2)$$

- By letting $N = 1$, expected queue length for the corresponding **Equivalent Single Queue** is simply

$$\epsilon \mathbb{E} \left[q^{(\epsilon)} \right] = \frac{\zeta^{(\epsilon)}}{2} + \mathbb{E} [q(t+1)u(t)] - \mathbb{E} [u^2] \quad (3)$$

Main Ideas (Cont'd)

- For the equivalent single queue, by exploiting the most important equation $q(t+1)u(t) = 0$ and the fact $\mathbb{E}[u^2]$ is $o(\epsilon)$, we have

$$\begin{aligned}\epsilon \mathbb{E}[q^{(\epsilon)}] &= \frac{\zeta^{(\epsilon)}}{2} + \mathbb{E}[q(t+1)u(t)] - \mathbb{E}[u^2] \\ &= \frac{\zeta^{(\epsilon)}}{2} - o(\epsilon)\end{aligned}$$

Main Ideas (Cont'd)

- For the equivalent single queue, by exploiting the most important equation $q(t+1)u(t) = 0$ and the fact $\mathbb{E}[u^2]$ is $o(\epsilon)$, we have

$$\begin{aligned}\epsilon \mathbb{E}[q^{(\epsilon)}] &= \frac{\zeta^{(\epsilon)}}{2} + \mathbb{E}[q(t+1)u(t)] - \mathbb{E}[u^2] \\ &= \frac{\zeta^{(\epsilon)}}{2} - o(\epsilon)\end{aligned}$$

- For the N server, by exploiting the equation $Q_n(t+1)U_n(t) = 0$ and the fact $\mathbb{E}[\|\bar{\mathbf{U}}(t)\|_1^2]$ is $o(\epsilon)$, we have

$$\begin{aligned}\epsilon \mathbb{E}\left[\sum_{n=1}^N \bar{Q}_n^{(\epsilon)}\right] &= \frac{\zeta^{(\epsilon)}}{2} + \mathbb{E}[\|\bar{\mathbf{Q}}(t+1)\|_1 \|\bar{\mathbf{U}}(t)\|_1] - \mathbb{E}[\|\bar{\mathbf{U}}(t)\|_1^2] \\ &= \frac{\zeta^{(\epsilon)}}{2} + N \mathbb{E}[\langle \bar{\mathbf{U}}, -\bar{\mathbf{Q}}_{\perp}(t+1) \rangle] - o(\epsilon) \\ &\leq \frac{\zeta^{(\epsilon)}}{2} + \sqrt{\mathbb{E}[\|\bar{\mathbf{Q}}_{\perp}\|^2]} o(\epsilon) - o(\epsilon)\end{aligned}$$

in which $\bar{\mathbf{Q}}_{\perp}$ is the perpendicular component of $\bar{\mathbf{Q}}$ with respect to $\mathbf{c} = \frac{1}{\sqrt{N}} \mathbf{1}$

Sufficient Conditions for Throughput Heavy-traffic Optimal

- (R1) If $\bar{\mathbf{Q}}$ exists, i.e., the underline Markov chain is positive recurrent.
- (R2) If $\mathbb{E} \left[\|\bar{\mathbf{Q}}\|^2 \right] \leq M$, i.e., the second moments is bounded by a constant.
- (R3) If $\mathbb{E} \left[\|\bar{\mathbf{Q}}_{\perp}\|^2 \right] \leq K$, which is independent of ϵ . This is often called steady-state collapse, which indicates that under heavy-traffic, queue-length vector concentrates around the line \mathbf{c} .

Sufficient Conditions for Throughput Heavy-traffic Optimal

- (R1) If $\bar{\mathbf{Q}}$ exists, i.e., the underline Markov chain is positive recurrent.
- (R2) If $\mathbb{E} \left[\|\bar{\mathbf{Q}}\|^2 \right] \leq M$, i.e., the second moments is bounded by a constant.
- (R3) If $\mathbb{E} \left[\|\bar{\mathbf{Q}}_{\perp}\|^2 \right] \leq K$, which is independent of ϵ . This is often called steady-state collapse, which indicates that under heavy-traffic, queue-length vector concentrates around the line \mathbf{c} .

How can we bound the moments?

A Very Useful Lemma

Lemma

For an irreducible aperiodic and positive Markov chain $\{X(t), t \geq 0\}$ over a countable state space \mathcal{X} , which converges in distribution to \bar{X} , and suppose $V : \mathcal{X} \rightarrow \mathbb{R}_+$ is a Lyapunov function. We define the T time slot drift of V at X as

$$\Delta V(X) := [V(X(t_0 + T)) - V(X(t_0))] \mathcal{I}(X(t_0) = X),$$

where $\mathcal{I}(\cdot)$ is the indicator function. Suppose for some positive finite integer T , the T time slot drift of V satisfies the following conditions:

- (C1) There exists an $\eta > 0$ and a $\gamma < \infty$ such that for any $t_0 = 1, 2, \dots$ and for all $X \in \mathcal{X}$ with $V(X) \geq \gamma$,

$$\mathbb{E}[\Delta V(X) | X(t_0) = X] \leq -\eta.$$

- (C2) There exists a constant $D < \infty$ such that for all $X \in \mathcal{X}$,

$$\mathbb{P}(|\Delta V(X)| \leq D) = 1.$$

Then, there exist finite constants $\{M_r, r \in \mathbb{N}\}$ such that for each positive r , $\mathbb{E}[V(\bar{X})^r] \leq M_r$, where M_r are fully determined by η , γ and D .

General Sufficient Conditions

Assuming bounded support of arrival and departure process, by exploiting the useful lemma and properties of projection to a convex cone, we are able to give sufficient condition to more general cases

General Sufficient Conditions

Assuming bounded support of arrival and departure process, by exploiting the useful lemma and properties of projection to a convex cone, we are able to give sufficient condition to more general cases

- (S1) If there exists a finite constant T_1 and $K_1 > 0$ and $\delta > 0$ such that for all t_0

$$\mathbb{E} \left[\sum_{t=t_0}^{t_0+T_1-1} \langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t_0) = \mathbf{Q} \right] \leq -\delta \|\mathbf{Q}\| + K_1 \quad (4)$$

holds, then the system is throughput optimal and has a stationary distribution with all moments bounded.

General Sufficient Conditions

Assuming bounded support of arrival and departure process, by exploiting the useful lemma and properties of projection to a convex cone, we are able to give sufficient condition to more general cases

- (S1) If there exists a finite constant T_1 and $K_1 > 0$ and $\delta > 0$ such that for all t_0

$$\mathbb{E} \left[\sum_{t=t_0}^{t_0+T_1-1} \langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}(t_0) = \mathbf{Q} \right] \leq -\delta \|\mathbf{Q}\| + K_1 \quad (4)$$

holds, then the system is throughput optimal and has a stationary distribution with all moments bounded.

- (S2) If there exists a finite constant T_2 , and constants $K_2 > 0$ and $\eta > 0$, both independent of ϵ such that for all t_0

$$\mathbb{E} \left[\sum_{t=t_0}^{t_0+T_2-1} \langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}(t_0) = \mathbf{Q} \right] \leq -\eta \|\mathbf{Q}_\perp\| + K_2 \quad (5)$$

holds, then the moments of perpendicular component with respect to any convex set \mathcal{C} is bounded. (steady-state collapse to a convex set \mathcal{C})

General Sufficient Conditions

Assuming bounded support of arrival and departure process, by exploiting the useful lemma and properties of projection to a convex cone, we are able to give sufficient condition to more general cases

- (S1) If there exists a finite constant T_1 and $K_1 > 0$ and $\delta > 0$ such that for all t_0

$$\mathbb{E} \left[\sum_{t=t_0}^{t_0+T_1-1} \langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t_0) = \mathbf{Q} \right] \leq -\delta \|\mathbf{Q}\| + K_1 \quad (4)$$

holds, then the system is throughput optimal and has a stationary distribution with all moments bounded.

- (S2) If there exists a finite constant T_2 , and constants $K_2 > 0$ and $\eta > 0$, both independent of ϵ such that for all t_0

$$\mathbb{E} \left[\sum_{t=t_0}^{t_0+T_2-1} \langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t_0) = \mathbf{Q} \right] \leq -\eta \|\mathbf{Q}_\perp\| + K_2 \quad (5)$$

holds, then the moments of perpendicular component with respect to any convex set \mathcal{C} is bounded. (steady-state collapse to a convex set \mathcal{C})

Clearly, if the convex set is $\mathbf{c} = \frac{1}{\sqrt{N}} \mathbf{1}$, (S1) implies (R1-R2) and (S2) implies (R3).

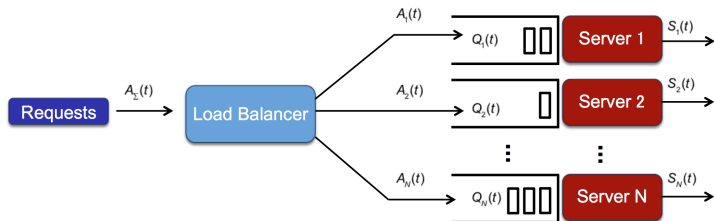
1 Introduction and Motivation

- Why do we need an effective and fast cloud?
- Load Balancing and Previous Works

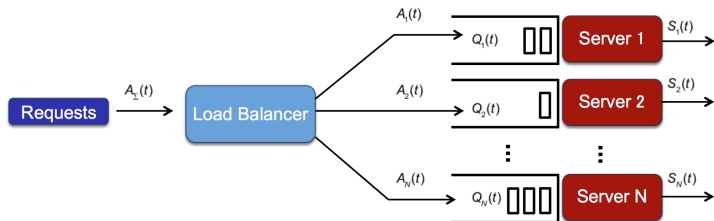
2 Throughput and Heavy-traffic optimality of General Load Balancing

- Model, Challenges, Contributions and Key insights
- Methodology and Sufficient Conditions
- **Homogeneous Servers**
- Heterogeneous Servers
- Conclusions

JBA in Homogeneous Servers

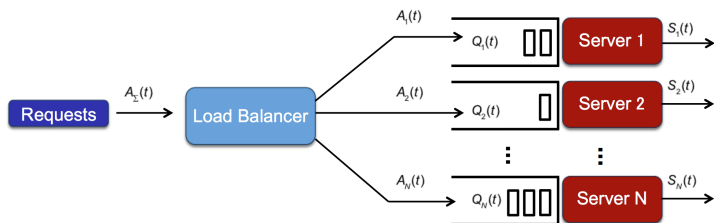


JBA in Homogeneous Servers



- All the N servers have the same average serve rate μ .

JBA in Homogeneous Servers



- All the N servers have the same average serve rate μ .
- The load balancer, under the JBA policy at each time-slot, randomly chose a queue among the queues that have workload less than the average workload at that time slot, and then forward all the incoming requests to that server.

S1 is Satisfied

- Let us first check (S1) by the choice $T = 1$:

$$\begin{aligned}\mathbb{E}[\langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}] &= \langle \mathbf{Q}, \mathbb{E}[\mathbf{A} | \mathbf{Q}] \rangle - \langle \mathbf{Q}, \boldsymbol{\mu} \rangle \\ &= \frac{\lambda_\Sigma}{L} \sum_{n=1}^L Q_n - \mu \sum_{n=1}^N Q_n \\ &= \left(\frac{\lambda_\Sigma}{N} - \mu\right) \sum_{n=1}^N Q_n - \left(\frac{\lambda_\Sigma}{N} - \frac{\lambda_\Sigma}{L}\right) \sum_{n=1}^N Q_n - \frac{\lambda}{N} \sum_{n=L+1}^N Q_n \\ &\leq -\frac{\epsilon}{N} \|\mathbf{Q}\|_1 - \lambda_\Sigma \frac{N-L}{N} (Q_{L+1} - Q_L) \\ &\leq -\frac{\epsilon}{N} \|\mathbf{Q}\|\end{aligned}\tag{6}$$

assume $Q_1(t) \leq Q_2(t) \leq \dots \leq Q_L(t) \leq Q^*(t) < Q_{L+1}(t) \leq \dots \leq Q_N(t)$,
and $Q^*(t) = \frac{1}{N} \sum Q_n(t)$ is the average queue length.

S1 is satisfied, and hence JBA is throughput optimal

S2 is Satisfied

- Let us turn to check (S2) with the line \mathbf{c} as the projection direction:

$$\begin{aligned}\mathbb{E}[\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}(t) = \mathbf{Q}] &\stackrel{(a)}{=} \langle \mathbf{Q}_\perp, \mathbb{E}[\mathbf{A} | \mathbf{Q}] \rangle - \langle \mathbf{Q}_\perp, \boldsymbol{\mu} \rangle \\ &= \langle \mathbf{Q}_\perp, \mathbb{E}[\mathbf{A} | \mathbf{Q}] \rangle \\ &= \frac{\lambda_\Sigma}{L} \sum_{n=1}^L (Q_n - Q^*) \\ &= -\frac{\lambda_\Sigma}{L} \sum_{n=1}^L |Q_n - Q^*| \\ &\leq -\frac{\lambda_\Sigma}{2N} \|\mathbf{Q}_\perp\| \\ &\leq -\frac{\mu}{2} \|\mathbf{Q}_\perp\|\end{aligned}\tag{7}$$

for all $0 < \epsilon \leq \frac{N\mu}{2}$. (a) follows $\langle \mathbf{Q}_\perp, \mathbf{1} \rangle = 0$

S2 is verified and hence JBA is heavy-traffic optimality

Theorems

Theorem

For any λ_{Σ} in the interior of \mathcal{R} , i.e., $\lambda_{\Sigma} < \mu_{\Sigma}$, the JBA routing policy stabilizes the system, and all the moments of the stationary distribution are bounded, i.e., there exist finite constants $\{M_r, r \in \mathbb{N}\}$ such that $\mathbb{E} \left[\|\bar{\mathbf{Q}}\|^r \right] \leq M_r$.

Theorem

Consider a set of load balancing system under JBA policy with the exogenous arrival process $\{A_{\Sigma}^{(\epsilon)}(t), t \geq 0\}$, parameterized by $\epsilon > 0$. Then, each of these systems, the expectation of the sum queue length in steady state is lower bounded by

$$\mathbb{E} \left[\sum_{n=1}^N \bar{Q}_n^{(\epsilon)} \right] \geq \frac{\zeta^{(\epsilon)}}{2\epsilon} - K \quad (8)$$

where $\zeta^{(\epsilon)} = (\sigma_{\Sigma}^{(\epsilon)})^2 + \nu_{\Sigma}^2 + \epsilon^2$, $K = \frac{NS_{\max}}{2}$. Therefore, in the heavy-traffic limit as $\epsilon \downarrow 0$, assuming the $(\sigma_{\Sigma}^{(\epsilon)})^2$ converges to a constant σ_{Σ}^2 , we have

$$\liminf_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[\sum_{n=1}^N \bar{Q}_n^{(\epsilon)} \right] \geq \frac{\zeta}{2}, \quad (9)$$

Theorems (Cont'd)

Theorem

Consider a set of load balancing system with the exogenous arrival process $\{A_{\Sigma}^{(\epsilon)}(t), t \geq 0\}$, parameterized by $\epsilon > 0$, such that the mean arrival rate is $\lambda_{\Sigma}^{(\epsilon)} = \mu_{\Sigma} - \epsilon$ and variance is denoted as $(\sigma_{\Sigma}^{(\epsilon)})^2$. Under the JBA algorithm, $\{\mathbf{Q}^{(\epsilon)}(t), t \geq 0\}$ converges in distribution to $\overline{\mathbf{Q}}^{(\epsilon)}$. Then, there exist finite constants $\{M_r, r \in \mathbb{N}\}$ which are independent of ϵ such that for all $r \in \mathbb{N}$,

$$\mathbb{E} \left[\left\| \overline{\mathbf{Q}}_{\perp}^{(\epsilon)} \right\| \right] \leq M_r, \quad (10)$$

for all system with $0 < \epsilon \leq \frac{N\mu}{2}$.

Theorems (Cont'd)

Theorem

Consider a set of load balancing system with the exogenous arrival process $\{A_{\Sigma}^{(\epsilon)}(t), t \geq 0\}$, parameterized by $\epsilon > 0$, such that the mean arrival rate is $\lambda_{\Sigma}^{(\epsilon)} = \mu_{\Sigma} - \epsilon$ and variance is denoted as $(\sigma_{\Sigma}^{(\epsilon)})^2$. Under the JBA algorithm, $\{Q^{(\epsilon)}(t), t \geq 0\}$ converges in distribution to $\bar{Q}^{(\epsilon)}$. For each system with $0 < \epsilon \leq \frac{N\mu}{2}$, the steady state average queue length satisfies

$$\mathbb{E} \left[\sum_{n=1}^N \bar{Q}_n^{(\epsilon)} \right] \leq \frac{\zeta^{(\epsilon)}}{2\epsilon} + \bar{B}^{(\epsilon)}, \quad (11)$$

where $\zeta^{(\epsilon)}$ is the same as in the lower bound, and $\bar{B}^{(\epsilon)}$ is $o(\frac{1}{\epsilon})$, i.e.,

$$\lim_{\epsilon \downarrow 0} \epsilon \bar{B}^{(\epsilon)} = 0.$$

Therefore, assuming the variance $(\sigma_{\Sigma}^{(\epsilon)})^2$ converges to a constant σ_{Σ}^2 , the upper bound becomes

$$\limsup_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[\sum_{n=1}^N \bar{Q}_n^{(\epsilon)} \right] \leq \frac{\zeta}{2} \quad (12)$$

How about Random Load Balancing

Under random load balancing, we have $\mathbf{A} - \mathbf{S} = -\frac{\epsilon}{N}\mathbf{1}$

- For (S1), we have $\mathbb{E}[\langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}] = -\frac{\epsilon}{N} \|\mathbf{Q}\| \leq -\frac{\epsilon}{2N} \|\mathbf{Q}\|$
- For (S2), we have $\mathbb{E}[\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}(t) = \mathbf{Q}] = 0$ for all t

Consider steady-state collapse, random load balancing actually does no harm in the sense that it would not incur any positive drift.

Can we utilize this fact to turn bad to good?

From 1 to any finite T

The load balancer uses JBA, Power of d , or JSQ every T time-slots, otherwise just random routing.

- For (S1), it is trivial to hold.
- For (S2), by letting $T_2 = T$, we have

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=t_0}^{t_0+T_2-1} \langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t_0) = \mathbf{Q} \right] \\ &= \sum_{t=t_0}^{t_0+T_2-1} \mathbb{E} [\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t_0) = \mathbf{Q}] \\ &\stackrel{(a)}{=} \sum_{t=t_0}^{t_0+T_2-1} \mathbb{E} [\mathbb{E} [\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t)] \mid \mathbf{Q}(t_0) = \mathbf{Q}] \\ &= \mathbb{E} [-\eta \|\mathbf{Q}_\perp(t^*)\| \mid \mathbf{Q}(t_0) = \mathbf{Q}] \\ &\leq -\eta \|\mathbf{Q}_\perp(t_0)\| + DT \end{aligned} \tag{13}$$

where (a) follows from tower property of conditional expectation.

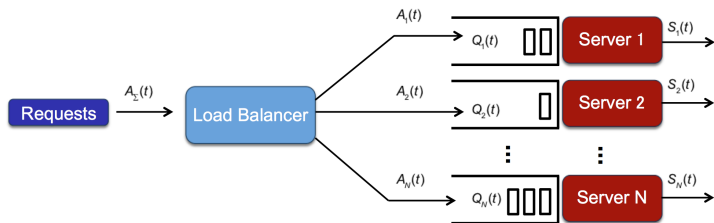
1 Introduction and Motivation

- Why do we need an effective and fast cloud?
- Load Balancing and Previous Works

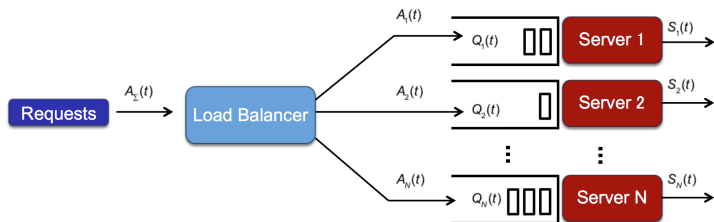
2 Throughput and Heavy-traffic optimality of General Load Balancing

- Model, Challenges, Contributions and Key insights
- Methodology and Sufficient Conditions
- Homogeneous Servers
- **Heterogeneous Servers**
- Conclusions

JBA in Heterogeneous Servers

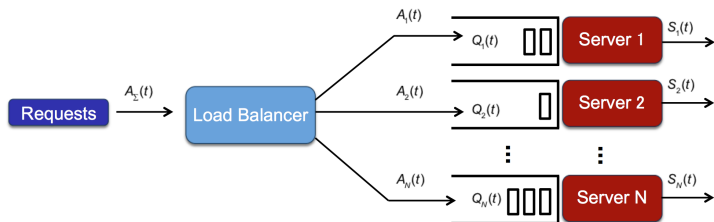


JBA in Heterogeneous Servers



- The N servers do not have the same average serve rate, each with μ_n .

JBA in Heterogeneous Servers



- The N servers do not have the same average serve rate, each with μ_n .
- The load balancer, under the JBA policy at each time-slot, randomly with probability **proportional to the service rate μ_n** to chose a queue among the queues that have workload less than the average workload at that time slot, and then forward all the incoming requests to that server.

$$\mathbb{P}(R_i) = \frac{\mu_i}{\sum_{i \leq L} \mu_i} \quad \text{if } i \leq L \quad (14)$$

JBA in Heterogeneous Servers (Cont'd)

Under the JBA load balancing, it can be shown that (S1) and (S2) still hold, i.e.,

- For (S1), we have $\mathbb{E}[\langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}] \leq -\delta \|\mathbf{Q}\|$ for some $\delta > 0$
- For (S2), we have $\mathbb{E}[\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}(t) = \mathbf{Q}] \leq -\eta \|\mathbf{Q}_\perp(t)\|$ for some $\eta > 0$ independent of ϵ

Under purely random load balancing with proportional probability, we have

- For (S1), we have $\mathbb{E}[\langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}] \leq -\delta \|\mathbf{Q}\|$ for some $\delta > 0$
- For (S2), we have $\mathbb{E}[\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle | \mathbf{Q}(t) = \mathbf{Q}] \leq \epsilon \|\mathbf{Q}_\perp(t)\|$ for all t

From 1 to any finite T

The load balancer uses JBA, Power of d , or JSQ every T time-slots, otherwise just random routing with proportional probability over N servers.

- For (S1), it is trivial to hold.
- For (S2), by letting $T_2 = T$, we have

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=t_0}^{t_0+T_2-1} \langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t_0) = \mathbf{Q} \right] \\ &= \sum_{t=t_0}^{t_0+T_2-1} \mathbb{E} [\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t_0) = \mathbf{Q}] \\ &= \sum_{t=t_0}^{t_0+T_2-1} \mathbb{E} [\mathbb{E} [\langle \mathbf{Q}_\perp(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \mid \mathbf{Q}(t)] \mid \mathbf{Q}(t_0) = \mathbf{Q}] \\ &\leq ((T-1)\epsilon - \eta) \|\mathbf{Q}_\perp(t_0)\| + DT^2 \\ &\leq -\frac{1}{2\eta} \|\mathbf{Q}_\perp(t_0)\| + DT^2 \end{aligned} \tag{15}$$

for all $0 < \epsilon \leq \frac{\eta}{2T}$, hence smaller ϵ means a larger T !

1 Introduction and Motivation

- Why do we need an effective and fast cloud?
- Load Balancing and Previous Works

2 Throughput and Heavy-traffic optimality of General Load Balancing

- Model, Challenges, Contributions and Key insights
- Methodology and Sufficient Conditions
- Homogeneous Servers
- Heterogeneous Servers
- **Conclusions**

Conclusions

- Can we generalize existing load balancing algorithm?
 - ▶ Yes, the proposed JBA policy
- Do we really need sampling for each time-slot for heavy-traffic optimality?
 - ▶ No, we can actually sampling every T slots whenever $T\epsilon \leq \alpha$ is satisfied.

Thank you!
Q & A