

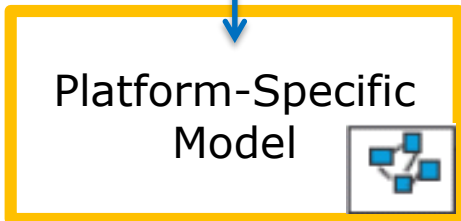
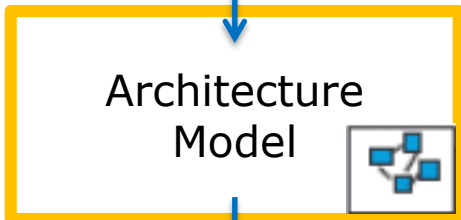
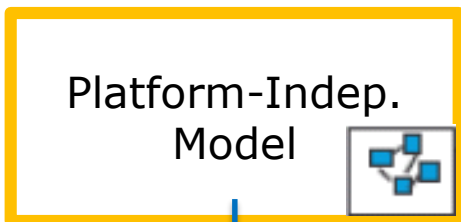
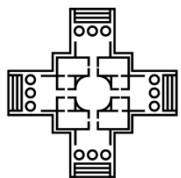
Maximum Likelihood Estimation of Closed Queueing Network Demands from Queue Length Data

Weikun Wang (*Imperial College London, UK*)

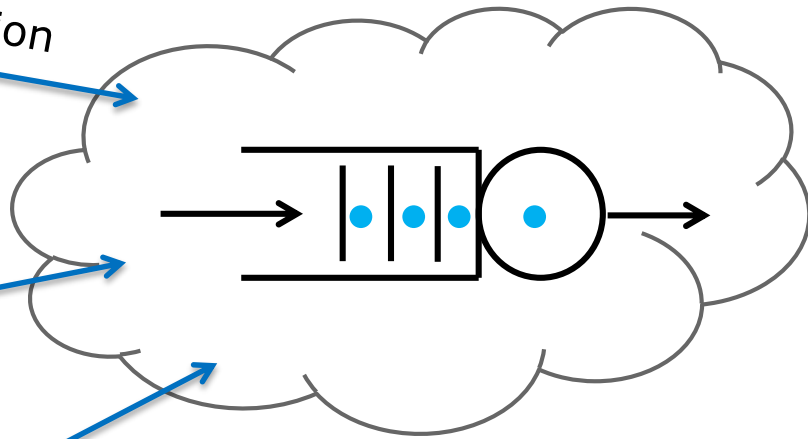
Giuliano Casale (*Imperial College London, UK*)

Ajay Kattepur (*TCS Innovation Labs, India*)

Manoj Nambiar (*TCS Innovation Labs, India*)



Performance
& Reliability Models

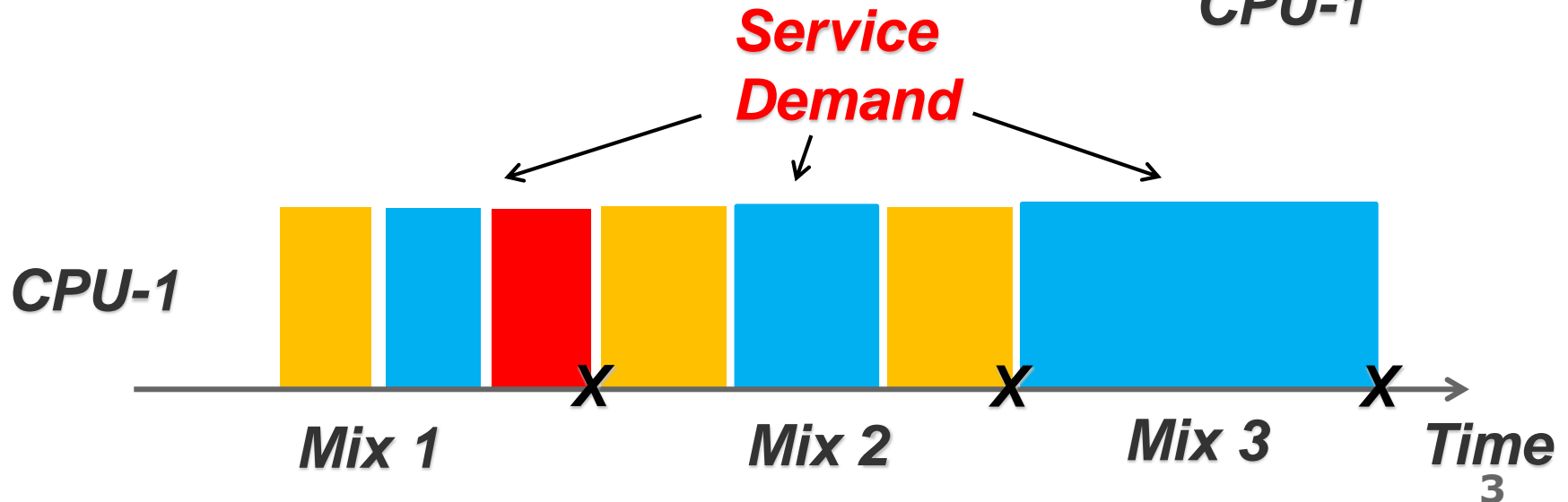
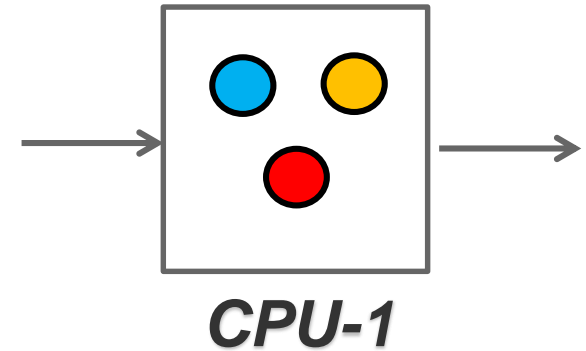


Optimization

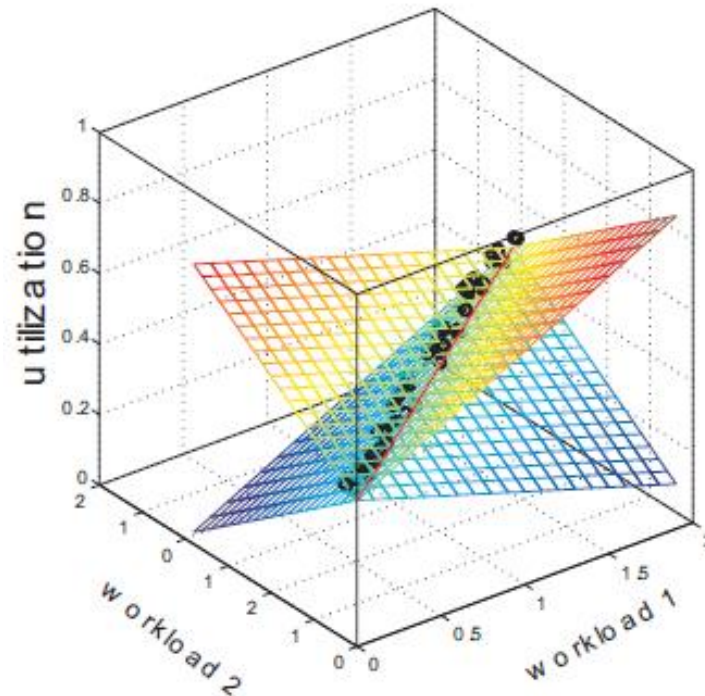
Initial
Configuration

Run-time
Management

- **Service demand** of a request
 - CPU time, bandwidth consumed, ...
- Multi-threaded software
 - *e.g.*, web servers



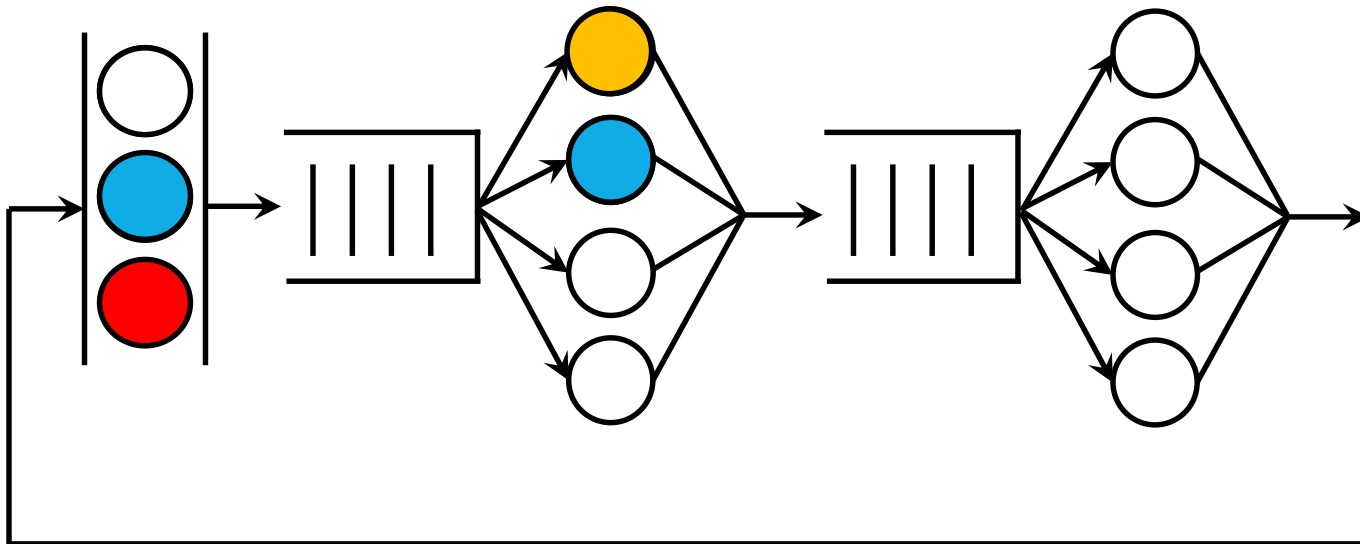
- Utilization based approaches
 - Regression based on utilization and throughput
 - Issues: collinearities, load-dependence, outliers, utilization unreliable/unavailable, ...



■ State observations

- Dataset (L points): $\mathbf{n}^l \in \mathcal{D}$.
- CQN State: $\mathbf{n} = (n_{01}, \dots, n_{0R}, n_{11}, \dots, n_{1R}, \dots, n_{MR})$

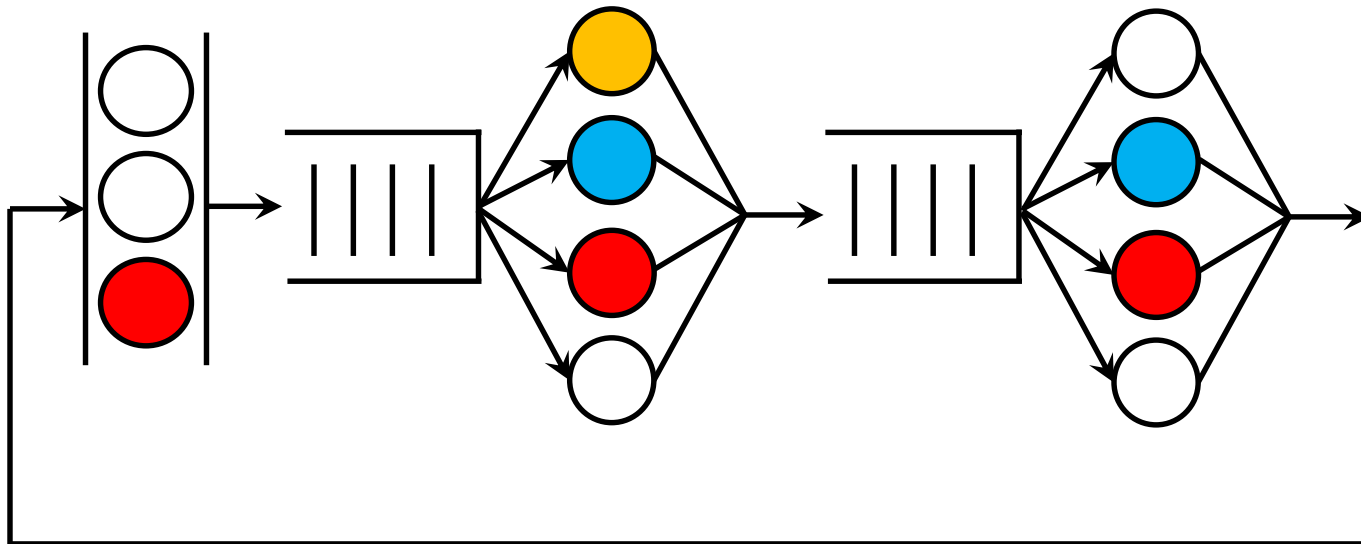
$\mathbf{n} = (\quad 0, 1, 3, \quad 1, 1, 0, \quad 0, 0, 0 \quad)$



■ State observations

- Dataset (L points): $\mathbf{n}^l \in \mathcal{D}$.
- CQN State: $\mathbf{n} = (n_{01}, \dots, n_{0R}, n_{11}, \dots, n_{1R}, \dots, n_{MR})$

$$\mathbf{n} = (\quad 0, 0, 1, \quad 1, 1, 1, \quad 0, 1, 1 \quad)$$



- Assume product-form state probabilities

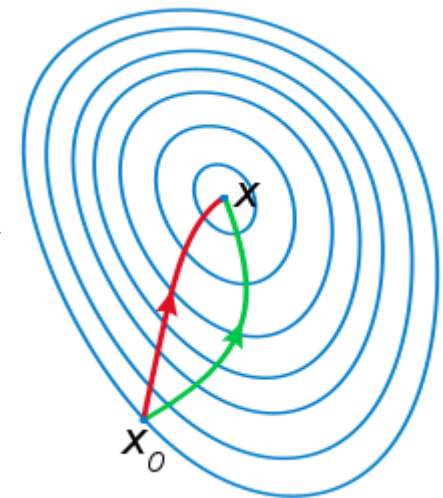
$$\mathbb{P}(\mathbf{n}|\boldsymbol{\theta}) = \prod_{j=1}^R \frac{\theta_{0j}^{n_{0j}}}{n_{0j}!} \prod_{i=1}^M n_i! \prod_{j=1}^R \frac{\theta_{ij}^{n_{ij}}}{n_{ij}! G(\boldsymbol{\theta})}$$

Service demand → $\theta_{ij}^{n_{ij}}$

Queue length → n_{ij}

Normalizing constant → $G(\boldsymbol{\theta})$

- Computationally challenging to evaluate $\mathbb{P}(\mathbf{n}|\boldsymbol{\theta})$
- Maximum likelihood estimation?
 - Infer demands with the probability



- Maximum likelihood estimation (MLE)

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \prod_{l=1}^L \mathbb{P}(\mathbf{n}^l | \theta)$$

parameter space \rightarrow $\underbrace{\hspace{10em}}$ *Likelihood $\mathcal{L}(\theta)$*

- Problem with direct computation

- Evaluation of $\mathbb{P}(\mathbf{n} | \theta)$ for each observation
- Slow due to the need for computing $G(\theta)$
- Very small probabilities when L is large

- Any other solution?

- A necessary condition for a point $\hat{\theta}$ inside Θ to be a MLE is that

$$Q_{ij}(\hat{\theta}) = \tilde{Q}_{ij}(D), \quad \forall i, j,$$

theoretical mean queue length *observed mean queue length*

Only **mean** queue length is required!

- How to find the MLE?
 - Change the value of θ , until the mean queue length predicted with MVA match $\tilde{Q}_{ij}(D)$
 - Fixed point iteration or an optimization program

- Assume the MLE to be asymptotically normal
- Confidence intervals for the MLE demands

$$\hat{\theta}_{ij} \pm c\sqrt{(\mathbf{I}(\hat{\boldsymbol{\theta}})^{-1})_{ij,ij}}$$

- $\mathbf{I}(\hat{\boldsymbol{\theta}}) = -\mathbf{H}(\hat{\boldsymbol{\theta}})$ is the Fisher Information matrix
- $\mathbf{H}(\hat{\boldsymbol{\theta}})$ is the Hessian matrix $\mathbf{H}(\boldsymbol{\theta})_{ij,kh} = \left. \frac{\partial^2 \log \mathcal{L}(\boldsymbol{\theta})}{\partial \theta_{ij} \partial \theta_{kh}} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$
- $\mathbf{H}(\hat{\boldsymbol{\theta}})$ works with mean queue length only!
 - Obtained by using standard MVA, no probabilities!

- Exact MLE can be found by direct search
 - Fixed-point iteration tends to be effective
- A simple approximation of the MLE:
 - Consider the demand vector θ^{bs} where

$$\theta_{ij}^{bs} = \frac{\tilde{Q}_{ij}(D)}{(N_j - \sum_{k=1}^M \tilde{Q}_{kj})} \frac{\theta_{0,j}}{(1 + \sum_{h=1}^R \tilde{Q}_{ih} - \tilde{Q}_{ij}(D)/N_j)}$$

- Then it must be

$$Q_{ij}^{bs}(\theta^{bs}) = \tilde{Q}_{ij}(D)$$

*observed mean
queue length*

■ CI: Complete Information

[J.F. Perez *et al.*, *IEEE Trans. Sw. Eng.*'15]

- Full knowledge of sample path
- Baseline approach

■ ERPS: Extended Regression for Processor Sharing

[J.F. Perez *et al.*, *IEEE Trans. Sw. Eng.*'15]

- Based on mean response time and arrival queue

■ GQL: Gibbs Sampling for Queue Lengths

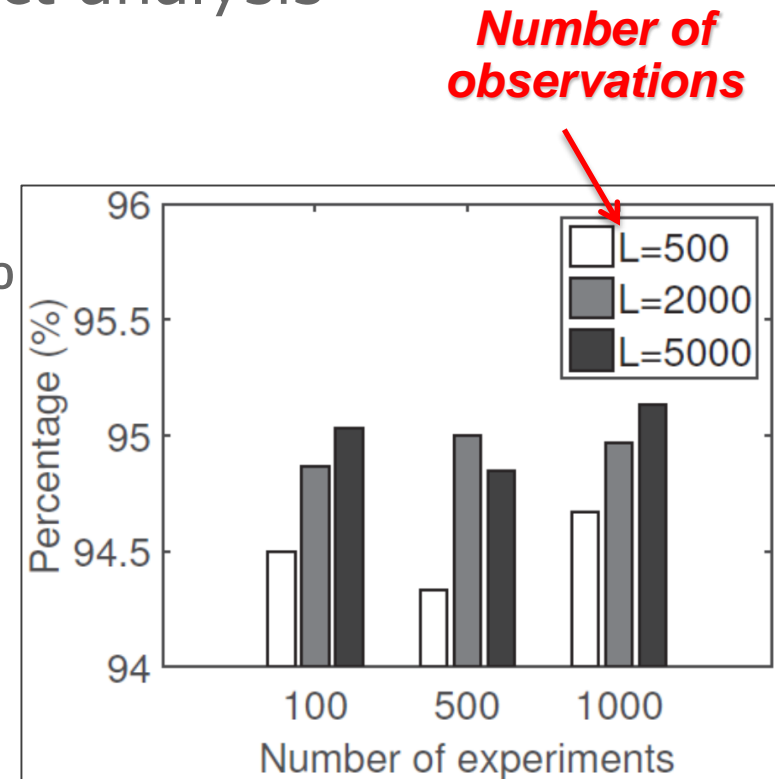
[W. Wang *et al.*, *Accepted to appear in ACM TOMACS*]

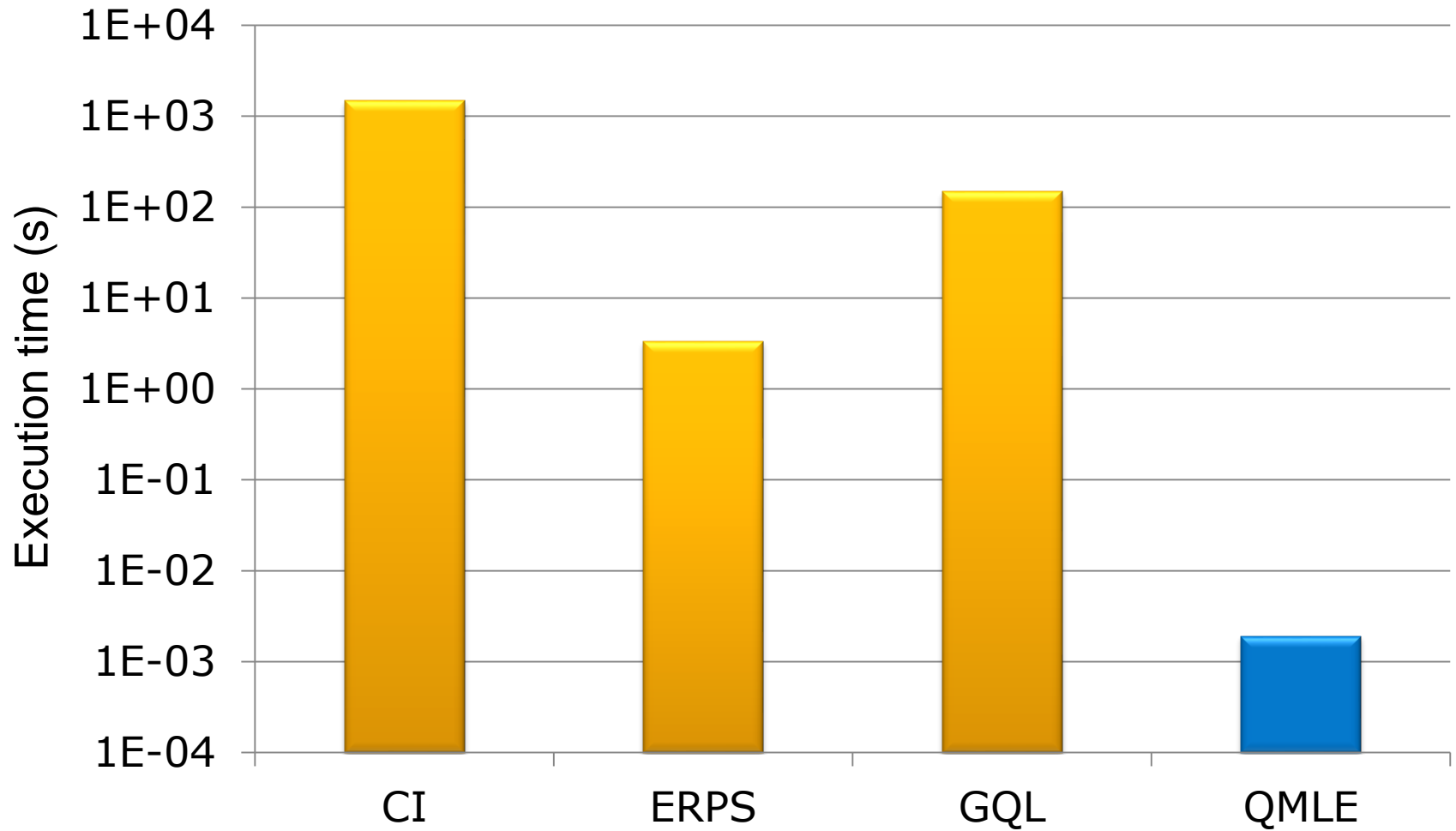
- Gibbs sampling based on queue length samples
- Many iterations until convergence

- ≈ 20000 random models
 - Randomized number of stations, classes, jobs
 - Focus on QMLE instead of exact analysis

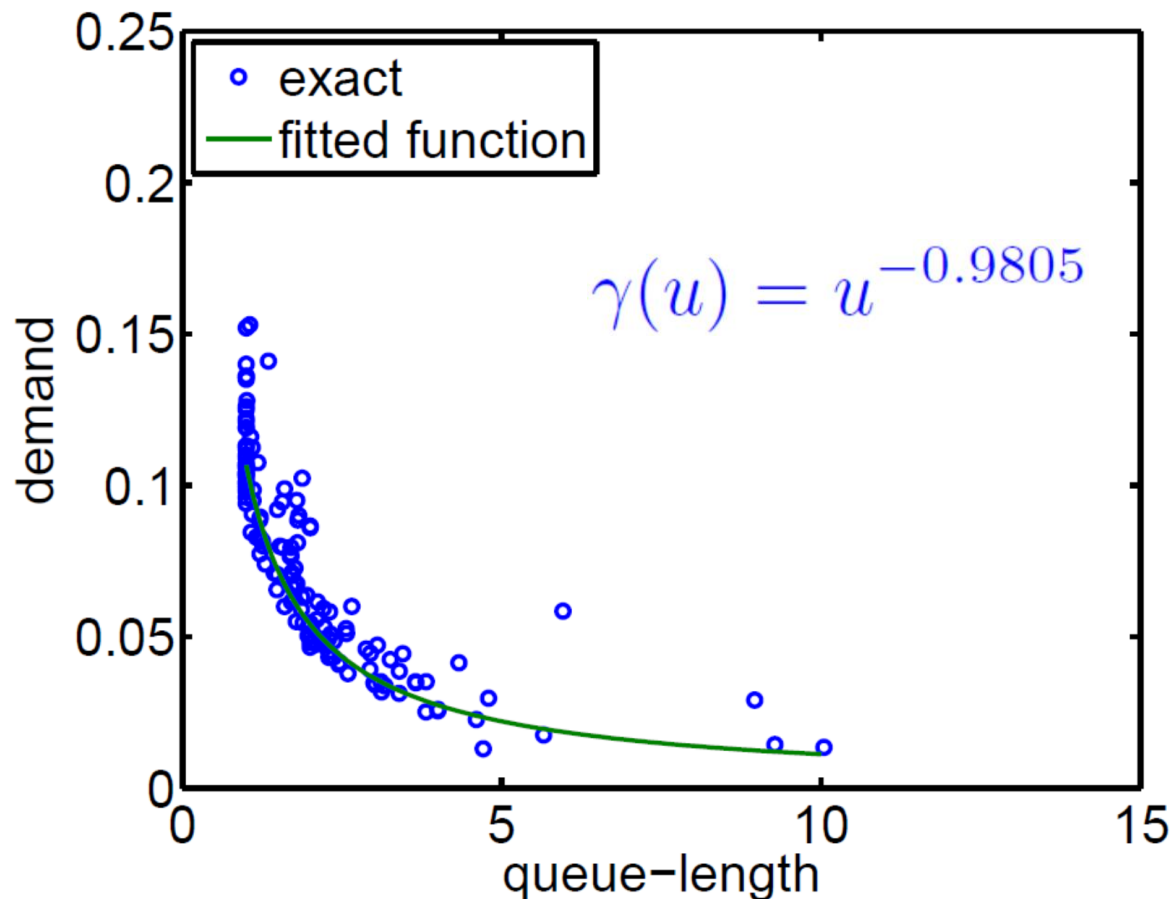
■ Results

- All the algorithms: below 10%
- QMLE has less than 4% error
- Confidence interval validated





- Mean demand varies under different load
- Real world system behavior
 - e.g. multi-core servers



- A scaling factor function $\gamma_i(u)$
 - load-independent : $\gamma_i(u) = 1, 1 \leq u \leq n_i$

- Product-form still holds

$$\mathbb{P}(\mathbf{n}|\boldsymbol{\theta}, \boldsymbol{\gamma}) = \left(\prod_{j=1}^R \frac{\theta_{0j}^{n_{0j}}}{n_{0j}!} \right) \prod_{i=1}^M n_i! \prod_{j=1}^R \frac{\theta_{ij}^{n_{ij}}}{n_{ij}! G(\boldsymbol{\theta})} \prod_{u=1}^{n_i} \gamma_i(u)$$

new term

- MLE

$$(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\gamma}}) = \arg \max_{(\boldsymbol{\theta}, \boldsymbol{\gamma}) \in \Theta} \prod_{l=1}^L \mathbb{P}(\mathbf{n}^l | \boldsymbol{\theta}, \boldsymbol{\gamma})$$

- Directly computation is infeasible
- A necessary condition for a point $(\hat{\theta}, \hat{\gamma})$ inside Θ to be a MLE is that

$$Q_{ij}(\hat{\theta}, \hat{\gamma}) = \tilde{Q}_{ij}(\mathbf{D}), \quad \forall i, j$$

and

$$\mathbb{P}(n_k = v | \hat{\theta}, \hat{\gamma}) = \mathbb{P}(\tilde{n}_k = v | \mathbf{D}), \quad \forall k, v$$

*Theoretical marginal
queue length probability*

*Empirical marginal
queue length probability*

- Works with marginal probability only!

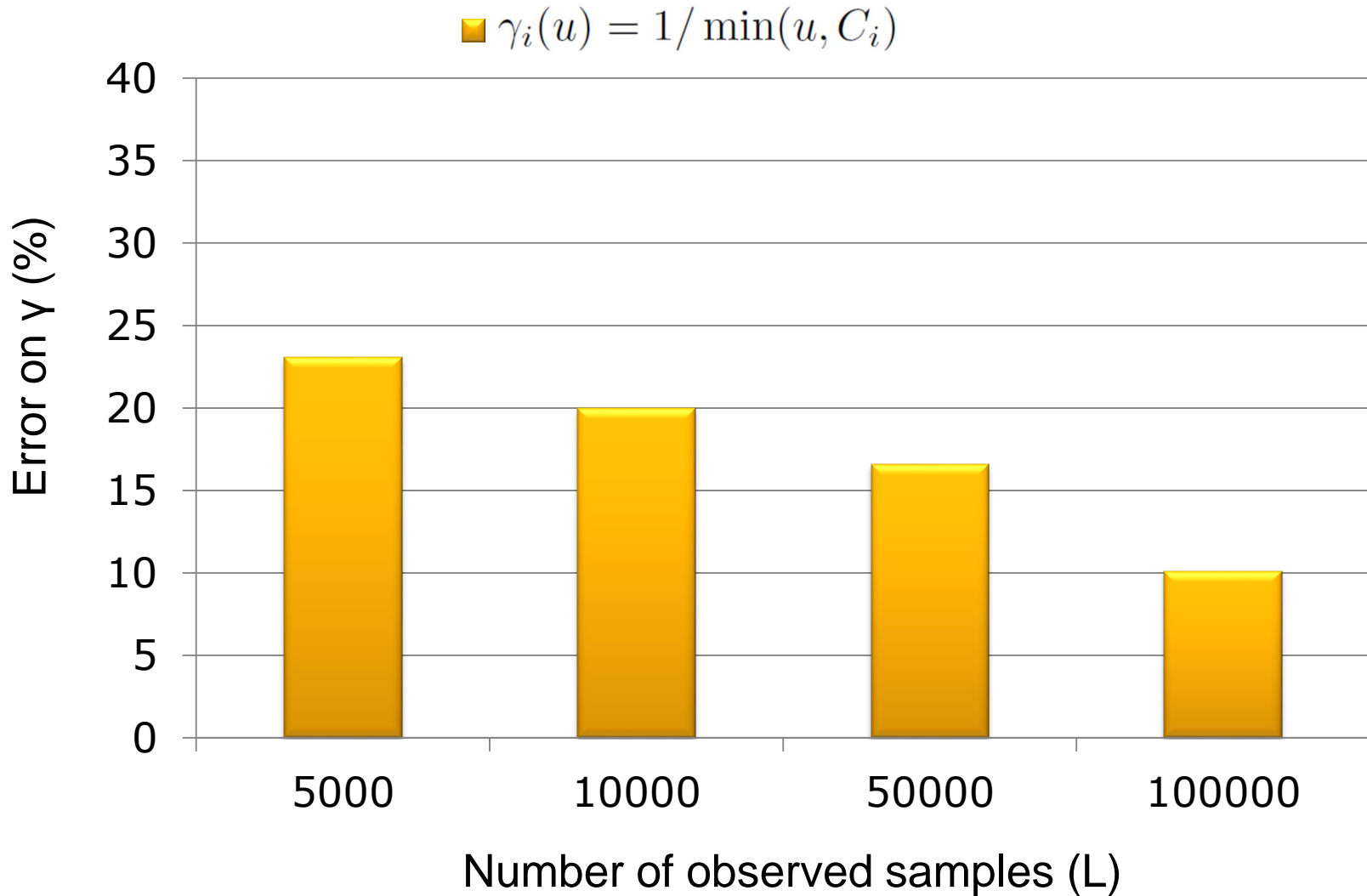
- How to find the MLE?
 - Solve by optimization program
- Confidence intervals
 - Hessian matrix can still be derived
 - Computation requires marginal probabilities and mean queue length only
- Drawback
 - Computationally expensive because of LD-MVA

■ Random models validation

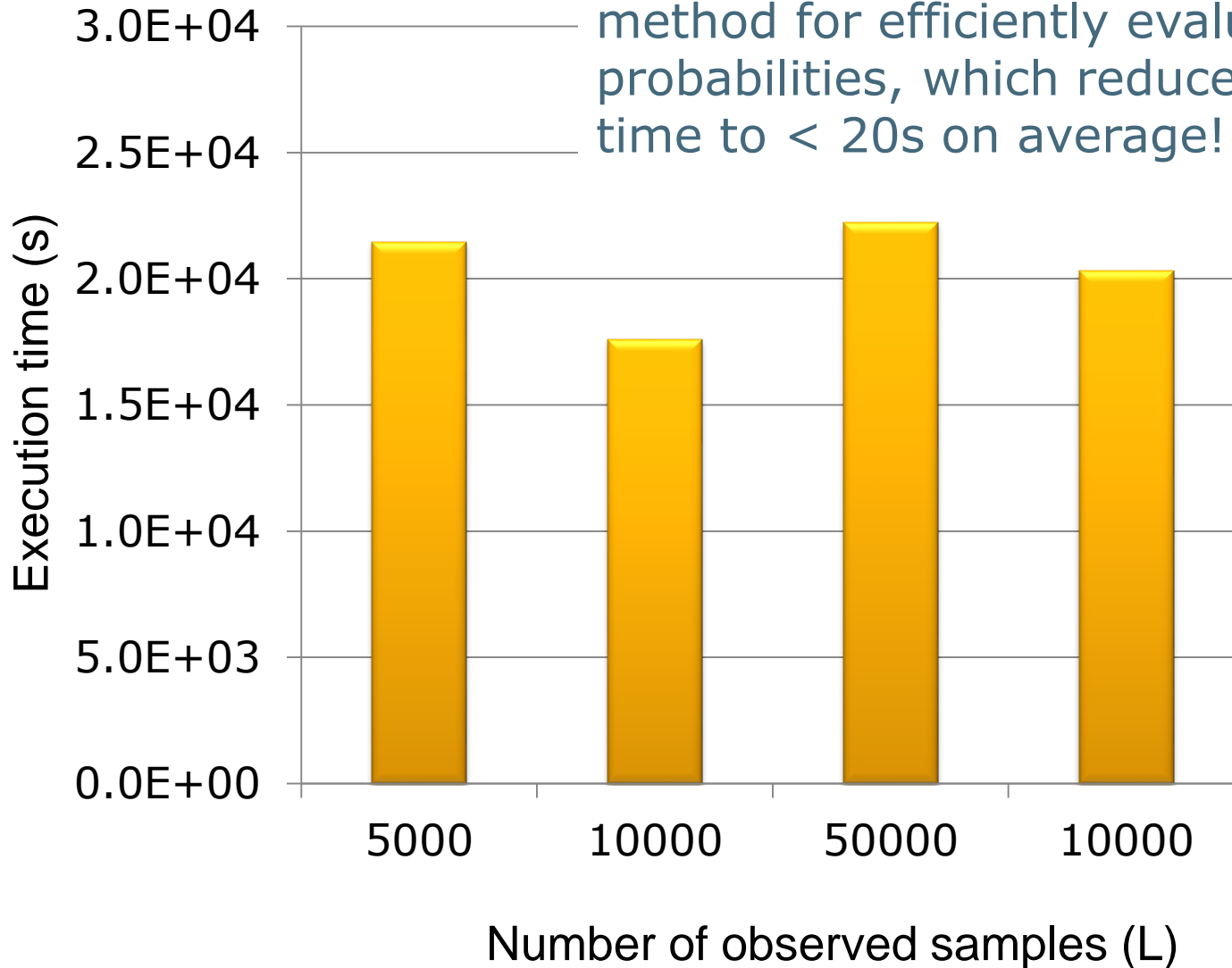
- 2 stations, 2 classes, 8 jobs, different think time
- MATLAB *fmincon* solver
- Compare the estimated $(\hat{\theta}, \hat{\gamma})$ against exact ones

■ Considered scaling factors

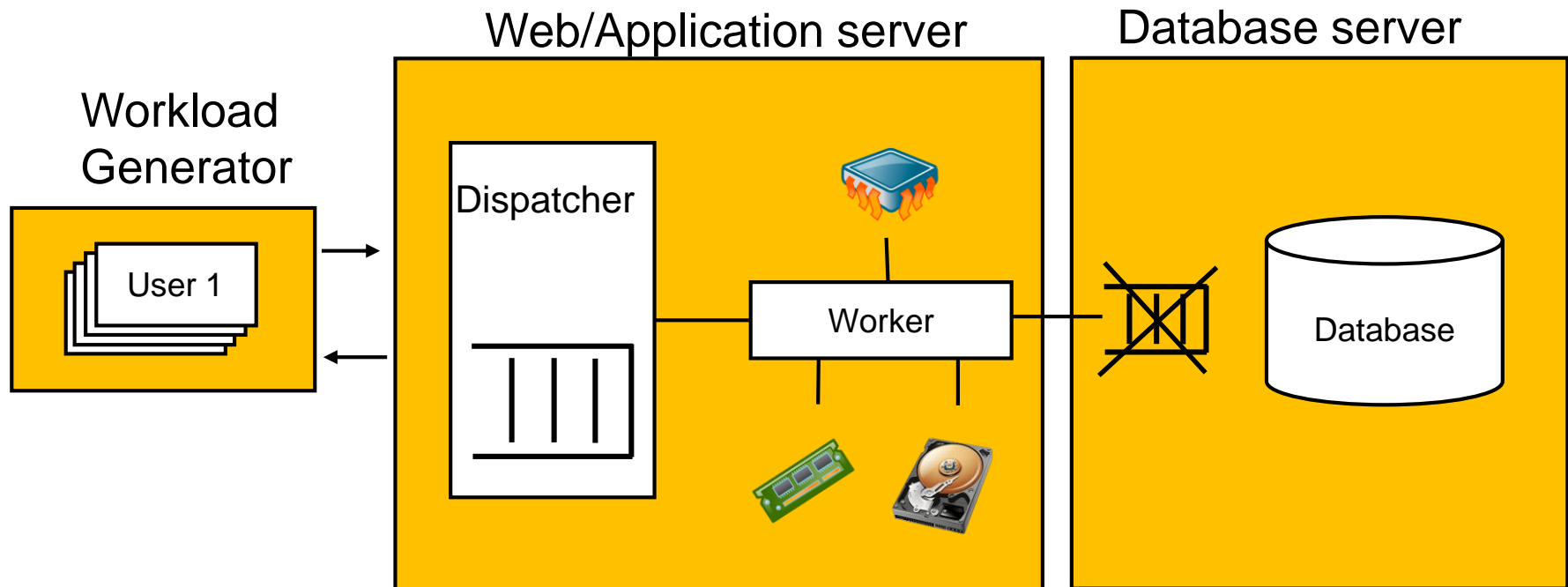
- $\gamma_i(u) = 1 / \min(u, C_i)$: resembles multi-core feature
 - C_i number of CPUs in queueing station i .



Progress: We found a new approximation method for efficiently evaluating marginal probabilities, which reduces the execution time to $< 20s$ on average!

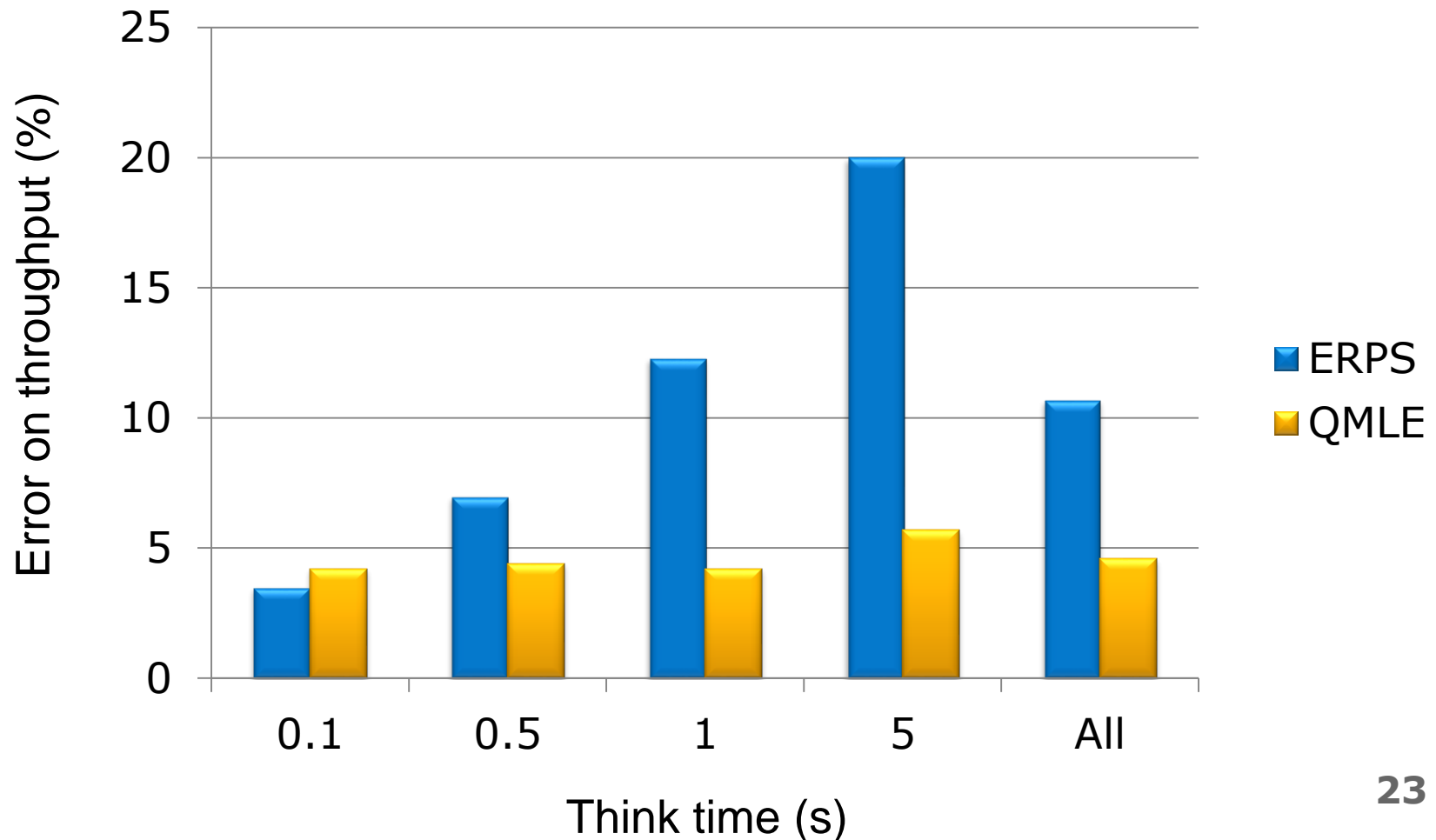


- 3-tier commercial application
 - Transactions grouped in R=1 class
 - 5 GB user data



■ Exact demand unknown

- Estimated demands using QMLE
- Validate observed throughput with estimated demands



- Demand estimation from queue length
 - Efficient
 - Confidence interval characterization
 - Load-dependent extension
- Ongoing work
 - Accelerate the load-dependent estimation
 - More experimental evaluations

Funded by FP7 MODAClouds, H2020 DICE, EPSRC OptiMAM

Thanks!



weikun.wang11@imperial.ac.uk