

Adaptive Control of Virtualized Resources in Utility Computing Environments

Pradeep Padala,
Kang G. Shin



Xiaoyun Zhu, Mustafa Uysal,
Zhikui Wang, Sharad Singhal,
Arif Merchant



Kenneth Salem



June 2007

EuroSys '07: Proceedings of the 2nd ACM SIGOPS/EuroSys
European Conference on Computer Systems 2007

Publisher: ACM



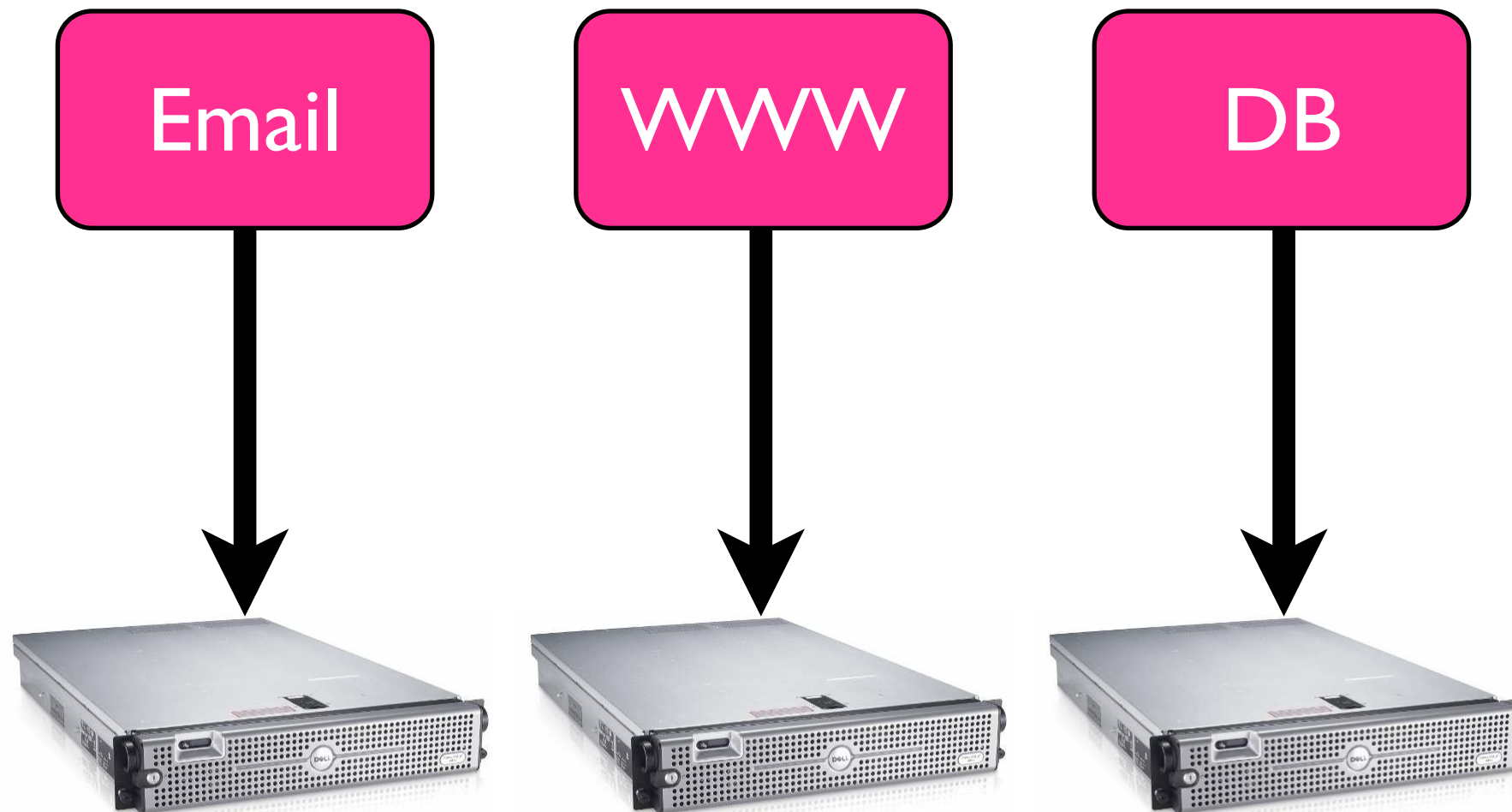
July 8, 2009

OUTLINE

- INTRODUCTION
- MODELING
- DESIGN
- EVALUATION
- CONCLUSIONS

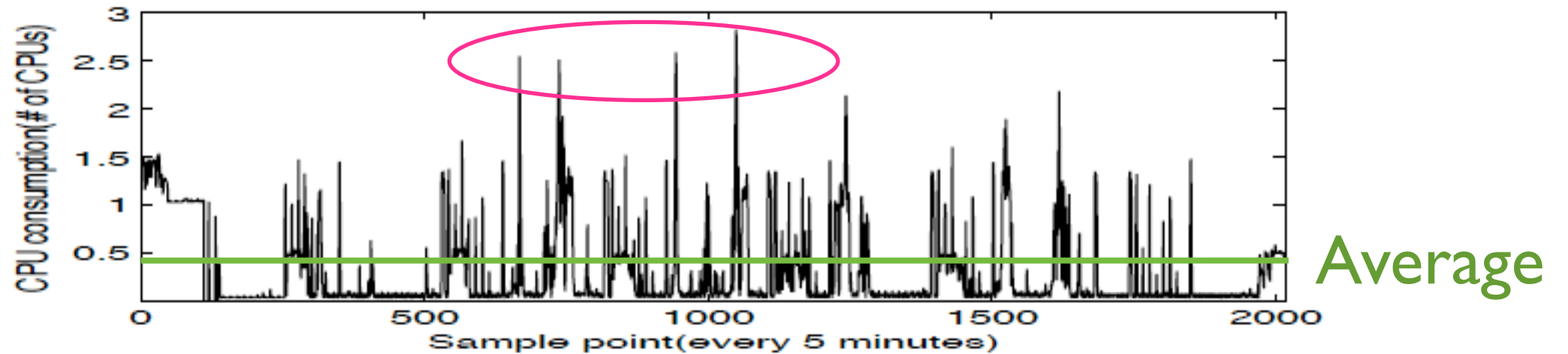
INTRODUCTION

Each application has its own dedicated servers in traditional way

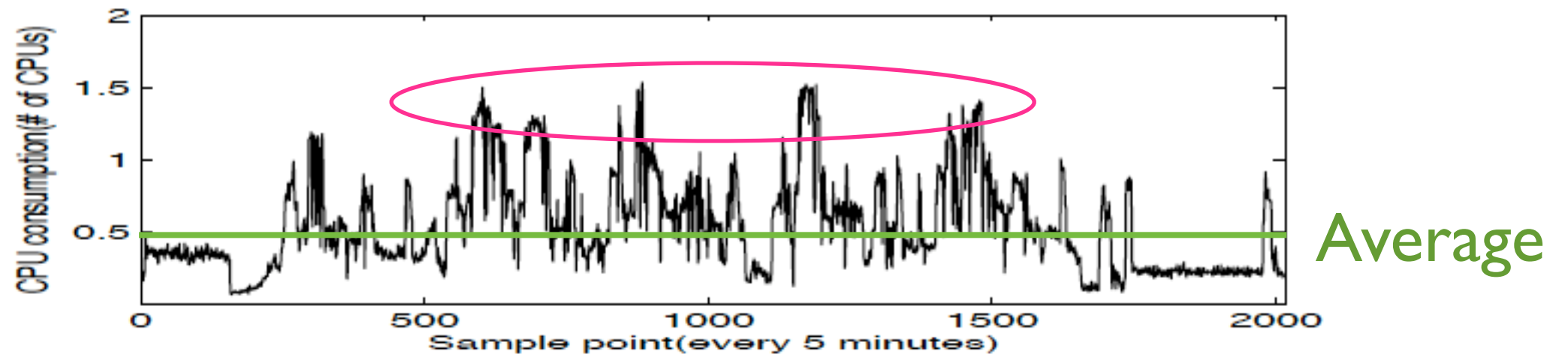


Data Center

INTRODUCTION



(a) CPU consumption of node 1

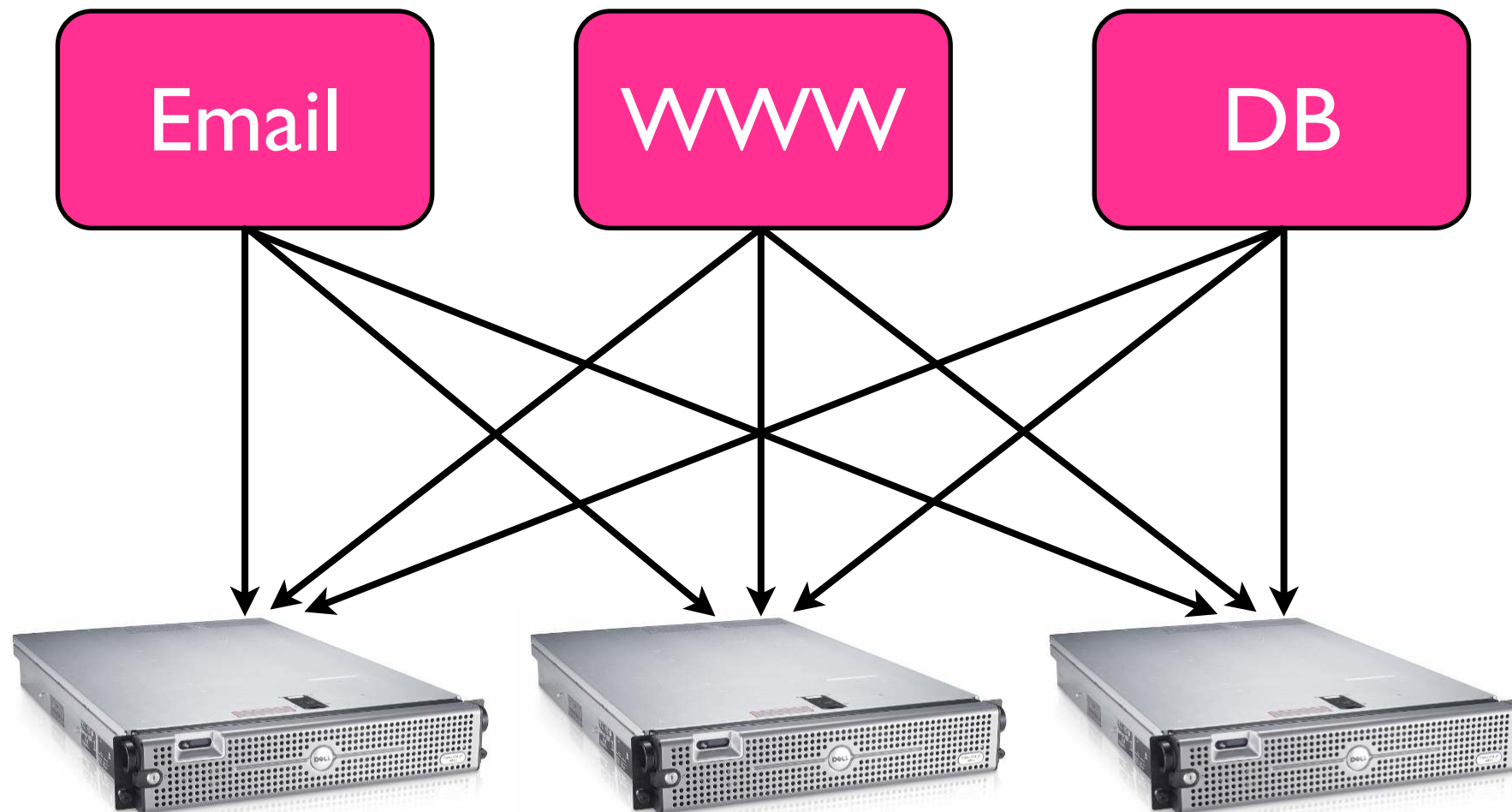


(b) CPU consumption of node 2

An example of data center server consumption
(each node has 6 CPUs)

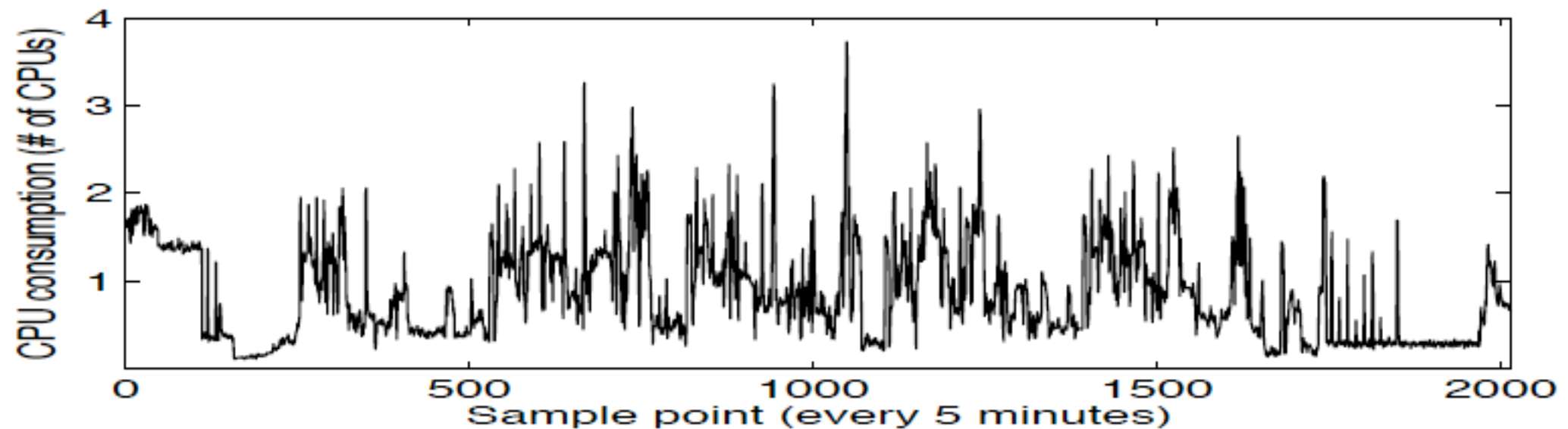
INTRODUCTION

Applications shares resources as their demands changes over time



Next-generation Data Center

INTRODUCTION



(c) Sum of CPU consumptions from both nodes

An example of data center server consumption
(each node has 6 CPUs)

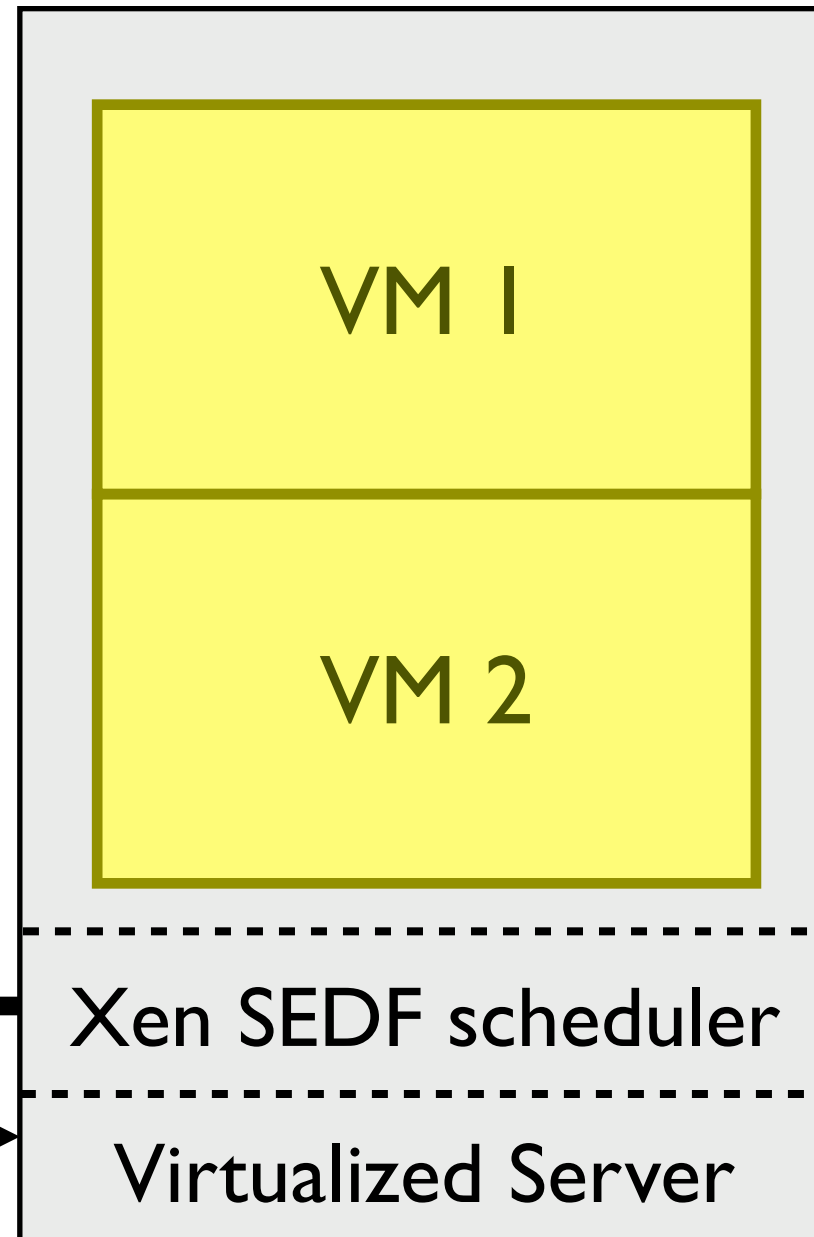
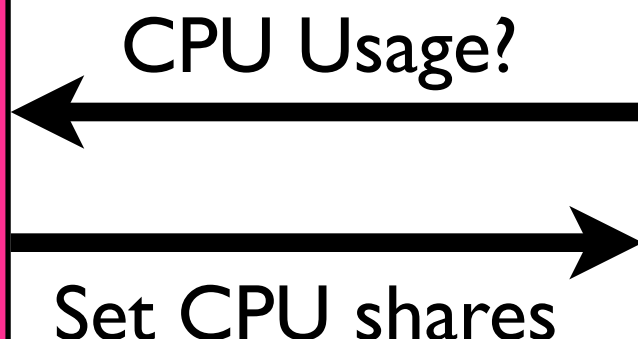
Solution:Adptive Controller

INTRODUCTION

Goals

- Good Utilization
- Good Performance
- QoS Differentiation

Goals met? **NO**



CPU
share

50%

50%

INTRODUCTION

Goals

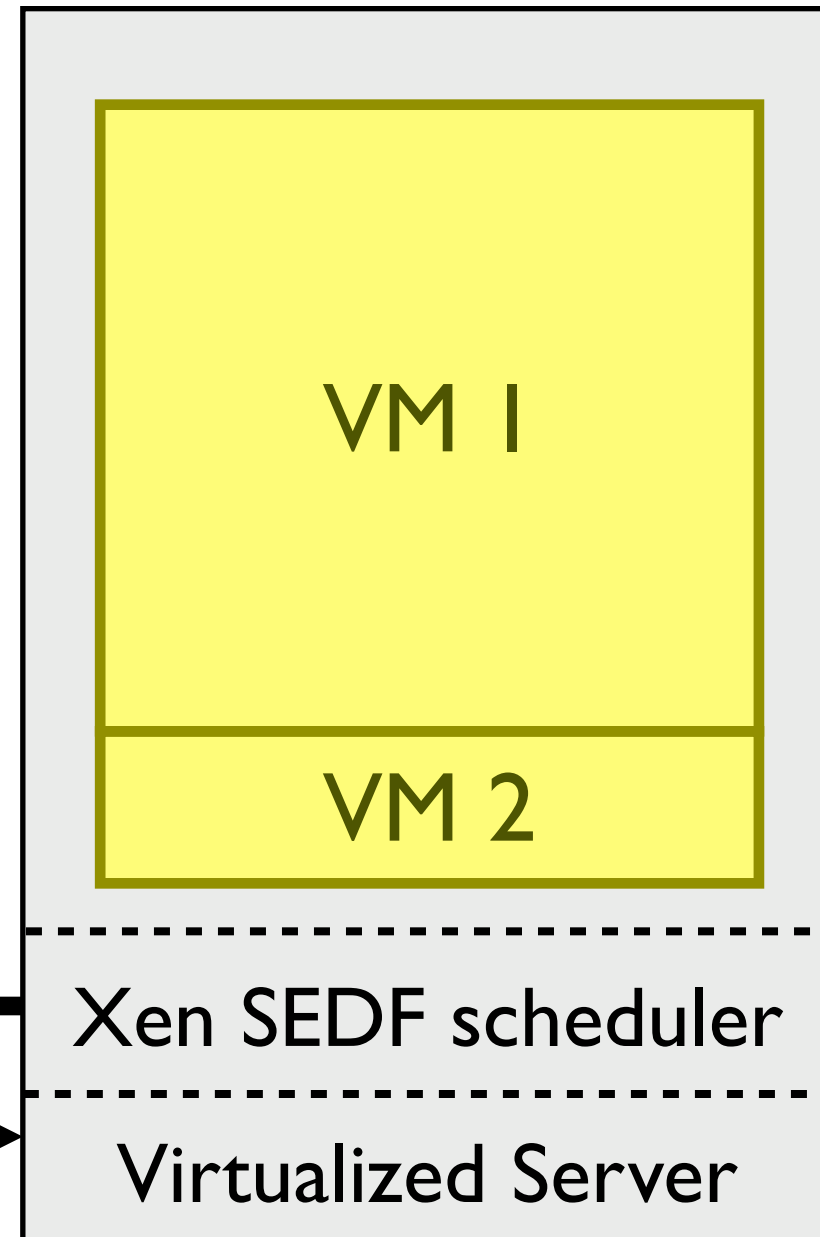
- Good Utilization
- Good Performance
- QoS Differentiation

Goals met? **NO**



CPU Usage?

Set CPU shares



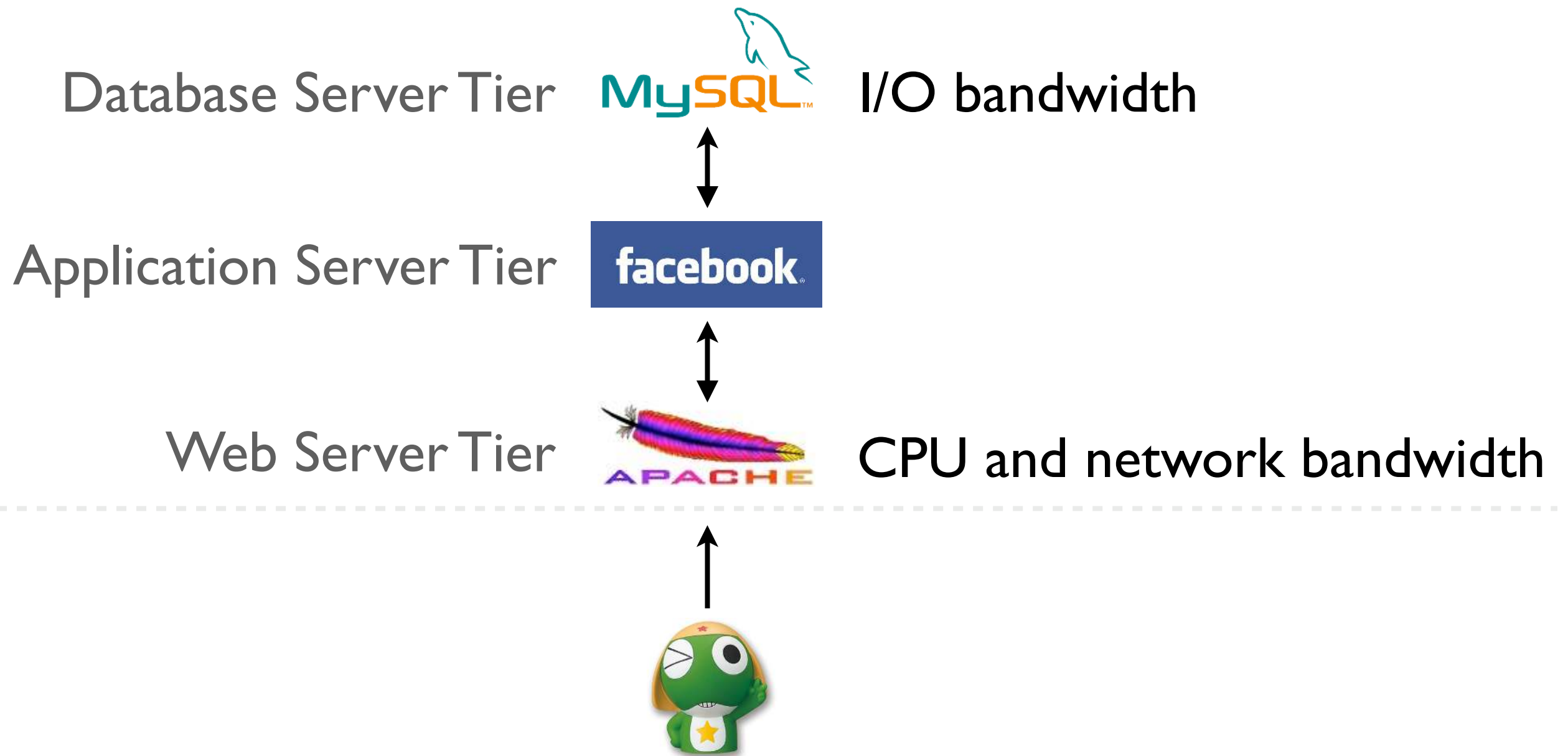
CPU
share

80%

20%

Control CPU shares based on goals

INTRODUCTION

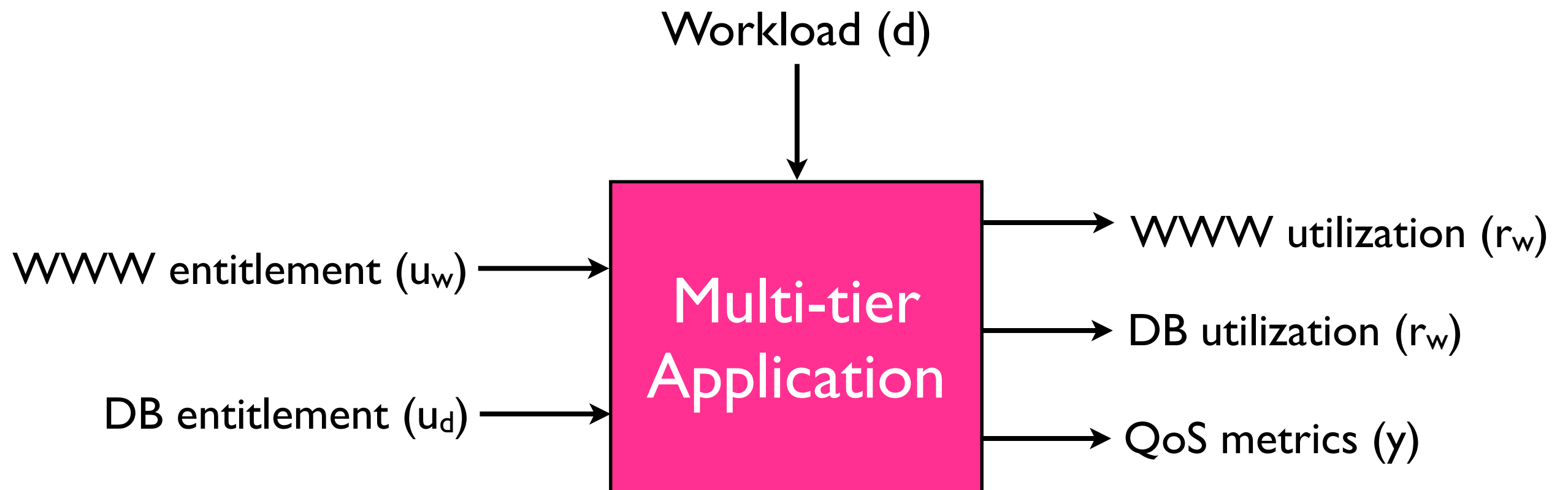


Enterprise applications typically employ a multi-tier architecture

MODELING

- First define some terminology
 - **entitlement (u)**: the percentage of CPU share allocated to a virtual machine
 - **consumption (v)**: the percentage of total CPU capacity actually used by the VM
 - **VM utilization (r)**: the ratio between consumption and entitlement ($r = v / u$)

MODELING



An input-output model for a multi-tier application

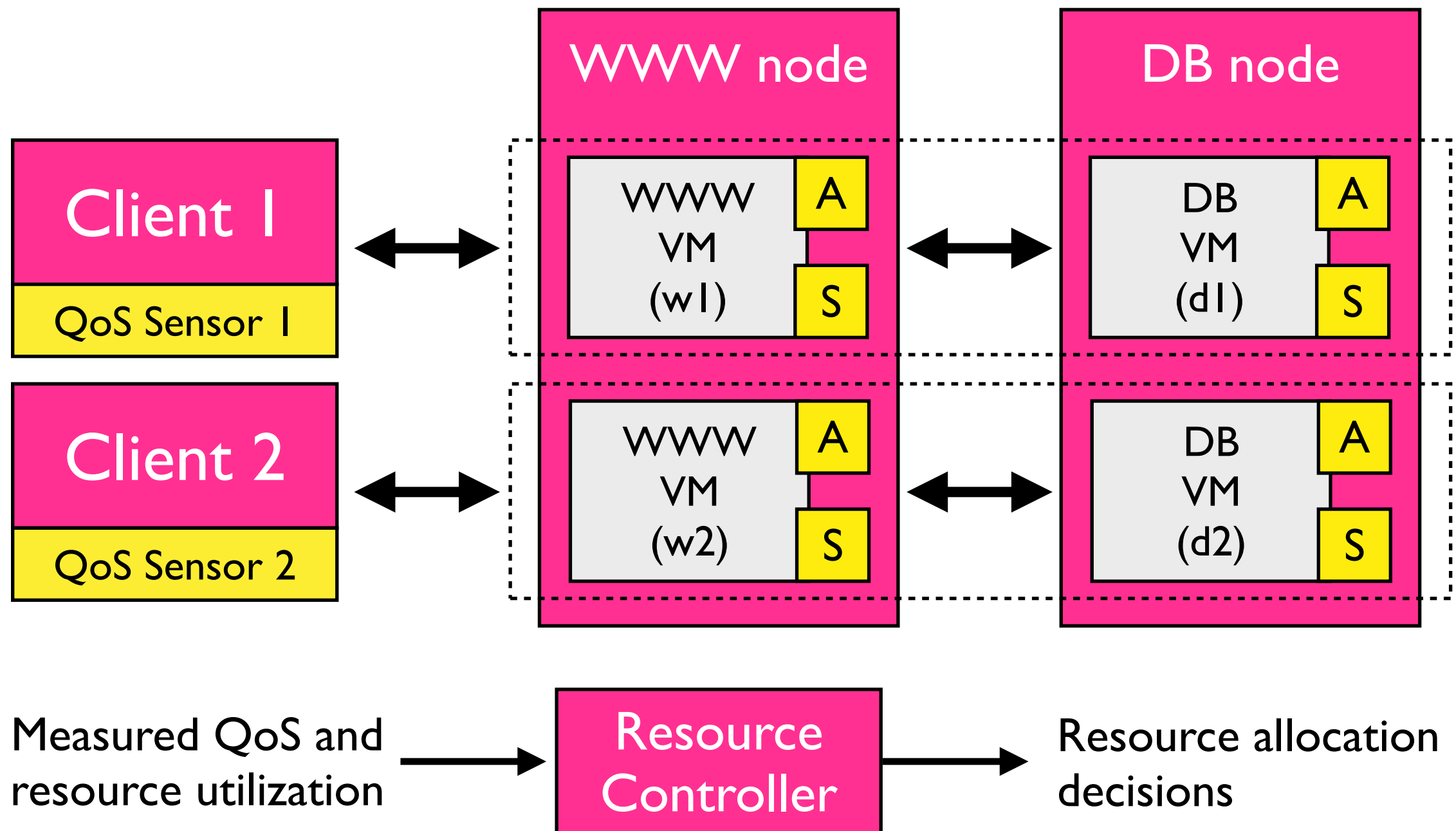
MODELING

Experimental Testbed

- 5 HP proliant servers, each has two processors, 4 GB of RAM, one Gigabit Ethernet interface, and two local SCSI disks
- Used two workload generators
 - RUBiS: an online auction site benchmark
 - TPC-W: a transactional e-Commerce benchmark

MODELING

Experimental Testbed



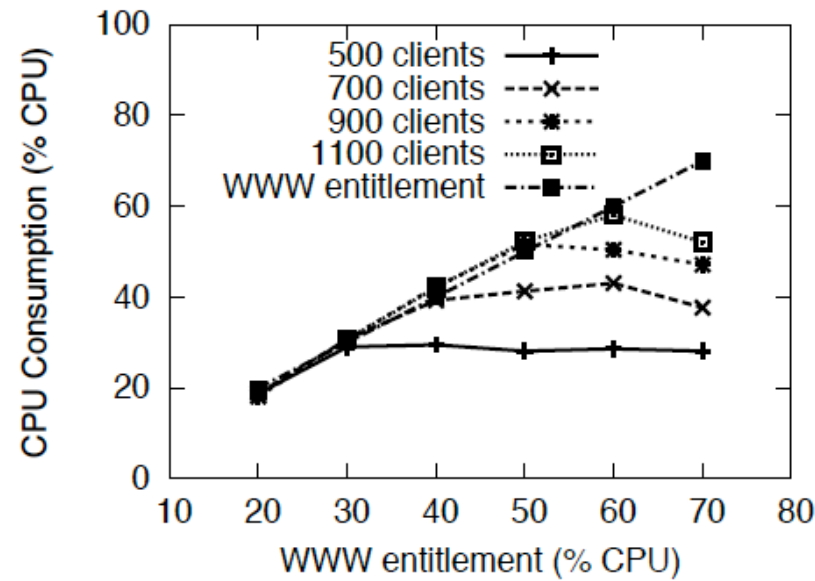
MODELING

Modeling single multi-tier application

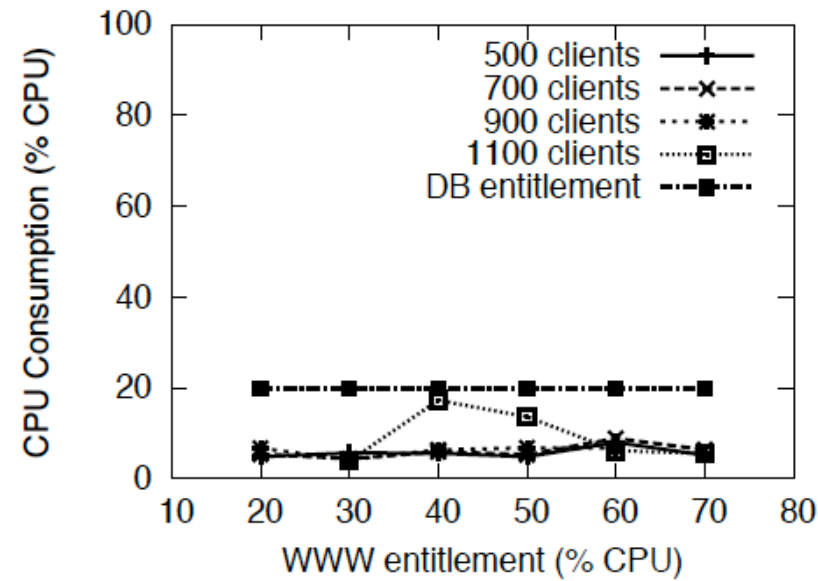
- To know how the changes in the WWW/DB entitlements impact the utilization of virtual machines and QoS metrics
- A single testbed node was used to host a two-tier implementation of RUBiS
- Pinned the WWW VM(20-70%), the DB VM(20%), as well as dom0(remaining) to one processor

MODELING

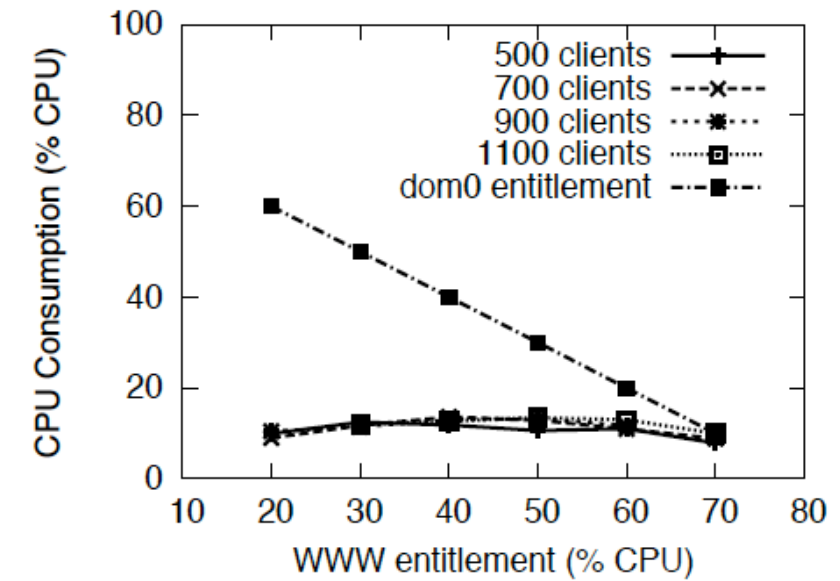
Modeling single multi-tier application



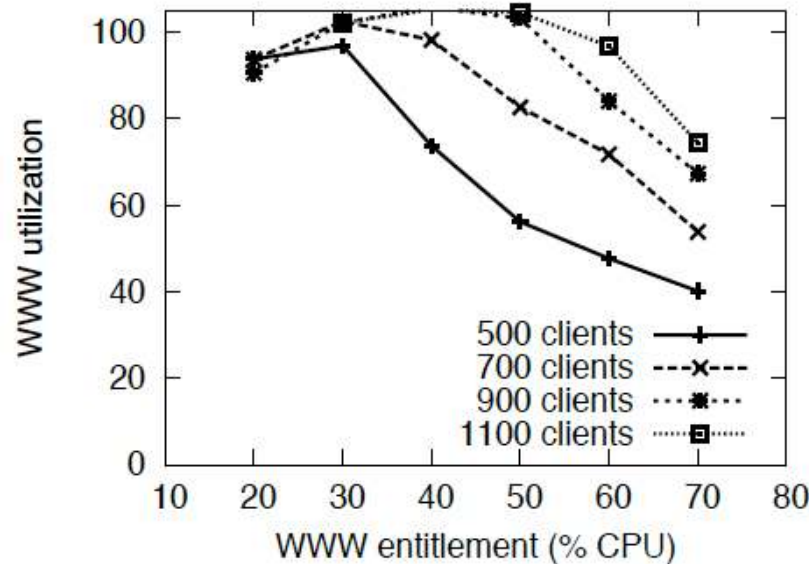
(a) WWW CPU consumption



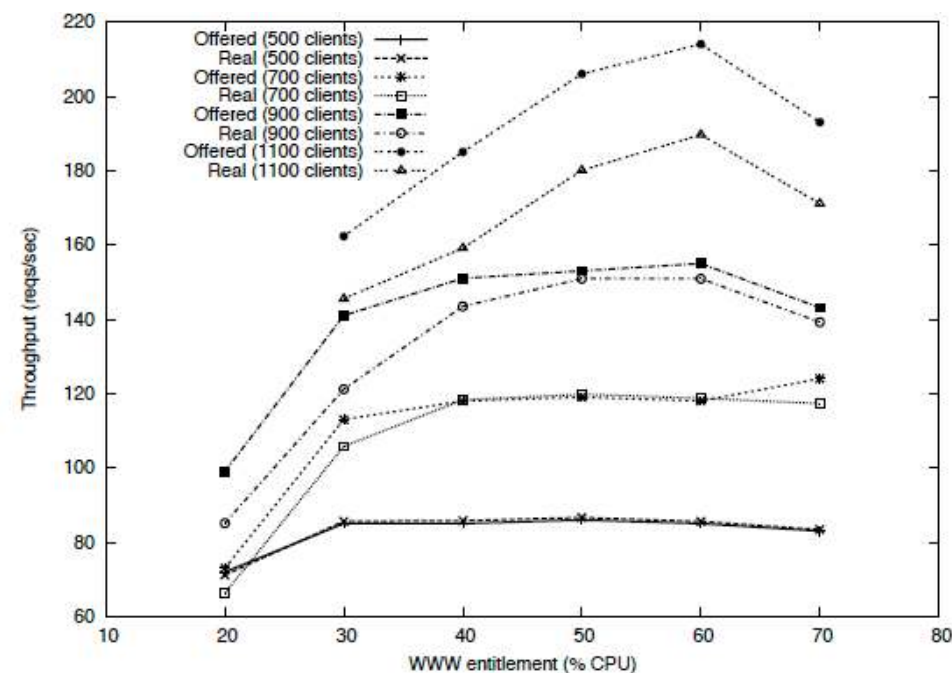
(b) DB CPU consumption



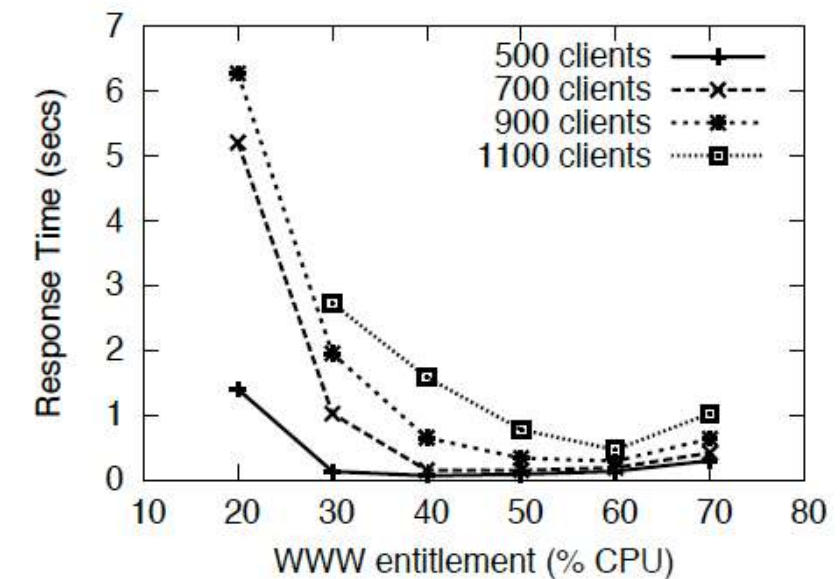
(c) dom0 CPU consumption



(d) WWW VM utilization



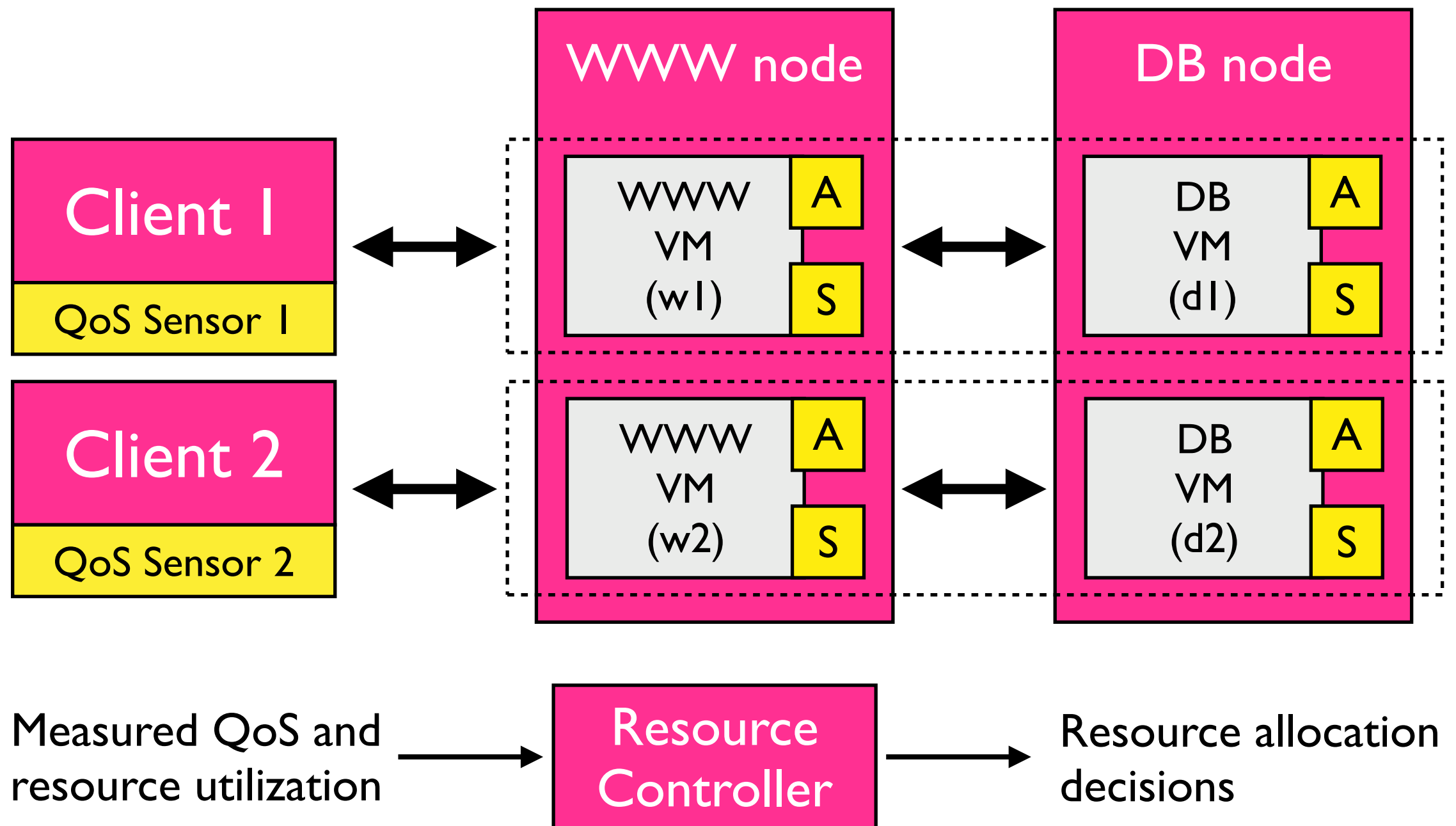
(e) Throughput



(f) Response time

MODELING

Modeling co-hosted multi-tier applications



MODELING

Modeling co-hosted multi-tier applications

- At any given time either the WWW node or the DB node may become saturated

	DB node unsat.	DB node sat.
WWW node unsat.	WU-DU	WU-DS
WWW node sat.	WS-DU	WS-DS

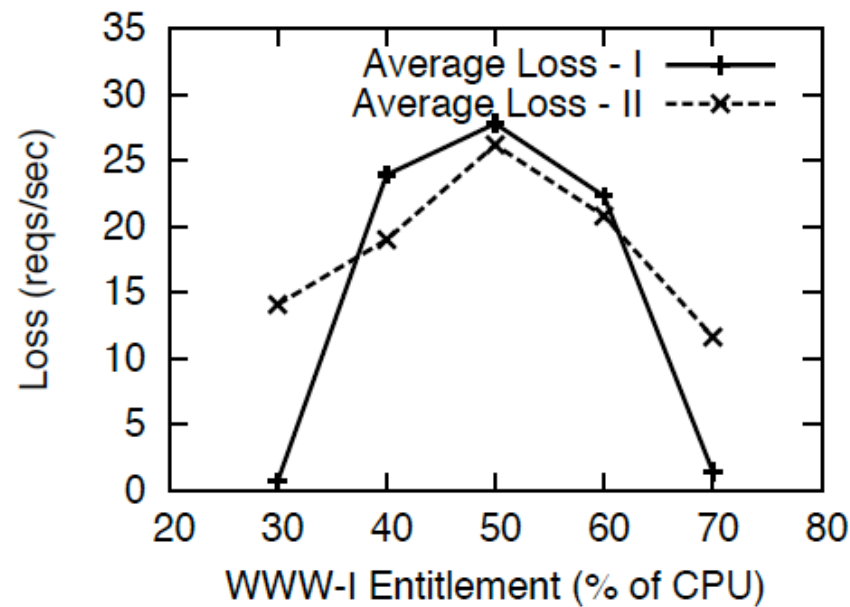
- Also need a QoS differentiation metric:

$$y_{ratio} = \frac{y_1}{y_1 + y_2}$$

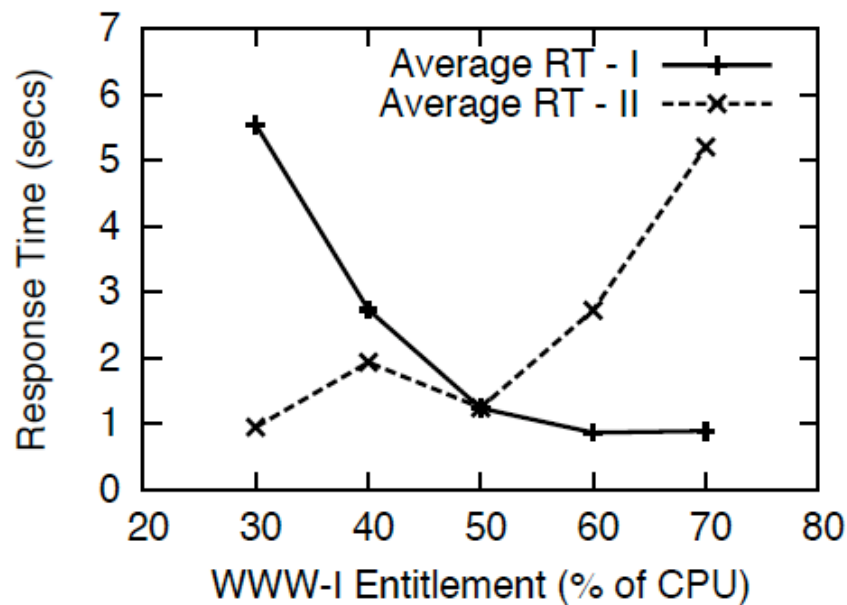
y can be average response time,
throughput, or request loss

MODELING

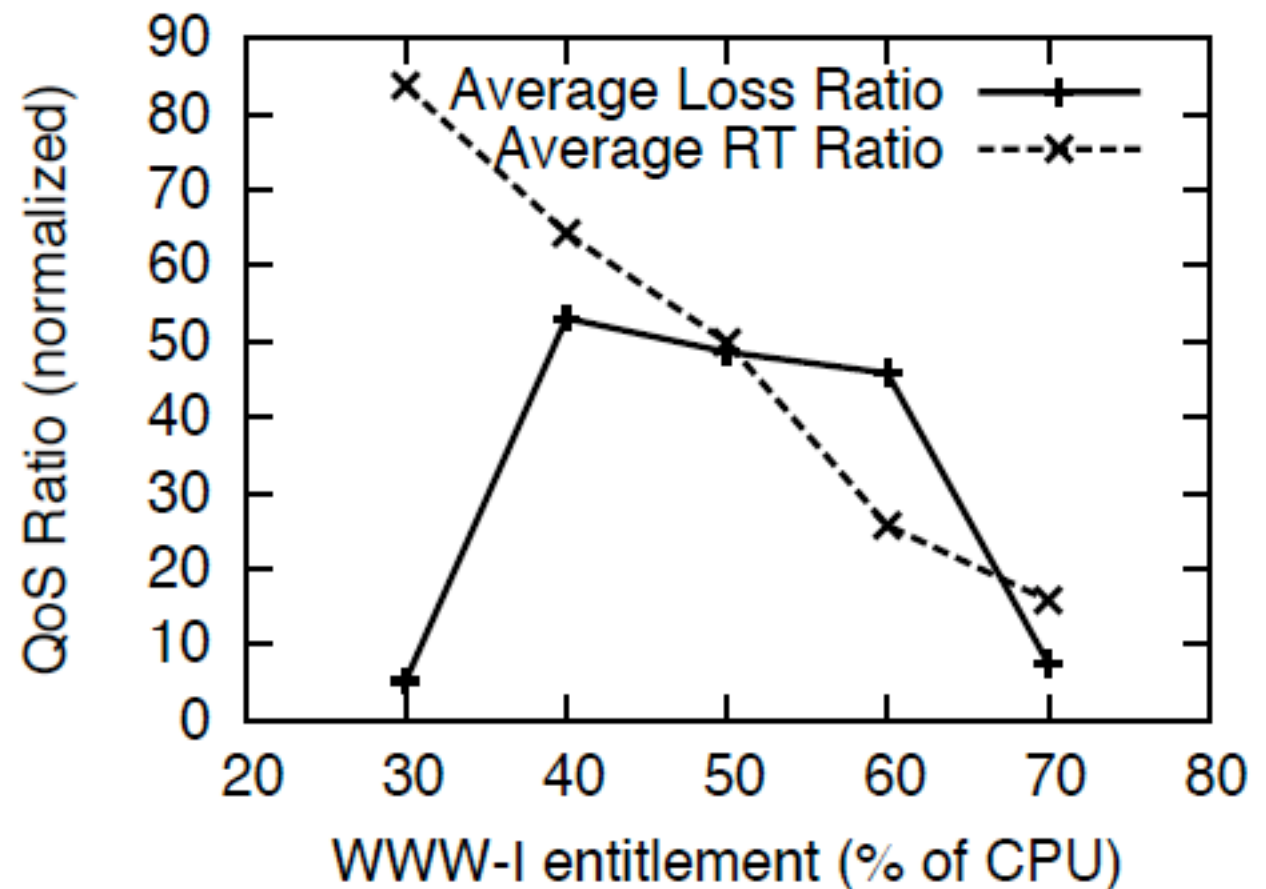
Modeling co-hosted multi-tier applications



(a) Loss



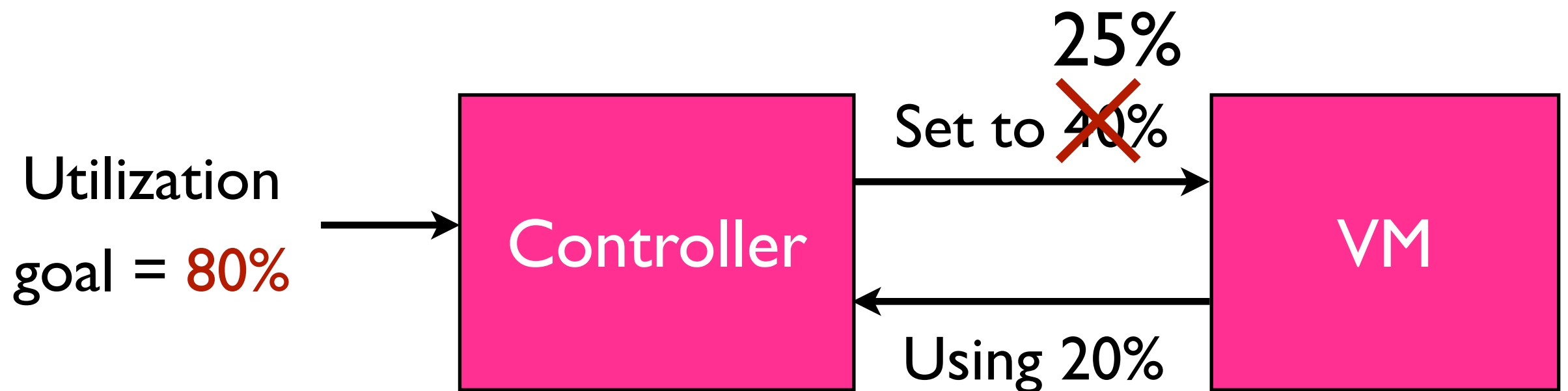
(b) Response time



(c) Loss ratio and RT ratio

Response time ratio is more **controllable** than loss ratio

DESIGN

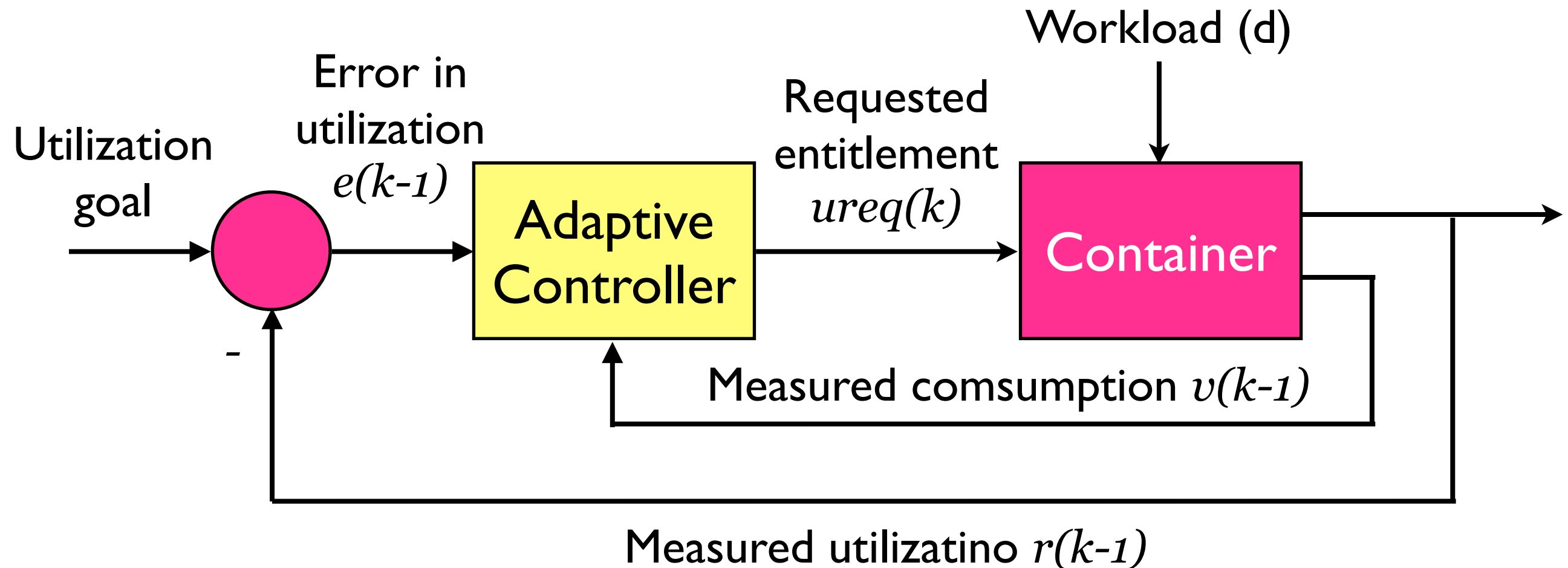


Utilization: $20/40 * 100 = 50\%$

New Utilization: $20/25 * 100 = 80\%$

An example of the controller

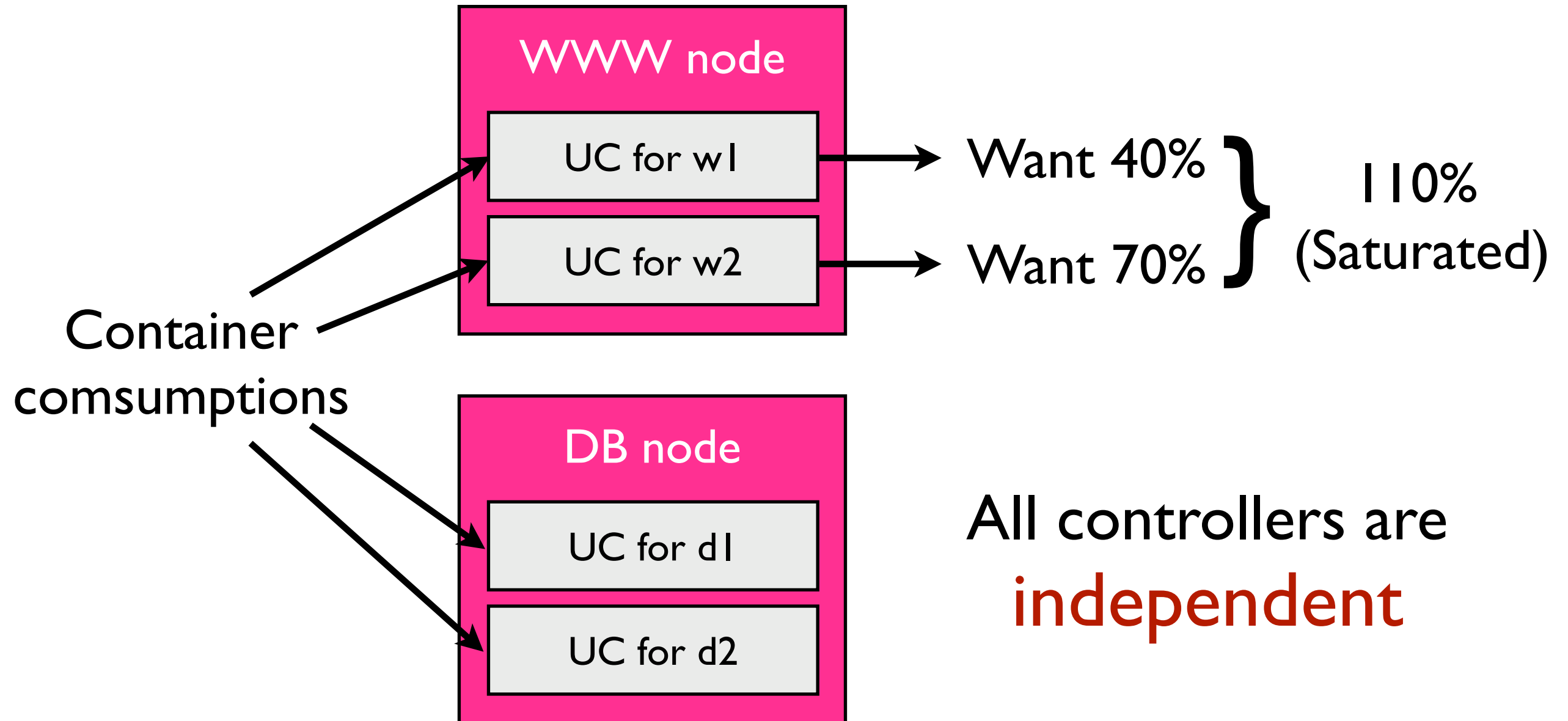
DESIGN



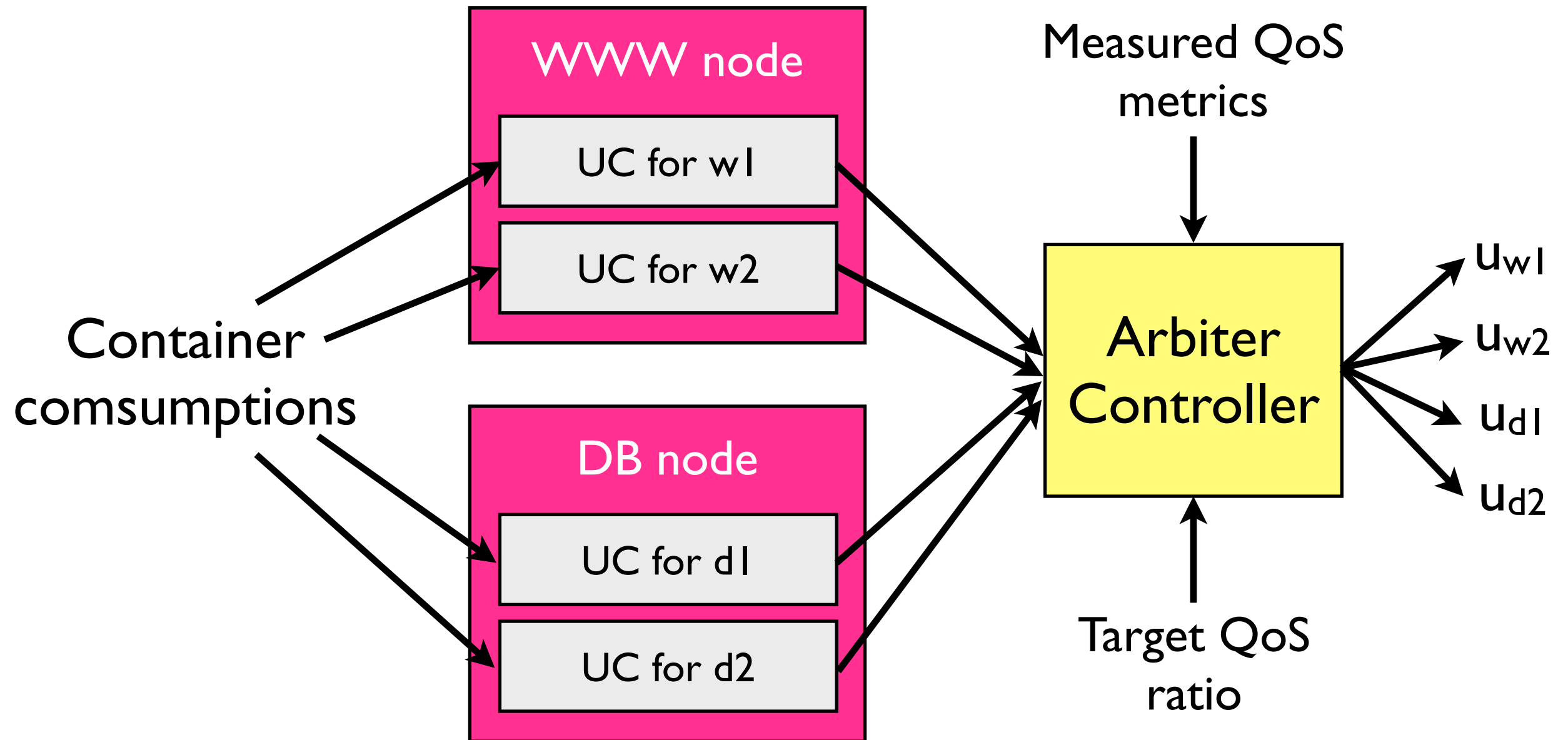
$$ureq(k) = ureq(k - 1) - K(k)e(k - 1)$$

Adaptive utilization controller

DESIGN

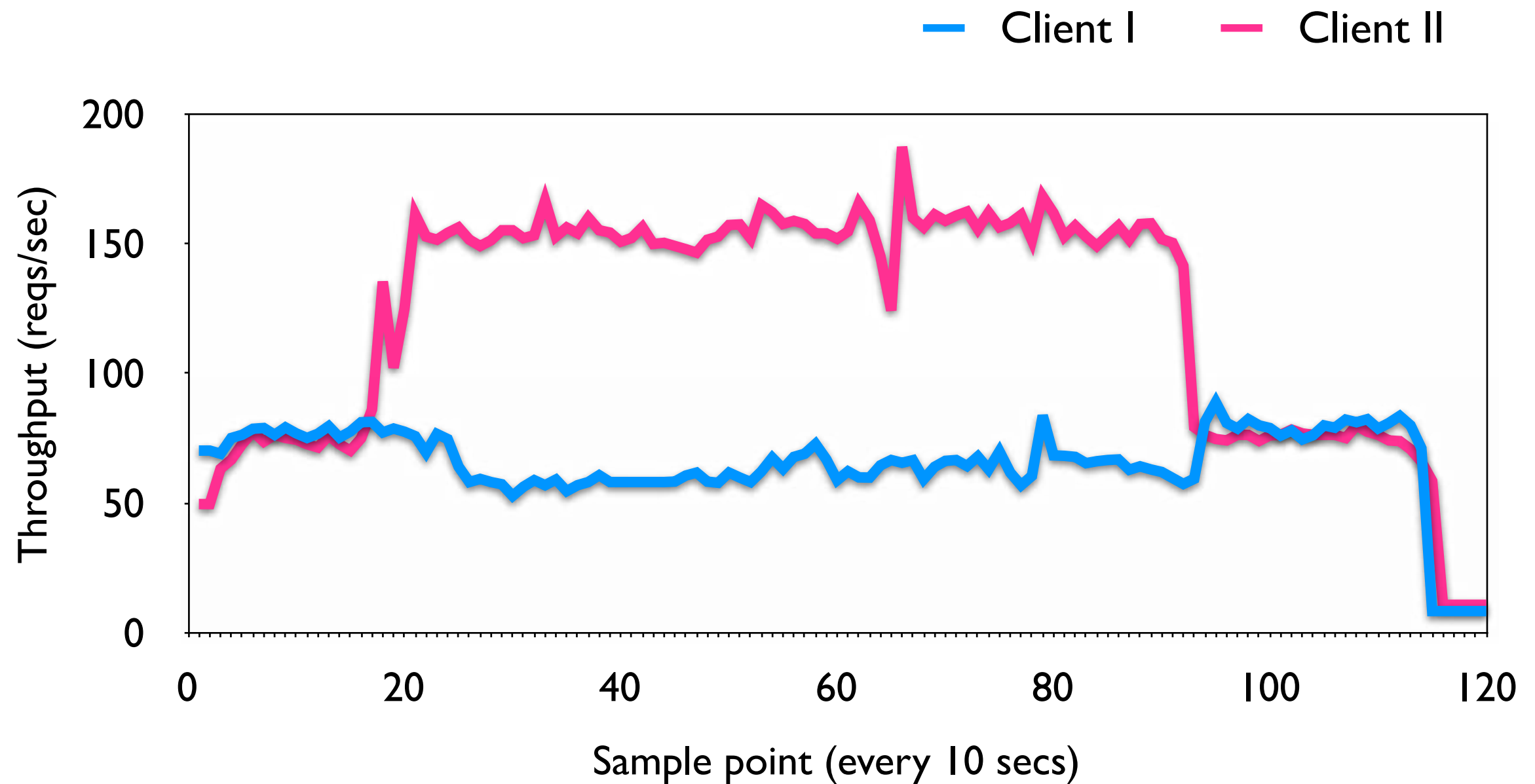


DESIGN



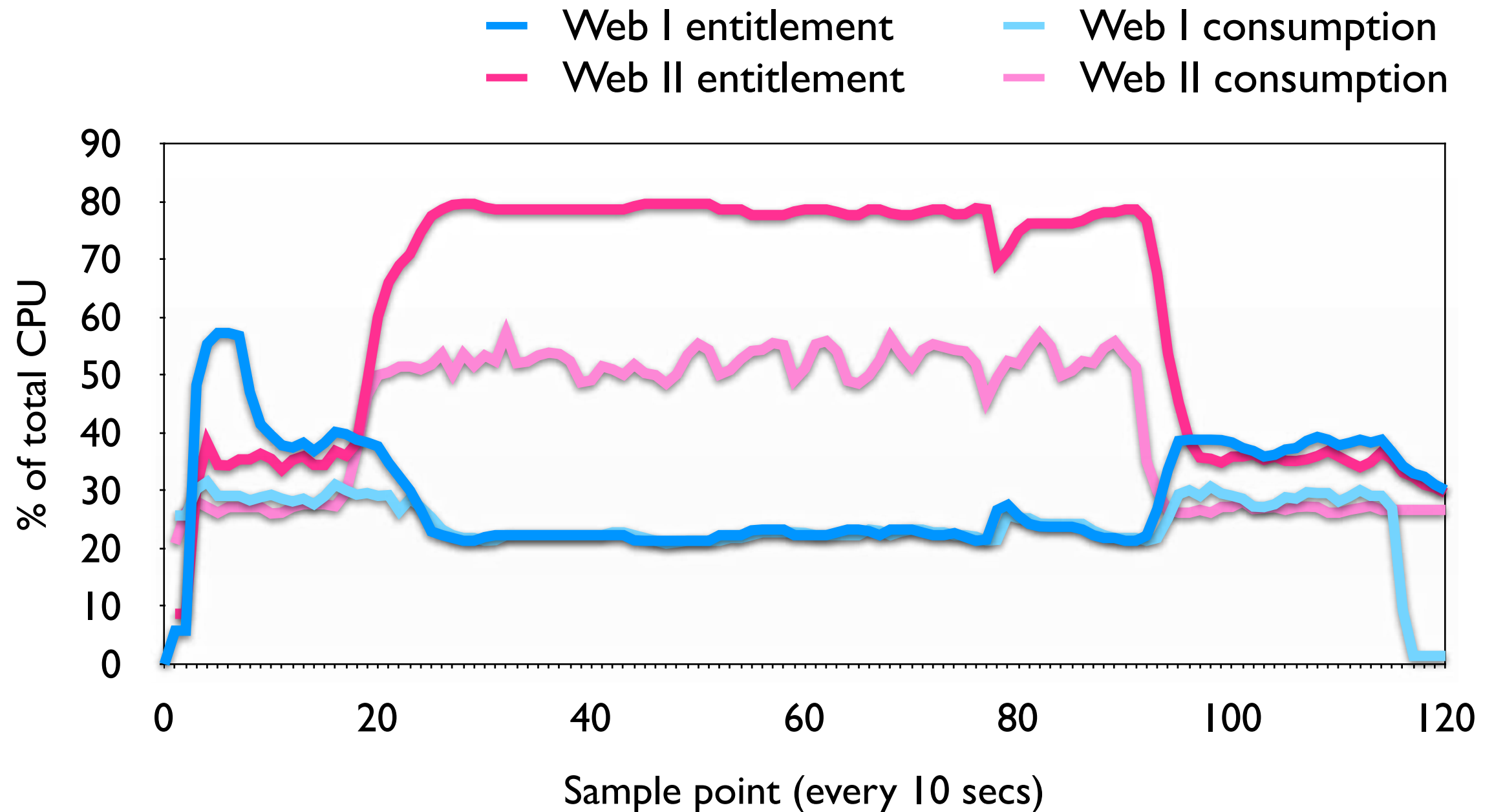
A two-layered controller architecture

EVALUATION



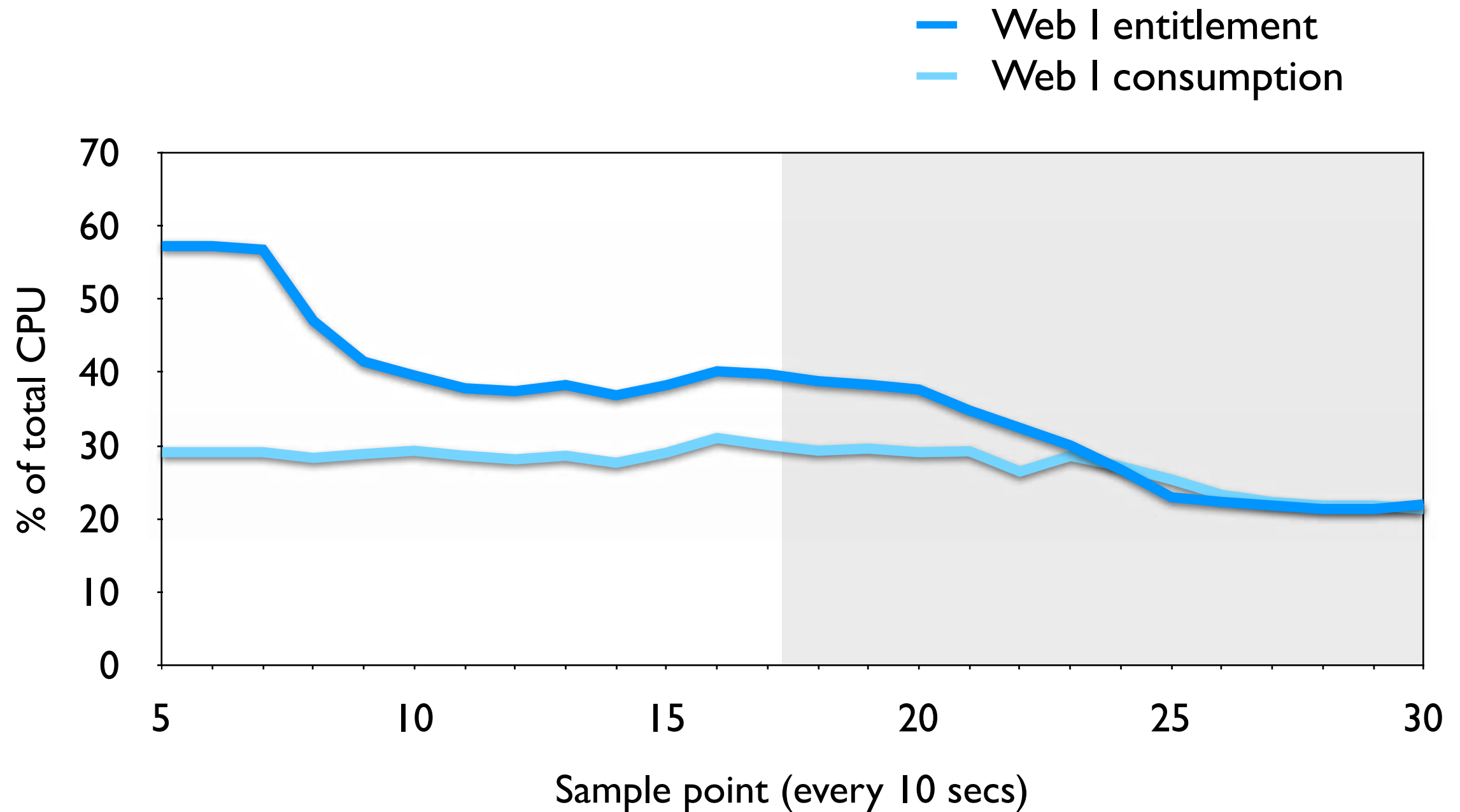
Varing load generated by clients

EVALUATION



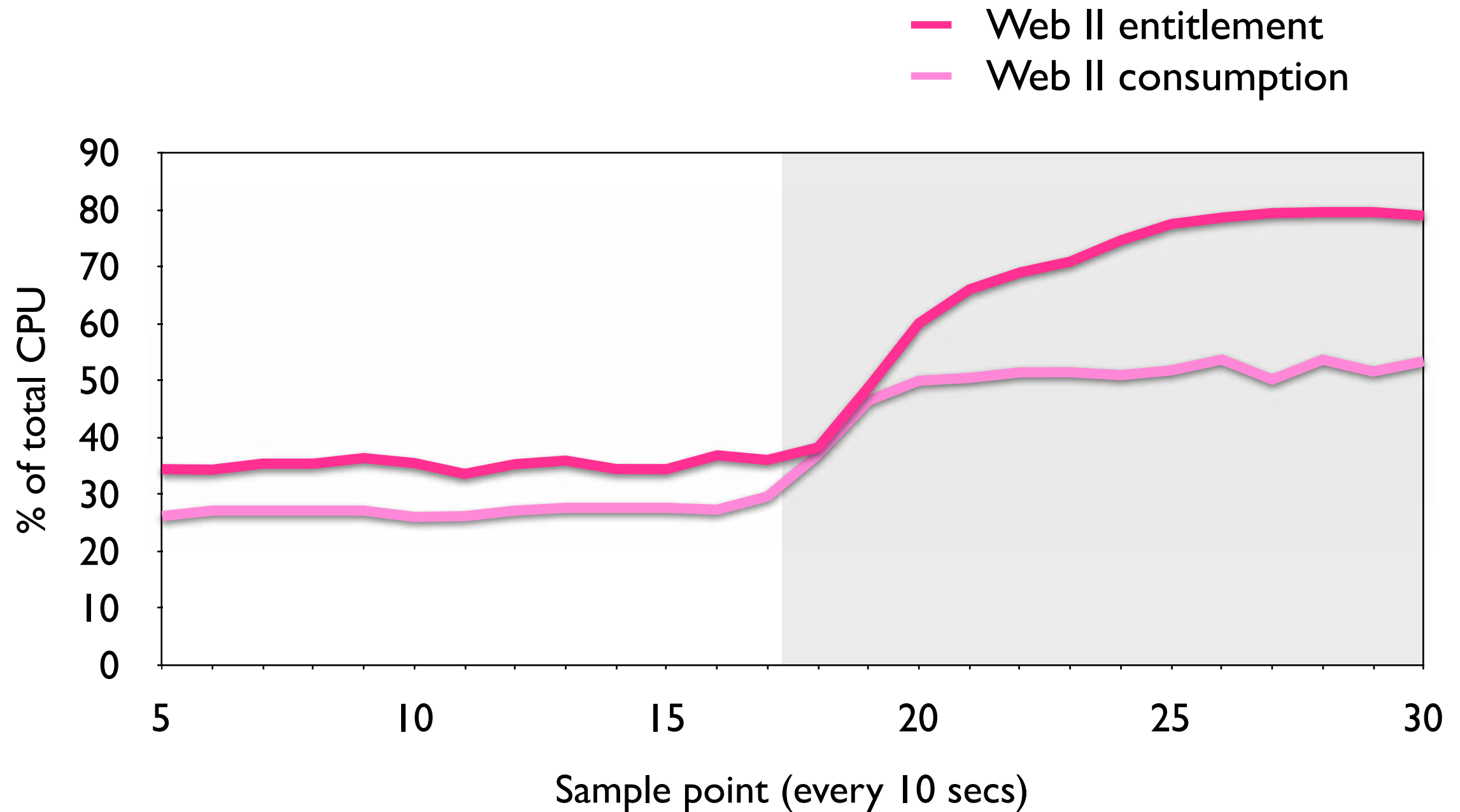
CPU entitlement and consumption

EVALUATION



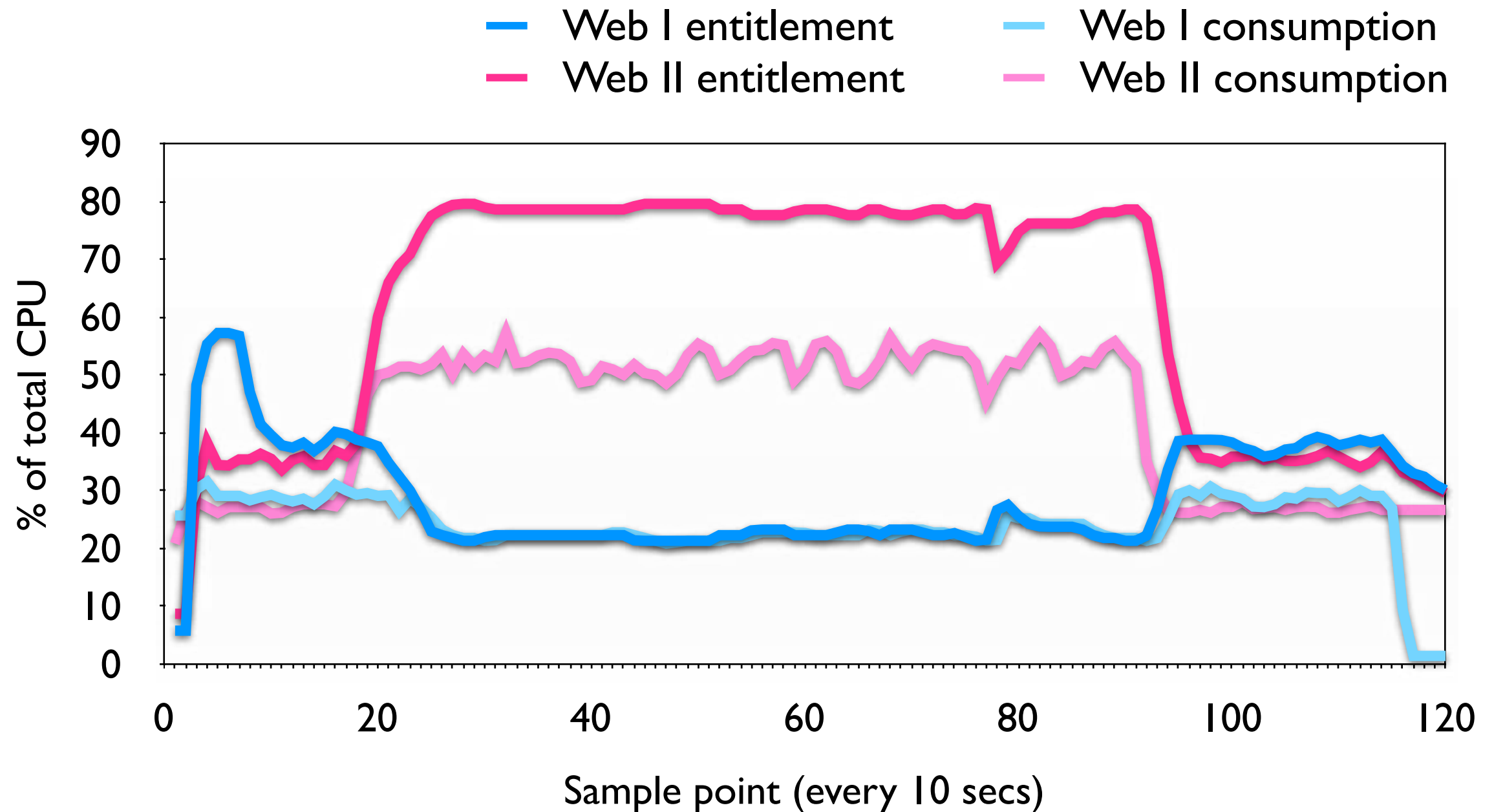
CPU entitlement and consumption

EVALUATION



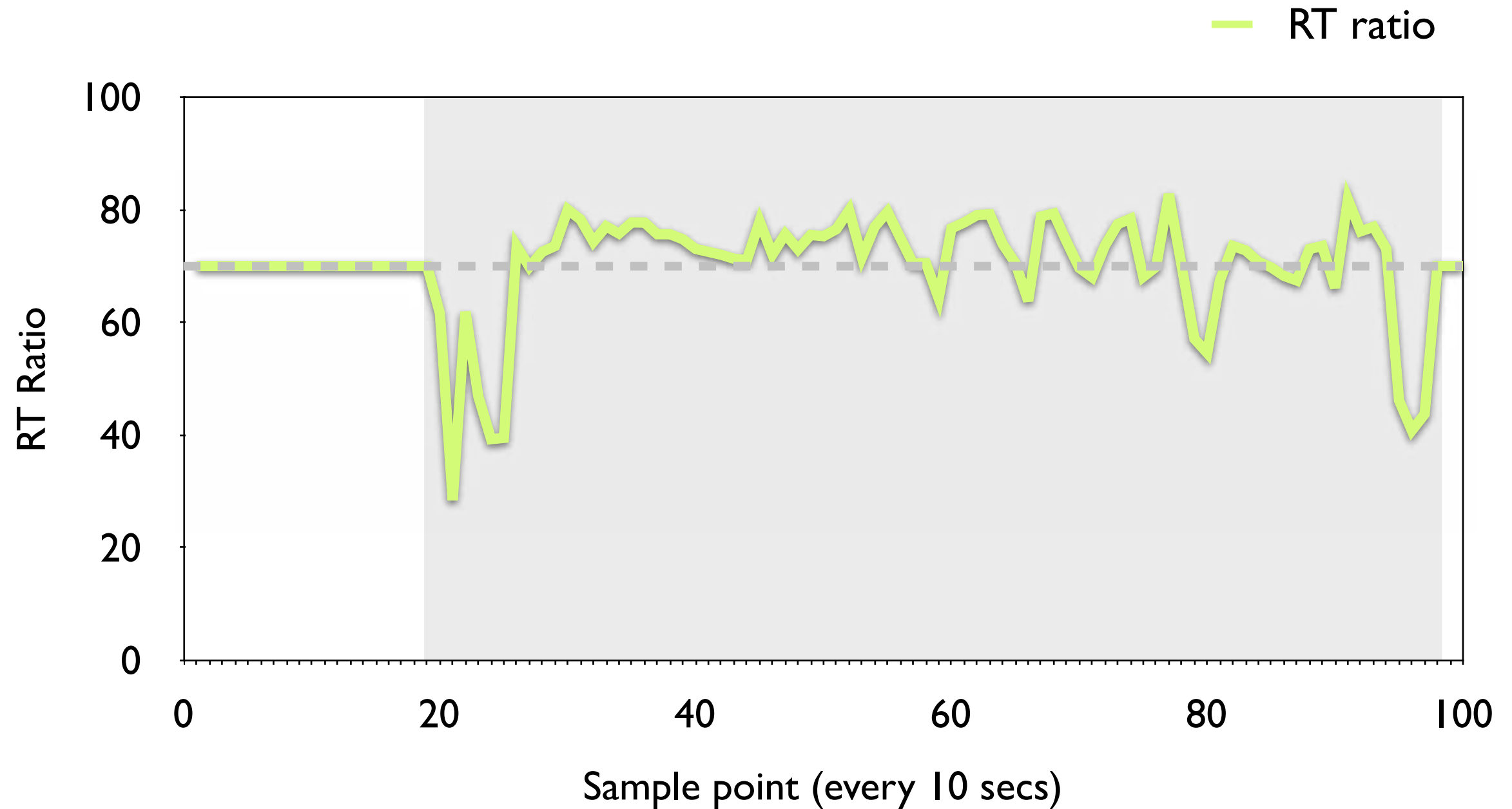
CPU entitlement and consumption

EVALUATION



CPU entitlement and consumption

EVALUATION



CPU entitlement and consumption

CONCLUSIONS

- Achieves high utilization of the data center while meeting application-level QoS goals
- Be able to provide a specified level of QoS differentiation between applications under overload conditions