

Model-based performance self-adaptation

Emilio Incerto and Mirco Tribastone

IMT School for Advanced Studies Lucca, Italy
{emilio.incerto,mirco.tribastone}@imtlucca.it

ICPE'19 Tutorial
Mumbai, India

- Performance is a fundamental aspect of computing systems:
*“[. . .] in many practical deployment scenarios, particularly mobile, **performance is the new correctness.**”¹*
- Reasoning about performance is hard:
 - Large input spaces, measurement and testing becomes difficult.
 - Behavior dependent on interplay between software and hardware, can be machine dependent.
 - Usage profiles may be uncertain, performance analysis done at design time **may not be sufficient.**

¹Mark Harman and Peter O’Hearn. “From Start-ups to Scale-ups: Opportunities and Open Problems for Static and Dynamic Program Analysis”. In: *SCAM*. 2018.

AWS Auto Scaling

[Overview](#) [Features](#) [Pricing](#) [Resources](#) [FAQs](#)

AWS Auto Scaling

Application scaling to optimize performance and costs

Get started with AWS Auto Scaling

AWS Auto Scaling monitors your applications and automatically adjusts capacity to maintain steady, predictable performance at the lowest possible cost. Using AWS Auto Scaling, it's easy to setup application scaling for multiple resources across multiple services in minutes. The service provides a simple, powerful user interface that lets you build scaling plans for resources including [Amazon EC2](#) instances and Spot Fleets, [Amazon ECS](#) tasks, Amazon DynamoDB tables and indexes, and [Amazon Aurora](#) Replicas. AWS Auto Scaling makes scaling simple with recommendations that allow you to optimize performance, costs, or balance between them. If you're already using [Amazon EC2 Auto Scaling](#) to dynamically scale your Amazon EC2 instances, you can now combine it with AWS Auto Scaling to scale additional resources for other AWS services. With AWS Auto Scaling, your applications always have the right resources at the right time.

It's easy to get started with AWS Auto Scaling using the [AWS Management Console](#), Command Line Interface (CLI), or SDK. AWS Auto Scaling is available at no additional charge. You pay only for the AWS resources needed to run your applications and [Amazon CloudWatch](#) monitoring fees.

Introducing Predictive Scaling

Automatically scale your compute capacity in advance of traffic changes using ML technology.

[Read the Blog>>](#)

Amazon EC2 Auto Scaling has moved

You have several options for scaling with AWS. Do you only need to manage Amazon EC2 instances?

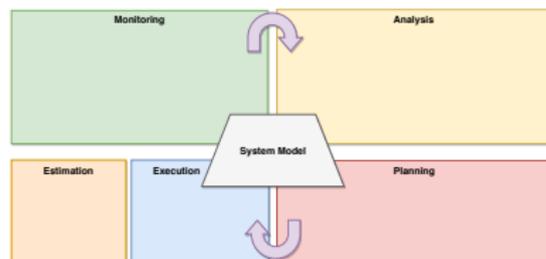
[Visit the Amazon EC2 Auto Scaling page »](#)

Amazon Web Services

Objective: Design systems that adapt themselves to changing environments while meeting desired performance-based service-level agreements (throughput, response time, utilization).

Three main activities:

- 1 monitor the system execution;
- 2 continuously update a model of the system;
- 3 trigger reconfigurations when required.



MAPE-K control loop

Objective: Design systems that adapt themselves to changing environments while meeting desired performance-based service-level agreements (throughput, response time, utilization)

Three main activities:

- 1 monitor the system execution;
- 2 continuously update a model of the system;
- 3 trigger reconfigurations when required.

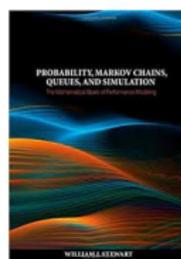
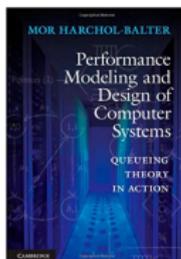
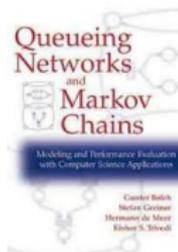
Three main difficulties:

- 1 non-intrusive monitoring infrastructure;
- 2 ***efficient and accurate predictive models;***
- 3 effective and robust planner.

- 1 Efficient performance models
- 2 Symbolic performance adaptation
- 3 Model-predictive control for performance-driven self-adaptation
- 4 Moving horizon estimation of service demands

- 1** Efficient performance models
- 2 Symbolic performance adaptation
- 3 Model-predictive control for performance-driven self-adaptation
- 4 Moving horizon estimation of service demands

- We build on queuing networks, an established class of performance models for computing systems.
- Need to convince ICPE audience?
 - Varsha Apte's workshop was yesterday...
 - Kishor Trivedi speaks tomorrow...



varsha apte
Professor, CSE, IIT Bombay

Home PANDA SYNERG Teaching Research For Students About Me Environment

Performance Engineering Training Workshop and Hackathon. April 7th, 2019

Co-located with ICPE 2019, JIT Bombay, 7th-11th April 2019
Available only with ICPE 2019 Conference Registration

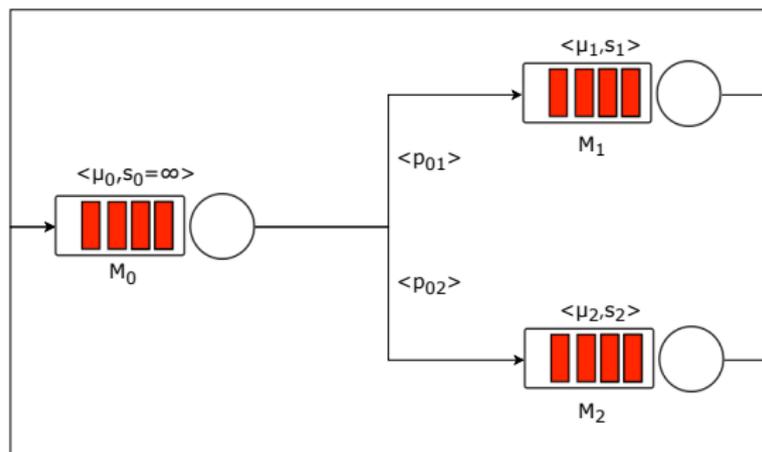
Response time vs number of users

How many metrics and parameters of the system can be predicted by just observing this graph? Turns out, you can deduce a lot if you know how to apply queuing theory to your load testing results. Knowing queuing theory can help make sense of load testing results, quickly spot experimental errors, and shorten the time to complete tests due to clever planning of tests. In this training event co-located with ICPE 2019, in the morning session, you will be introduced to the fundamental results from queuing theory, and will learn how to apply them during load testing. Numerous exercises will be given to reinforce your learning. In the afternoon

Performance Sherlock (Morning Session)

Outline

- Introduction to performance metrics and parameters
- Introduction to Open Queuing systems
- Asymptotic Analysis of Open queuing systems
- Basic Laws of queuing systems (Little's Law, etc)
- Introduction to Closed Queuing systems
- Asymptotic Analysis of Closed queuing systems
- Mean Value Analysis
- Practice exercises



A load balancing system

- Exponentially distributed service times—can be relaxed;²
- single class of users—multiple classes later;
- multiple servers.

²Matthias Kowal et al. “Scaling Size and Parameter Spaces in Variability-Aware Software Performance Models”. In: *ASE*. 2015.

- Queuing networks enjoy efficient analysis techniques (e.g., BCMP theorem), **but**:
 - In general, each parameterization requires a distinct analysis — large cost when exploring large parameter spaces.

Our contribution:

Efficient parametric analysis of queuing networks.

- Analytical results assume time-homogeneous networks, i.e., parameters do not vary with time — adaptation requires changing the system, hence the parameters of its model.
- Most analytical results available for the steady-state regime — adaptation might need to change parameters at all time points.

Our contribution:

Approximate queuing network analysis based on fluid models.

- Consider a queuing network with probability matrix P , vector of exogenous arrivals λ , and service rates μ .
- Consider concrete as well as symbolic parameters.
- Solve for stationary average performance indices using a computer algebra system.³

$$P = \begin{bmatrix} 0.0 & p_{1,2}^* & p_{1,3}^* & 0.0 \\ 0.0 & 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 0.0 & 1.0 \\ 0.5 & 0.0 & 0.0 & 0.0 \end{bmatrix} \quad \lambda = \begin{bmatrix} \lambda_1^* \\ 0.0 \\ 0.0 \\ 0.1 \end{bmatrix} \quad \mu = \begin{bmatrix} 1.5 \\ 2.0 \\ 3.0 \\ \mu_4^* \end{bmatrix}$$

³Matthias Kowal, Ina Schaefer, and Mirco Tribastone. “Family-Based Performance Analysis of Variant-Rich Software Systems”. In: *FASE*. 2014.

- Consider a queuing network with probability matrix P , vector of exogenous arrivals λ , and service rates μ .
- Consider concrete as well as symbolic parameters.
- Solve for stationary average performance indices using a computer algebra system.⁴

$$\begin{aligned}
 \text{AQL} = & -\frac{20\lambda_1^* + 1}{80\lambda_1^* + 60\rho_{1,2}^* + 60\rho_{1,3}^* - 116} \\
 & - \frac{\rho_{1,2}^*(20\lambda_1^* + 1)}{4(21\rho_{1,2}^* + 20\rho_{1,3}^* + 20\lambda_1^*\rho_{1,2}^* - 40)} \\
 & - \frac{10\lambda_1^*\rho_{1,2}^* + 10\lambda_1^*\rho_{1,3}^* + 1}{20\mu_4^*\rho_{1,2}^* - 40\mu_4^* + 20\mu_4^*\rho_{1,3}^* + 40\lambda_1^*\rho_{1,2}^* + 40\lambda_1^*\rho_{1,3}^* + 4} \\
 & - \frac{\rho_{1,3}^*(20\lambda_1^* + 1)}{4(30\rho_{1,2}^* + 31\rho_{1,3}^* + 20\lambda_1^*\rho_{1,2}^* - 60)}
 \end{aligned}$$

⁴Matthias Kowal, Ina Schaefer, and Mirco Tribastone. “Family-Based Performance Analysis of Variant-Rich Software Systems”. In: *FASE*. 2014.

<i>Variables</i>				<i>Runtimes (s)</i>			
<i>Nodes</i>	<i>p</i>	<i>m</i>	<i>g</i>	<i>SYM</i>	<i>NUM</i>	<i>NUM/SYM</i>	<i>PC</i>
4	1	0	0	0.011	0.049	4.47	0.545
4	0	1	0	0.004	0.043	10.78	0.185
4	0	0	1	0.009	0.045	4.92	0.190
4	1	1	1	0.011	0.046	4.13	0.214
4	2	2	2	0.014	0.049	3.60	0.319
4	4	4	4	0.020	0.049	2.43	0.310
24	1	0	0	0.011	0.066	5.96	1.141
24	0	1	0	0.004	0.063	14.60	1.124
24	0	0	1	0.011	0.067	6.29	1.152
24	1	1	1	0.013	0.068	5.04	1.244
24	2	2	2	0.016	0.068	4.31	1.100
24	0	6	0	0.007	0.065	9.94	1.004
24	4	4	4	0.099	0.070	0.70	1.904

<i>Variables</i>				<i>Runtimes (s)</i>			
<i>Nodes</i>	<i>p</i>	<i>m</i>	<i>g</i>	<i>SYM</i>	<i>NUM</i>	<i>NUM/SYM</i>	<i>PC</i>
142	1	0	0	0.016	6.540	397.34	5.425
142	0	1	0	0.005	8.107	1664.28	4.908
142	0	0	1	0.015	10.808	726.41	4.836
142	1	1	1	0.017	10.069	601.51	5.113
142	2	2	2	0.019	10.052	526.50	4.936
142	0	6	0	0.007	7.802	1191.79	5.137
142	4	4	4	0.026	10.985	429.83	4.942
302	1	0	0	0.019	19.186	1007.95	11.026
302	0	1	0	0.006	13.680	2292.46	11.192
302	0	0	1	0.018	13.399	728.73	11.050
302	1	1	1	0.021	19.520	909.12	11.591
302	2	2	2	0.024	19.459	820.30	11.214
302	0	6	0	0.006	13.896	2258.21	11.089
302	4	4	4	0.030	14.879	495.06	11.718

- Queuing networks enjoy efficient analysis techniques (e.g., BCMP theorem), **but**:
 - In general, each parameterization requires a distinct analysis — large cost when exploring large parameter spaces.

Our contribution:

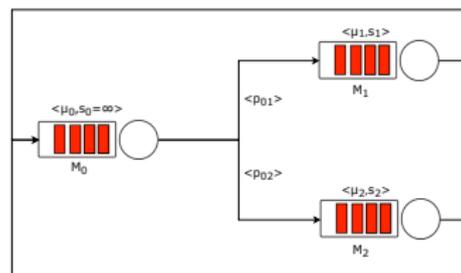
Efficient parametric analysis of queuing networks.

- Analytical results assume time-homogeneous networks, i.e., parameters do not vary with time — adaptation requires changing the system, hence the parameters of its model.
- Most analytical results available for the steady-state regime — adaptation might need to change parameters at all time points.

Our contribution:

Approximate queuing network analysis based on fluid models.

Queuing network



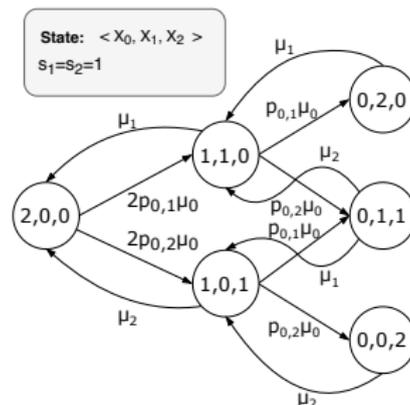
Forward (master) equations

$$\dot{\pi}(2,0,0) = -2\mu_0\pi(2,0,0) + \mu_1\pi(1,1,0) + \mu_2\pi(1,0,1)$$

⋮

$$\dot{\pi}(0,0,2) = -\mu_2\pi(0,0,2) + p_{0,2}\mu_0\pi(1,0,1)$$

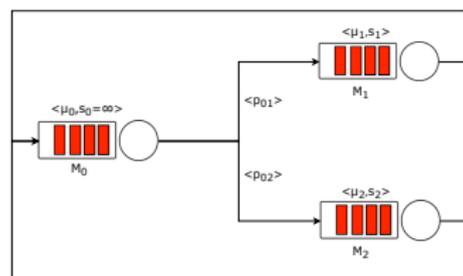
Underlying continuous-time Markov chain



State explosion

Number of states is exponential!

Queuing network



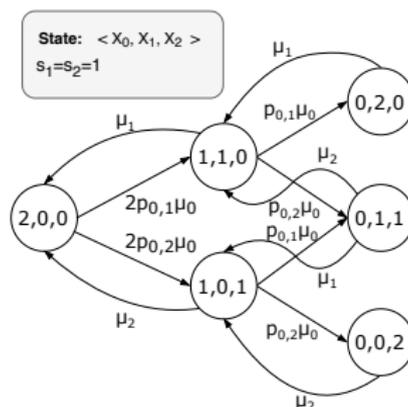
Forward (master) equations

$$\dot{\pi}(2,0,0) = -2\mu_0\pi(2,0,0) + \mu_1\pi(1,1,0) + \mu_2\pi(1,0,1)$$

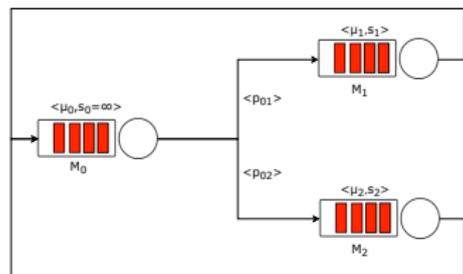
$$\vdots$$

$$\dot{\pi}(0,0,2) = -\mu_2\pi(0,0,2) + p_{0,2}\mu_0\pi(1,0,1)$$

Underlying continuous-time Markov chain



<i>Clients</i>	<i>States</i>	<i>Transitions</i>
2	6	71
10	66	220
20	231	840
40	861	3280
80	3321	12960

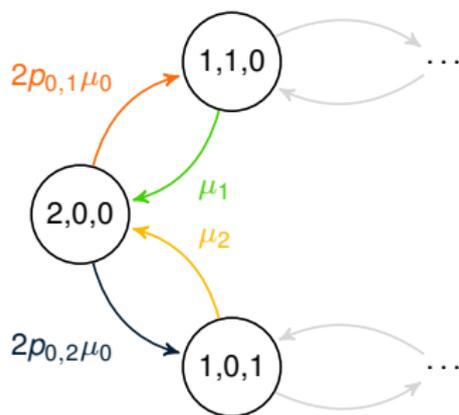


$$Q_0 \rightarrow Q_1, \quad \rho_{0,1} \mu_0 Q_0$$

$$Q_0 \rightarrow Q_2, \quad \rho_{0,2} \mu_0 Q_0$$

$$Q_1 \rightarrow Q_0, \quad \mu_1 \min\{Q_1, s_1\}$$

$$Q_2 \rightarrow Q_0, \quad \mu_2 \min\{Q_2, s_2\}$$



$$\Downarrow$$

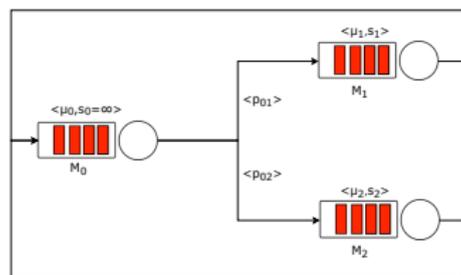
$$l_1 = (-1, +1, 0), \quad f_1(x) = \rho_{0,1} \mu_0 x_0$$

$$l_2 = (-1, 0, -1), \quad f_2(x) = \rho_{0,2} \mu_0 x_0$$

$$l_3 = (+1, -1, 0), \quad f_3(x) = \mu_1 \min\{x_1, s_1\}$$

$$l_4 = (+1, 0, -1), \quad f_4(x) = \mu_2 \min\{x_2, s_2\}$$

Jumps and rate functions



$$Q_0 \rightarrow Q_1, \quad p_{0,1} \mu_0 Q_0$$

$$Q_0 \rightarrow Q_2, \quad p_{0,2} \mu_0 Q_0$$

$$Q_1 \rightarrow Q_0, \quad \mu_1 \min\{Q_1, s_1\}$$

$$Q_2 \rightarrow Q_0, \quad \mu_2 \min\{Q_2, s_2\}$$

$$\Downarrow$$

Master Equation

$$\dot{\pi}_n(t) = - \sum_i f_i(n) \pi_n(t) + \sum_i f_i(n - l_i) \pi_{n-l_i}(t)$$

for all states $n = (n_0, n_1, n_2)$.

$$l_1 = (-1, +1, 0), \quad f_1(x) = p_{0,1} \mu_0 x_0$$

$$l_2 = (-1, 0, -1), \quad f_2(x) = p_{0,2} \mu_0 x_0$$

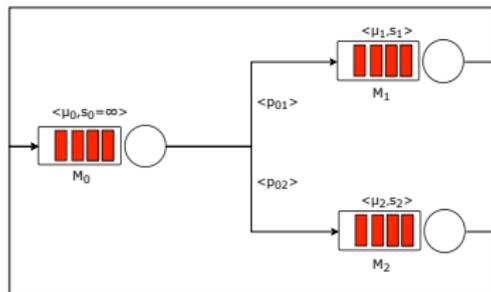
$$l_3 = (+1, -1, 0), \quad f_3(x) = \mu_1 \min\{x_1, s_1\}$$

$$l_4 = (+1, 0, -1), \quad f_4(x) = \mu_2 \min\{x_2, s_2\}$$

Jumps and rate functions

Master equation

$$\dot{\pi}_n(t) = - \sum_i f_i(n) \pi_n(t) + \sum_i f_i(n - l_i) \pi_{n-l_i}(t), \quad \text{for all states } n.$$



- Call $X(t) = (X_0(t), X_1(t), X_2(t))$ the queue length process.
- The true equations for the average queue lengths are:

$$\mathbb{E}[\dot{X}](t) = \sum_i l_i \mathbb{E} [f_i(X(t))]$$

- **Problem: equations are not closed!**

$$\begin{aligned} l_1 &= (-1, +1, 0), & f_1(x) &= p_{0,1} \mu_0 x_0 \\ l_2 &= (-1, 0, -1), & f_2(x) &= p_{0,2} \mu_0 x_0 \\ l_3 &= (+1, -1, 0), & f_3(x) &= \mu_1 \min\{x_1, s_1\} \\ l_4 &= (+1, 0, -1), & f_4(x) &= \mu_2 \min\{x_2, s_2\} \end{aligned}$$

True equations

$$\mathbb{E}[\dot{X}](t) = \sum_i l_i \mathbb{E}[f_i(X(t))]$$

$$l_1 = (-1, +1, 0), \quad f_1(x) = \rho_{0,1} \mu_0 x_0$$

$$l_2 = (-1, 0, -1), \quad f_2(x) = \rho_{0,2} \mu_0 x_0$$

$$l_3 = (+1, -1, 0), \quad f_3(x) = \mu_1 \min\{x_1, s_1\}$$

$$l_4 = (+1, 0, -1), \quad f_4(x) = \mu_2 \min\{x_2, s_2\}$$

$$\begin{aligned} \mathbb{E}[\dot{X}_0](t) &= -\rho_{0,1} \mu_0 \mathbb{E}[X_0](t) - \rho_{0,2} \mu_0 \mathbb{E}[X_0](t) \\ &\quad + \mathbb{E}[\mu_1 \min\{X_1(t), s_1\}] + \mathbb{E}[\mu_2 \min\{X_2(t), s_2\}] \end{aligned}$$

$$\mathbb{E}[\dot{X}_1](t) = +\rho_{0,1} \mu_0 \mathbb{E}[X_0](t) - \mathbb{E}[\mu_1 \min\{X_1(t), s_1\}]$$

$$\mathbb{E}[\dot{X}_2](t) = +\rho_{0,2} \mu_0 \mathbb{E}[X_0](t) - \mathbb{E}[\mu_2 \min\{X_2(t), s_2\}]$$

How to close the equations? $\mathbb{E}[\min\{X, Y\}] \approx \min\{\mathbb{E}[X], \mathbb{E}[Y]\}$

True equations

$$\dot{\mathbb{E}}[X](t) = \sum_i l_i \mathbb{E} [f_i(X(t))]$$

$$l_1 = (-1, +1, 0), \quad f_1(x) = \rho_{0,1} \mu_0 x_0$$

$$l_2 = (-1, 0, -1), \quad f_2(x) = \rho_{0,2} \mu_0 x_0$$

$$l_3 = (+1, -1, 0), \quad f_3(x) = \mu_1 \min\{x_1, s_1\}$$

$$l_4 = (+1, 0, -1), \quad f_4(x) = \mu_2 \min\{x_2, s_2\}$$

Fluid Approximation

$$\begin{aligned} \dot{X}_0(t) = & -\rho_{0,1} \mu_0 X_0(t) - \rho_{0,2} \mu_0 X_0(t) \\ & + \mu_1 \min\{X_1(t), s_1\} + \mu_2 \min\{X_2(t), s_2\} \end{aligned}$$

$$\dot{X}_1(t) = +\rho_{0,1} \mu_0 X_0(t) - \mu_1 \min\{X_1(t), s_1\}$$

$$\dot{X}_2(t) = +\rho_{0,2} \mu_0 X_0(t) - \mu_2 \min\{X_2(t), s_2\}$$

How to close the equations?

$$\mathbb{E}[\min\{X, Y\}] \approx \min\{\mathbb{E}[X], \mathbb{E}[Y]\}$$

Markov population process

$$l_1 = (-1, +1, 0), \quad f_1(x) = \rho_{0,1} \mu_0 x_0$$

$$l_2 = (-1, 0, -1), \quad f_2(x) = \rho_{0,2} \mu_0 x_0$$

$$l_3 = (+1, -1, 0), \quad f_3(x) = \mu_1 \min\{x_1, K s_1\}$$

$$l_4 = (+1, 0, -1), \quad f_4(x) = \mu_2 \min\{x_2, K s_2\}$$

CTMC family

$$X^1(0) = (Q_0, Q_1, Q_2)$$

$$X^2(0) = 2(Q_0, Q_1, Q_2)$$

...

$$X^K(0) = K(Q_0, Q_1, Q_2)$$

- Fix some initial condition (Q_0, Q_1, Q_2)
- Consider Markov chains with increasing clients as well as servers
- In the limit as K goes to infinity **one sample path** of $X^K(t)/K$ will converge to the fluid equations⁵

⁵T. G. Kurtz. "Solutions of ordinary differential equations as limits of pure Markov processes". In: *J. Appl. Prob.* (1970).

- Let $1, \dots, N$ be the queuing stations;
- let $P = (p_{i,j})_{1 \leq i,j \leq N}$ be the routing probability matrix;
- let s_1, \dots, s_N be the server multiplicities;
- let μ_1, \dots, μ_N be the service rates.

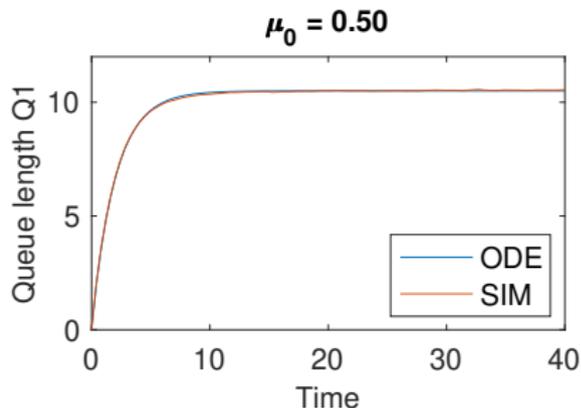
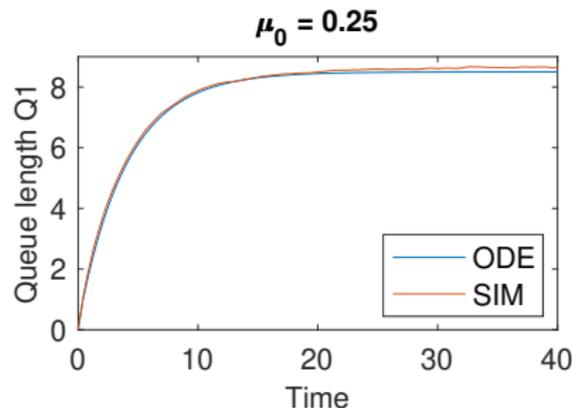
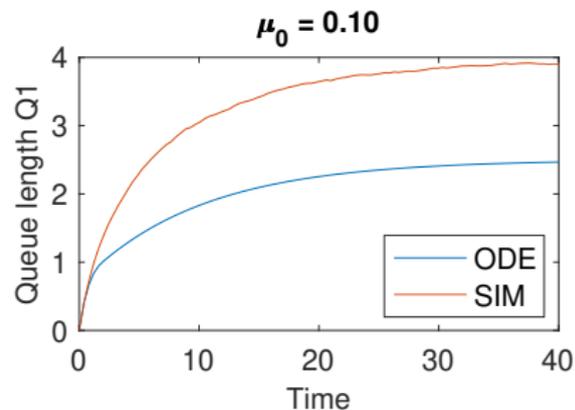
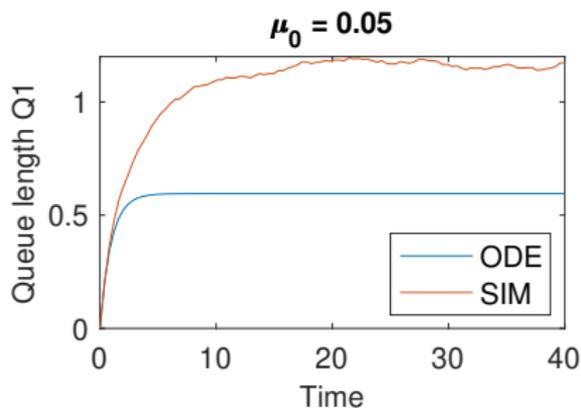
The fluid approximation is given by⁶

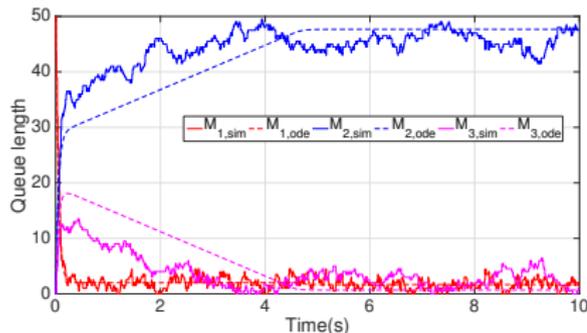
$$\dot{X}_i(t) = -\mu_i \min\{X_i(t), s_i\} + \sum_j p_{j,i} \mu_j \min\{X_j(t), s_j\}, \quad i = 1, \dots, N.$$

It can be also defined for time-varying parameters:

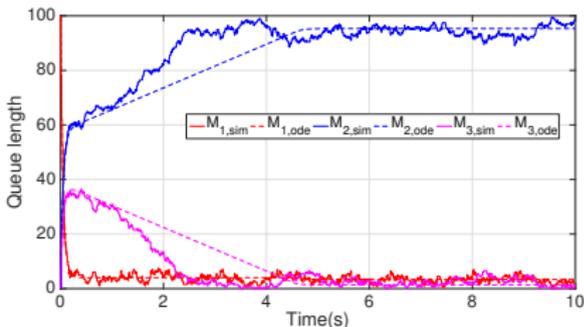
$$\dot{X}_i(t) = -\mu_i(t) \min\{X_i(t), s_i(t)\} + \sum_j p_{j,i}(t) \mu_j(t) \min\{X_j(t), s_j(t)\}$$

⁶Mirco Tribastone. "A Fluid Model for Layered Queueing Networks". In: *IEEE Trans. Software Eng.* (2013).

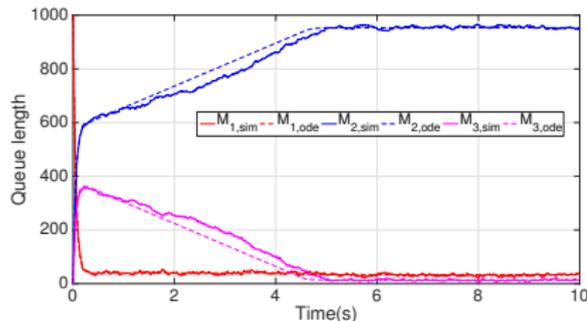




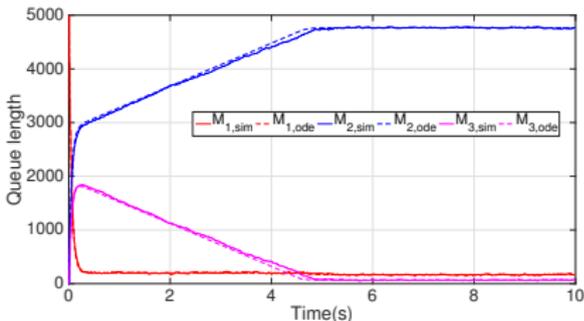
(a) $K=1$



(b) $K=2$



(c) $K=20$



(d) $K=100$

Comparison between a sample path of the CTMC and the fluid model

The screenshot displays the ERODE software interface with the following components:

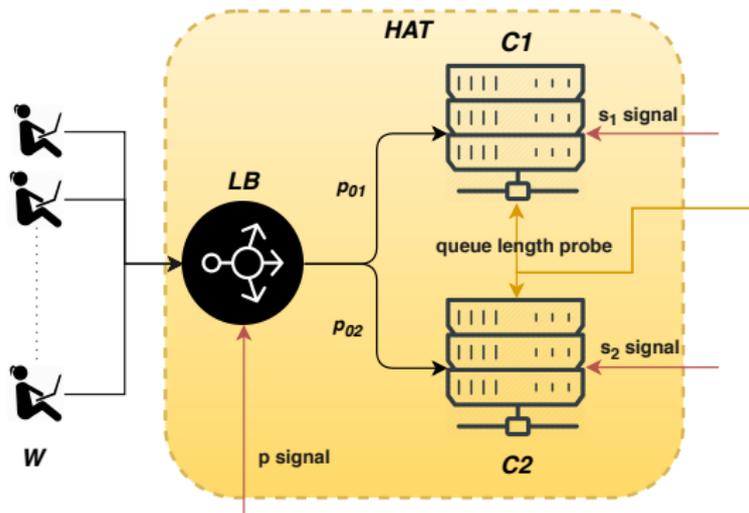
- Project Explorer:** Shows a tree view of the project files, including 'Examples', 'ExampleODE ode', 'ExampleRN ode', and 'InfluenceNetworks'.
- Outline:** Lists the structure of the 'ExampleODE' model, including parameters, initial conditions, ODEs, and views.
- Code Editor (Left):** Contains the source code for 'ExampleODE.ode', including parameters ($r1 = 1.0, r2 = 2.0$), initial conditions ($Au = 1.0, Ap = 2.0, B = 3.0, AuB = 4.0$), and a system of ODEs:

$$\begin{aligned} d(Au) &= -r1 \cdot Au + r2 \cdot Ap - 3 \cdot Au \cdot B + 4 \cdot AuB \\ d(Ap) &= r1 \cdot Au - r2 \cdot Ap - 3 \cdot Ap \cdot B + 4 \cdot ApB \\ d(B) &= -3 \cdot Au \cdot B + 4 \cdot AuB - 3 \cdot Ap \cdot B + 4 \cdot ApB \\ d(AuB) &= 3 \cdot Au \cdot B - 4 \cdot AuB \\ d(ApB) &= 3 \cdot Ap \cdot B - 4 \cdot ApB \end{aligned}$$
- Code Editor (Right):** Contains the source code for 'ExampleRN.ode', which defines a reduced model with similar parameters and initial conditions, and a set of reactions:

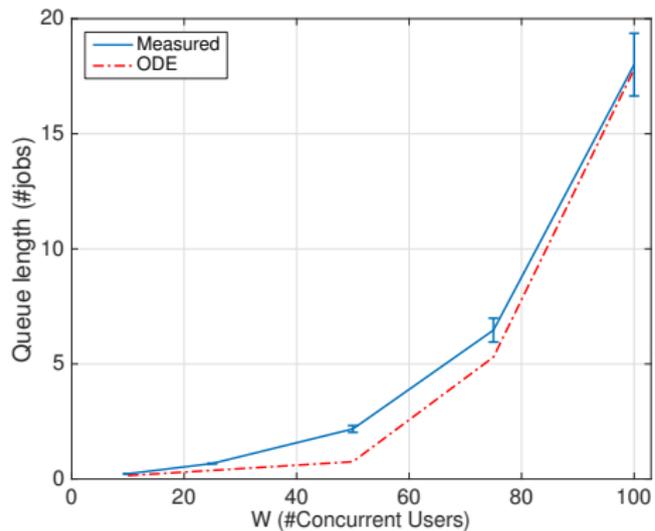
$$\begin{aligned} Au &\rightarrow Ap & r1 \\ Ap &\rightarrow Au & r2 \\ Au + B &\rightarrow AuB & 3.0 \\ AuB &\rightarrow Au + B & 4.0 \\ Ap + B &\rightarrow ApB & 3.0 \\ ApB &\rightarrow Ap + B & 4.0 \end{aligned}$$
- Plot:** A line graph titled 'simulateODE(End=1.0) ExampleRN - ODE solutions - All species/variables'. The y-axis is 'Species/variable concentrations' ranging from -0.3 to 3.3, and the x-axis is 'Time' ranging from -0.01 to 1.01. The plot shows the time evolution of Au (blue), Ap (red), B (green), AuB (purple), and ApB (magenta).
- Console:** Displays the execution log, including the time taken to solve the ODEs (0.086 s) and the number of parameters, species, and reactions.

<https://sysma.imtlucca.it/tools/erode/>

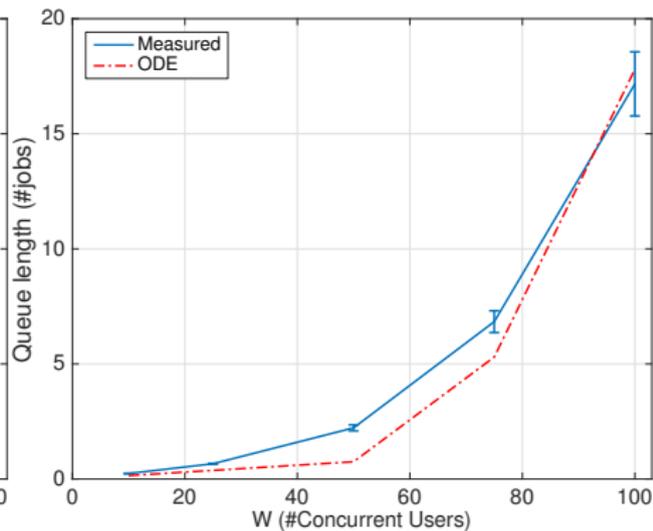
L. Cardelli, M. Tschaikowski, M. Tribastone, and A. Vandin. “ERODE: A tool for the evaluation and reduction of ordinary differential equations,” *TACAS’17*.



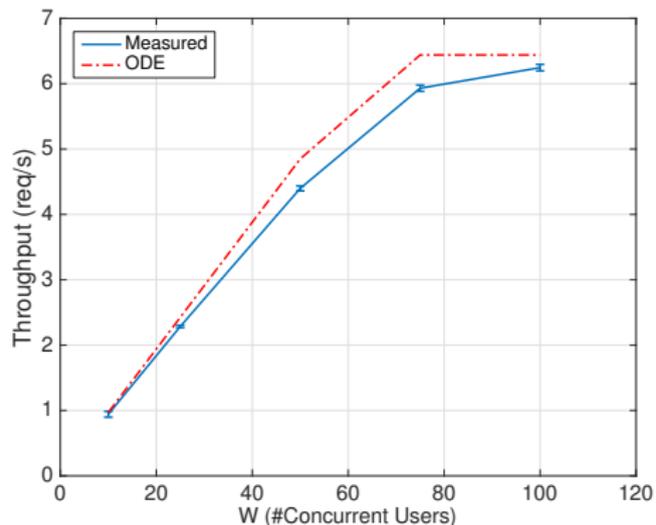
- In house developed web application
- It resembles CPU intensive behaviour
- Dispatcher (i.e., LB) and Computation Node (i.e., C1, C2)



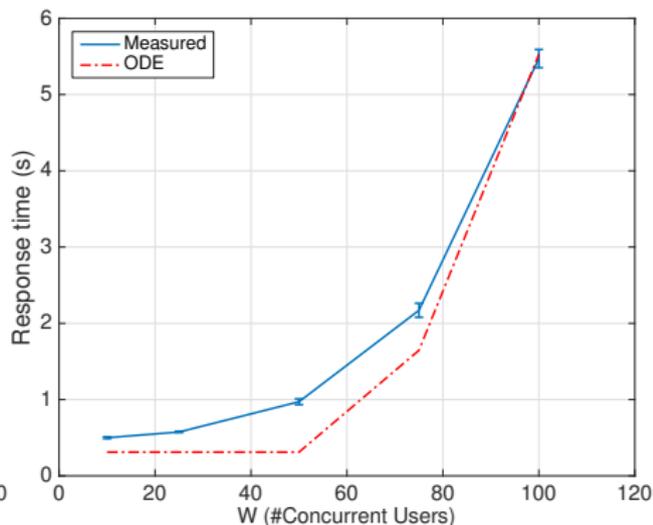
(a) N_1 queue length



(b) N_2 queue length



(a) System throughput



(b) System response time

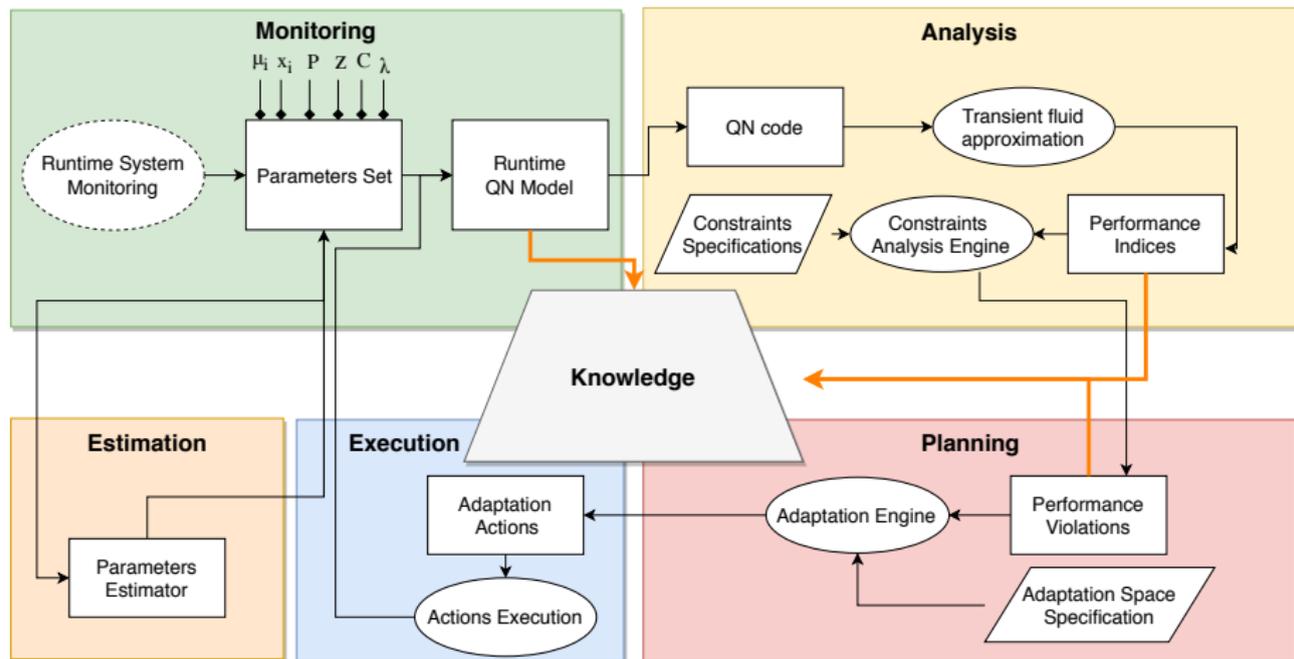
Objective: Design systems that adapt themselves to changing environments while meeting desired performance-based service-level agreements (throughput, response time, utilization)

Three main activities:

- 1 monitor the system execution;
- 2 continuously update a model of the system;
- 3 trigger reconfigurations when required.

Three main difficulties:

- 1 non-intrusive monitoring infrastructure;
- 2 efficient and accurate predictive models;
- 3 ***effective and robust planner.***



An MAPE-K framework based on fluid QN

- 1 Efficient performance models
- 2 Symbolic performance adaptation**
- 3 Model-predictive control for performance-driven self-adaptation
- 4 Moving horizon estimation of service demands

- 💡 **Idea:** Encode performance-driven adaptation as a satisfiability modulo theories (SMT) problem.⁷
- **Goal:** Query an SMT solver for a feasible assignment of the system parameters needed to satisfy QoS requirements or getting a formal proof of its non-existence.
- We rely on a combination of:
 - **queuing networks:** quantitative model to represent QoS attributes of the system;
 - **symbolic analysis:** represent all possible system configuration as a set of nonlinear real constraints;
 - **SMT:** devise feasible system configurations.

⁷Emilio Incerto, Mirco Tribastone, and Catia Trubiani. “Symbolic Performance Adaptation”. In: *SEAMS*. 2016.

- Satisfiability Modulo Theory⁸ checks the satisfiability of logical formulas over one or more theories, i.e., if there exists an assignment of real numbers satisfying

$$3x + 2y - z \geq 4 \wedge x \geq 0 \wedge y \geq 0$$

- We interpret the symbolic expressions as a set of constraints that the parametric values have to satisfy
- We augment them with further constraints representing the QoS requirements \mathcal{Q}



Exploiting the Non-Linear Real Arithmetic theory \mathcal{NRA} we encode this problem as an SMT problem suitable to devise the set of feasible parameters satisfying \mathcal{Q}

⁸Clark Barrett et al. “Handbook of Satisfiability: Volume 185 Frontiers in Artificial Intelligence and Applications”. In: ed. by A. Biere et al. IOS Press, 2009. Chap. Satisfiability Modulo Theories.

- \mathcal{M} be the symbolic solution of the QN
- \mathcal{Q} be a set of QoS requirements
- \mathcal{D} be a set of domain assumptions
- \mathcal{R} be a set of resource constraints

The QoS-based adaptation is turned into the satisfiability problem:

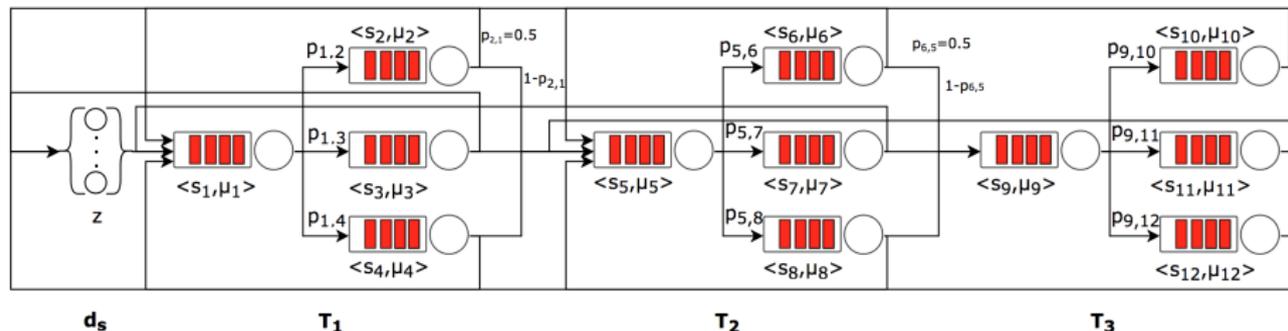
Find an assignment of the variables for model \mathcal{M} that ensures \mathcal{Q} subjected to constraints $\mathcal{D} \wedge \mathcal{R}$.

$$\mathcal{Q} := \{T_0 \geq C/2\}$$

$$\mathcal{D} := \{C = 350\}$$

$$\mathcal{R} := \left\{ \forall i, j \in \mathcal{S}. \sum_{k \in \mathcal{S}} p_{i,k} = 1.0 \wedge \right.$$

$$\left. 0 \leq p_{i,j} \leq 1.0 \wedge 1 \leq s_i \leq 40 \wedge 0.02 \leq \frac{1}{\mu_i} \leq 10 \right\}$$

QN model for a three-tier system⁹

Control objective

$$\mathcal{R}_s = C/T_1 \leq 1$$

$$C \in [1, 304],$$

$$Z = 0.$$

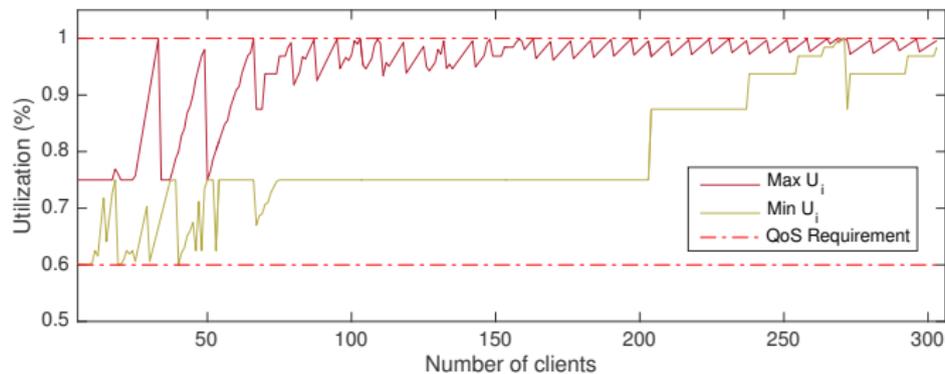
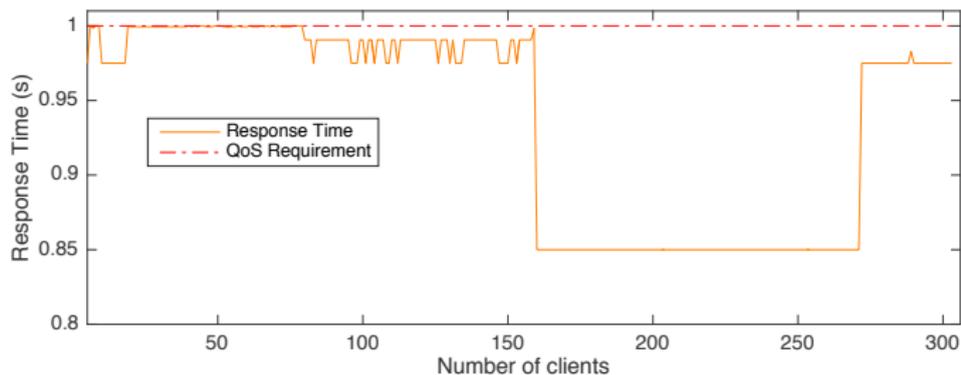
⁹Bhuvan Urgaonkar et al. "An analytical model for multi-tier internet services and its applications". In: *SIGMETRICS*. 2005.

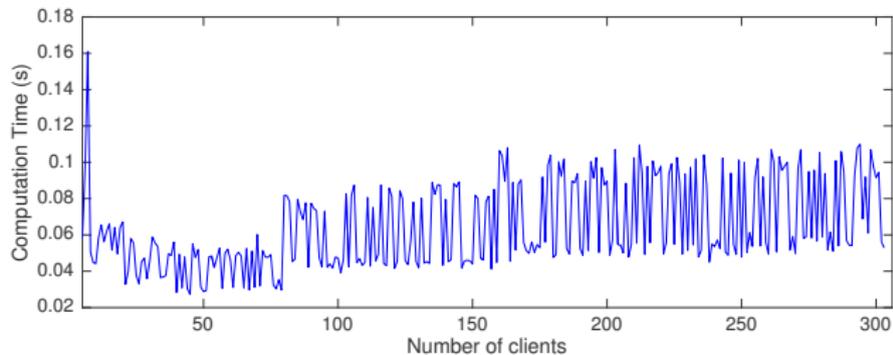
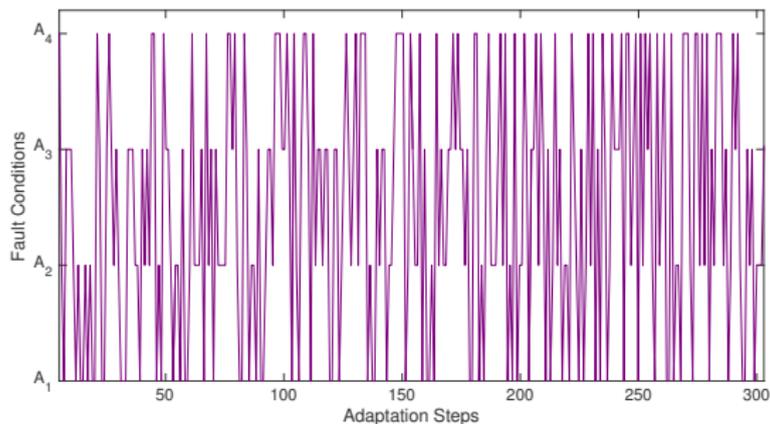
$$Q := (\mathcal{T}_1 \geq C) \wedge \forall i \in \mathcal{S}. ((\mathcal{U}_i \leq 1.0 \wedge \mathcal{U}_i \geq 0.6) \vee \mathcal{U}_i = 0.0)$$

$$\mathcal{D} \wedge \mathcal{R} := (\mathcal{Z} = 0) \wedge \mathcal{R}_1 \wedge \mathcal{R}_2$$

$$\mathcal{R}_1 := \forall i, j \in \mathcal{S}. (0.0 \leq p_{i,j} \leq 1.0) \wedge \forall i \in \mathcal{S}. \sum_{j \in \mathcal{S}} p_{i,j} = 1.0$$

$$\mathcal{R}_2 := \forall i \in \mathcal{S}. ((1.0 \leq s_i \leq 40.0) \wedge (0.06 \leq \frac{1}{\mu_i} \leq 10.0))$$





■ Advantages:

- Computation pushed at design-time phase for the symbolic solution
- Adaptation step is performed in a few ms even for realistic fully parametric three-tier models
- The solver may return an unsatisfiability result, meaning that there is no configurations satisfying the requirements

■ Limitations:

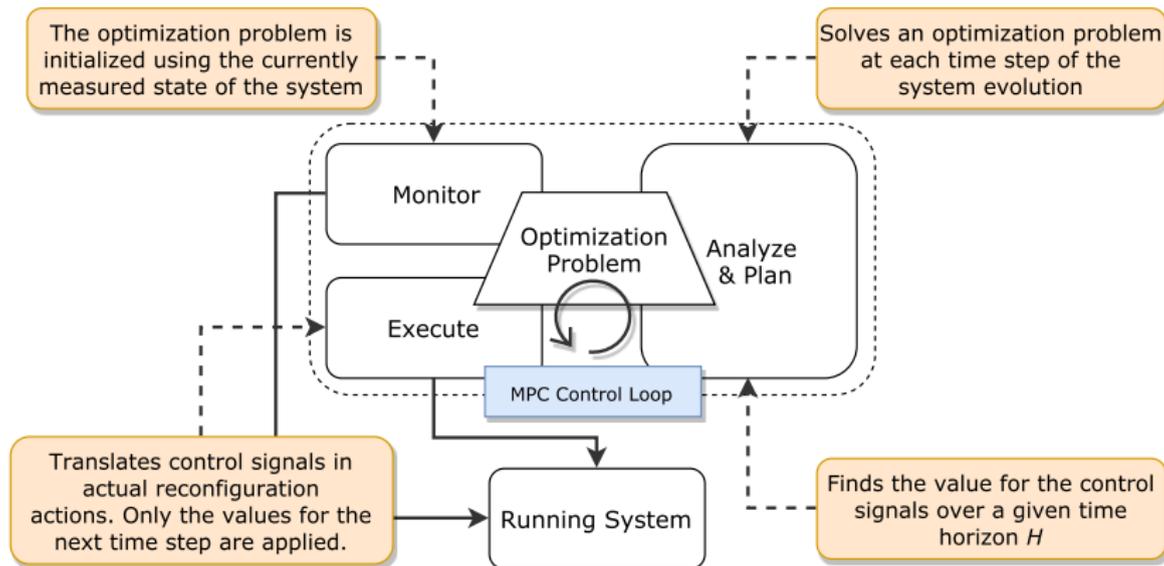
- Single class QNs model
- Closed form QNs model for the symbolic solution
- Steady state analysis
- Scalability of the SMT solver

- 1 Efficient performance models
- 2 Symbolic performance adaptation
- 3 Model-predictive control for performance-driven self-adaptation**
- 4 Moving horizon estimation of service demands

- 💡 **Idea:** Encode performance-driven self-adaptation as a model predictive control (MPC) problem¹⁰
- **Goals:** Get optimal assignments of the system parameters needed for steering an application toward a desired operating point (e.g., throughput)
 - fully automated
 - multiple adaptation knobs
 - considers actual run-time conditions
- 🤖 involves the solution of mixed integer nonlinear programs (MINLPs)

💡 As a main technical result we formally translate the naive MINLP MPC formulation in a mixed integer quadratic programming (MIQP) one

¹⁰Emilio Incerto, Mirco Tribastone, and Catia Trubiani. “Software Performance Self-Adaptation through Efficient Model Predictive Control”. In: *ASE'17*.



The main idea is to encode the discrete time version the ODE model as constraints of the optimization problem

$$\min_{\mu, \mathbf{s}, \mathbf{p}, \mathbf{x}} \sum_{k=0}^{H-1} \sum_{i=1}^M w_{i,k} (m_i(k) - \hat{r}_i(k))^2 + w_{s_i,k} \Delta \mathbf{s}_{i,k}^2$$

s.t.

$$x_i(k+1) = (-\mu_i(k) \min\{x_i(k), s_i(k)\} + \sum_{j \in S} p_{j,i}(k) \mu_j(k) \min\{x_j(k), s_j(k)\}) \Delta t + x_i(k)$$

$$s_i(k) \in \{\underline{s}_i, \underline{s}_i + 1, \dots, \bar{s}_i\}$$

$$0 \leq p_{i,j}(k) \leq 1, \sum_{j \in S} p_{i,j}(k) = 1$$

$$\text{for } 1 \leq i \leq M, 0 \leq k \leq H-1,$$

$\hat{r}_i(k)$ set point

💡 We reformulate the system relying on “virtual” adaptation knobs which are related to the original ones

- $\min\{x_i(k), s_i(k)\} = \alpha_i(k)$ can be encoded through standard MIP techniques¹¹
- $\mu_i(k)\alpha_i(k)$ can be encoded by adding slack variables $\gamma_i(k)$ and proper bound constraints

$$\mu_i(k)\alpha_i(k) = \hat{\mu}_i\alpha_i(k) + \gamma_i(k) \leftrightarrow \mu_i(k) = \hat{\mu}_i + \frac{\gamma_i(k)}{\alpha_i(k)}$$

$$\begin{cases} \mu_i(k) \leq \bar{\mu}_i \leftrightarrow \gamma_i(k) \leq (\bar{\mu}_i - \hat{\mu}_i)\alpha_i(k) \\ \mu_i(k) \geq \underline{\mu}_i \leftrightarrow \gamma_i(k) \geq (\underline{\mu}_i - \hat{\mu}_i)\alpha_i(k) \end{cases}$$

- The new linear time system is

$$\dot{x}_i(t) = \gamma_i(t) + \sum_{j \in \mathcal{S}} (-\gamma_j(t) + \zeta_{j,i}(t)), \quad i, j \in \mathcal{S}$$

¹¹www.gurobi.com/documentation/8.1/refman/constraints.html

Theorem

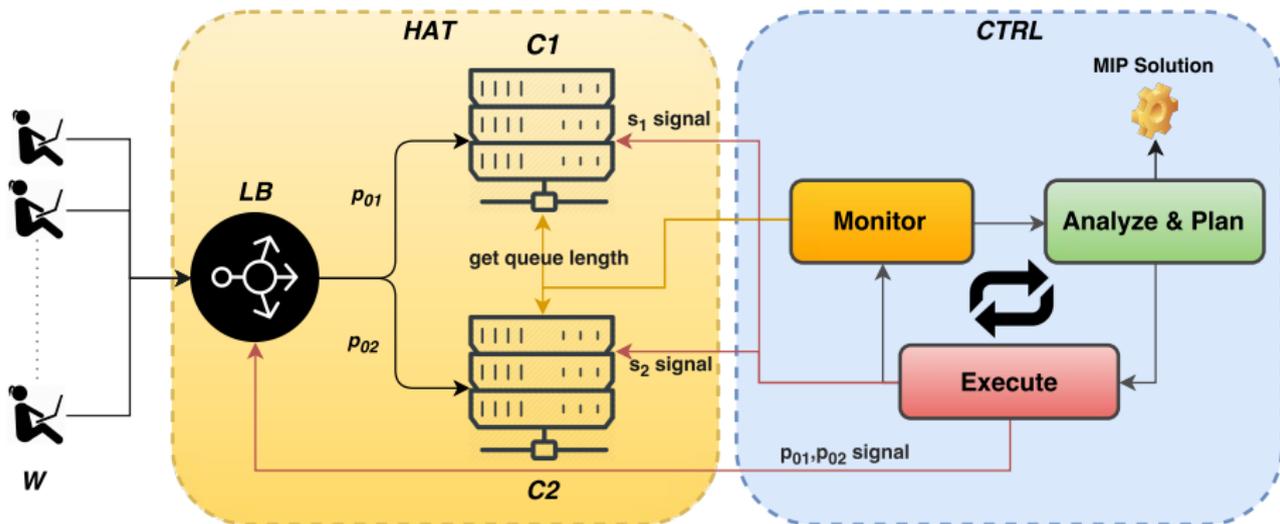
Denoting by $S = \{\mu_i^*(k), p_{i,j}^*(k), s_i^*(k), x_i^*(k)\}$ an optimal solution of the non-linear adaptation problem, there exists an MPC problem based on the MIP formulation with linear dynamics such that its optimal solution $S' = \{\gamma'_i(k), x'_i(k), \zeta'_{i,j}(k), s'_i(k)\}$ satisfies:

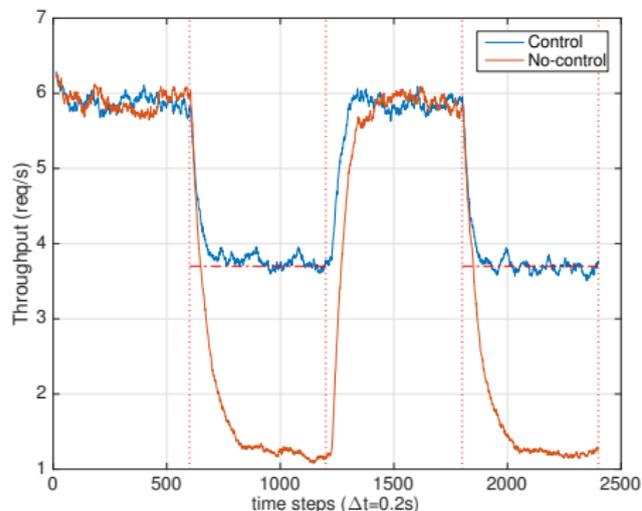
$$\mu_i^*(k) = \begin{cases} -\frac{\gamma'_i(k)}{x'_i(k)\Delta t} & \text{if } x'_i(k) \leq s'_i(k) \\ -\frac{\gamma'_i(k)}{s'_i(k)\Delta t} & \text{if } x'_i(k) > s'_i(k) \end{cases}$$

$$p_{i,j}^*(k) = \frac{\gamma'_i(k) - \zeta'_{i,j}(k)}{\gamma'_i(k)}, \quad s_i^*(k) = s'_i(k),$$

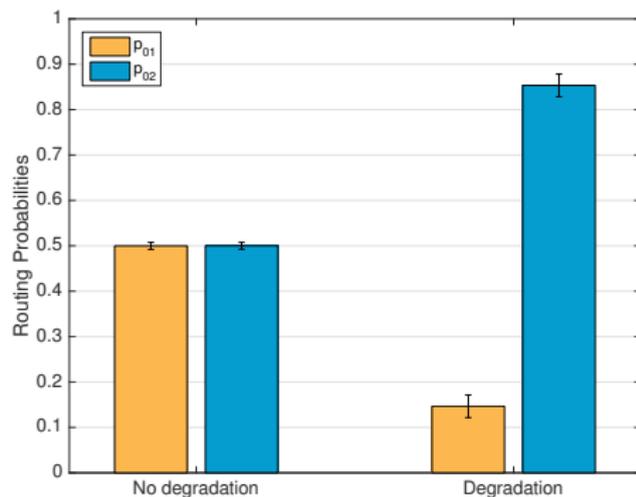
for all $k = 0, \dots, H - 1$.

- We evaluate the effectiveness and the scalability of the MPC approach on a real system
 - an in-house developed web application
- By studying two non trivial adaptation scenarios
 - hardware degradation
 - workload fluctuation
- Scalability comparison with probabilistic model checking

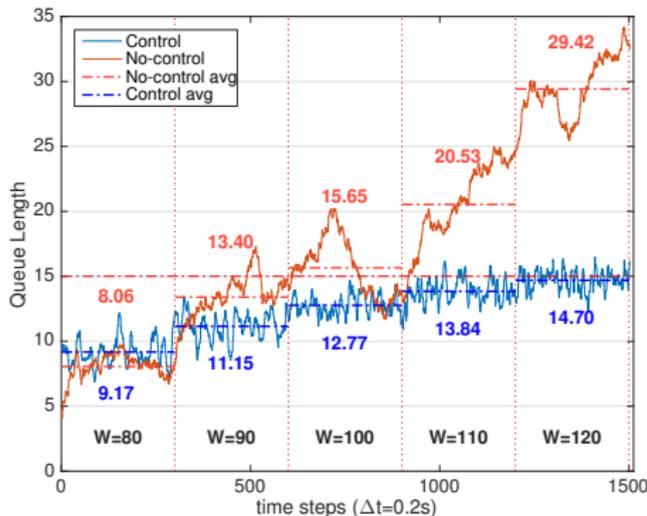




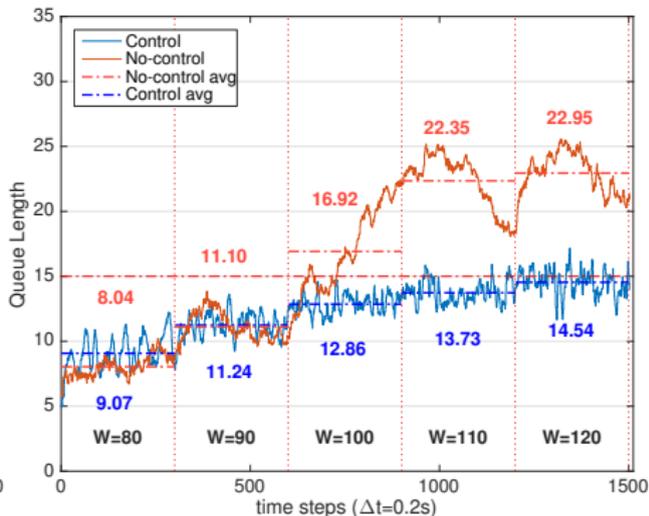
(a) System throughput



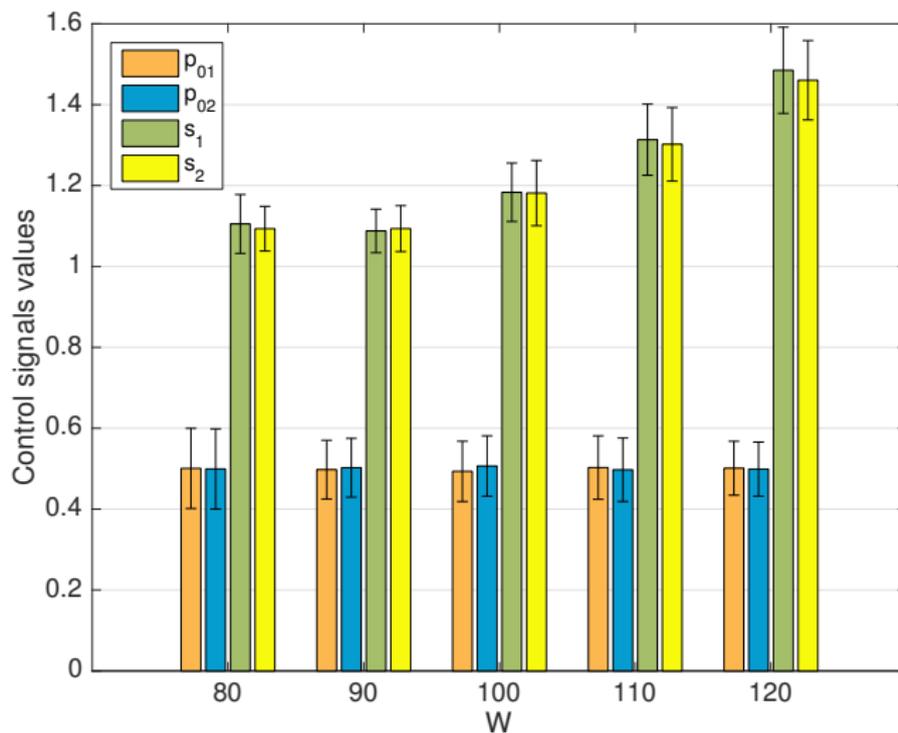
(b) Optimal control signals



(a) C_1 queue length



(b) C_2 queue length



Statistics of optimal control signals

Comparison against MDP (TO: timeout after 120 s)¹²

<i>W</i>	<i>MIP</i>	<i>Markov Decision Processes</i>		
	Runtime(s)	Runtime(s)	# States	# Transitions
80	0.0037	71	3 018 789	334 732 743
90	0.0036	87	3 805 074	421 958 628
100	0.0040	TO	4 682 259	519 272 613
110	0.0038	TO	5 650 344	626 674 698
120	0.0041	TO	6 709 329	744 164 883

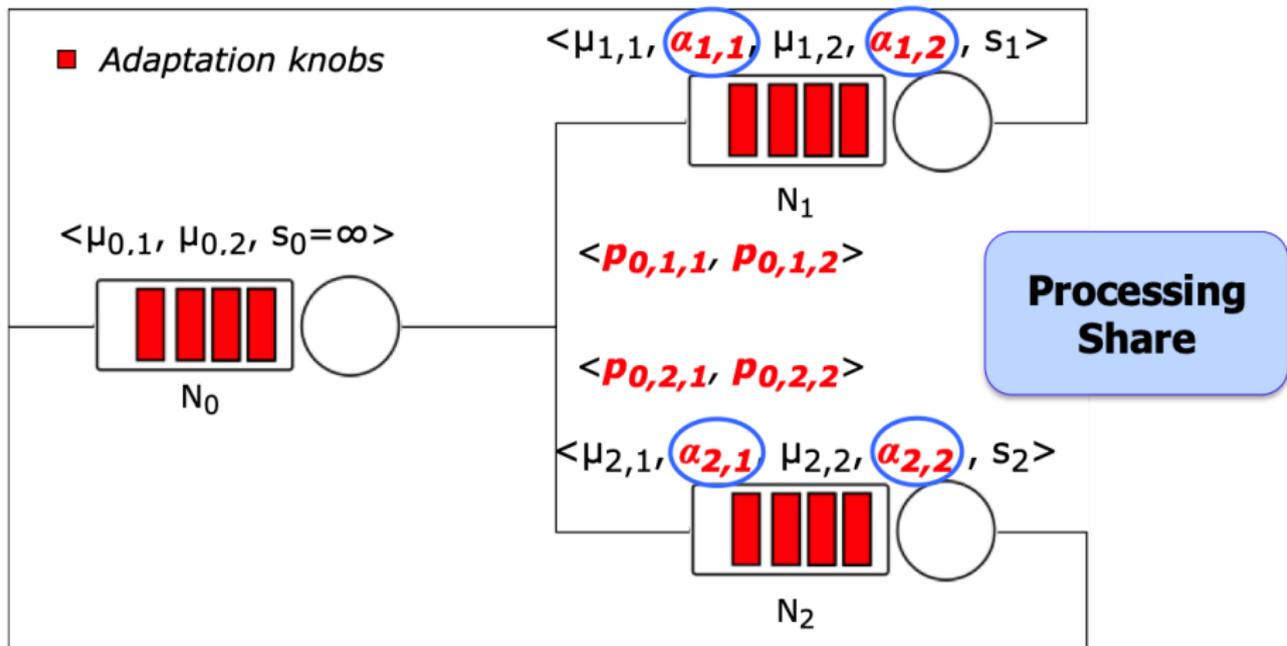
¹²Marta Kwiatkowska, Gethin Norman, and David Parker. “PRISM 4.0: Verification of Probabilistic Real-time Systems”. In: *CAV*. 2011.

- Contribution: a model predictive control based self-adaptive approach to continuously meet performance requirements
- Advantages:
 - fully automated
 - efficient
 - it works during the transient regime
 - the proposed linearization technique can be straightforwardly applied to more expressive models (e.g., Petri Nets)
- Limitations:
 - single-class model
 - QN parametrization
 - throughput and queue length QoS requirements only

- Meeting performance targets of co-located applications (e.g., virtualized cloud environments) is challenging
- Scaling techniques:
 - *Vertical scaling*: assigns resource shares on each individual machine
 - *Horizontal scaling*: chooses the number of virtual machines employed
- State-of-the-art approaches apply vertical and horizontal scaling in an isolated fashion (i.e., symmetric load balancing)
- Ineffective when machines have different hardware characteristics (e.g., software aging or hardware degradation)

- A model based approach:
 - Multi-class QN as the quantitative model
 - Analysed by means of QNs fluid approximation
 - Combined horizontal and vertical scaling formulated as an MPC quadratic programming problem
- Main technical results:¹³
 - A multi-class model that enables an accurate representation of the **capped allocation paradigm**
 - The specification of latency-based requirements
 - Extension of [ASE'17]

¹³Emilio Incerto, Mirco Tribastone, and Catia Trubiani. “Combined Vertical and Horizontal Autoscaling Through Model Predictive Control”. In: *EURO-PAR*. 2018.



Multi class QN model of load balancer with two co-located applications

- For each station i and each class c we define the ODE:

$$\dot{x}_{i,c}(t) = -\mu_{i,c} \min\{x_{i,c}(t), \alpha_{i,c}(t)s_i\} + \sum_{j \in S} \rho_{j,i,c}(t) \mu_{j,c} \min\{x_{j,c}(t), \alpha_{j,c}(t)s_j\}$$

with $\alpha_{i,c}(t) \geq 0$, $\sum_{c \in C} \alpha_{i,c}(t) \leq 1$.

- We define the following metrics:
 - Class- c throughput at station i

$$T_{i,c}(t) = \mu_{i,c} \min\{x_{i,c}(t), \alpha_{i,c}(t)s_i\}$$

- Instantaneous response time

$$R_{i,c}(t) = \frac{x_{i,c}(t)}{T_{i,c}(t)}$$

- In [ASE'17] we employed MPC for performance runtime adaptation for single class queuing networks
- We showed how it could be formulated as a series of mixed-integer quadratic program (MIQP)
- Here we extend this formulation for controlling the multi-class QN under the cap allocation sharing
 - A positive side effect: quadratic programming (QP) formulation

- The QP formulation is possible due to the following:

$$\min\{x_{i,c}(t), \alpha_{i,c}(t)s_i\} = \gamma_{i,c}(t)$$



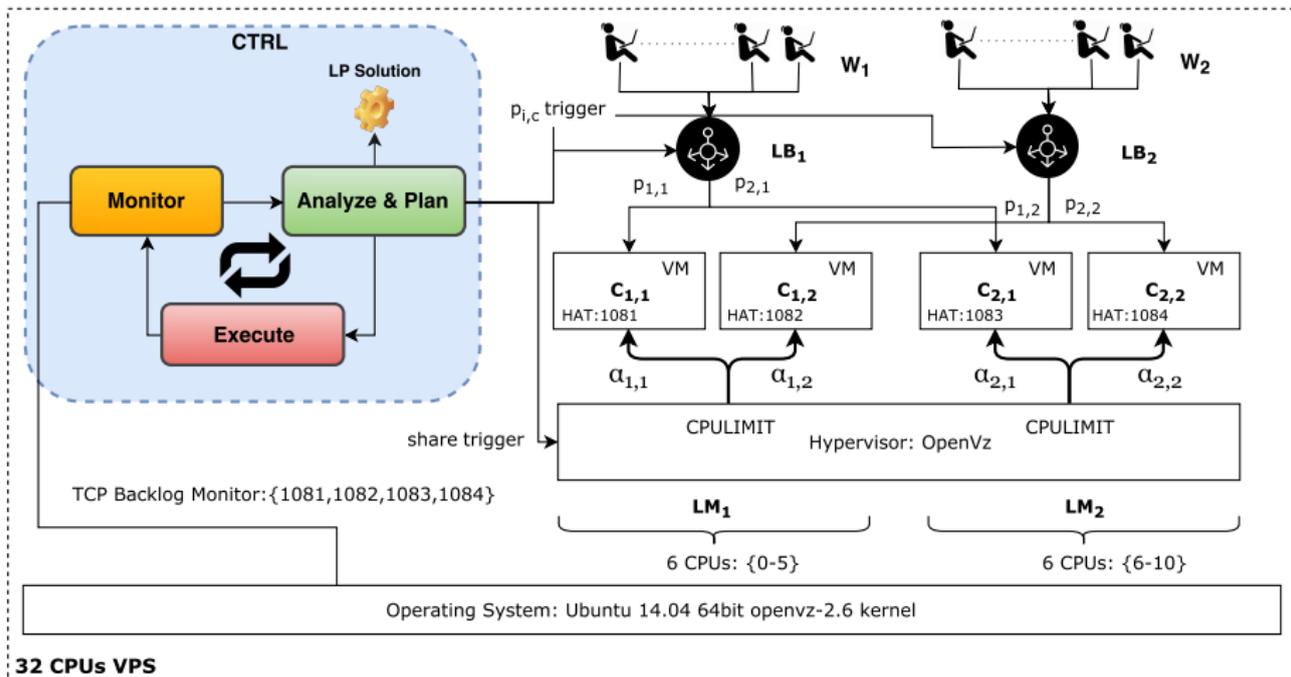
$$\begin{array}{ll} \gamma_{i,c}(t) \geq 0, & \gamma_{i,c}(t) \leq s_i \\ \gamma_{i,c}(t) \leq x_{i,c}(t) & \sum_{c \in C} \gamma_{i,c}(t) \leq s_i \end{array}$$

with $i \in S, c \in C$

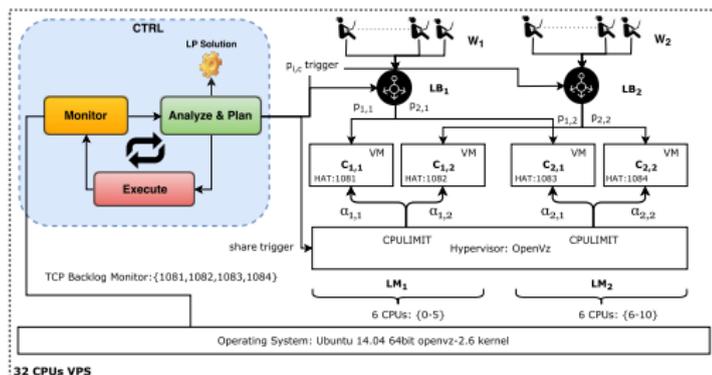
- then if $\gamma_{i,c}^*(t) \neq x_{i,c}^*(t)$ the real actuator can be computed as

$$\alpha_{i,c}^*(t) = \frac{\gamma_{i,c}^*(t)}{s_i}$$

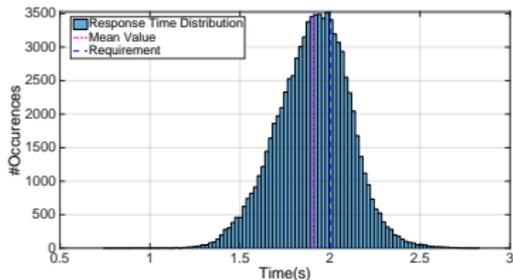
- We show that the combined vertical and horizontal adaptation can efficiently meet performance targets when either of the two techniques alone cannot
- We used the OpenVz hypervisor while horizontal scaling has been enabled through a NodeJS-based load balancer
- For emulating a multi class scenario, we ran two instances of the same load balanced HAT deployment consisting of two OpenVz virtual machines each



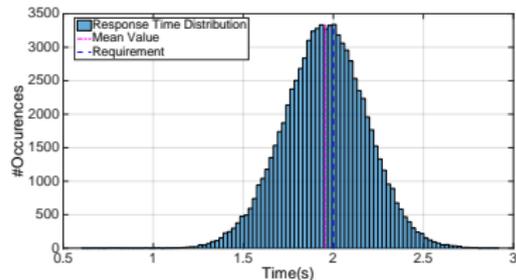
From a symmetric set-up, inject a degradation event such that service rate at node LM1 becomes 3 times smaller



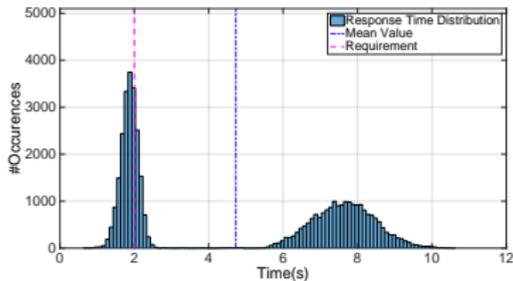
- Objective: set points $R_1 = 2 s$ and $T_2 = 50 r/s$
- Workload of 200 users and think times with average 1 s
- Control approaches evaluated in two separate 20-minute-long sessions



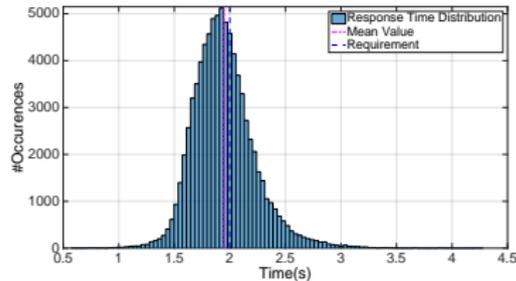
(a) Vertical scaling only



(b) Vertical & horizontal scaling

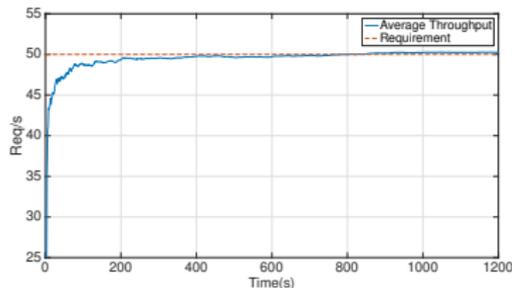


(c) Vertical scaling only

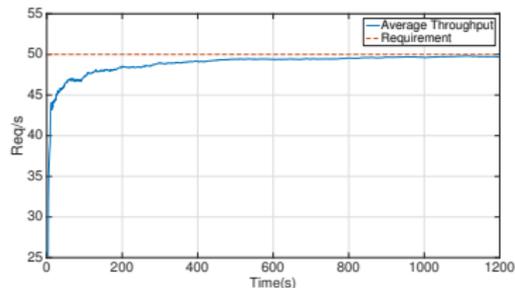


(d) Vertical & horizontal scaling

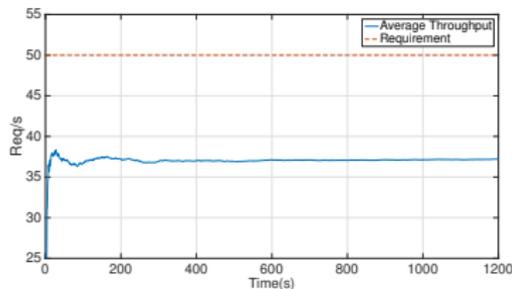
Response time distribution without (a,b) and with (c,d) degradation



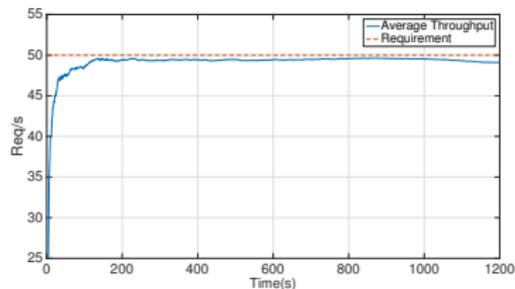
(a) Vertical scaling only



(b) Vertical & horizontal scaling



(c) Vertical scaling only



(d) Vertical & horizontal scaling

Class-2 average throughput with (a,b) and without (c,d) degradation

- **Contribution:** an efficient approach for performance adaptation of distributed co-located applications.
- The main novelties lay:
 - The combined usage of vertical and horizontal scaling techniques
 - A multi-class fluid model for co-located applications under a capped resources allocation scheduler.
- Future works:
 - modeling response time distribution instead of its average only
 - include resource contention policies for network, memory, I/O
 - consider more expressive resource schedulers and system performance interactions

- 1 Efficient performance models
- 2 Symbolic performance adaptation
- 3 Model-predictive control for performance-driven self-adaptation
- 4 Moving horizon estimation of service demands**

- Well calibrated model parameters are necessary for computing accurate predictions
- When dealing with queuing networks service demands are fundamental
- The estimation need to be performed:
 - continuously
 - non intrusively



MHE: We formulate the estimation problem as a quadratic program solved according to the moving horizon paradigm.¹⁴

¹⁴Emilio Incerto, Annalisa Napolitano, and Mirco Tribastone. “Moving Horizon Estimation of Service Demands in Queuing Networks”. In: *MASCOTS*. 2018.

$$\underset{\bar{x}, T}{\text{minimize}} \sum_{k=1}^H \sum_{i=1}^M (\bar{x}_i(k) - \tilde{x}_i(k))^2,$$

subject to:

$$\bar{x}_i(k+1) = \bar{x}_i(k) - T_i(k) + \sum_{j=1}^M p_{j,i} T_j(k)$$

$$\bar{x}_i(0) = \tilde{x}_i(0), \quad 1 \leq i \leq M, 0 \leq k \leq H-1.$$

$$\mu_i^* := \frac{\sum_{k=0}^{H-1} T_i^*(k)}{\sum_{k=0}^{H-1} \min \{ \bar{X}_i^*(k), s_i \}}, \quad 1 \leq i \leq M.$$

Accuracy comparison between the queue length maximum likelihood estimation (QMLE)¹⁶ and our approach (MHE).

K	$x(0) = (3, 0, 0)$ $H = 2347, U_2 \approx 0.10$		$x(0) = (9, 0, 0)$ $H = 688, U_2 \approx 0.30$		$x(0) = (12, 0, 0)$ $H = 521, U_2 \approx 0.40$		$x(0) = (19, 0, 0)$ $H = 353, U_2 \approx 0.60$		$x(0) = (26, 0, 0)$ $H = 262, U_2 \approx 0.80$	
	QMLE	MHE	QMLE	MHE	QMLE	MHE	QMLE	MHE	QMLE	MHE
1	0.52	9.25 ± 1.03	1.37	9.63 ± 1.06	2.07	7.90 ± 1.01	3.40	6.58 ± 0.81	5.15	4.89 ± 0.69
2	448.30	4.13 ± 0.62	126.54	3.93 ± 0.58	67.18	4.20 ± 0.63	5.46	3.90 ± 0.56	2.33	3.59 ± 0.54
5	184.02	2.26 ± 0.33	60.41	3.02 ± 0.43	42.09	2.76 ± 0.38	8.78	2.07 ± 0.33	1.65	2.06 ± 0.34
10	92.29	1.65 ± 0.27	30.53	1.99 ± 0.31	23.18	1.82 ± 0.31	9.50	2.09 ± 0.30	3.89	1.50 ± 0.24
20	45.18	1.37 ± 0.21	15.01	1.13 ± 0.19	11.32	1.36 ± 0.18	6.41	1.36 ± 0.19	5.81	1.17 ± 0.18
50	18.67	0.74 ± 0.10	6.08	0.81 ± 0.14	4.57	0.78 ± 0.11	2.72	0.81 ± 0.12	5.17	0.73 ± 0.10

¹⁶Weikun Wang et al. “Maximum likelihood estimation of closed queueing network demands from queue length data”. In: *ICPE*. 2016.

MHE scalability analysis

M	<i>Errors</i>				<i>Runtimes (s)</i>			
	min	avg	95-th	max	min	avg	95-th	max
5	1.60	2.53	4.12	4.50	0.03	0.03	0.03	0.04
10	1.63	2.46	3.28	3.56	0.08	0.08	0.09	0.09
15	1.59	2.63	3.48	4.56	0.18	0.18	0.19	0.19
20	1.62	2.52	3.19	3.82	0.34	0.38	0.48	0.52

- We defined optimization methods based on fluid queuing networks enabling fast adaptation reactions under strict time constraints.
- Through an extensive numerical validation we have shown how complex multi-dimensional adaptation actions can be computed at a small computational cost in real-world scenarios.
- We developed the first transient estimation technique for service demands of QNs enabling effective and efficient parameter estimation in a minimally intrusive manner.

- The definition of model predictive control problem based on the fluid interpretation of layered queuing networks.
- The definition of performance-driven self-adaptation approaches considering higher-order statistics of the controlled quantities, e.g., variance
- The development of minimally intrusive estimation techniques for multi-class QNs with non-exponentially distributed service times

- A number of results on self-adaptation by means of formal quantitative verification.^{17 18}
- Runtime use of discrete time Markov chains for reliability of self-adaptive systems.¹⁹
- Control-theoretic approaches for learning and updating a linear model from measurements and controlling it via a single parameter.²⁰
- Control of software performance via service-rate adaption using a queuing model.²¹

¹⁷Carlo Ghezzi et al. “Managing non-functional uncertainty via model-driven adaptivity”. In: *ICSE’13*.

¹⁸A. Filieri, G. Tamburrelli, and C. Ghezzi. “Supporting Self-adaptation via Quantitative Verification and Sensitivity Analysis at Run Time”. In: *IEEE Trans. Software Eng.* (2016).

¹⁹Antonio Filieri et al. “Self-adaptive software meets control theory: A preliminary approach supporting reliability requirements”. In: *ASE’11*.

²⁰Antonio Filieri, Henry Hoffmann, and Martina Maggio. “Automated Design of Self-adaptive Software with Control-theoretical Formal Guarantees”. In: *ICSE’14*.

²¹Davide Arcelli et al. “Control theory for model-based performance-driven software adaptation”. In: *QoSA’15*.