# Maximizing Resource Efficiency for Next Generation Cloud Platforms

**Jashwant Raj Gunasekaran**

*Advisors*: *Dr. Mahmut T. Kandemir & Dr. Chita R. Das*

*High Performance Computing Lab*

Dissertation Defense

May 6, 2021

PennState
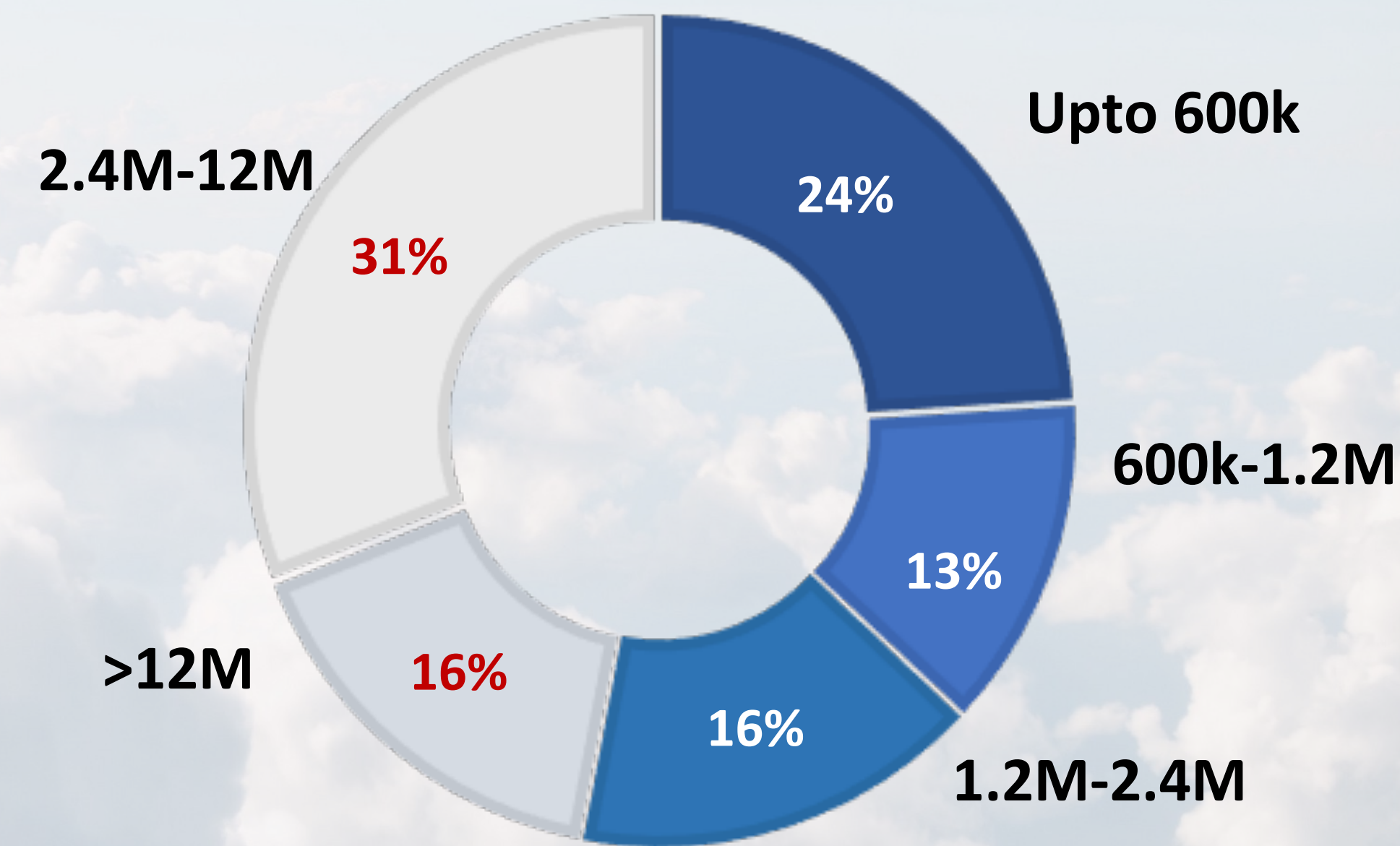College of Engineering

PennState
High Performance
Computing Lab

# RESEARCH PHILOSOPHY

Cloud is about *how* you do computing, not *where* you do computing!

Paul Maritz, Former CEO, Vmware

PennState
College of Engineering
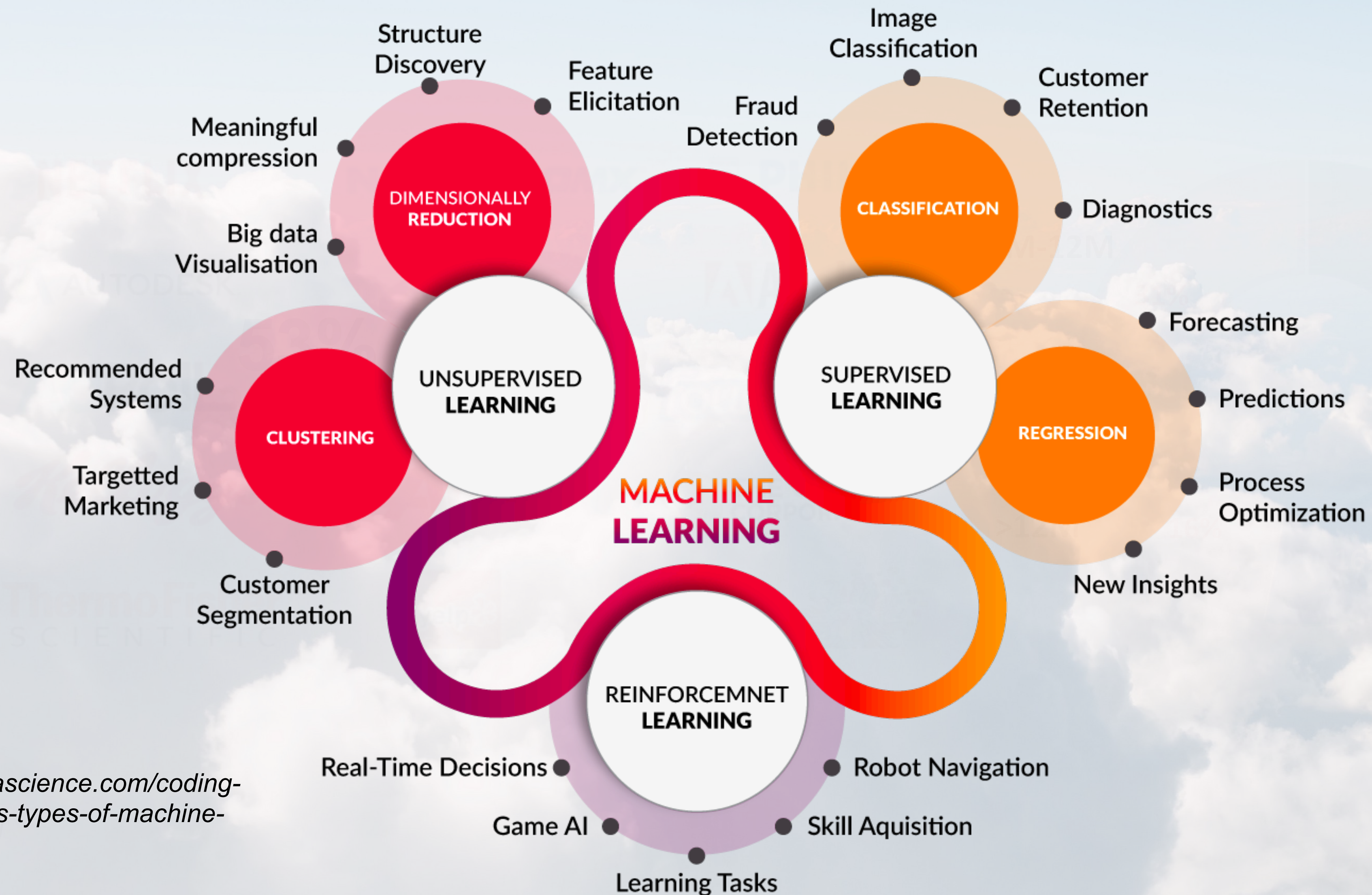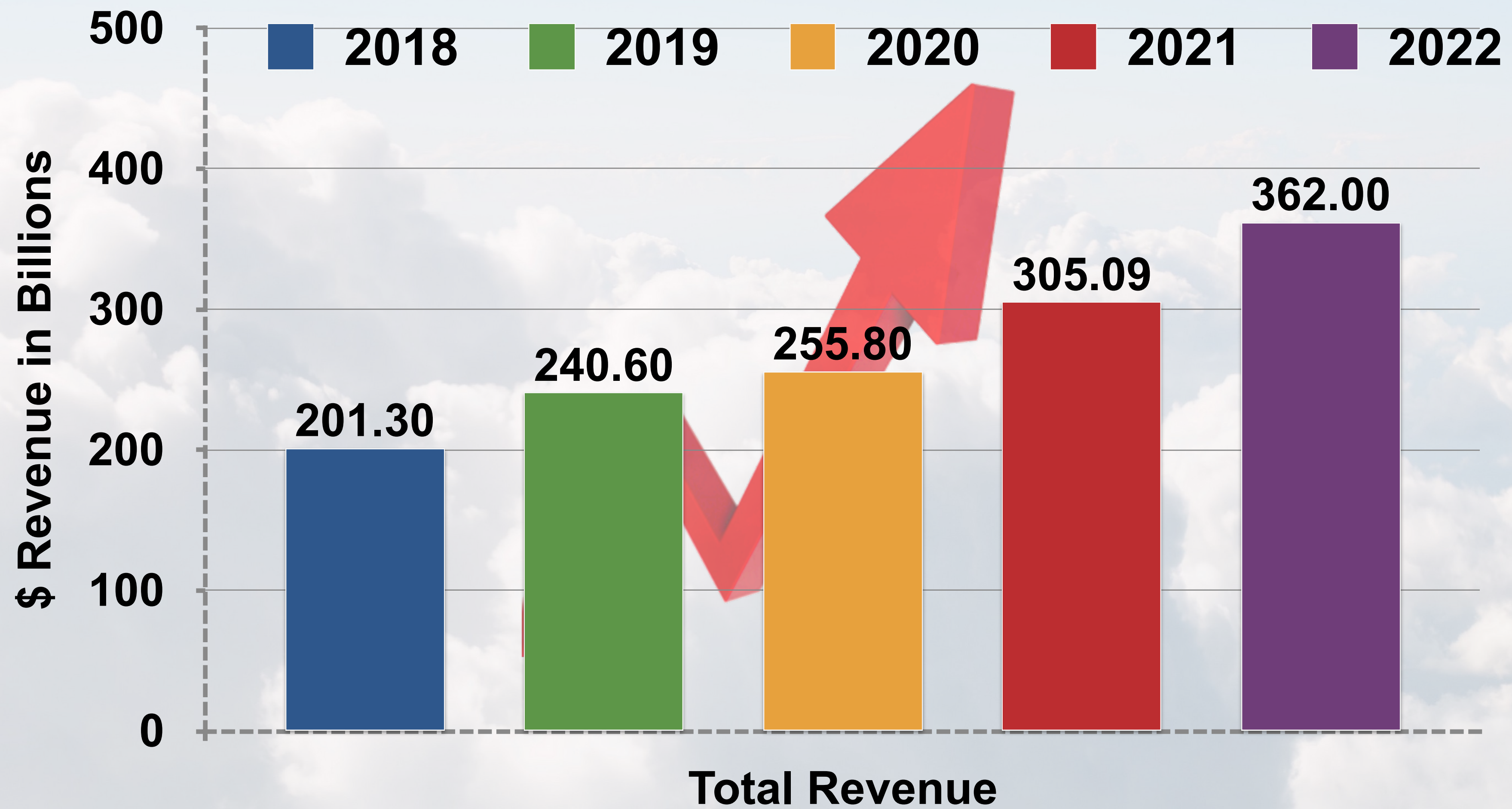
PennState
High Performance
Computing Lab

# PUSH FOR MORE CLOUD ADOPTION

**53%**

Upto 600k — 24%

600k-1.2M — 13%

1.2M-2.4M — 16%

>12M — 16%

2.4M-12M — 31%

*Source: Flexera 2020 Cloud*
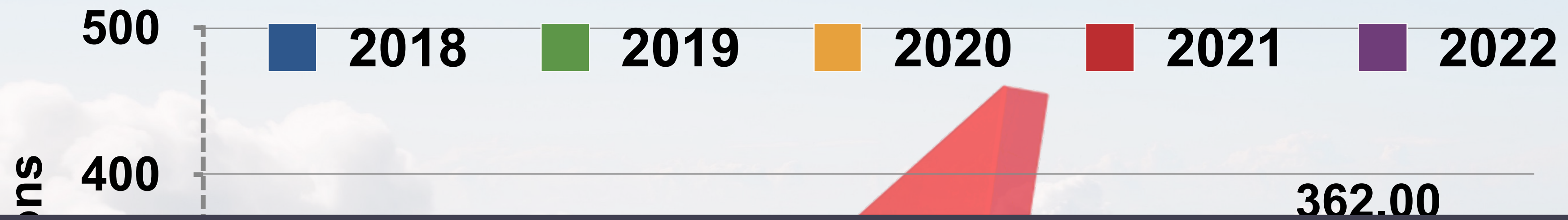
PennState College of Engineering

PennState High Performance Computing Lab

Source: https://towardsdatascience.com/coding-deep-learning-for-beginners-types-of-machine-
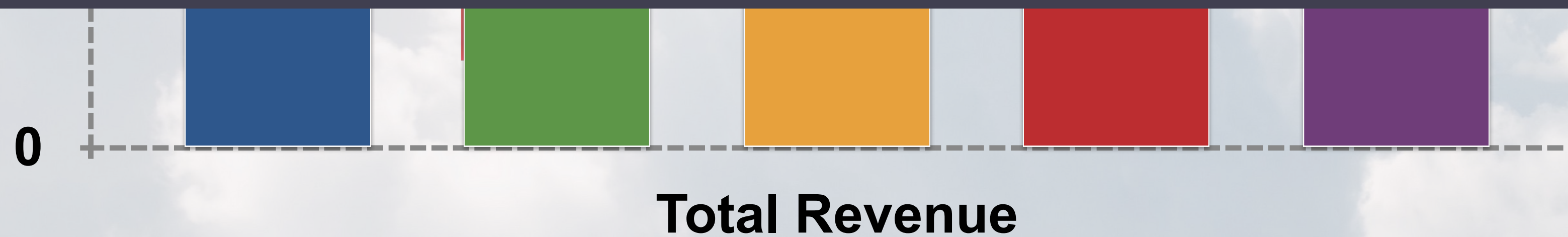
# PUBLIC CLOUD REVENUE

# TENANT-SIDE PROBLEMS

~35%        ~77%                    ~73%

Resource Selection        AutoScaling

PennState
College of Engineering

PennState
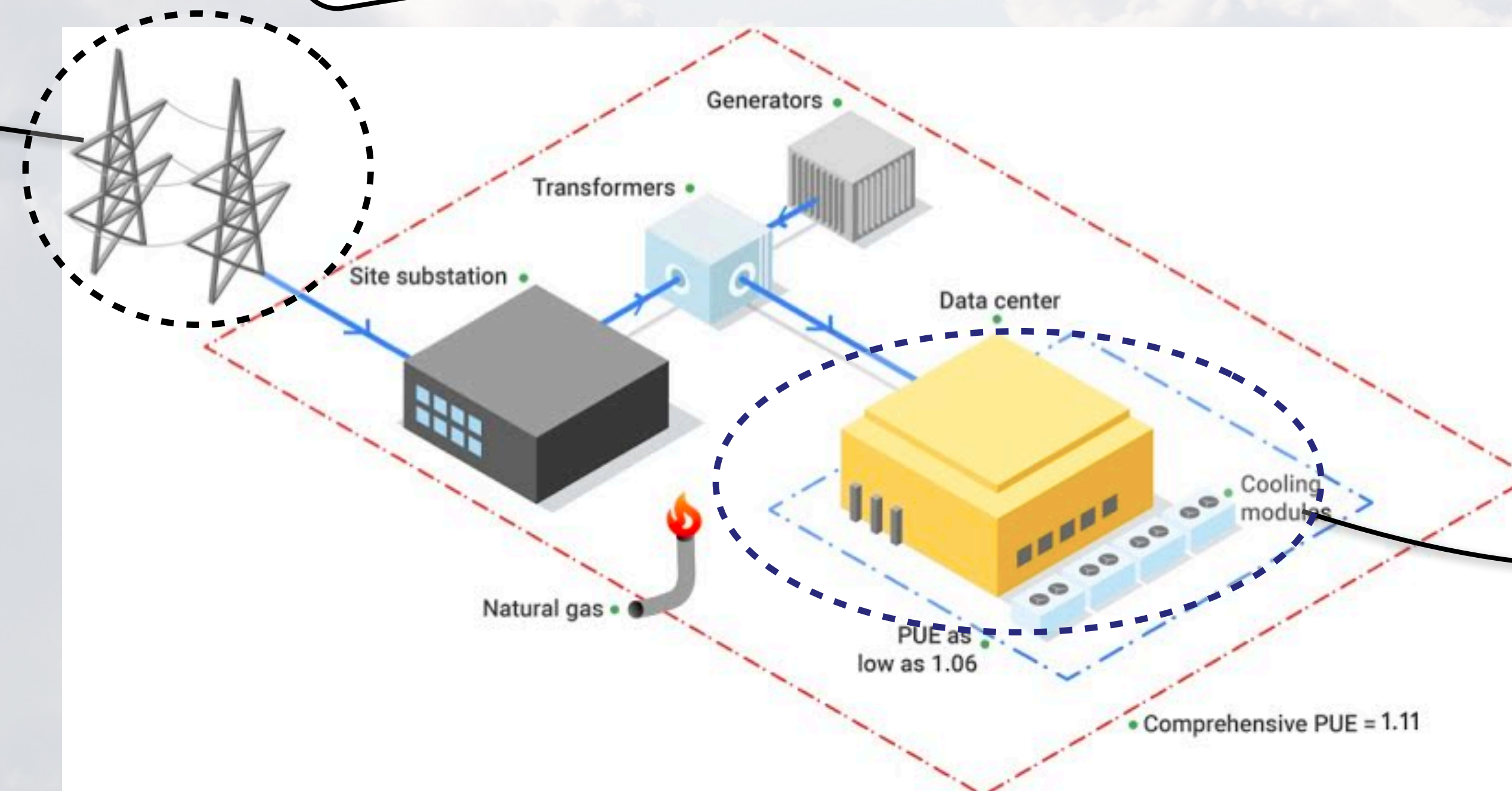High Performance
Computing Lab

# PROVIDER EXPENDITURE



**OpEx**

**CapEx > ~5x OpEx !**

**CapEx**

*Amortized cost of 50000 servers in Microsoft datacenter. Source: The Cost of Cloud, ACM SIGCOMM'09*
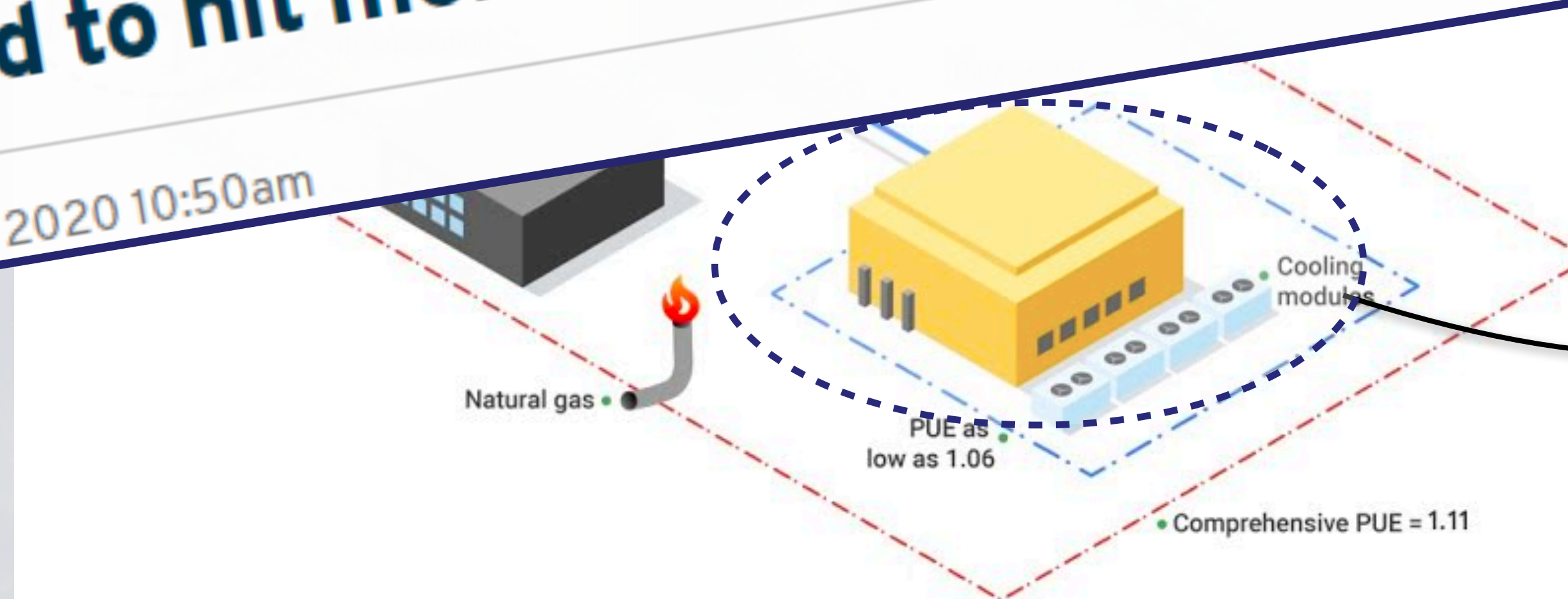
PennState
College of Engineering

PennState
High Performance
Computing Lab

7

**OpEx**

**CapEx > ~5x OpEx !**

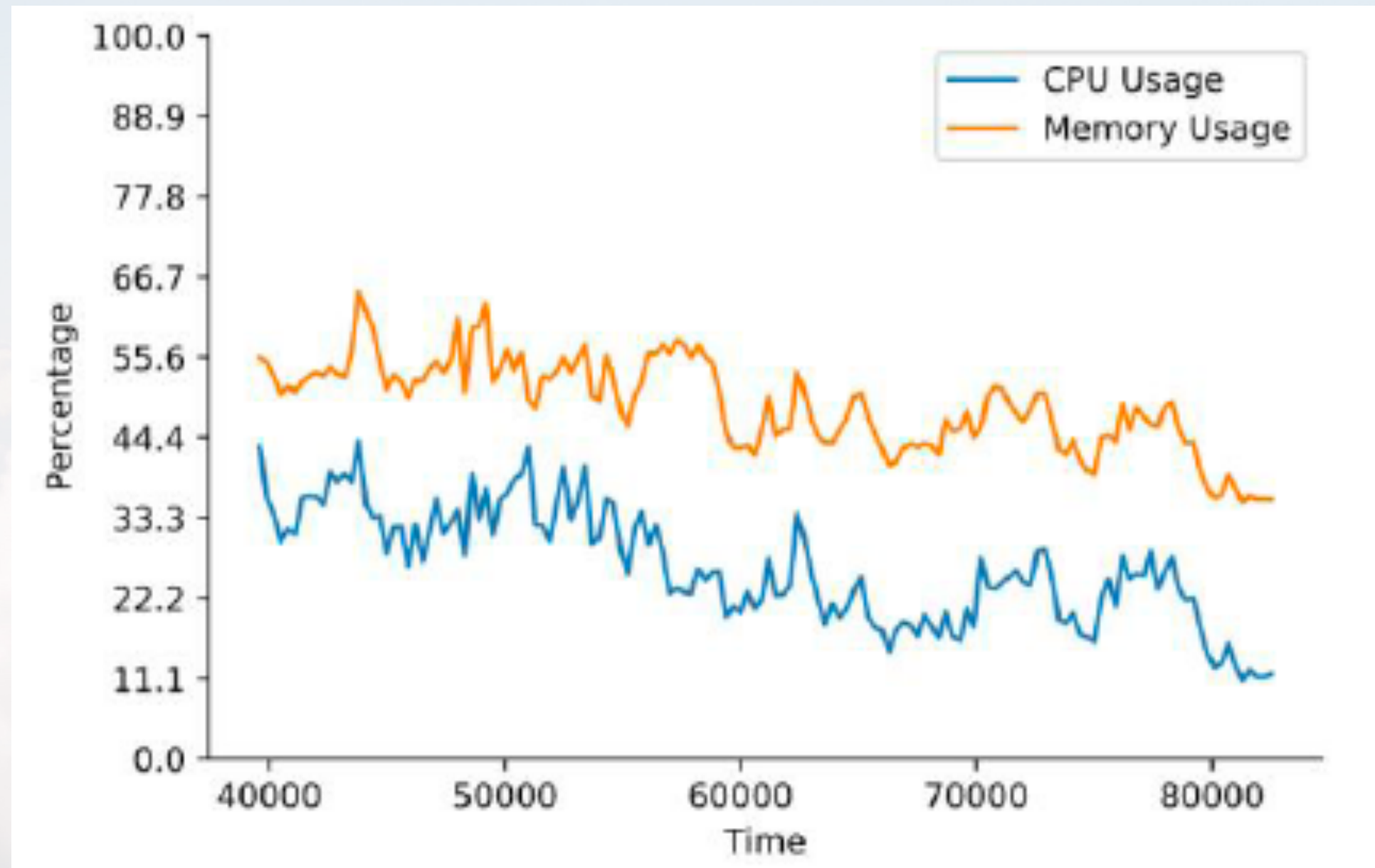Report: Despite Covid-19 disruption in 2020, data center capex poised to hit more than $200B over next five years

by Mike Robuck | Jul 24, 2020 10:50am

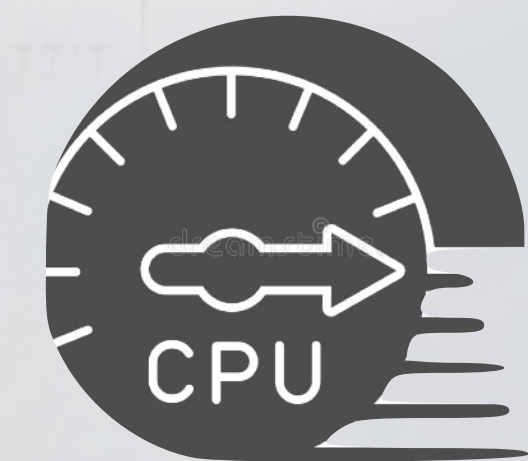Natural gas

PUE as low as 1.06

Cooling modules

Comprehensive PUE = 1.11

*Amortized cost of 50000 servers in Microsoft datacenter. Source: The Cost of Cloud, ACM SIGCOMM'09*

PennState College of Engineering

PennState High Performance Computing Lab

# PROVIDER SIDE PROBLEMS



Source: Alibaba Datacenter Case Study, IEEE Access'19

Overstated Requirements

Communication

Tenants

Providers

Blackbox Applications

Overprovisioning

**~13-40%**   **~42-65%**

PennState
College of Engineering

PennState
High Performance
Computing Lab

# SERVERLESS COMPUTING



"...Distributed Event-based programming Service..." - **OpenWhisk**

"Run code without thinking about servers.Pay for only the compute time you consume" - **AWS Lambda**

"...logic can be spun up on-demand in response to events originating from anywhere...." - **Google Cloud Functions**

**PennState**
College of Engineering

**PennState**
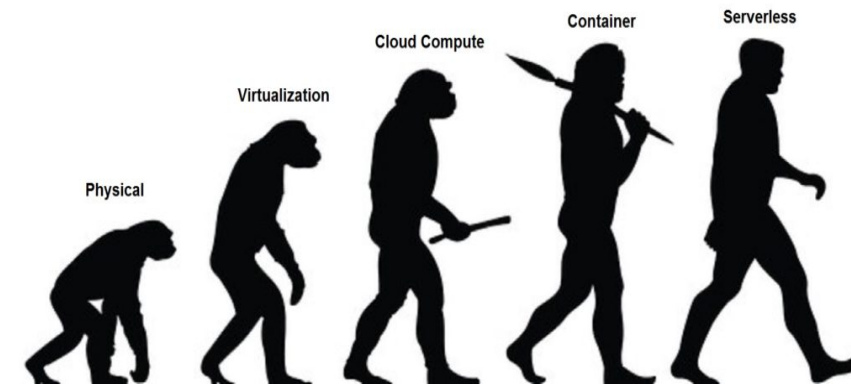High Performance
Computing Lab

# SERVERLESS COMPUTING



"…Distributed Event-based programming Service…" - **OpenWhisk**

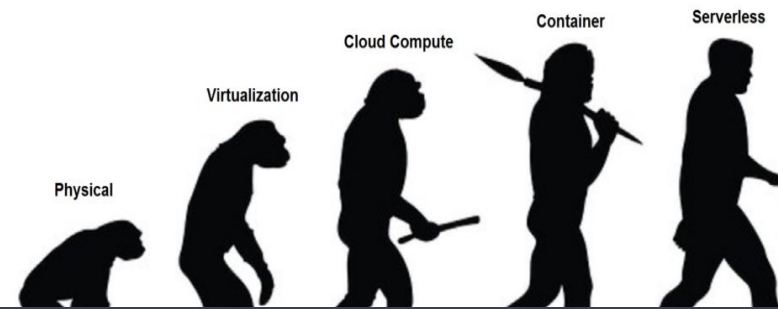"Run code without thinking about servers.Pay for only the compute time you consume" - **AWS Lambda**

"…logic can be spun up on-demand in response to events originating from anywhere…." - **Google Cloud Functions**

**Very Fast Startup**

**58%**

# WHAT WE NEED ?

How to solve?

# DISSERTATION CONTRIBUTIONS



**Applications**

**Tenant Resource Management**

Autoscaling | Application Modes
Resource Procurement

- Cost
- Latency
- Performance

Req- Resp

**Public Cloud Resources**

IaaS — EC2, Spot | PaaS | FaaS

**Provider Resource Management**

OPEX | Bin-Packing | Load Balancing
CAPEX | CPU CPU CPU GPU GPU FPGAs

- Power
- Energy
- Utilization

**Physical Servers**

*Spock*- Cost Efficient and Latency Aware Autoscaling, **IEEE CLOUD' 2019**

*Cocktail*- Improving Machine Learning Performance at Low Cost, **NSDI' 2021 (Under-Revision), WoSC'2021**

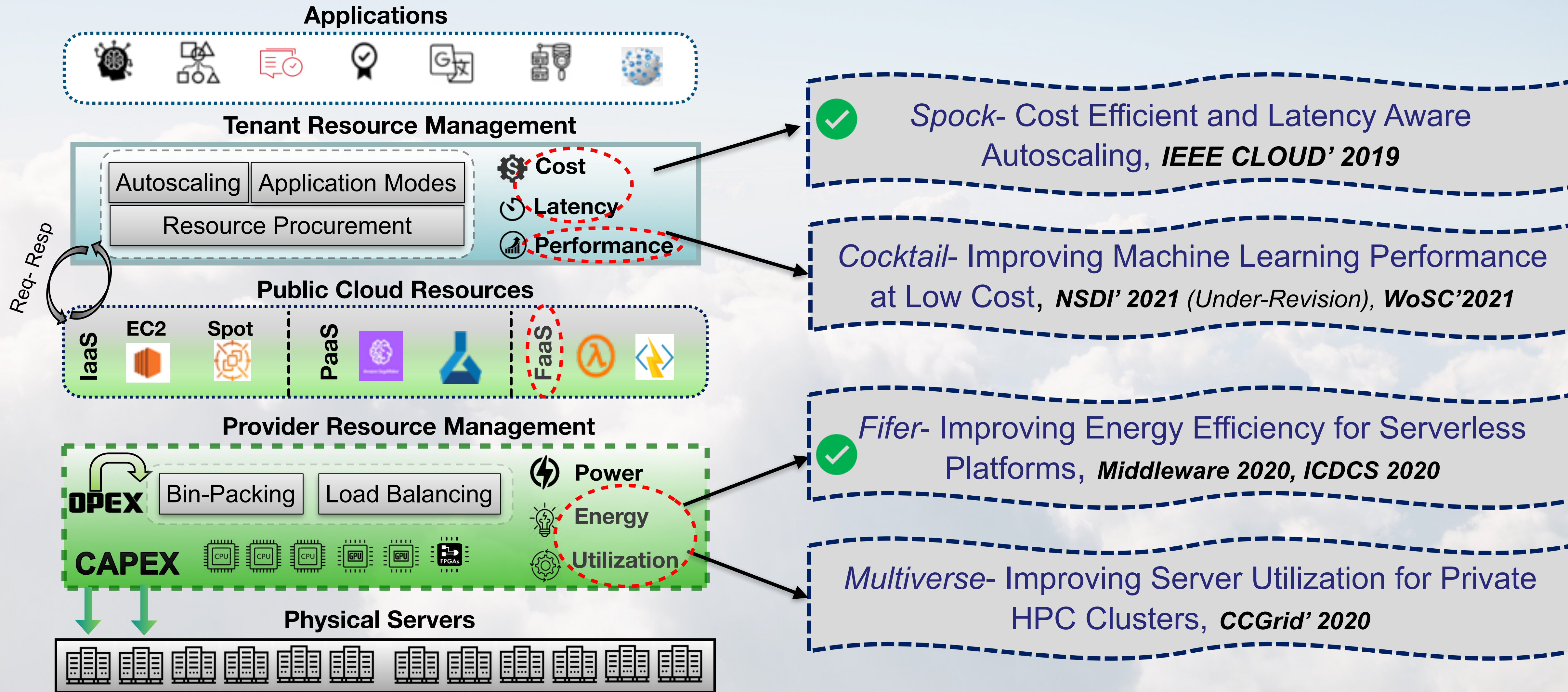*Fifer*- Improving Energy Efficiency for Serverless Platforms, **Middleware 2020, ICDCS 2020**

*Multiverse*- Improving Server Utilization for Private HPC Clusters, **CCGrid' 2020**

PennState College of Engineering

11

PennState High Performance Computing Lab

# MODEL SERVING HOSTED ON CLOUD

# PRIOR WORKS

- Utilization based autoscaling- *Urgaonkar et al PODC'03*

  ➡ Not suitable for millisecond scale applications

- Relaxed VM scale down - *Gandhi et al SC'12, TOCS'12*

  ➡ Intermittent over-provisioning

- Exploiting different VM instance types *Wang et al. Eurosys'17,*

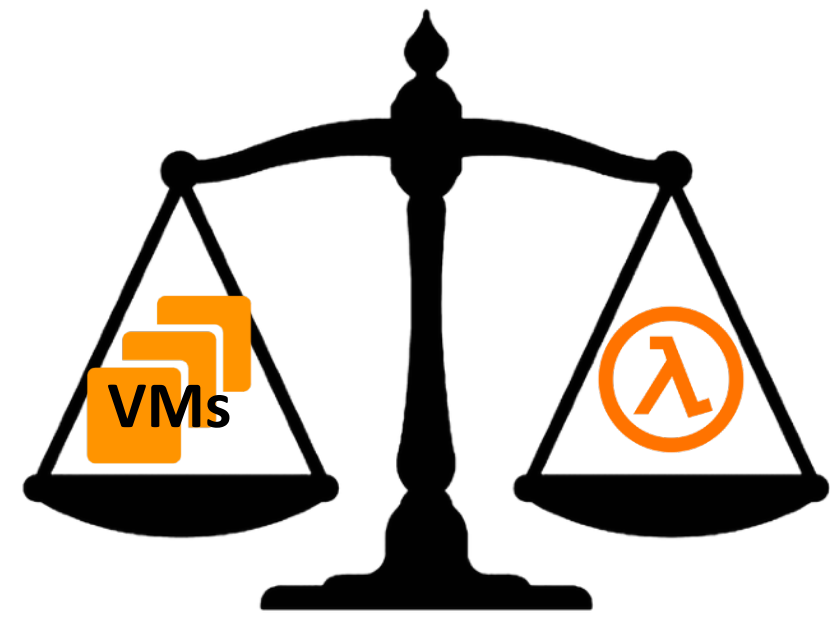  ➡ They are complementary to our proposal.

# PRIOR WORKS

- Utilization based autoscaling- *Urgaonkar et al PODC'03*
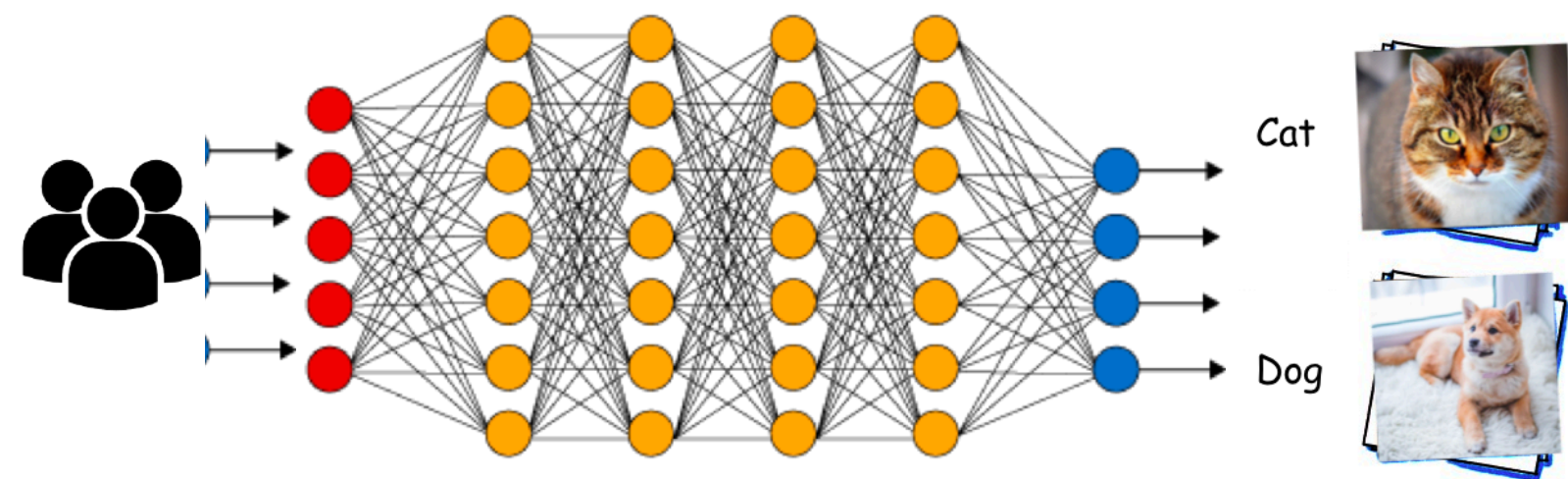
  ➡ Not suitable for millisecond scale applications
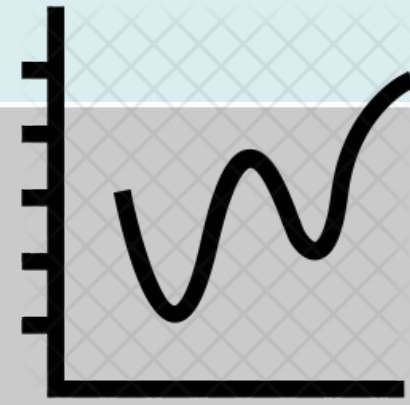
## Only VM based solutions are largely expensive

- Exploiting different VM instance types *Wang et al. Eurosys'17,*

  ➡ They are complementary to our proposal.

PennState
College of Engineering

PennState
High Performance
Computing Lab

# KEY FINDINGS



Deep Learning Inferences

| Arrival | Resource | 💰 Cost | ⏱ SLO |
|---|---|---|---|
| Bursty | λ | Pay per use 🙂 | Pre warmed 😊 |
| | VMs | Over provisioned 🙁 | Too much Scaling 🙁 |
| Predictable | λ | Per-unit Cost 🙁 | Pre warmed 😊 |
| | VMs | Known Demand 🙂 | Reduced Scaling 🙂 |

PennState College of Engineering

PennState High Performance Computing Lab

Can we multiplex both?

# SPOCK: EXPLOITING SERVERLESS FUNCTIONS FOR SLO AND COST AWARE AUTOSCALING



➤ Offload queries to lambdas while starting new VMs.

➤ Reduces SLO violations during request surge.

➤ Reduce intermittent over-provisioning VMs

# SPOCK: EXPLOITING SERVERLESS FUNCTIONS FOR SLO AND COST AWARE AUTOSCALING

➤ Offload queries to lambdas while starting new VMs.

➤ Reduces SLO violations during request surge.

➤ Reduce intermittent over-provisioning VMs

# SPOCK: EXPLOITING SERVERLESS FUNCTIONS FOR SLO AND COST AWARE AUTOSCALING

- - - VM   - - - SLA   → Lambda   —— Arrival rate

Scale Up

Scale Down

➤ Offload queries to lambdas while starting new VMs.

Spock reduces SLO violations by **~74%** with **~33%** cost savings

Time (Sec)

➤ Reduce intermittent over-provisioning VMs

PennState College of Engineering

15

PennState High Performance Computing Lab

**Provisioning Latency**
**Resource** ✅ **(Spock)**

How to improve accuracy with low latency and low cost?

Metrics

Cost   Latency   Accuracy

- **InFaas** uses different resource types to ensure low latency at low cost.

-  **Clipper** uses model **ensembling** to achieve higher accuracy.

Crankshaw et al  CIDR'15, NSDI'17, SoCC'20
Yadawkar et al Arxiv'19

**Metrics**

**Cost**  **Latency**  **Accuracy**

**Large Model**  **Small Models**



- **InFaas** uses different resource types to ensure low latency at low cost.

-  **Clipper** uses model **ensembling** to achieve higher accuracy.

Crankshaw et al  CIDR'15, NSDI'17, SoCC'20
Yadawkar et al Arxiv'19

Cost

0.75
0.5
0.25

Accuracy     Latency

······ InFaas   ······ Clipper   —— Cocktail

PennState
College of Engineering

PennState
High Performance
Computing Lab

**Metrics**

**Cost**

**Latency**

**Accuracy**

**Large Model**

**Small Models**

# How to do ensembling?

– **Clipper** uses model **ensembling** to achieve higher accuracy.

Accuracy

*Latency*

......... InFaas  ......... Clipper  —— Cocktail

Crankshaw et al CIDR'15, NSDI'17, SoCC'20
Yadawkar et al Arxiv'19

PennState
College of Engineering

PennState
High Performance
Computing Lab

# Model Ensembling Framework

Majority Voting

**Model Selection**
- ■ NasNetMobile
- ■ MobileNetV2
- ■ InceptionV3

Requests

Host Server

Aggregator

Model-1  Model-2  Model-N

Model-1  Model-2  Model-N

**Cloud Resources for Individual Models (Virtual Machines)**

Model Ensembling Framework

Model Selection

NasNetMobile
MobileNetV2

Requests

Host Server

Aggregator

Majority Voting

High Resource Footprint
What about Model Selection?

Model-1    Model-2    . . .    Model-N

Cloud Resources for Individual Models (Virtual Machines)

18

**Most accurate model**
- ✳ **~2x** parameters, latency
- ✳ **~2%** more accuracy

- How to bridge the 2% accuracy gap?
- What about cost?

IEEE Access'18 Benchmark Analysis of Representative Deep Neural Network Architectures

**Most accurate model**

✳ **~2x** parameters, latency

✳ **~2%** more accuracy

- What about cost?

## How to ensemble?

IEEE Access'18 Benchmark Analysis of Representative Deep Neural Network Architectures

# FULL ENSEMBLE

Model Set: Top 12 frequently used models from Keras Tensorflow

Choose baseline models in decreasing order of accuracy

Combine all models which are under the latency of baseline model.

# FULL ENSEMBLE



Latency Comparison

Accuracy Comparison

# FULL ENSEMBLE

Latency Comparison

■ Single    ■ Ensemble

Accuracy Comparison

■ Single    ■ Ensemble

## What about Cost?

Lat 100

0

**NasNetLarge** **IncepResV2**    **Xception**

Top 79.5

78

**NasNetLarge** **IncepResV2**    **Xception**

PennState
College of Engineering

PennState
High Performance
Computing Lab

# FULL ENSEMBLING COST



Ensembling is up-to **2x** expensive.

Spot instances can potentially reduce cost.

PennState
College of Engineering

PennState
High Performance
Computing Lab

# What can we do?

| Baseline(BL) | NASLarge | IRV2 | Xception | DNet121 | NASMob |
|---|---|---|---|---|---|
| #Models | 10 | 8 | 7 | 5 | 2 |

✦ Do we need so many models?

✦ How to autoscale resources for each model?

✦ How to handle instance failures?

PennState
College of Engineering

PennState
High Performance
Computing Lab

**Compared to Full-Ensemble (N models)**



**Most accurate N/2 models**

**Accuracy** 🎯

# STATIC ENSEMBLING

**Compared to Full-Ensemble (N models)**



How to dynamically select the models?

mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

# DYNAMIC MODEL SELECTION

## Leverage Class-wise Accuracy

Mobilenet (MNet) ⟹ Slug ☺

Mobilenet (MNet) ⟹ Quill ☹

# COCKTAIL- MULTIDIMENSIONAL OPTIMIZATION FOR ENSEMBLE LEARNING IN CLOUD

User Requests

Dynamic Model Selection

w1 w2 w3 w4 . . . wk

Aggregator

Weight Matrix

L

N

output

Queries

Model-1  Model-2  Model-3  Model-4  Model-n
         GPU       CPU       CPU       GPU

Prediction Policy

Importance Sampling

Autoscaler

**Class-wise dictionary**

**Weighted Selection**

**Dedicated Pools**

**Per model Scaling**

**Fault tolerant**

PennState
College of Engineering

PennState
High Performance
Computing Lab

25

## Experiment Setup

- *40 EC2 CPU/GPU VMs*
- *Wiki Twitter Traces*

| Dataset | Application | Classes | Train-set | Test-set |
|---------|-------------|---------|-----------|----------|
| ImageNet [56] | Image | 1000 | 1.2M | 50K |
| CIFAR-100 [116] | Image | 100 | 50K | 10K |
| SST-2 [117] | Text | 2 | 9.6K | 1.8K |
| SemEval [118] | Text | 3 | 50.3K | 12.2K |

# MAJOR RESULTS



Cocktail incurs ~**32%** lower cost

Cocktail reduces #models by ~**50%** on average

Cocktail yields ~**2x** lower latency

Cocktail gains upto ~**1.25%** more accuracy

# DISSERTATION CONTRIBUTIONS



**Applications**

**Tenant Resource Management**

Autoscaling | Application Modes
Resource Procurement

Cost
Latency
Performance

*Spock*- Cost Efficient and Latency Aware Autoscaling, *IEEE CLOUD' 2019*

*Cocktail*- Improving Machine Learning Performance at Low Cost, *NSDI' 2021 (Under-Revision)*

**Public Cloud Resources**

IaaS | EC2 | Spot | PaaS | FaaS

Req- Resp

**Provider Resource Management**

OPEX | Bin-Packing | Load Balancing
CAPEX | CPU CPU CPU GPU GPU FPGAs

Power
Energy
Utilization

*Fifer*- Improving Energy Efficiency for Serverless Platforms, *Middleware, ICDCS 2020*

*Multiverse*- Improving Server Utilization for Private HPC Clusters, *CCGrid' 2020*

**Physical Servers**

PennState College of Engineering

PennState High Performance Computing Lab

# RECAP



**58% use Serverless to reduce cost and accelerate development.**

# Provider Challenges?

**58%** use Serverless to reduce cost and accelerate development.

# SERVERLESS FUNCTION CHAINS



Stages in Intelligent Personal Assistant (IPA)

Request

1500ms

Users

Response

Image Classification

Natural Language Processing

Question Answering

Container Cold Starts

Containers for Each Microservice

PennState
College of Engineering

30

PennState
High Performance
Computing Lab

# SERVERLESS FUNCTION CHAINS

**Stages in Intelligent Personal Assistant (IPA)**

Request

Cold-starts contribute ~2000 to 7500 ms overheads to overall latency

Cold Starts

**Containers for Each Microservice**

30

PennState
College of Engineering

PennState
High Performance
Computing Lab

# CURRENT SERVERLESS PLATFORMS

- Spawn new containers if existing containers are busy.
  - ➡ Leads to SLO violations due to cold-starts.
  - ➡ Many idle containers. Wasted power and energy.

- Employing static queuing of requests on fixed pool of containers
  - ➡ Leads to SLO violations due to queuing.

- Not aware of application execution times and response latency requirements.
  - ➡ Colossal container overprovisioning.

*Wang et al, Peeking behind the curtains of Serverless Platforms in ATC'18*

*Shahrad et al, Serverless in the Wild, in ATC'21*

# CURRENT SERVERLESS PLATFORMS

- Spawn new containers if existing containers are busy.
  - ➡ Leads to SLO violations due to cold-starts.
  - ➡ Many idle containers. Wasted power and energy.

# How can we do better?

- Not aware of application execution times and response latency requirements.
  - ➡ Colossal container overprovisioning.

*Wang et al, Peeking behind the curtains of Serverless Platforms in ATC'18*

*Shahrad et al, Serverless in the Wild, in ATC'21*

Slack = Response Latency ⊖ Execution Time (ET)

Multi-staged applications have ample slack
(200-700ms)

Execution times of each function is predictable-
(20-100ms)

Slack > ~7x ET !

Slack Aware
Provisioning

PennState
College of Engineering

PennState
High Performance
Computing Lab

**Slack-aware batching and queuing**

**Queuing delay-based Reactive container scaling**
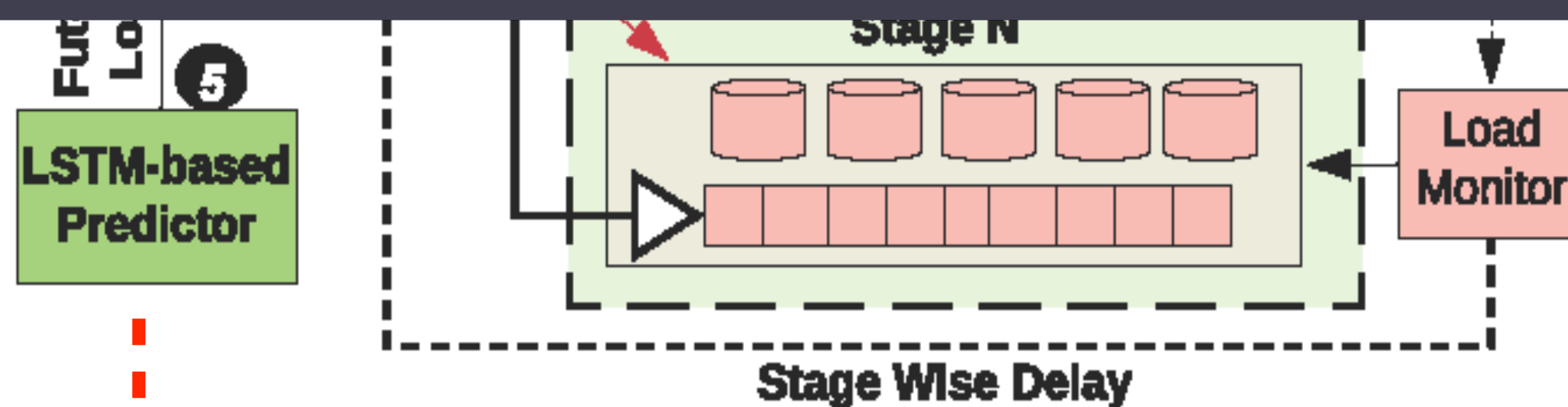
**Proactive container scaling**

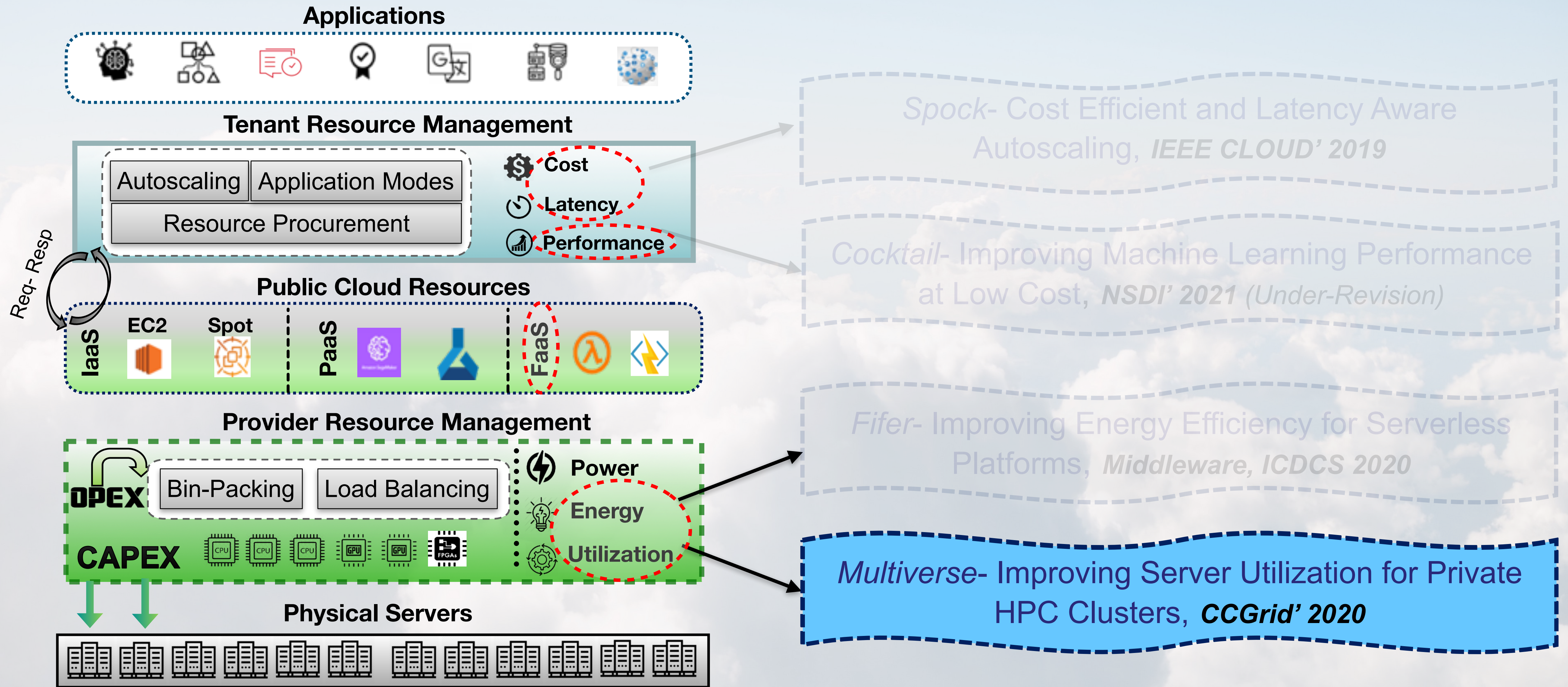# FIFER: STAGE-AWARE PROACTIVE CONTAINER PROVISIONING AND MANAGEMENT



**Slack-aware batching**

## Fifer spawns ~60% less containers.
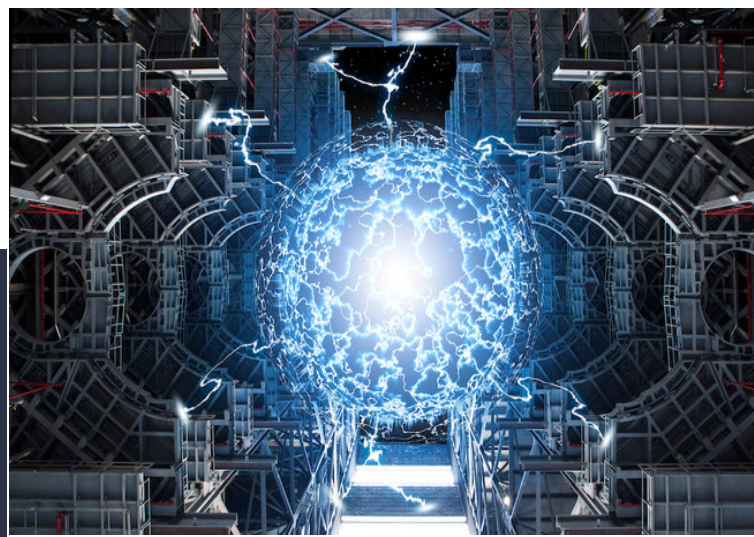## Fifer is ~31% more energy efficient.

**Proactive container scaling**

PennState
College of Engineering

PennState
High Performance
Computing Lab

# DISSERTATION CONTRIBUTIONS

# HIGH PERFORMANCE COMPUTING



US government awards millions to HPE, Intel, and others in hopes they'll build next-gen

**THE VERGE**

*Department of Energy*

Secretary of Energy Rick Perry Announces $1.8 Billion Initiative for New Supercomputers

The Worldwide HPC Server Market: $6.7 Billion in First Half 2019

Hyperion: AI-driven HPC Industry Continues to Push Growth Projections

By Doug Black

High-Performance Computing as a Service Market is Expected to Reach $17.00 Billion by 2026, Says Allied Market Research

PennState
College of Engineering

PennState
High Performance
Computing Lab

# VIRTUALIZED HPC



**Heterogeneous Compute**

**Flexibility**

**Isolation and Security**

https://blogs.vmware.com/apps/2018/09/vhpc-ra-part1.html

# CHALLENGES WITH HPC

**HPC Schedulers**



- Focus on throughput and utilization.

- Batch Jobs are usually long running.

- Fair sharing and fixed node reservations.

**HPC Schedulers**

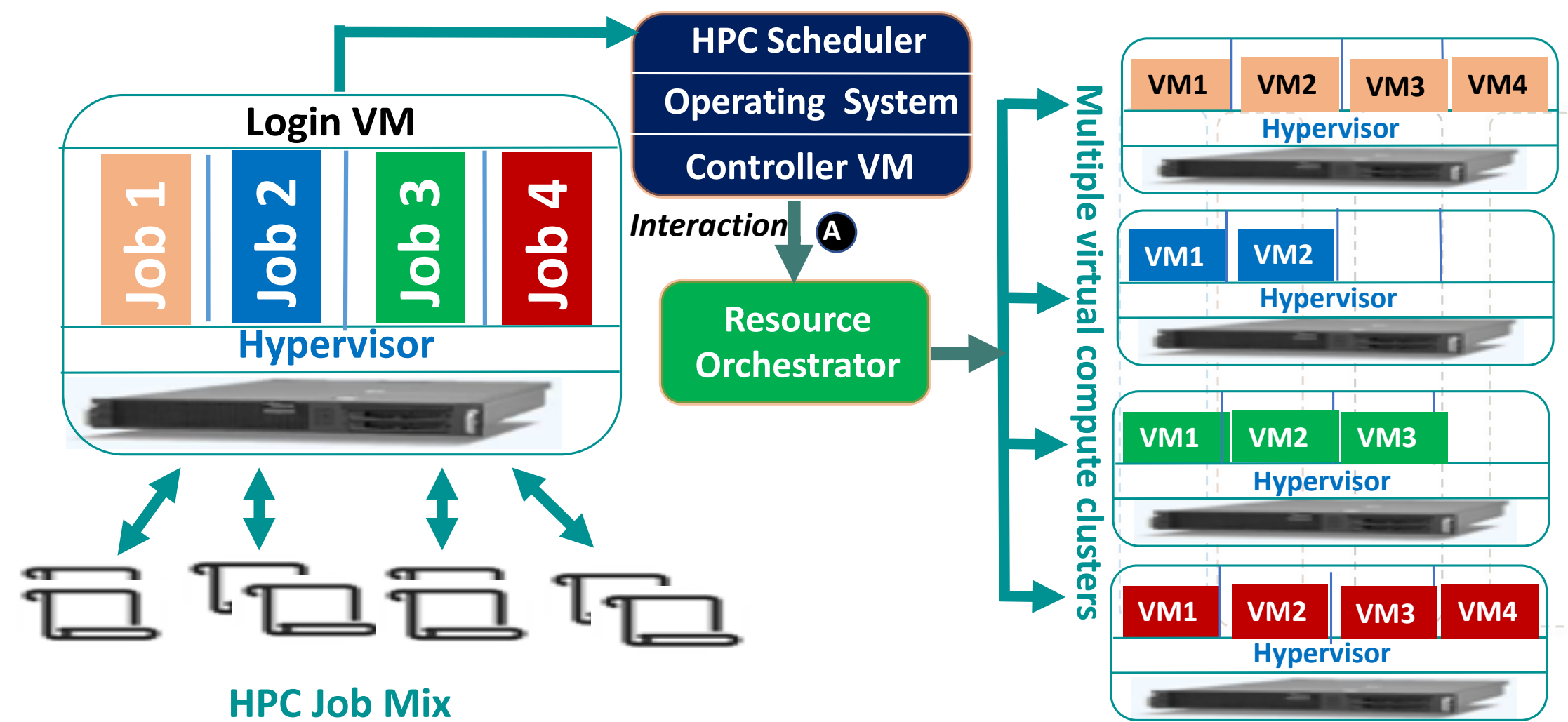- Focus on throughput and utilization.

No interaction with VM orchestrators
Results in Underutilization

reservations.

High Throughput Computing

PennState
College of Engineering

PennState
High Performance
Computing Lab

- Static Provisioning

HPC Scheduler
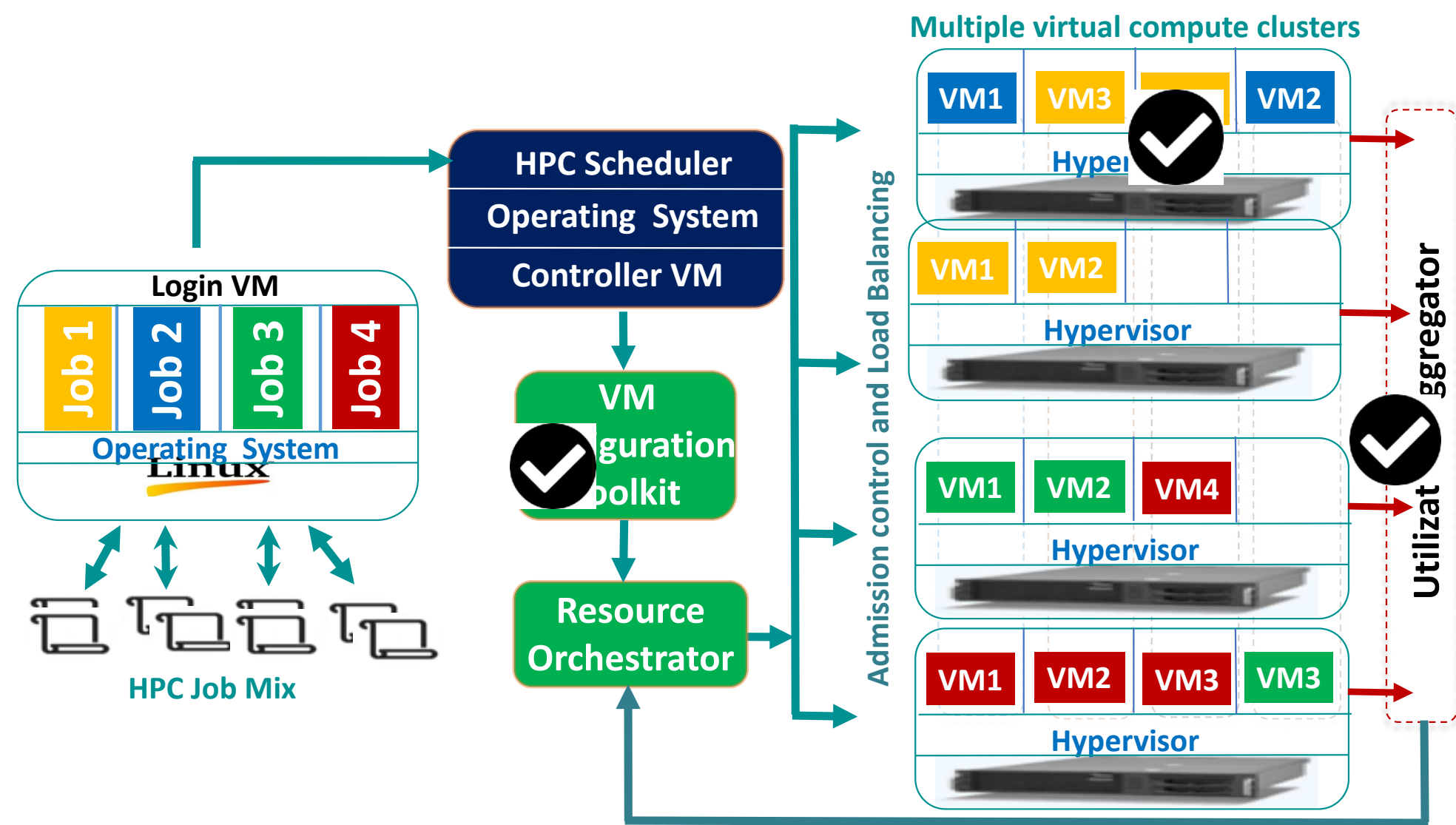
Operating System

| VM1 | VM2 | VM3 | VM4 |

# How to solve this problem?

cluster resources

# MULTIVERSE- DYNAMIC VM PROVISIONING FOR HIGH PERFORMANCE COMPUTING CLUSTERS



**Seamless interaction with integration**

**Dynamic VM Provisioning**

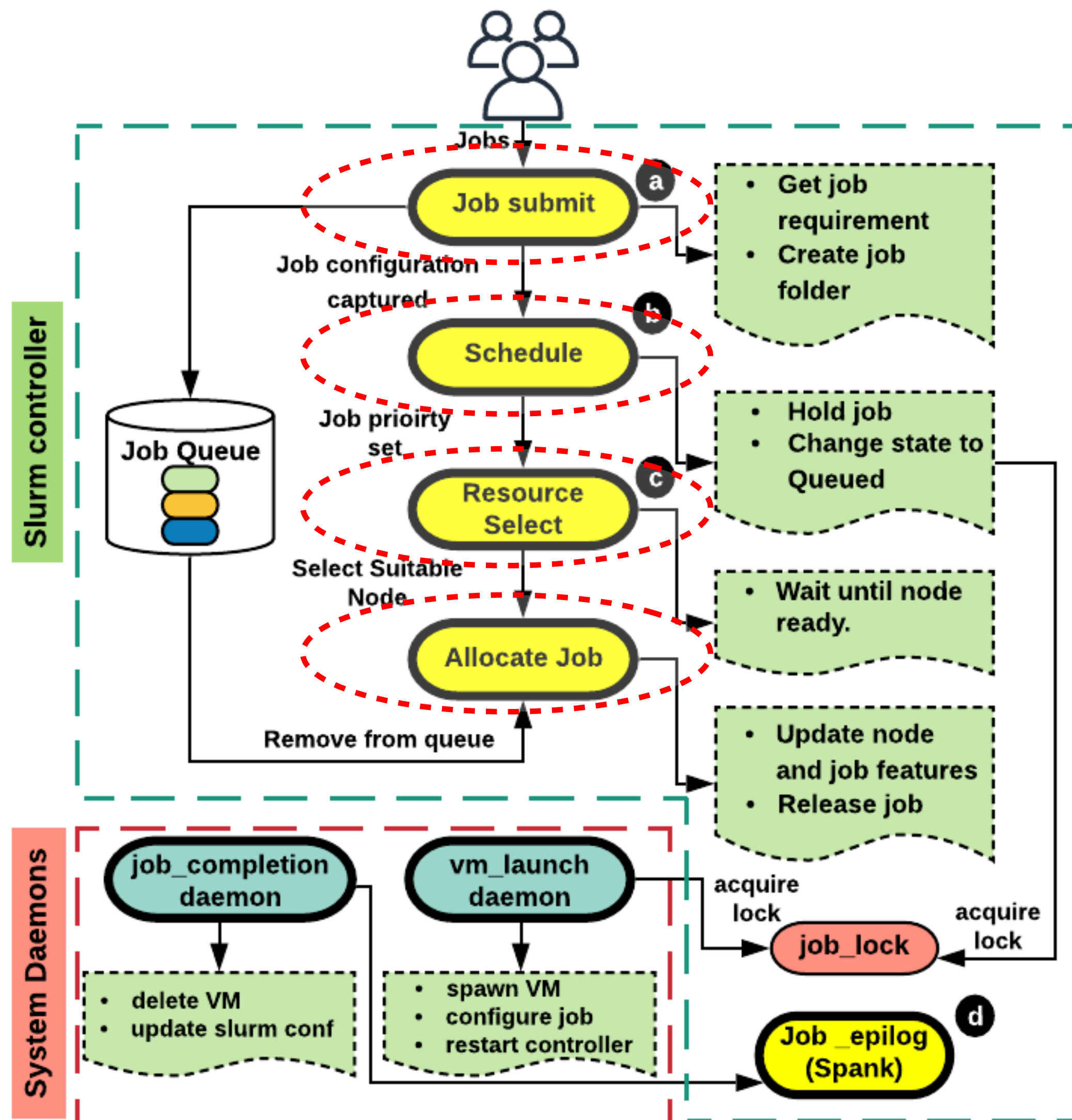**Leverage Instant Clone**

**Expose Real-time Cluster Statistics**

# MULTIVERSE DESIGN

- Parse Job Requirements
- Customized VM launch
- Map Jobs to VMs (concurrency)

- Need to be thread-safe
- Schedulers are multi-threaded and are thread-safe.

**We built a thread safe finite-state machine using linux flock utility.**

PennState
College of Engineering

PennState
High Performance
Computing Lab

**Each phase corresponds to a plugin**

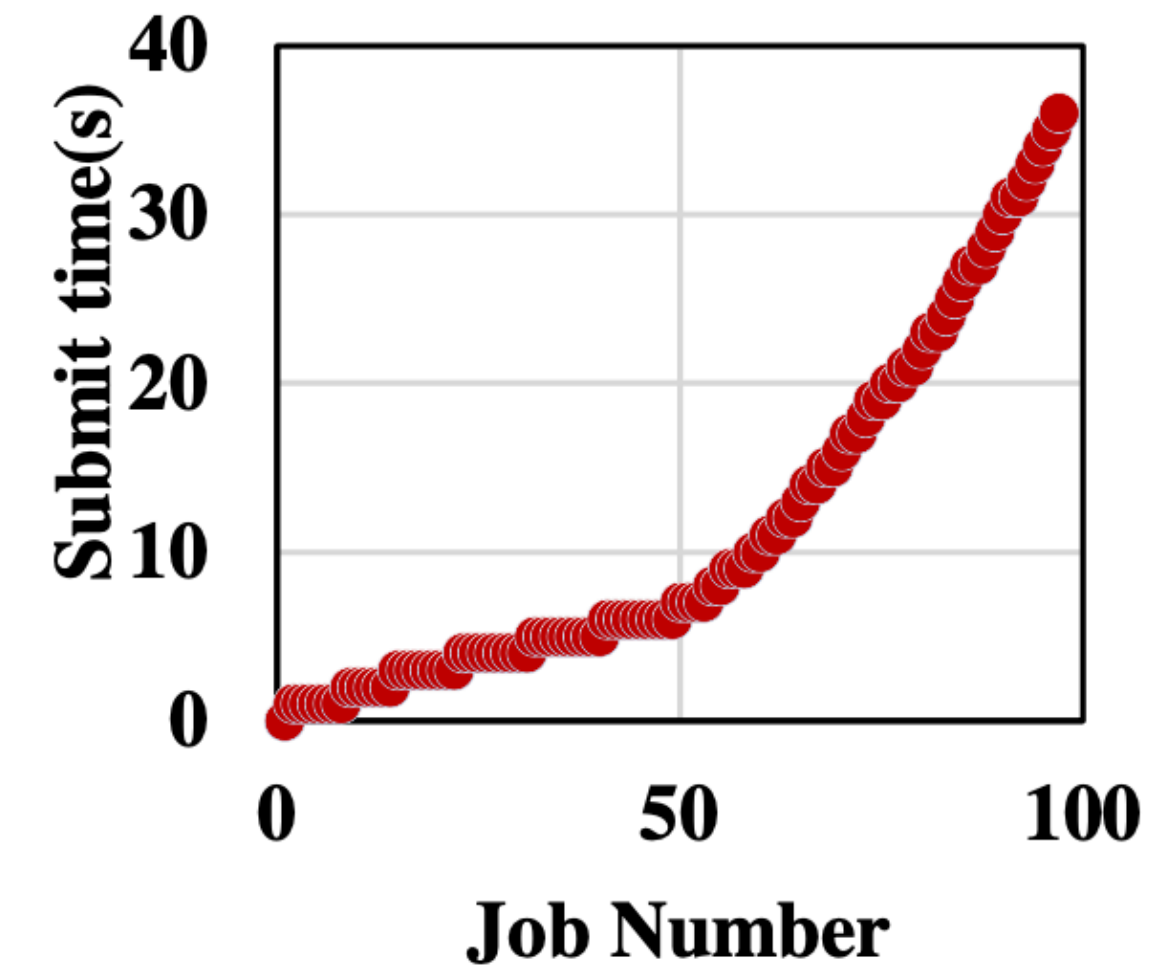**System Daemons ensure concurrency**

**Spank Plugins for VM Cleanup**

# EVALUATION SETUP



## Experiment Setup

- *220 core* HPC cluster.
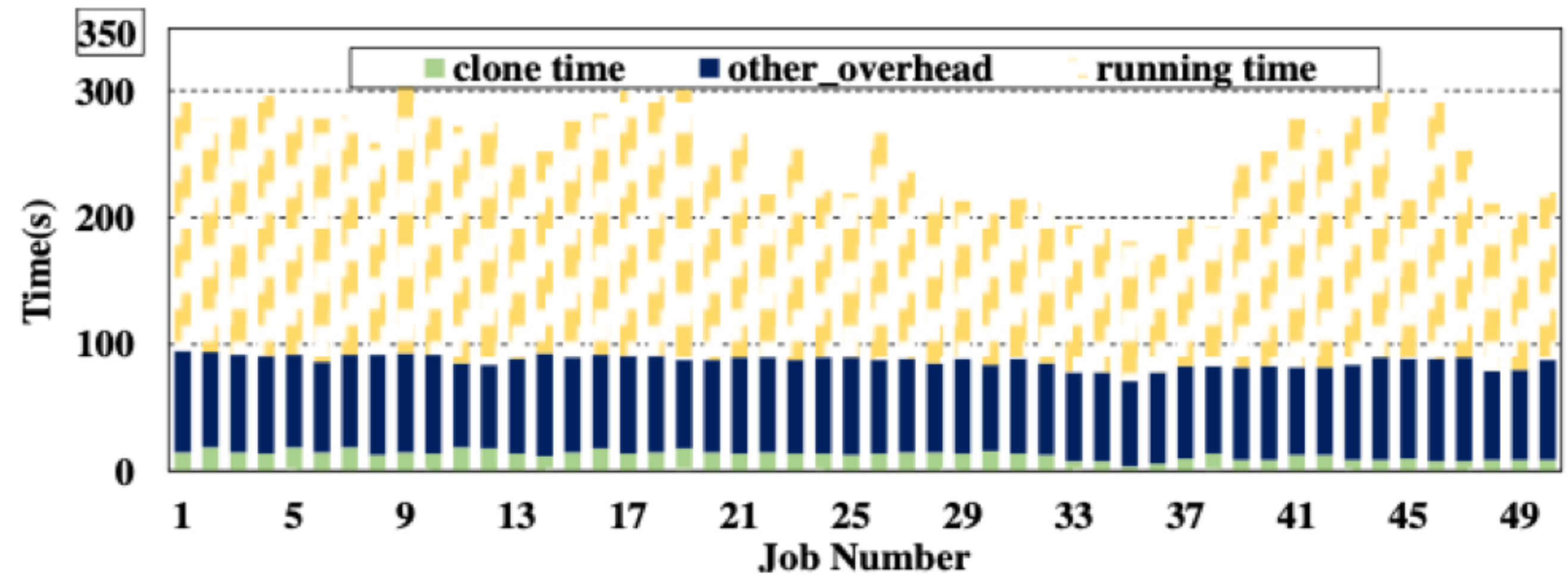
- *1TB* Memory

- *72TB* shared datastore

## Workload
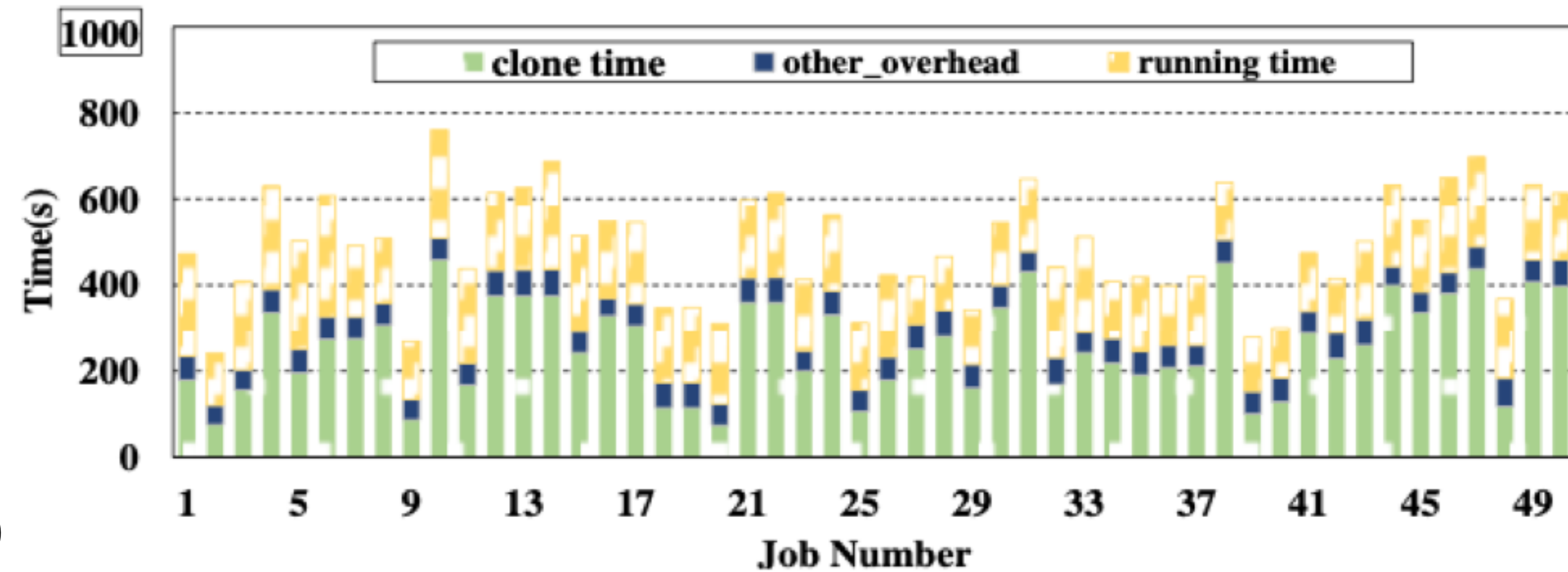
- HPCC, HPL, RandomAccess.
- Small (2vCPU, 4GB), Large (8vCPU, 16GB)
- 50 job/s, 100jobs/s

PennState
College of Engineering
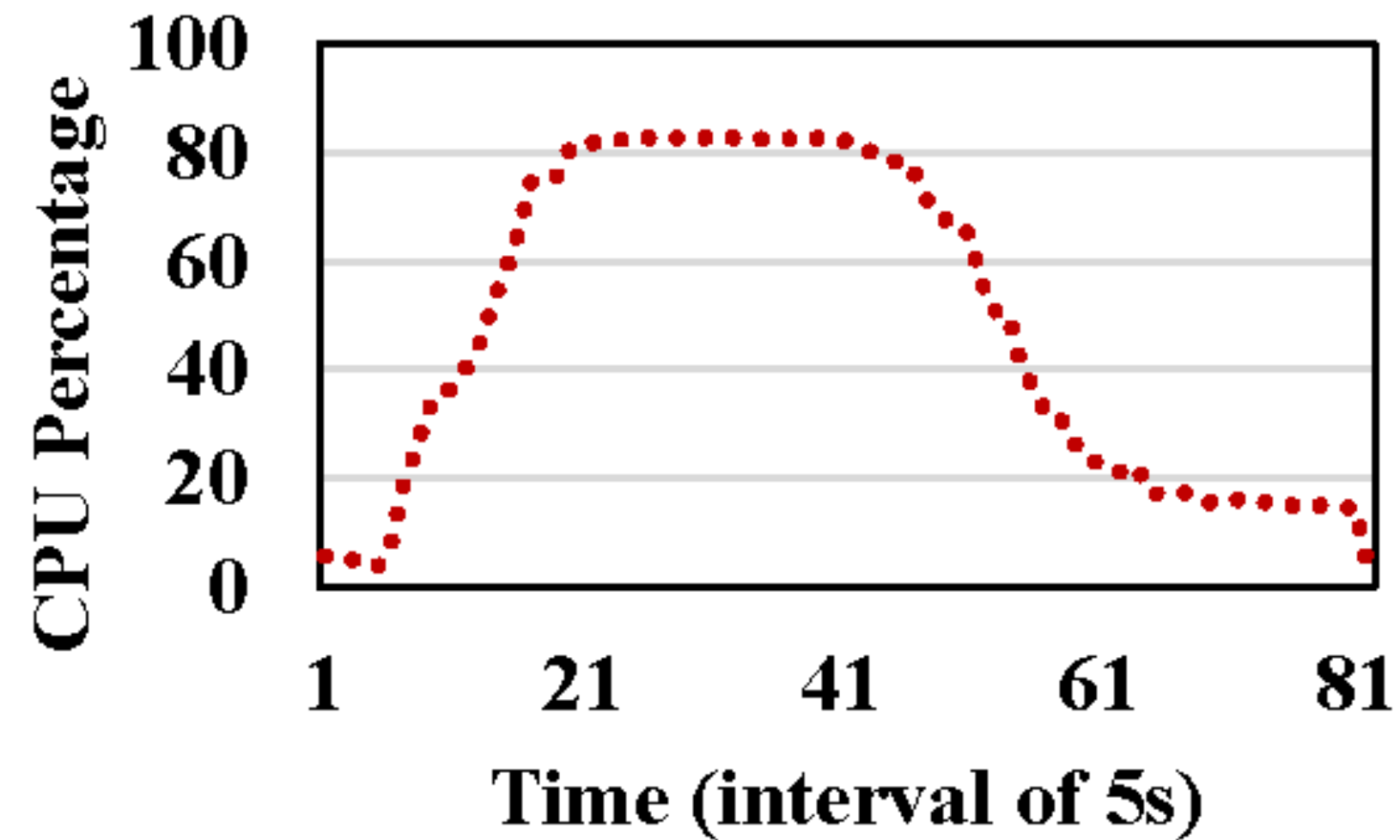
42

PennState
High Performance
Computing Lab

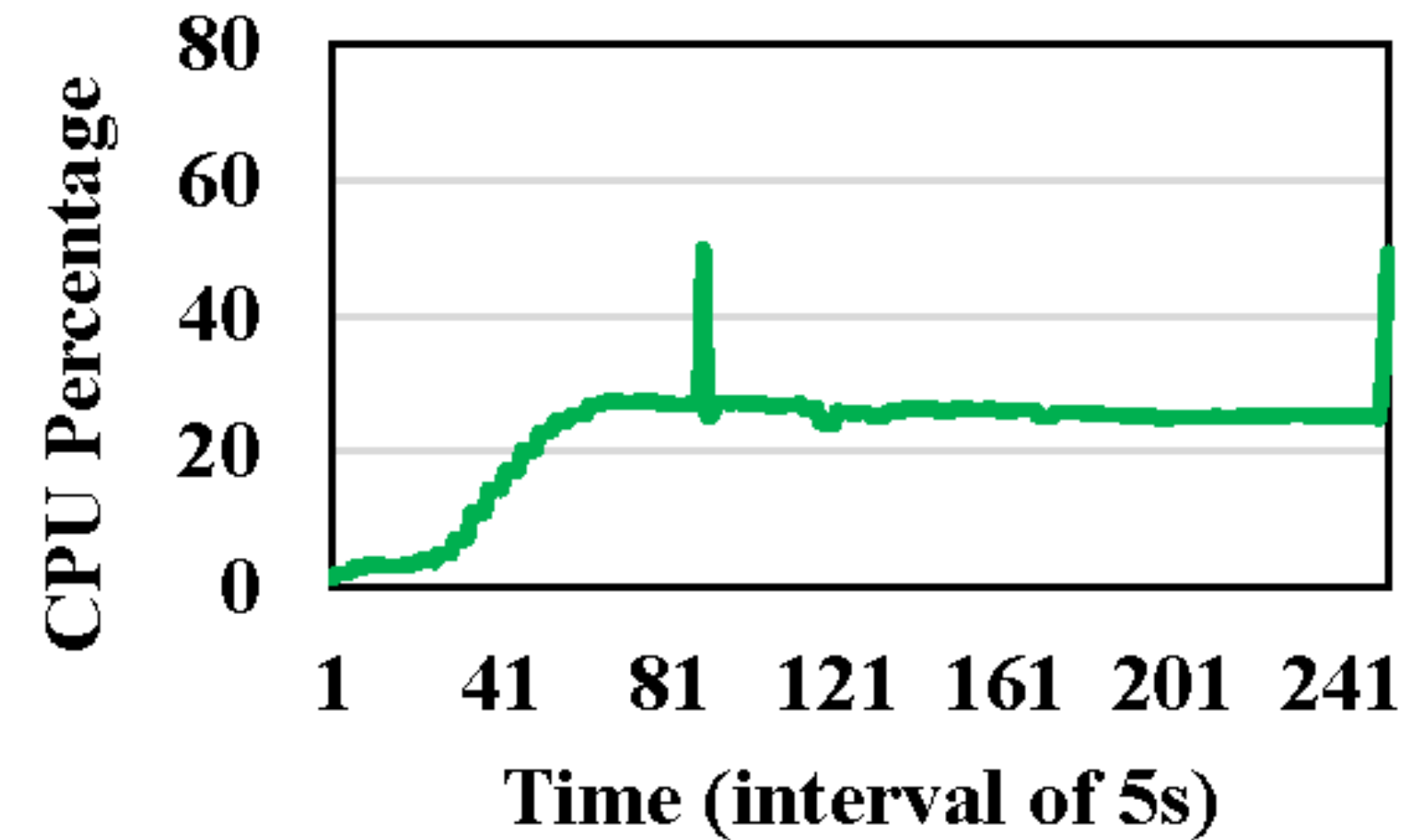**Full Clone**

**~3x Fast!**

**Instant Clone**

# MAJOR RESULTS

**Instant Clone**

**Full Clone**

~**1.5x** more throughput.

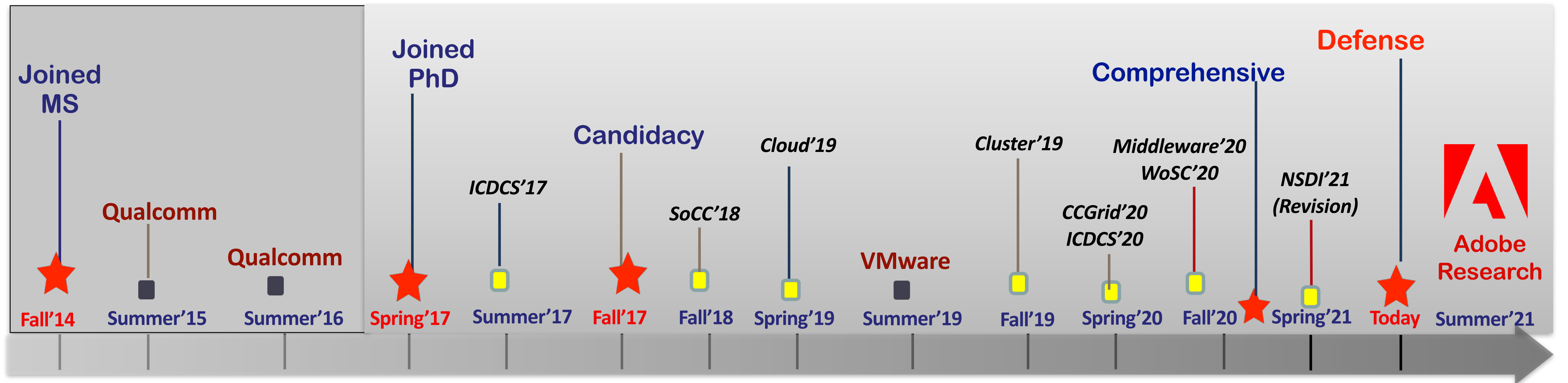~**40%** higher CPU utilization.

# FUTURE RESEARCH DIRECTIONS

## SHORT TERM

- Dynamic DAGs in Serverless

- Stateful Serverless Storage Costs

- Machine Learning Training Costs

## LONG TERM

- Federated learning in Public Cloud

- Online Real-time training using serverless

- HPC in public cloud

PennState
College of Engineering

PennState
High Performance
Computing Lab

# Doctoral Committee



**Dr. Mahmut Kandemir**
**Dissertation Co-Advisor**
Professor
Department of CSE



**Dr. Chita Das**
**Dissertation Co-Advisor**
Distinguished Professor
**Head** of CSE Department



**Dr. George Kesidis**
**Committee Member**
Professor
Department of CSE



**Dr. Bhuvan Urgaonkar**
**Committee Member**
Associate Professor
Department of CSE



**Dr. Anton Nekrutenko**
**Committee Member**
Professor
Department of BME



**Dr. Bikash Sharma**
**Special Member**
Infrastructure Engineer
Facebook

PennState
College of Engineering

PennState
High Performance
Computing Lab

# ACKNOWLEDGEMENTS



**Nachiappan**

**Prashanth**

**Prasanna**

**Adhi (My Wife)**

**Ria (My Kid)**

# ACKNOWLEDGEMENTS



All other fellow lab mates

# Thank You