



LUNDS
UNIVERSITET

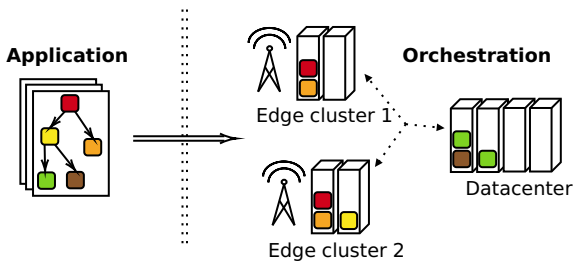
Cloud application modeling using mean-field fluid models

Johan Ruuskanen





Introduction, managing cloud applications



Trends in modern cloud computing

- Applications split into graphs of smaller services
- Clouds of multiple clusters

Complex service graphs and dynamic environments



Introduction, managing cloud applications

Problem, how to deploy/manage an application such that

- a) users receive a good QoS (e.g. low latency, robustness)
- b) the costs are minimized (e.g. allocated resources)

Automatic adaption of resources and scheduling

Popular research topic considering single service application, and recently more considering service-graph applications.

Good decisions necessitates good models



Introduction, managing cloud applications

Problem, how to deploy/manage an application such that

- a) users receive a good QoS (e.g. low latency, robustness)
- b) the costs are minimized (e.g. allocated resources)

Automatic adaption of resources and scheduling

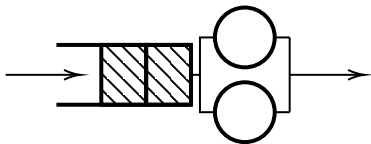
Popular research topic considering single service application, and recently more considering service-graph applications.

Good decisions necessitates good models



Applications as queuing networks

Common to model a service¹ as a queue.



Lifetime of a request: (i) arrives, (ii) is assigned a service time from G_S , (iii) processed according to *discipline* and (iv) departs once completed.

Queuing disciplines

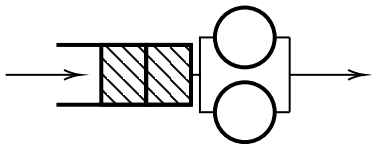
- First come, first served (FCFS)
- Processor sharing (PS)
- Pure delay (INF)

¹i.e. server, but not necessarily a physical computer



Applications as queuing networks

Common to model a service¹ as a queue.



Lifetime of a request: (i) arrives, (ii) is assigned a service time from G_S , (iii) processed according to *discipline* and (iv) departs once completed.

Queuing disciplines

- First come, first served (FCFS)
- Processor sharing (PS)
- Pure delay (INF)

¹i.e. server, but not necessarily a physical computer



Applications as queuing networks

Applications of many stages, use many queues in a network

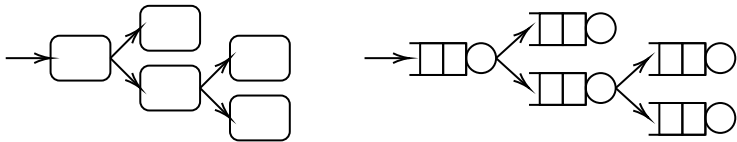


Figure: Simple example, where each stage is a service in a service graph.

Exists many extensions, one important is

Multi-class queues; Each queue has a set of *classes*, each request is assigned to one. Each class has its own G_s , and destination once completed.

$P_{i,j}^{r,s}$ - the probability that a completed request of class r in queue i gets routed to class s in queue j .



Applications as queuing networks

Applications of many stages, use many queues in a network

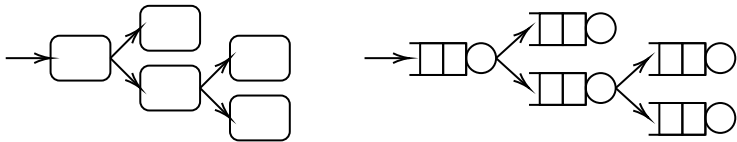


Figure: Simple example, where each stage is a service in a service graph.

Exists many extensions, one important is

Multi-class queues; Each queue has a set of *classes*, each request is assigned to one. Each class has its own G_s , and destination once completed.

$P_{i,j}^{r,s}$ - the probability that a completed request of class r in queue i gets routed to class s in queue j .



Evaluating a queuing network

$X_{i,r}(t)$ - population of requests of class r in queue i at time t .

Always possible to estimate the PMF of $X_{i,r}(t) \forall i, r, t \geq 0$ using MC simulations.

Very computationally intensive, not suitable for most cases.

Instead, approximate important metrics (e.g. mean queue length, response time)

Exists many methods

Stationary, product-form -> methods utilizing the BCMP theorem

Transient, non-product-form -> fluid models



Evaluating a queuing network

$X_{i,r}(t)$ - population of requests of class r in queue i at time t .

Always possible to estimate the PMF of $X_{i,r}(t) \forall i, r, t \geq 0$ using MC simulations.

Very computationally intensive, not suitable for most cases.

Instead, approximate important metrics (e.g. mean queue length, response time)

Exists many methods

Stationary, product-form -> methods utilizing the BCMP theorem

Transient, non-product-form -> fluid models



Fluid model of a queuing network

Model $\mathbb{E}[X_{i,r}(t)]$ as $x_{i,r}(t)$, where $\mathbf{X}(0) = \mathbf{x}(0)$ and

$$\dot{x}_{i,r}(t) = f_{i,r}^{in}(\mathbf{x}(t)) - f_{i,r}^{out}(\mathbf{x}(t))$$

Difficult to find f^{in}, f^{out} such that $\mathbf{x}(t)$ is a good approximation.

Much research has been done considering the single-queue/single-class case.

Queuing networks trickier, for some types the *mean-field approximation* gives one way.



Fluid model of a queuing network

Model $\mathbb{E}[X_{i,r}(t)]$ as $x_{i,r}(t)$, where $\mathbf{X}(0) = \mathbf{x}(0)$ and

$$\dot{x}_{i,r}(t) = f_{i,r}^{in}(\mathbf{x}(t)) - f_{i,r}^{out}(\mathbf{x}(t))$$

Difficult to find f^{in}, f^{out} such that $\mathbf{x}(t)$ is a good approximation.

Much research has been done considering the single-queue/single-class case.

Queuing networks trickier, for some types the *mean-field approximation* gives one way.



Fluid model of a queuing network

Mean-field approximation

Let \mathbf{X} be a vector of populations in a *density-dependent population process* (special type of CTMC).

Transition l such that at an event $\mathbf{X}(t^+) = \mathbf{X}(t) + l$ with rate function $f(\mathbf{X}, l)$. The drift then becomes $F(\mathbf{X}) = \sum_{l \in \mathcal{L}} l f(\mathbf{X}, l)$

Mean-field approximation; $\dot{\mathbf{x}} = F(\mathbf{x})$, certain conditions $\nu^{-1} \mathbf{X}^{(\nu)} \rightarrow \mathbf{x}$ at all t when $\nu \rightarrow \infty$ (Kurtz's theorem).

Mean-field fluid model

For some queuing networks, possible to translate to such a CTMC.

Applies to multi-class queuing networks of PS and INF queues where G_s has a *Phase-type* distribution ² ³.

²**Closed networks:** F. Pérez and G. Casale, *Line: Evaluating Software Applications in Unreliable Environments*, IEEE Transactions on Reliability (2017)

³**Open/mixed networks:** (allowing arrivals/departures) pre-print available



Fluid model of a queuing network

Mean-field approximation

Let \mathbf{X} be a vector of populations in a *density-dependent population process* (special type of CTMC).

Transition l such that at an event $\mathbf{X}(t^+) = \mathbf{X}(t) + l$ with rate function $f(\mathbf{X}, l)$. The drift then becomes $F(\mathbf{X}) = \sum_{l \in \mathcal{L}} l f(\mathbf{X}, l)$

Mean-field approximation; $\dot{\mathbf{x}} = F(\mathbf{x})$, certain conditions $\nu^{-1} \mathbf{X}^{(\nu)} \rightarrow \mathbf{x}$ at all t when $\nu \rightarrow \infty$ (Kurtz's theorem).

Mean-field fluid model

For some queuing networks, possible to translate to such a CTMC.

Applies to multi-class queuing networks of PS and INF queues where G_s has a *Phase-type* distribution ^{2 3}.

²**Closed networks:** F. Pérez and G. Casale, *Line: Evaluating Software Applications in Unreliable Environments*, IEEE Transactions on Reliability (2017)

³**Open/mixed networks:** (allowing arrivals/departures) pre-print available



Fluid model of a queuing network

Mean-field approximation

Let \mathbf{X} be a vector of populations in a *density-dependent population process* (special type of CTMC).

Transition l such that at an event $\mathbf{X}(t^+) = \mathbf{X}(t) + l$ with rate function $f(\mathbf{X}, l)$. The drift then becomes $F(\mathbf{X}) = \sum_{l \in \mathcal{L}} l f(\mathbf{X}, l)$

Mean-field approximation; $\dot{\mathbf{x}} = F(\mathbf{x})$, certain conditions $\nu^{-1} \mathbf{X}^{(\nu)} \rightarrow \mathbf{x}$ at all t when $\nu \rightarrow \infty$ (Kurtz's theorem).

Mean-field fluid model

For some queuing networks, possible to translate to such a CTMC.

Applies to multi-class queuing networks of PS and INF queues where G_S has a *Phase-type* distribution ² ³.

²**Closed networks:** F. Pérez and G. Casale, *Line: Evaluating Software Applications in Unreliable Environments*, IEEE Transactions on Reliability (2017)

³**Open/mixed networks:** (allowing arrivals/departures) pre-print available



Mean-field fluid model

Phase-type distribution

Represent a distribution as the *time to absorption* in a single-sink CT Markov random walk across some graph.

Parameterized (for every class r in every queue i)

- $\alpha \in \mathbb{R}^{S_{i,r}}$, prob. vector of starting transient state
- $\Psi \in \mathbb{R}^{S_{i,r} \times S_{i,r}}$, matrix of transition rates between transient states
- $\psi \in \mathbb{R}^{S_{i,r}}$, transition rates between transient states and the sink

We can now introduce $X_{i,r,a}$

- population of requests in state a , in class r , in queue i .



Mean-field fluid model

Phase-type distribution

Represent a distribution as the *time to absorption* in a single-sink CT Markov random walk across some graph.

Parameterized (for every class r in every queue i)

- $\alpha \in \mathbb{R}^{S_{i,r}}$, prob. vector of starting transient state
- $\Psi \in \mathbb{R}^{S_{i,r} \times S_{i,r}}$, matrix of transition rates between transient states
- $\psi \in \mathbb{R}^{S_{i,r}}$, transition rates between transient states and the sink

We can now introduce $X_{i,r,a}$

- population of requests in state a , in class r , in queue i .



Mean-field fluid model

Assume a multi-class queuing network of PS and INF queues under Poisson arrivals

Nice thing with PS and INF queues, order does not matter.

$$- \theta_{i,r,a}(\mathbf{X}) = X_{i,r,a} \frac{\min(k_i, \sum_{s,b} X_{i,s,b})}{\sum_{s,b} X_{i,s,b}}$$

Requests in i, r, a times the share of each request in queue i

then with PH distributions, the evolution of \mathbf{X} is a CTMC.



Mean-field fluid model

Exists four types of transitions (l_1 and l_2 from Peréz & Casale)

- $e_{i,r,a}$, zero vector with 1 on position i, r, a .

between phases: $l_1 = e_{i,r,b} - e_{i,r,a}$

$$f^n(\mathbf{X}, l_1) = \Psi_{a,b}^{i,r} \theta_{i,r,a}(\mathbf{X})$$

between classes: $l_2 = e_{j,s,b} - e_{i,r,a}$

$$f^c(\mathbf{X}, l_2) = \psi_a^{i,r} \alpha_b^{j,s} P_{i,j}^{r,s} \theta_{i,r,a}(\mathbf{X})$$

arrivals: $l_3 = e_{i,r,a}$

$$f^a(\mathbf{X}, l_3) = \alpha_a^{i,r} \lambda^{i,r}$$

departures: $l_4 = -e_{i,r,a}$

$$f^d(\mathbf{X}, l_4) = \psi_a^{i,r} \left(1 - \sum_{j,s} P_{i,j}^{r,s} \right) \theta_{i,r,a}(\mathbf{X})$$



Mean-field fluid model

Drift in each i, r, a

$$F_{i,r,a}(\mathbf{X}) = \sum_b \Psi_{b,a}^{i,r} \theta_{i,r,b}(\mathbf{X}) + \alpha_a^{i,r} \sum_{j,s,b} \psi_b^{j,s} P_{j,i}^{s,r} \theta_{j,s,b}(\mathbf{X}) + \alpha_a^{i,r} \lambda^{i,r}$$

Assuming \mathbf{X} subsequently ordered in phases/classes/queues

$$\Psi = \text{diag}(\Psi^{1,1}, \Psi^{1,2}, \Psi^{1,3}, \dots)$$

$$A = \text{diag}(\alpha^{1,1}, \alpha^{1,2}, \alpha^{1,3}, \dots)$$

$$B = \text{diag}(\psi^{1,1}, \psi^{1,2}, \psi^{1,3}, \dots)$$

$$P = \begin{bmatrix} P_{1,1}^{\cdot} & P_{1,2}^{\cdot} & \dots \\ P_{2,2}^{\cdot} & P_{2,2}^{\cdot} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

We can create $W = \Psi + BPA^T$ and

$$F(\mathbf{X}) = W^T \theta(\mathbf{X}) + A\lambda$$



Mean-field fluid model

Drift in each i, r, a

$$F_{i,r,a}(\mathbf{X}) = \sum_b \Psi_{b,a}^{i,r} \theta_{i,r,b}(\mathbf{X}) + \alpha_a^{i,r} \sum_{j,s,b} \psi_b^{j,s} P_{j,i}^{s,r} \theta_{j,s,b}(\mathbf{X}) + \alpha_a^{i,r} \lambda^{i,r}$$

Assuming \mathbf{X} subsequently ordered in phases/classes/queues

$$\mathbf{\Psi} = \text{diag}(\Psi^{1,1}, \Psi^{1,2}, \Psi^{1,3}, \dots)$$

$$\mathbf{A} = \text{diag}(\alpha^{1,1}, \alpha^{1,2}, \alpha^{1,3}, \dots)$$

$$\mathbf{B} = \text{diag}(\psi^{1,1}, \psi^{1,2}, \psi^{1,3}, \dots)$$

$$\mathbf{P} = \begin{bmatrix} P_{1,1} & P_{1,2} & \dots \\ P_{2,1} & P_{2,2} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

We can create $\mathbf{W} = \mathbf{\Psi} + \mathbf{BPA}^T$ and

$$F(\mathbf{X}) = \mathbf{W}^T \theta(\mathbf{X}) + \mathbf{A}\lambda$$



Mean-field fluid model

The entire mean-field fluid model can then be expressed as

$$\dot{\mathbf{x}} = \mathbf{W}^T \theta(\mathbf{x}) + \mathbf{A} \boldsymbol{\lambda}$$

$$\mathbf{x}(0) = \mathbf{X}(0)$$

then $\lim_{\nu \rightarrow \infty} \nu^{-1} \mathbf{X}^{(\nu)} = \mathbf{x}$ at all t ,

where $\mathbf{X}^{(\nu)}$ is \mathbf{X} with k , $\mathbf{X}(0)$ and $\boldsymbol{\lambda}$ scaled with ν .

However, can give poor performance for smaller system sizes



Mean-field fluid model

The entire mean-field fluid model can then be expressed as

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{W}^T \theta(\mathbf{x}) + \mathbf{A} \boldsymbol{\lambda} \\ \mathbf{x}(0) &= \mathbf{X}(0)\end{aligned}$$

then $\lim_{\nu \rightarrow \infty} \nu^{-1} \mathbf{X}^{(\nu)} = \mathbf{x}$ at all t ,

where $\mathbf{X}^{(\nu)}$ is \mathbf{X} with k , $\mathbf{X}(0)$ and $\boldsymbol{\lambda}$ scaled with ν .

However, can give poor performance for smaller system sizes



Improving the mean-field fluid model

Why is this?

Want \mathbf{x} to approximate $\mathbb{E}(\mathbf{X})$ but

$$\frac{d}{dt} \mathbb{E}[\mathbf{X}] = \mathbb{E}[F(\mathbf{X})] \neq F(\mathbb{E}[\mathbf{X}]) = \frac{d}{dt} \mathbf{x}$$

the queuing network case

$$\mathbb{E}[\mathbf{W}^T \theta(\mathbf{X}) + \mathbf{A}\boldsymbol{\lambda}] = \mathbf{W}^T \mathbb{E}[\theta(\mathbf{X})] + \mathbf{A}\boldsymbol{\lambda} \neq \mathbf{W}^T \theta(\mathbb{E}[\mathbf{X}]) + \mathbf{A}\boldsymbol{\lambda}$$

Can we find another $\hat{\theta}(\mathbb{E}[\mathbf{X}])$ that improves accuracy?



Improving the mean-field fluid model

Why is this?

Want \mathbf{x} to approximate $\mathbb{E}(\mathbf{X})$ but

$$\frac{d}{dt} \mathbb{E}[\mathbf{X}] = \mathbb{E}[F(\mathbf{X})] \neq F(\mathbb{E}[\mathbf{X}]) = \frac{d}{dt} \mathbf{x}$$

the queuing network case

$$\mathbb{E}[\mathbf{W}^T \theta(\mathbf{X}) + \mathbf{A}\boldsymbol{\lambda}] = \mathbf{W}^T \mathbb{E}[\theta(\mathbf{X})] + \mathbf{A}\boldsymbol{\lambda} \neq \mathbf{W}^T \theta(\mathbb{E}[\mathbf{X}]) + \mathbf{A}\boldsymbol{\lambda}$$

Can we find another $\hat{\theta}(\mathbb{E}[\mathbf{X}])$ that improves accuracy?



Improving the mean-field fluid model

Problem,

$$\mathbb{E} [\theta_{i,r,a}(\mathbf{X})] = \sum_{\mathbf{z}} \mathbb{P}(\mathbf{X} = \mathbf{z}) z_{i,r,a} \frac{\min(k_i, \sum_{s,b} z_{i,s,b})}{\sum_{s,b} z_{i,r,a}}$$
$$\theta_{i,r,a}(\mathbb{E}[\mathbf{X}]) = \mathbb{E}[X_{i,r,a}] \frac{\min(k_i, \sum_{s,b} \mathbb{E}[X_{i,s,b}])}{\sum_{s,b} \mathbb{E}[X_{i,s,b}]}$$

First, let $\theta_{i,r,a}(\mathbf{X}) = X_{i,r,a} g_{i,r,a}(\mathbf{X})$,

$g_{i,r,a}(\mathbf{X})$ is the processor share of queue i and $g_{i,r,a}(\mathbf{X}) = g_{i,s,b}(\mathbf{X})$

Let $\hat{\theta}_{i,r,a}(\mathbb{E}[\mathbf{X}]) = \mathbb{E}[X_{i,r,a}] \hat{g}_{i,r,a}(\mathbb{E}[\mathbf{X}])$, then by summing over all states/classes in queue i

$$\hat{g}_{i,r,a}(\mathbb{E}[\mathbf{X}]) = \frac{\sum_c \mathbb{P}(\sum_{s,b} X_{i,r,a} = c) \min(k_i, c)}{\sum_{s,b} \mathbb{E}[X_{i,s,b}]} = \frac{k_i \rho_i(\mathbf{X})}{\sum_{s,b} \mathbb{E}[X_{i,s,b}]}$$

Dependence on the PMF of \mathbf{X} , we need to allow \hat{g} to change



Improving the mean-field fluid model

Problem,

$$\mathbb{E} [\theta_{i,r,a}(\mathbf{X})] = \sum_{\mathbf{z}} \mathbb{P}(\mathbf{X} = \mathbf{z}) z_{i,r,a} \frac{\min(k_i, \sum_{s,b} z_{i,s,b})}{\sum_{s,b} z_{i,r,a}}$$
$$\theta_{i,r,a}(\mathbb{E}[\mathbf{X}]) = \mathbb{E} [X_{i,r,a}] \frac{\min(k_i, \sum_{s,b} \mathbb{E}[X_{i,s,b}])}{\sum_{s,b} \mathbb{E}[X_{i,s,b}]}$$

First, let $\theta_{i,r,a}(\mathbf{X}) = X_{i,r,a} g_{i,r,a}(\mathbf{X})$,

$g_{i,r,a}(\mathbf{X})$ is the processor share of queue i and $g_{i,r,a}(\mathbf{X}) = g_{i,s,b}(\mathbf{X})$

Let $\hat{\theta}_{i,r,a}(\mathbb{E}[\mathbf{X}]) = \mathbb{E} [X_{i,r,a}] \hat{g}_{i,r,a}(\mathbb{E}[\mathbf{X}])$, then by summing over all states/classes in queue i

$$\hat{g}_{i,r,a}(\mathbb{E}[\mathbf{X}]) = \frac{\sum_c \mathbb{P}(\sum_{s,b} X_{i,r,a} = c) \min(k_i, c)}{\sum_{s,b} \mathbb{E}[X_{i,s,b}]} = \frac{k_i \rho_i(\mathbf{X})}{\sum_{s,b} \mathbb{E}[X_{i,s,b}]}$$

Dependence on the PMF of \mathbf{X} , we need to allow \hat{g} to change



Improving the mean-field fluid model

One such possible function is

$$\hat{g}_{i,r,a}(\mathbf{x} | p_i) = \frac{1}{(1 + (k_i^{-1} \sum_{s,b} \mathbf{x}_{i,s,b})^{p_i})^{1/p_i}}$$

The inverse p-norm, can be seen as a smoothing of $g_{i,r,a}(\mathbf{X})$ with parameter p_i .

$p_i \rightarrow \infty$ gives back $g_{i,r,a}(\mathbf{X})$.

Nice because of monotonicity:

- given data at stationarity, "optimal" \mathbf{p}^* can be found



Evaluation, M/M/1 queue

First considering the most simplistic queuing network,
- a single queue with 1 server, 1 class and 1 phase.

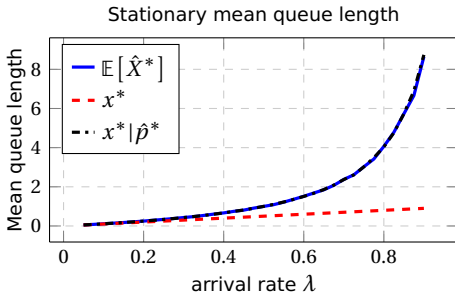
The mean field model then gives $\dot{x} = -\mu \min(1, x) + \lambda$

stationary point: $x = \lambda/\mu = \rho \leq 1$,

However, true mean is well-known: $\mathbb{E}[X] = \frac{\rho}{1-\rho}$.



Evaluation, M/M/1 queue



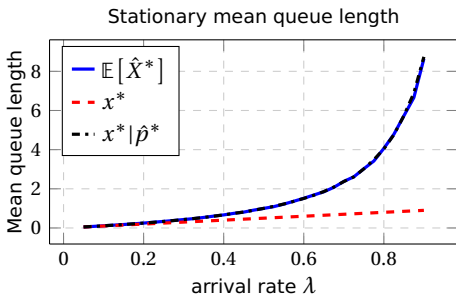
The p^* found is consistently around 1, which gives

$$\dot{x} = x \cdot \hat{g}(x | p = 1) + \lambda = \frac{x}{x+1} + \lambda$$

known as the Tipper model



Evaluation, M/M/1 queue



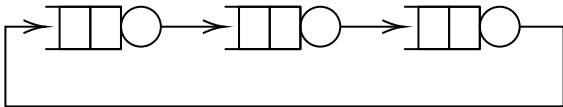
The p^* found is consistently around 1, which gives

$$\dot{x} = x \cdot \hat{g}(x | p = 1) + \lambda = \frac{x}{x+1} + \lambda$$

known as the Tipper model



Evaluation, three tandem queues



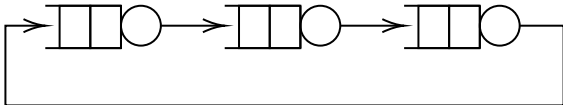
Three single class queues, queue 1 (INF) and queue 2 & 3 (PS)

$$W = \begin{bmatrix} -\mu_1 & \mu_1 & 0 & 0 & 0 \\ 0 & -4.0 & 4.0 & 0 & 0 \\ 0 & 0 & -4.0 & 4.0 & 0 \\ 1.9 & 0 & 0 & -2.0 & 0.1 \\ 0.1 & 0 & 0 & 0 & -0.1 \end{bmatrix}$$

$$\theta(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \cdot \min(4, x_2 + x_3) / (x_2 + x_3) \\ x_3 \cdot \min(4, x_2 + x_3) / (x_2 + x_3) \\ x_4 \cdot \min(8, x_4 + x_5) / (x_4 + x_5) \\ x_5 \cdot \min(8, x_4 + x_5) / (x_4 + x_5) \end{bmatrix}$$



Evaluation, three tandem queues



Three single class queues, queue 1 (INF) and queue 2 & 3 (PS)

$$W = \begin{bmatrix} -\mu_1 & \mu_1 & 0 & 0 & 0 \\ 0 & -4.0 & 4.0 & 0 & 0 \\ 0 & 0 & -4.0 & 4.0 & 0 \\ 1.9 & 0 & 0 & -2.0 & 0.1 \\ 0.1 & 0 & 0 & 0 & -0.1 \end{bmatrix}$$

$$\theta(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \cdot \min(4, x_2 + x_3) / (x_2 + x_3) \\ x_3 \cdot \min(4, x_2 + x_3) / (x_2 + x_3) \\ x_4 \cdot \min(8, x_4 + x_5) / (x_4 + x_5) \\ x_5 \cdot \min(8, x_4 + x_5) / (x_4 + x_5) \end{bmatrix}$$



Evaluation, three tandem queues

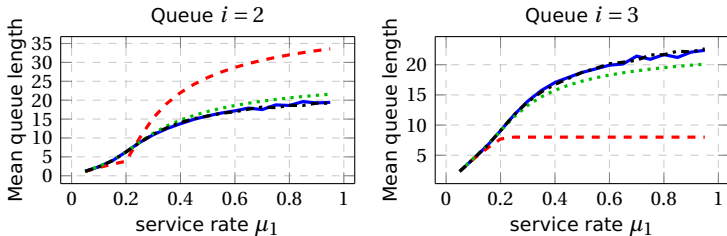


Figure: (Blue) queue length from simulation, (Red) mean-field model, (Black) smoothed model with p^* estimated at every μ_1 , (Green) smoothed model with p estimated at $\mu_1 = 0.2$



Conclusion & future work

Conclusion

- Managing applications in the cloud is tricky
- Model using queuing networks, evaluate using fluid models
- Mean-field approximation for networks of PS queues
- Not necessarily good, consider using smoothed model

Next steps

- Test on a real system.
- How to construct a fluid model that tracks a running application.



Conclusion & future work

Conclusion

- Managing applications in the cloud is tricky
- Model using queuing networks, evaluate using fluid models
- Mean-field approximation for networks of PS queues
- Not necessarily good, consider using smoothed model

Next steps

- Test on a real system.
- How to construct a fluid model that tracks a running application.