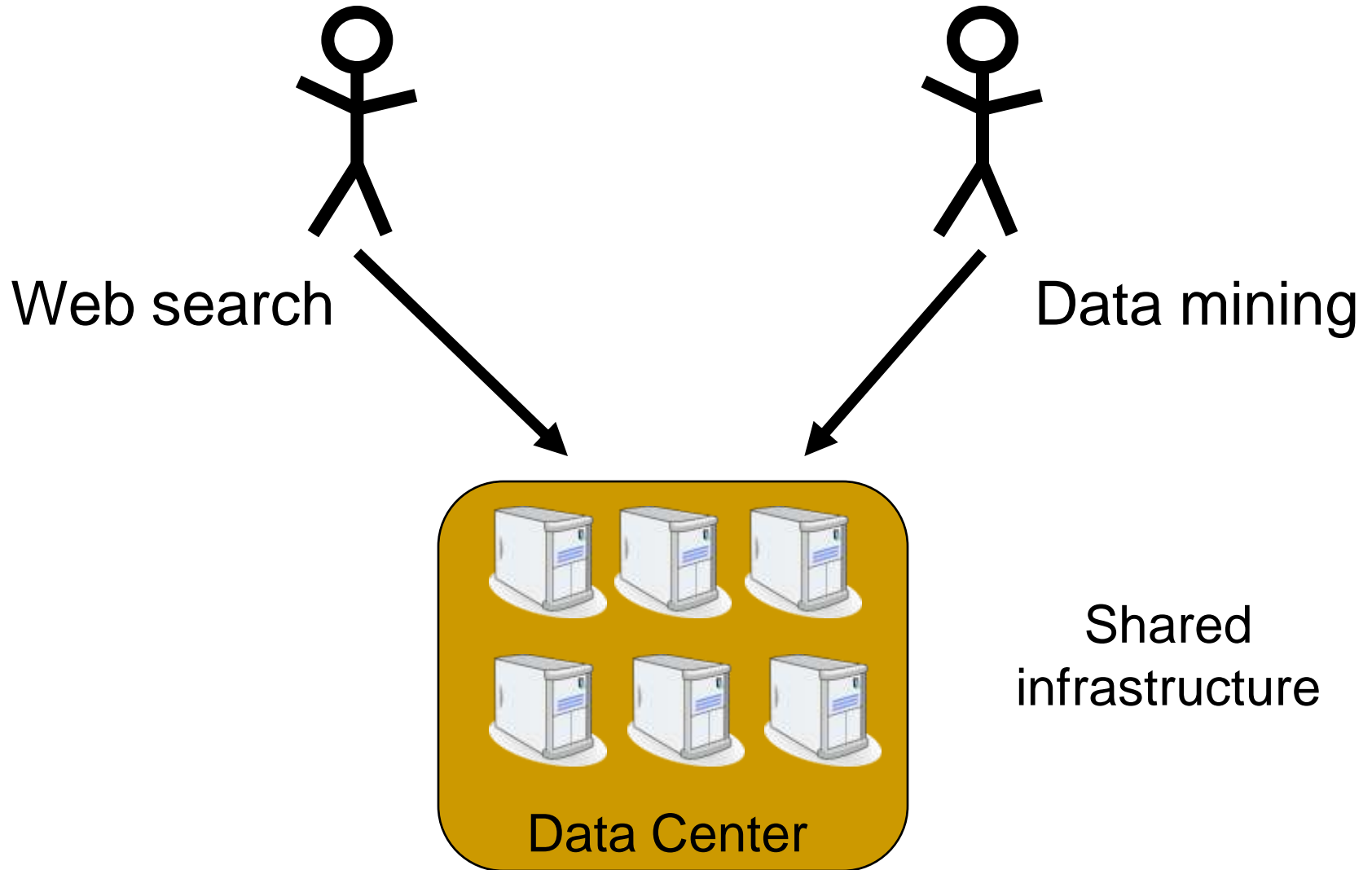# Got control ?

## AutoControl: Automated Control of MultipleVirtualized Resources

Pradeep Padala, Karen Hou, Xiaoyun Zhu*, Mustfa Uysal[†], Zhikui Wang[†], Sharad Singhal[†], Arif Merchant[†], Kang G. Shin

University of Michigan, VMware* and HP Labs[†]

# Typical scenario in shared infrastructures



Web search

Data mining

Shared
infrastructure

Data Center

# Application requirements

|  |  |
|---|---|
| Web search | Data mining |
| Fast searches | Analyze large data |
| ✓ Low response time | ✓ High throughput |

✓ QoS differentiation 3:1

# Hosting applications

**Physical partitioning**

app1 web
Node I

app1 db
Node II

app2
Node III

app3
Node IV

× Wasteful
× Difficult to manage

**Virtual data center**

app1 web | app2
Virtualization
Node I

app1 db | app3
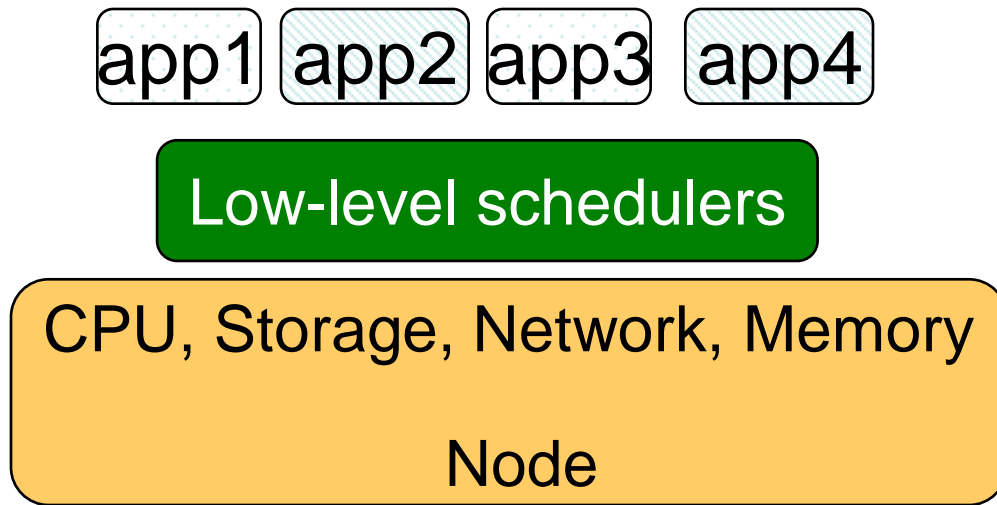Virtualization
Node II

✓ Improved utilization
✓ Reduced costs
✓ High flexibility

Problem: How to allocate resources?

# Approach I: Work-conserving mode

| app1 | app2 | app3 | app4 |

**Low-level schedulers**

CPU, Storage, Network, Memory

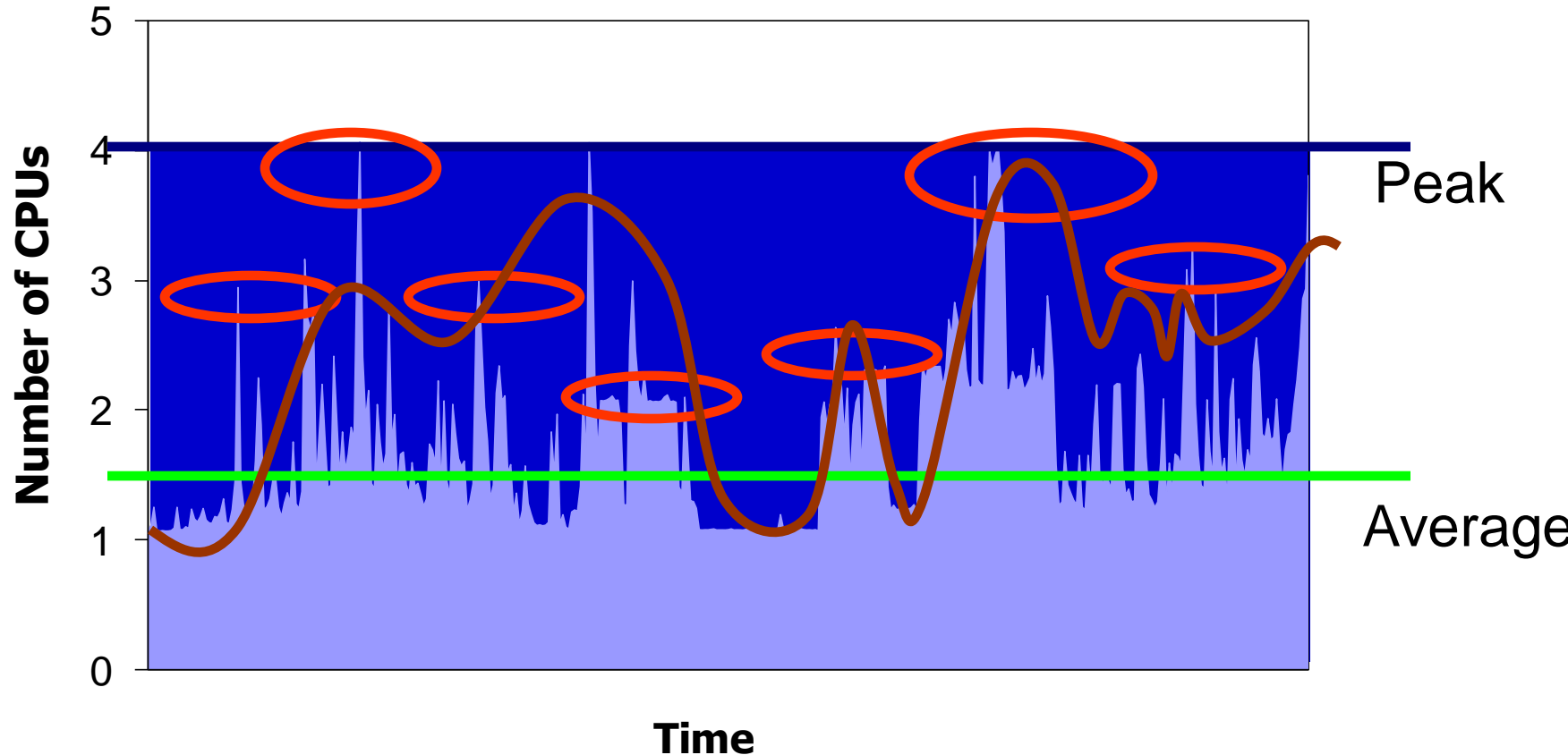Node

All applications can use as much resources as they require

- Greedy applications cause SLO (service level objective) violations

- How to prioritize? – no differentiation

- How to use scheduler mechanisms to meet targets?

# Approach II: Static allocation
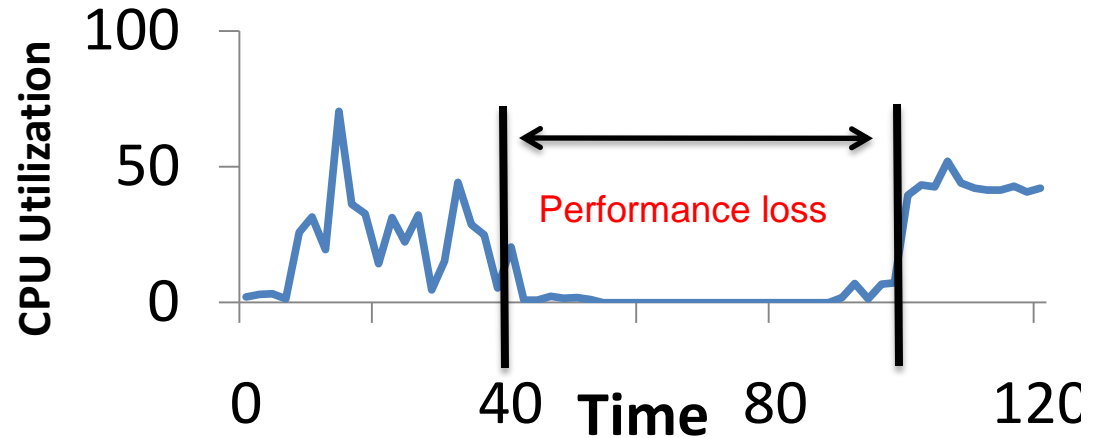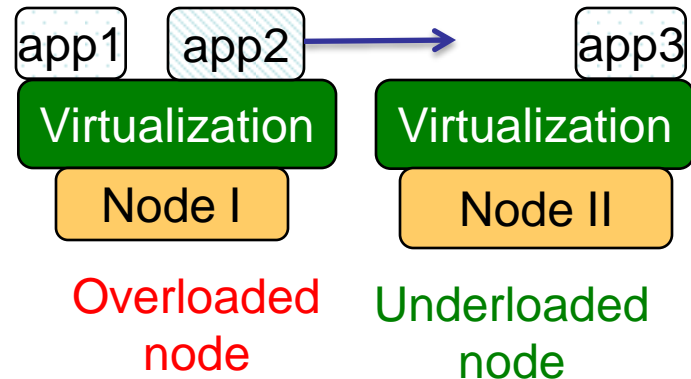
Bursty load     Poor respo Wasted resources
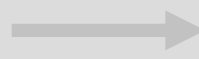
Finding the right share is hard!

# Approach III: Migration



- **Good** choice for long term overload

- **Poor** choice for bursty loads

- Adds **overhead** to an already overloaded node

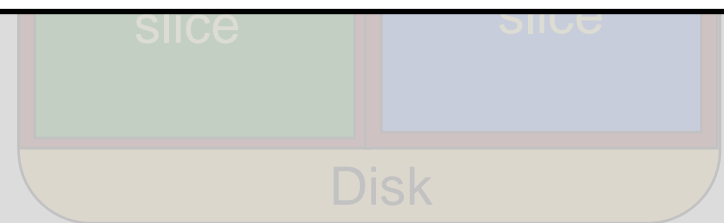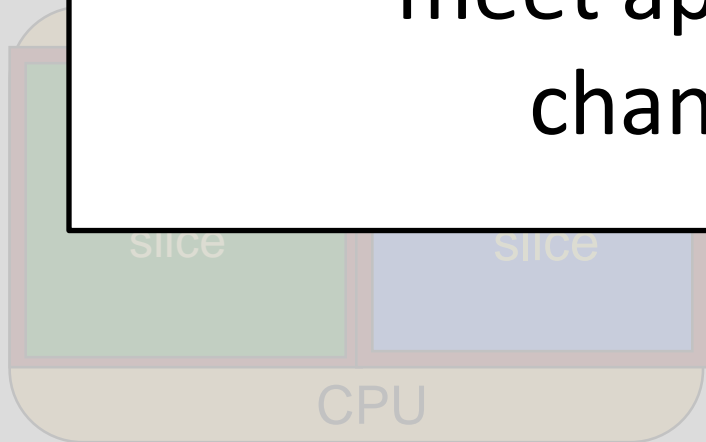- **SLO violations** while being migrated

App goals → Controller    Goals met ?    Yes

**Automatically** set resource shares to meet application targets in changing conditions

slice    slice    slice    slice

CPU    Disk

20%    50%    70%    50%

# Previous work

- Distributed resource allocation
  - AronSIGMETRICS00, ChaseSOSP01, ShenOSDI02
  - Orthogonal to our approach

- Low-level schedulers
  - Credit, CFQ, SFQ, WaldspurgerOSDI02, GulatiTR07
  - Policy vs. Mechanism

- QoS mechanisms
  - ChandraIWQOS03, UrgaonkarICAC05
  - Developed for a single resource or application

- Control theory based
  - AbdelZaherTPDS02, HellersteinBook04, KarlssonIWQOS04
  - Applied in other scenarios

# Outline
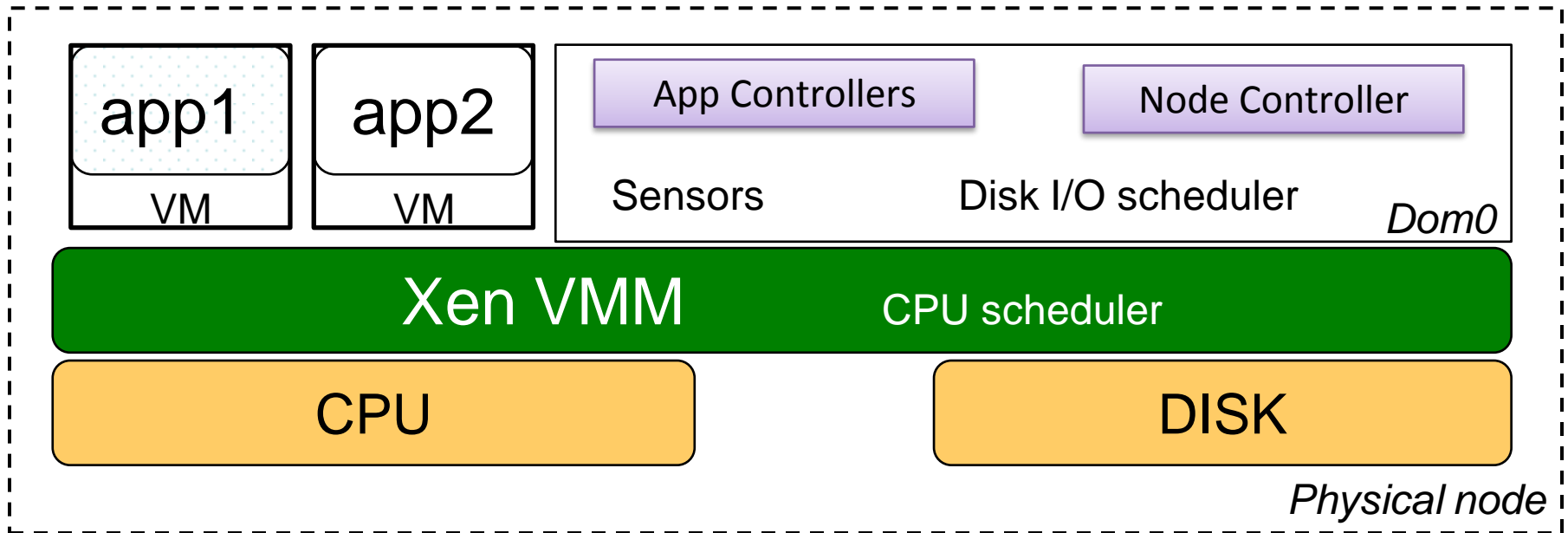
- *Motivation*
- *Background*
- *Idea*

- Modeling
- Controller Design
  - Application Controller
  - Node Controller

- Evaluation
  - Synthetic workloads
  - CPU and Disk bottlenecks

# AutoControl system – 1,000ft view

app1
VM

app2
VM

App Controllers

Node Controller

Sensors

Disk I/O scheduler

*Dom0*

Xen VMM          CPU scheduler

CPU

DISK

*Physical node*

Every control interval

App Controller — Figures out the resource share required for a single app to meet its targets
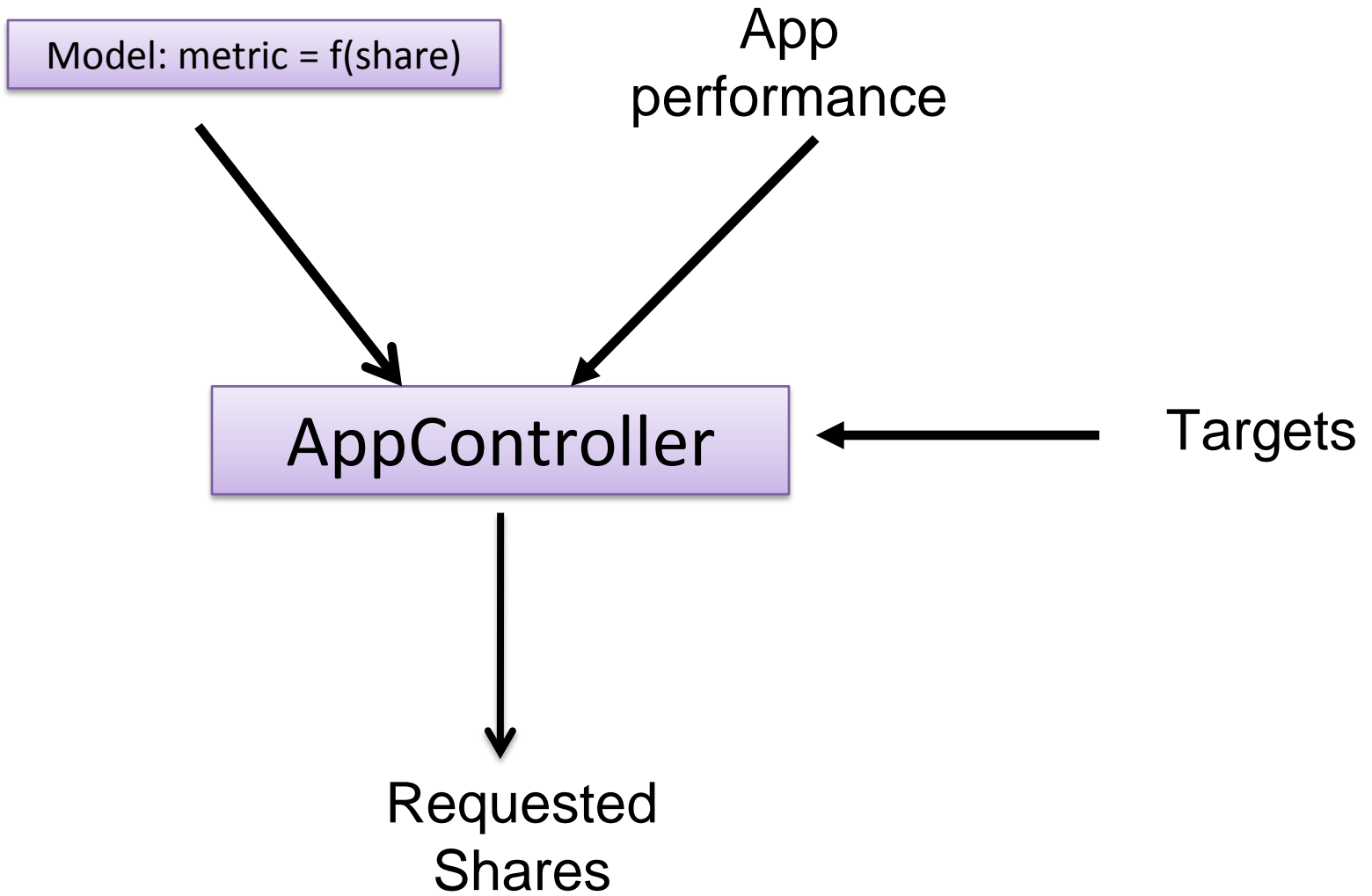
Node Controller — Arbitrates among multiple applications
All node controllers are independent

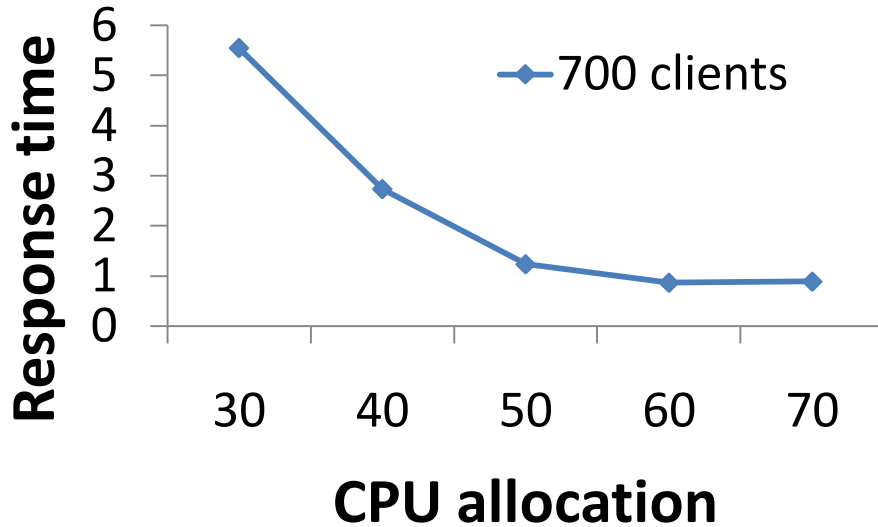CPU and disk schedulers — Final shares are set
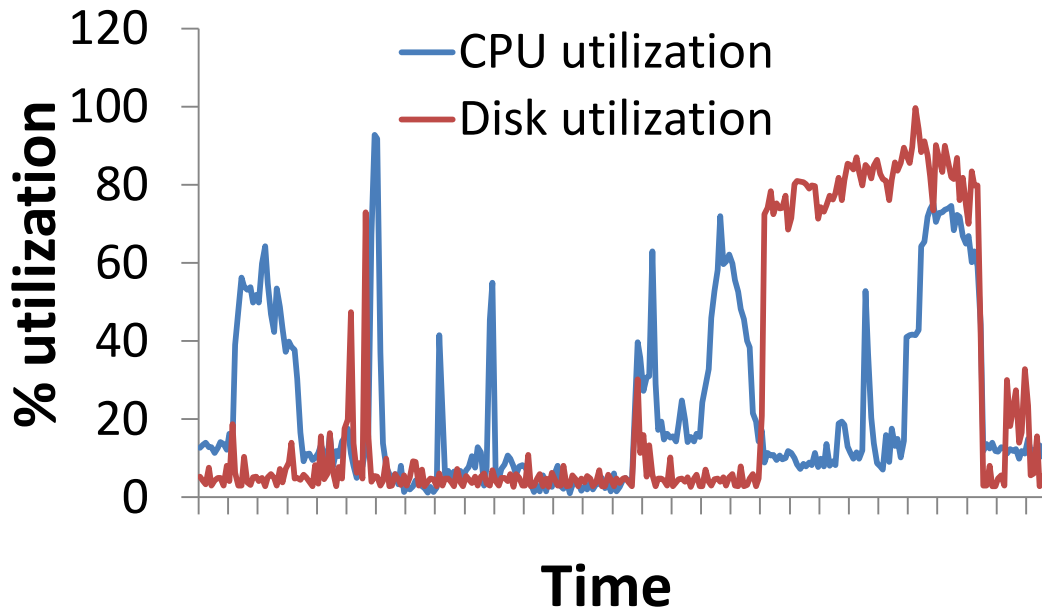
# Application controller

Model: metric = f(share)

App performance

**AppController**

Targets

Requested Shares
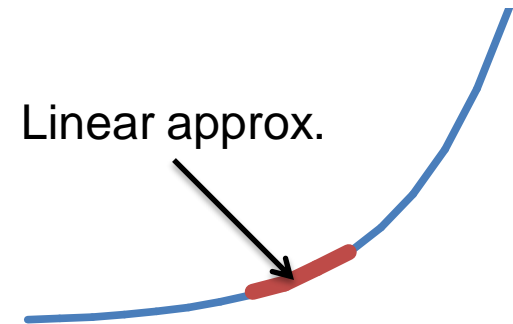
# Why is modeling hard?



Non linear relationships

Multiple resources

Bottleneck shifts

Multi-tier applications

# Solution: Dynamic black box modeler

✓ Nonlinearity approximated using linear equations

Linear approx.

✓ Multiple resources & multi-tier apps modeled with Multi Input Multi Output (MIMO) model

$$y_k = (a1)(y_{k-1}) + (b1\ b2)\begin{pmatrix} webcpu_{k-1} \\ webdsk_{k-1} \end{pmatrix} + (b3\ b4)\begin{pmatrix} dbcpu_{k-1} \\ dbdsk_{k-1} \end{pmatrix}$$

Current performance            Prev performance            Resource shares

*First order*

✓ Parameters (a1… b1 …) updated recursively (Recursive Least Squares RLS)

Aggression factors

$$Cost = W\|y - y_{ref}\|^2 + Q\|u_k - u_{k-1}\|^2$$

$$y = f(u_k)$$

$$y_{ref} = \text{target}$$

$$u_k = \begin{pmatrix} cpu \\ dsk \end{pmatrix}$$

K$^{\text{th}}$ time interval

Track target          Don't go wild
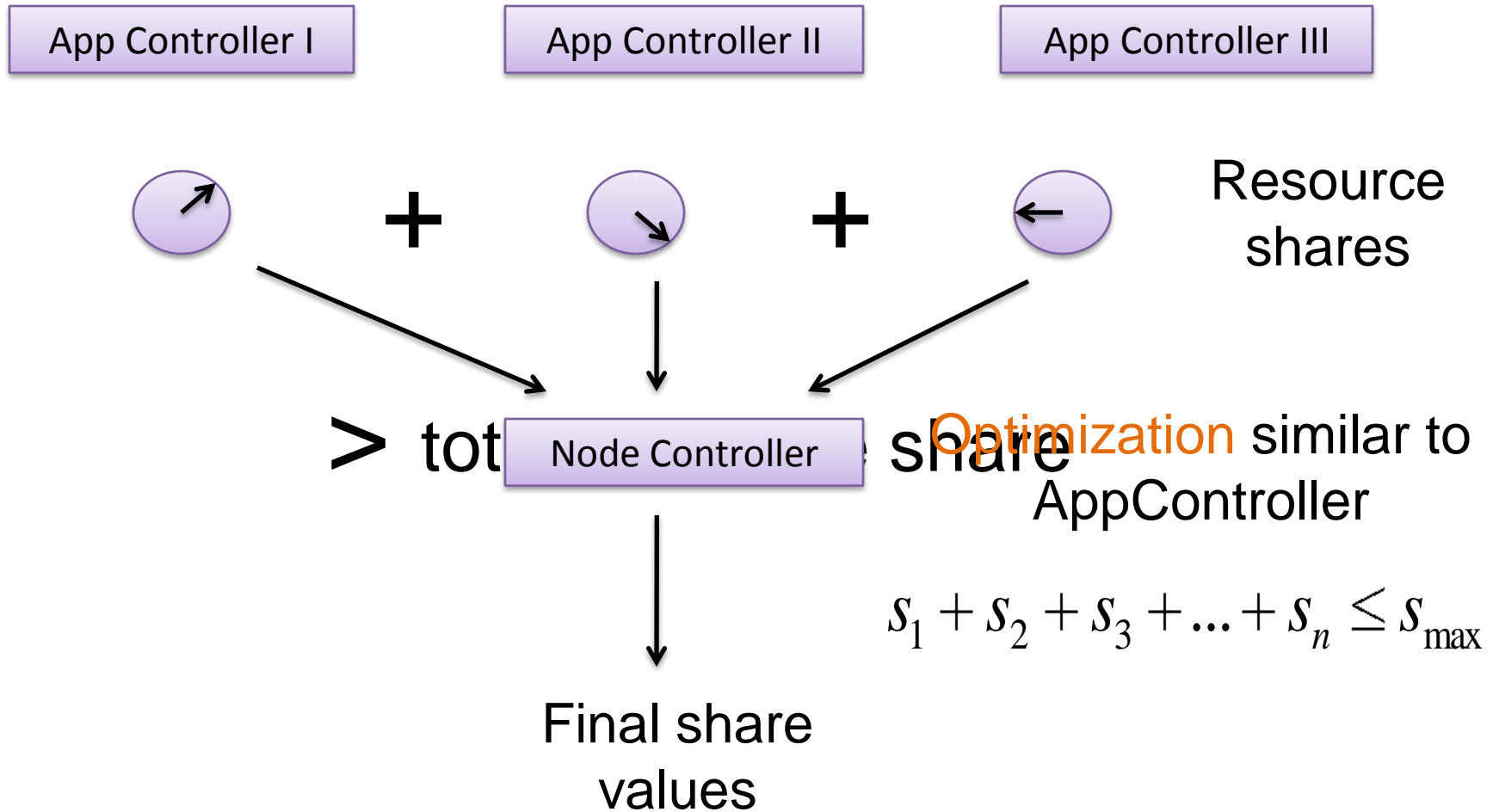
✓ Minimize cost

✓ Quadratic solvers to find $u_k$

Simplified Linear Quadratic Regulator formulation

Gory details: [CDC'07]

# Control with consolidation
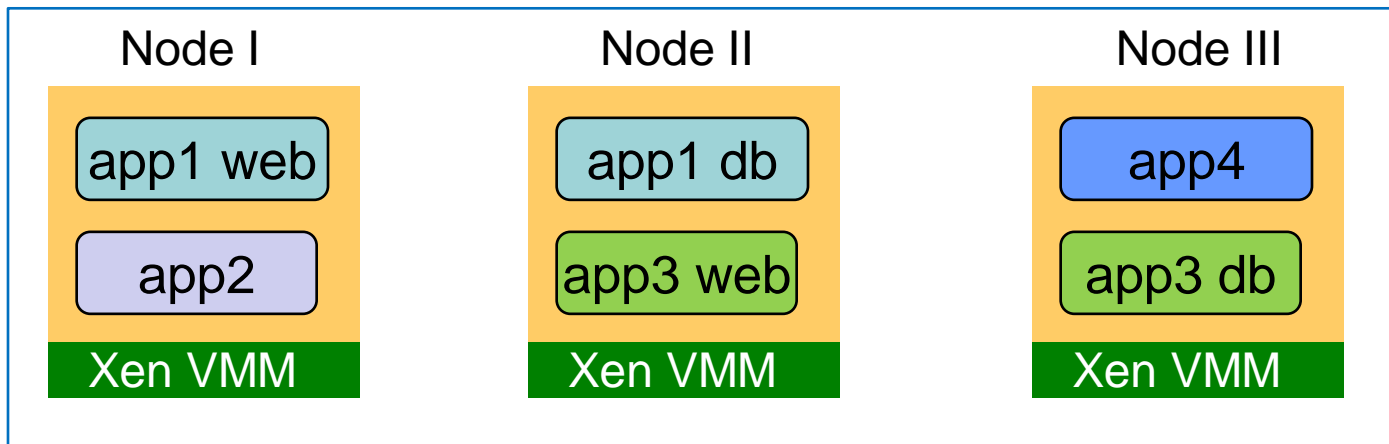
App Controller I    App Controller II    App Controller III



Resource shares

$+$    $+$

$>$ total node share

Node Controller

Optimization similar to AppController

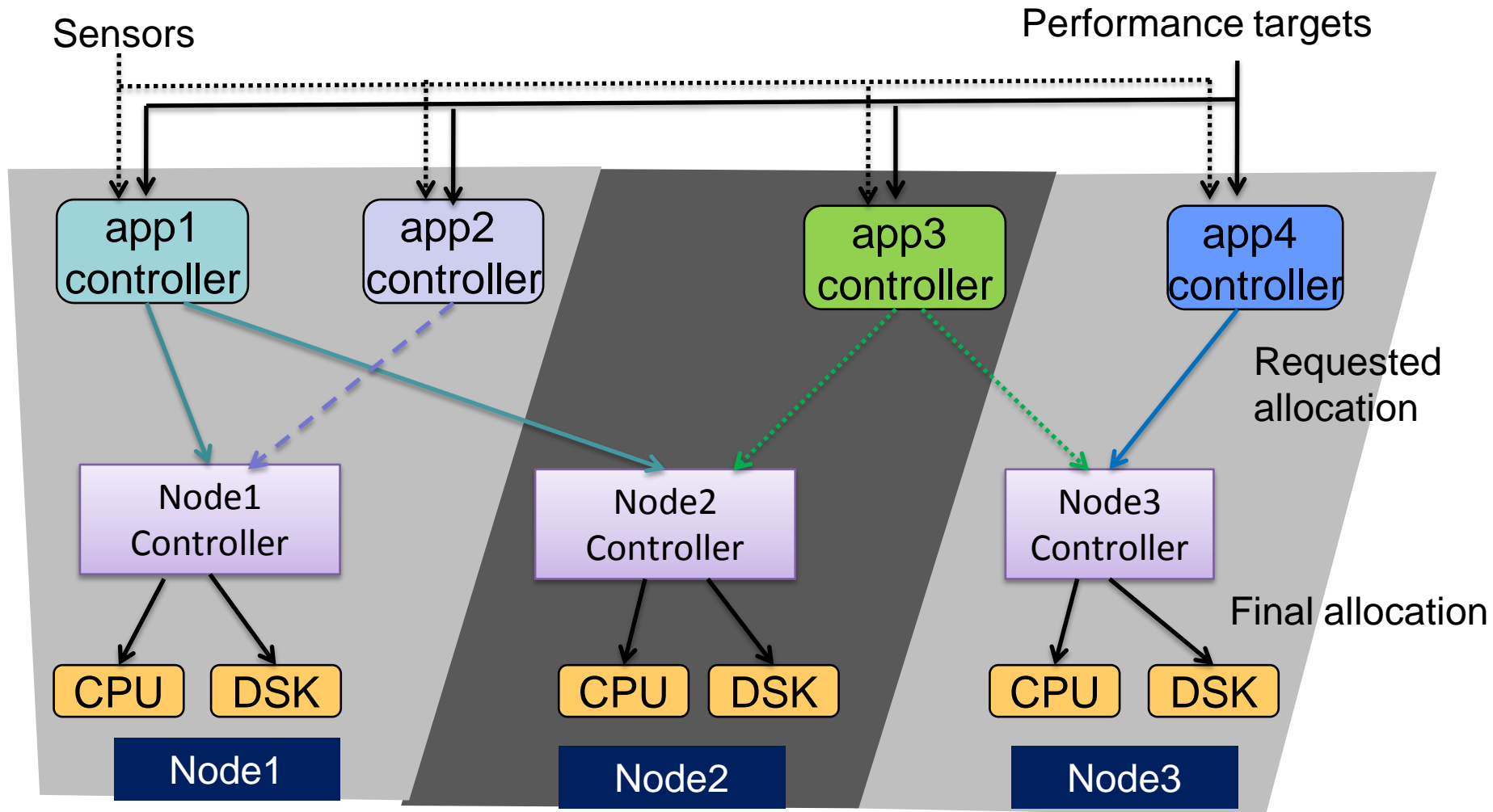$$s_1 + s_2 + s_3 + ... + s_n \leq s_{max}$$

Final share values

# Control for data center scale

- Why not centralized controller? – variable explosion

- Why not combine app and node controllers?
  - Applications may span multiple nodes



- One AppController for application

- One NodeController per node

# Distributed control



Experiments with 40 nodes on Emulab are successful

# Outline



- *Motivation*
- *Background*
- *Idea*



- Modeling
- Controller Design
  – Application Controller
  – Node Controller



- Evaluation
  – Synthetic workloads
  – CPU and Disk bottlenecks

# Evaluation

- Applications
  - RUBiS: eBay style auction benchmark
  - TPC-W: Transactional web e-Commerce benchmark
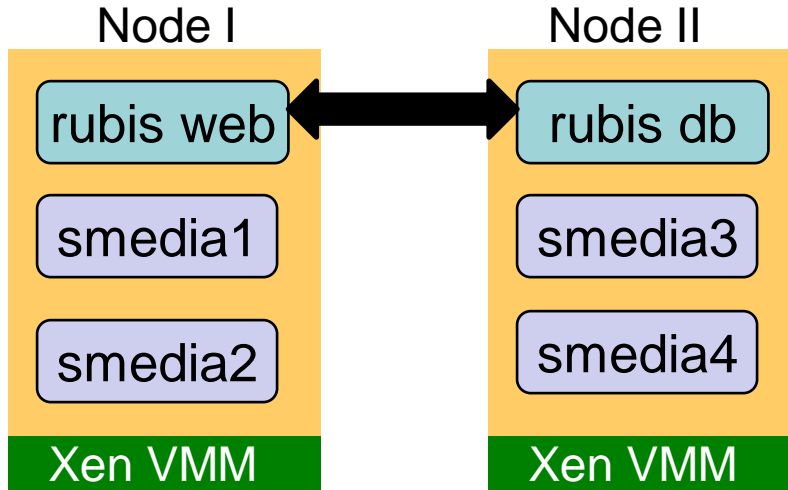  - Smedia: Custom built secure media server
- Workloads
  - Synthetic
  - CPU and disk bottlenecks
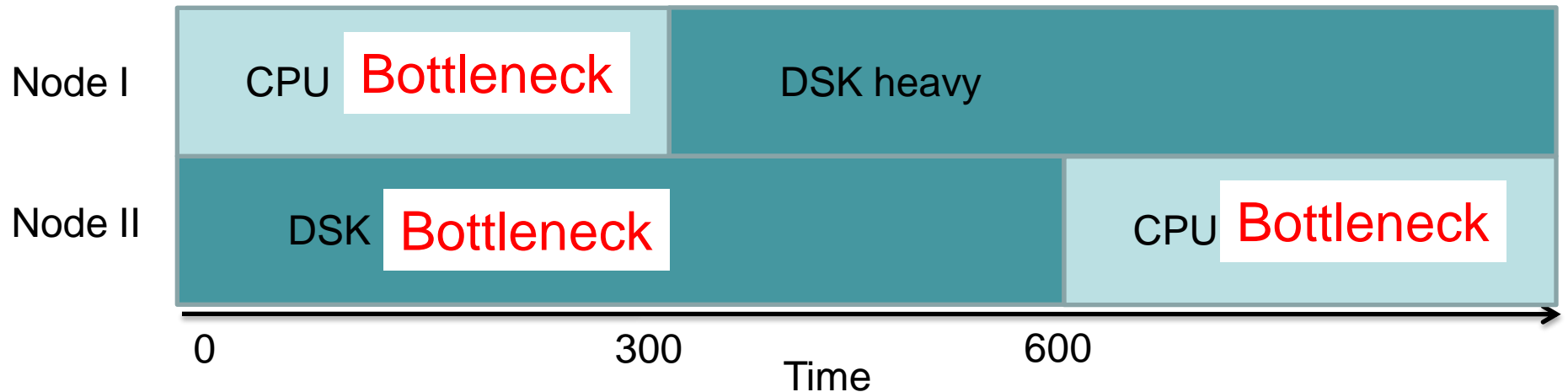- Evaluation questions
  - Can the controller meet targets?
  - Can it identify bottlenecks over time in different tiers and fix them?
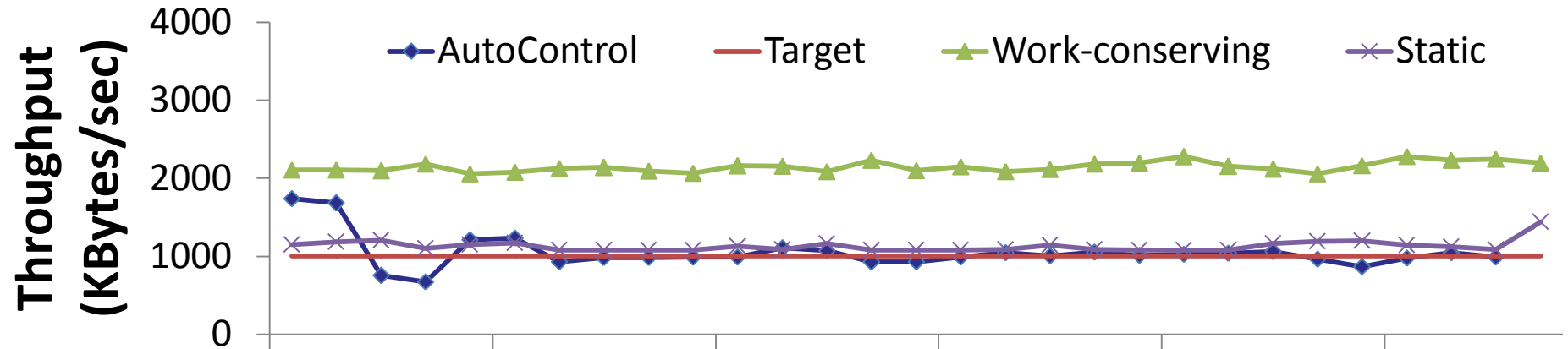  - Can it identify bottlenecks of different resources? (ex. CPU/Disk)

# Experimental setup

Node I

| rubis web |
|---|
| smedia1 |
| smedia2 |
| Xen VMM |

Node II
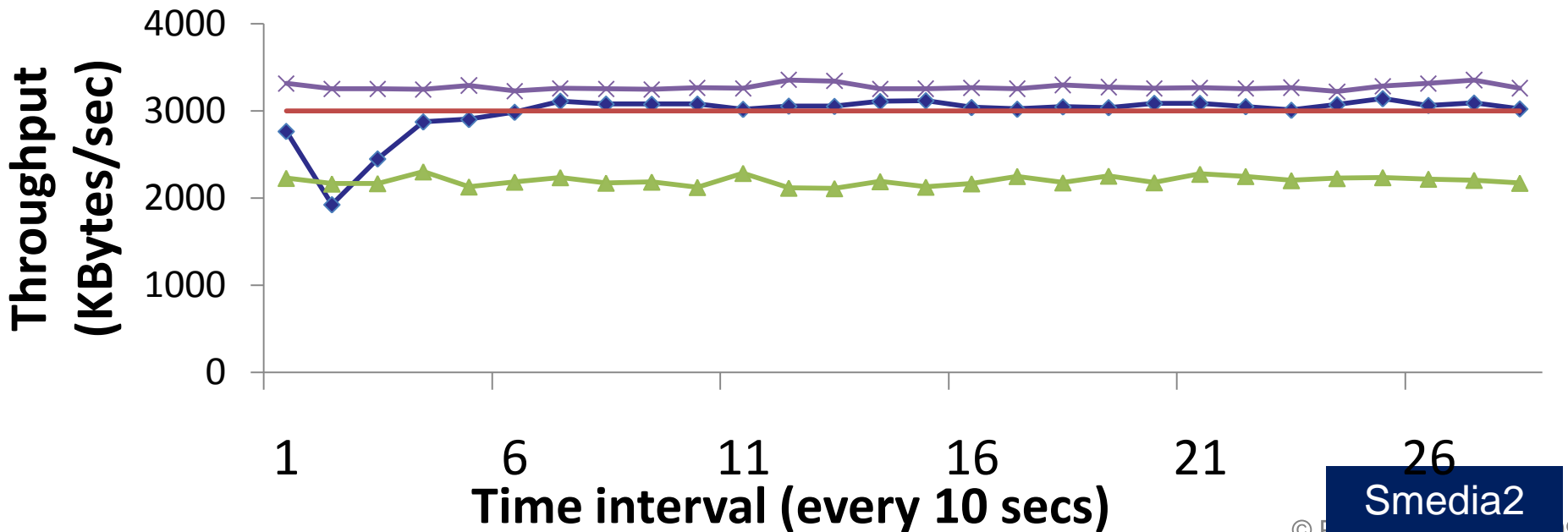
| rubis db |
|---|
| smedia3 |
| smedia4 |
| Xen VMM |

Targets
RUBiS – 100 req/sec
Smedia1 – 1000 Kbytes/sec
Smedia2 – 3000 Kbytes/sec
Smedia3 – 15000 Kbytes/sec
Smedia4 – 10000 Kbytes/sec

**Node I**
CPU **Bottleneck**    DSK heavy

**Node II**
DSK **Bottleneck**    CPU **Bottleneck**

0        300        600
Time

# Node 1: CPU bottleneck

# Node 2: DISK -> CPU bottleneck



**Time interval (every 10 secs)**
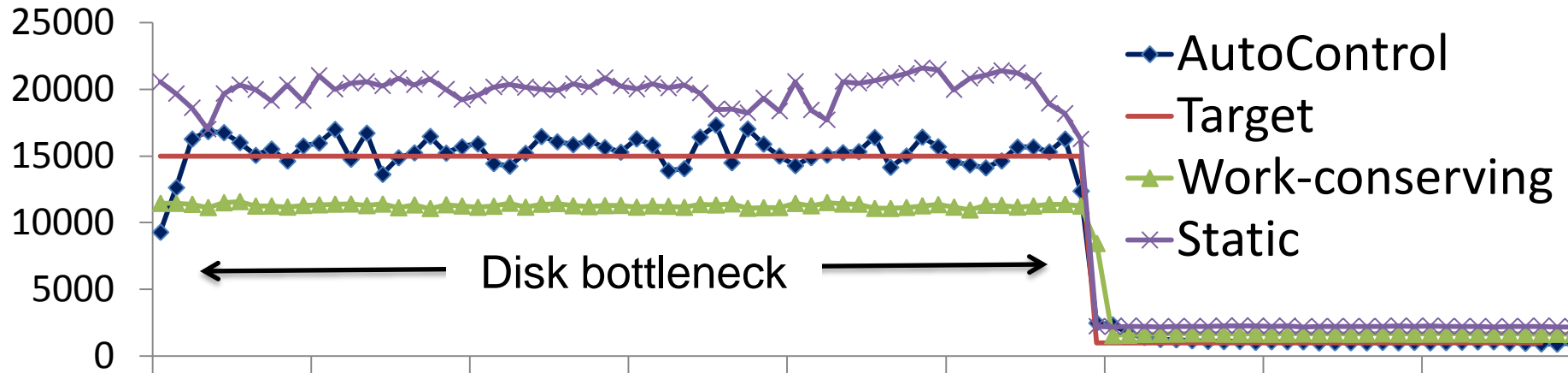
Legend (top chart):
- AutoControl
- Target
- Work-conserving
- Static

Disk bottleneck
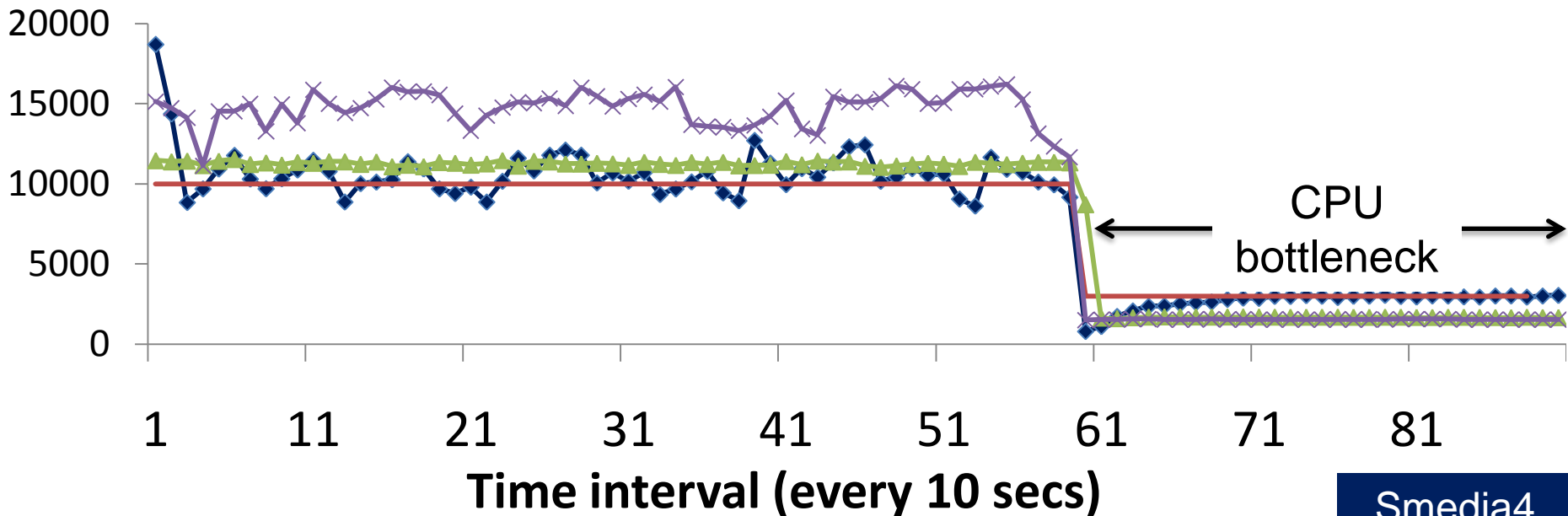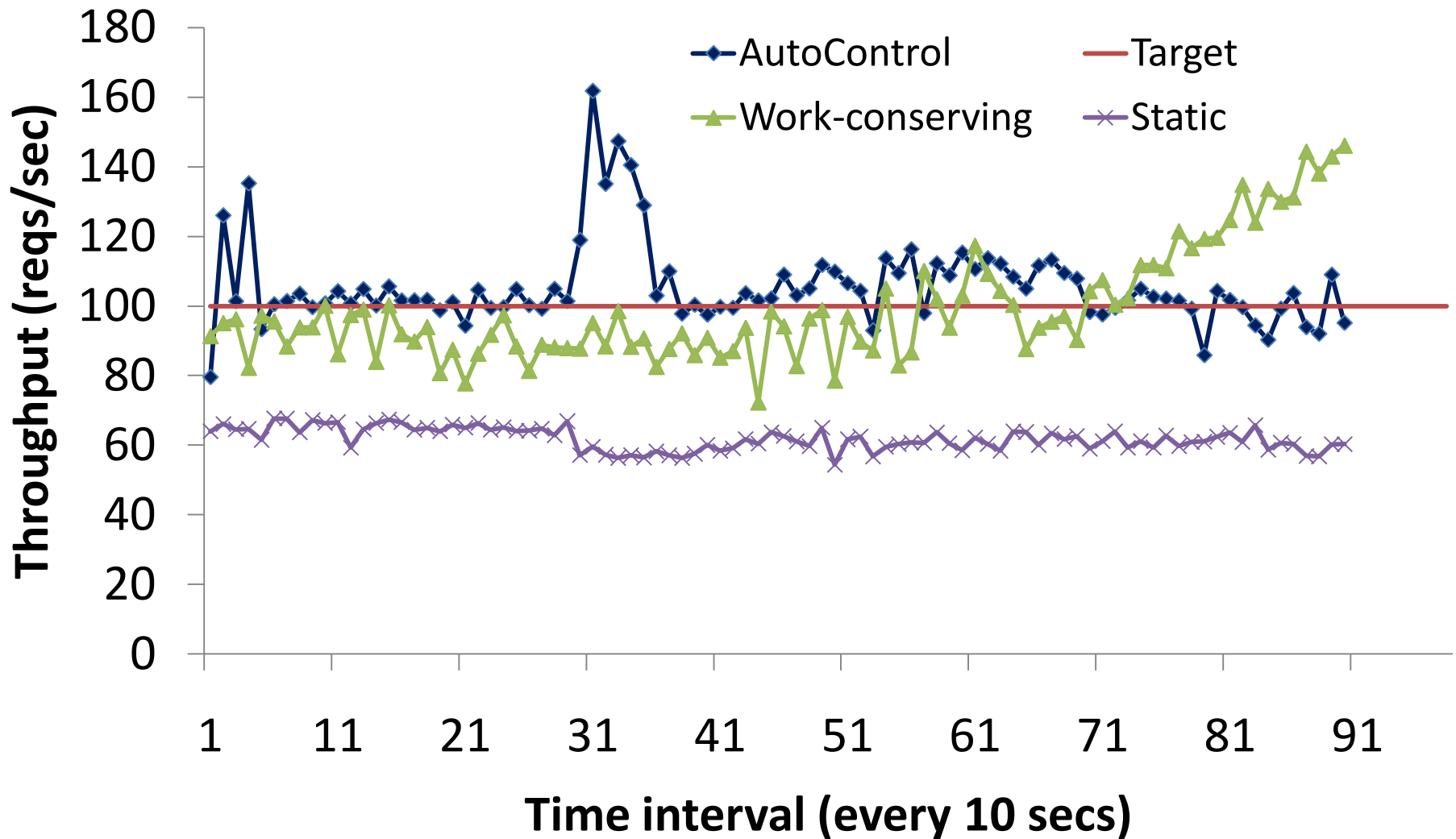
CPU bottleneck

Smedia3

Smedia4

# Node 1 & 2: RUBiS performarnce

# Experiment summary – average error

$$\text{Error} = \frac{\|y - y_{ref}\|}{yref} * 100$$

| App | Work-conserving | Static | AutoControl |
|---|---|---|---|
| RUBiS | 13.8% | 38.3% | 8.2% |
| Smedia1 | 100% | 12.1% | 4.3% |
| Smedia2 | 26.2% | 9.6% | 2.2% |
| Smedia3 | 44.3% | 61.4% | 10.1% |
| Smedia4 | 24.6% | 47.5% | 9.3% |

AutoControl achieves <=10% error

# Limitations and future work

- Modeling challenges

  o Non linear, fast changing workloads create problems

  o Combining white-box and black-box models

- Actuator and sensor behavior

  o Inaccuracies in measurements may lead to inaccurate models.

  o We are limited by what the scheduler can do

- Network and memory control

  o Earlier efforts with network control were unsuccessful

  o Preliminary memory control [IM'09 min-conference]

- Integrating with migration

# Summary

## Automated Control of Multiple Virtualized Resources

✓Feedback control can be successfully applied to computer systems

  o Dynamic black box modeler captures complex dynamics

  o AppController can compute shares to meet targets for a single app

  o NodeController arbitrates among competing apps

✓Distributed architecture that scales well



ppadala@umich.edu