

# Adaptive Resource Views for Containers

**Hang Huang**<sup>1</sup>, Jia Rao<sup>2</sup>, Song Wu<sup>1</sup>, Hai Jin<sup>1</sup>, Kun Suo<sup>2</sup>, Xiaofeng Wu<sup>2</sup>

<sup>1</sup> Huazhong University of Science and Technology

<sup>2</sup> The University of Texas at Arlington



HUAZHONG UNIVERSITY  
OF SCIENCE & TECHNOLOGY



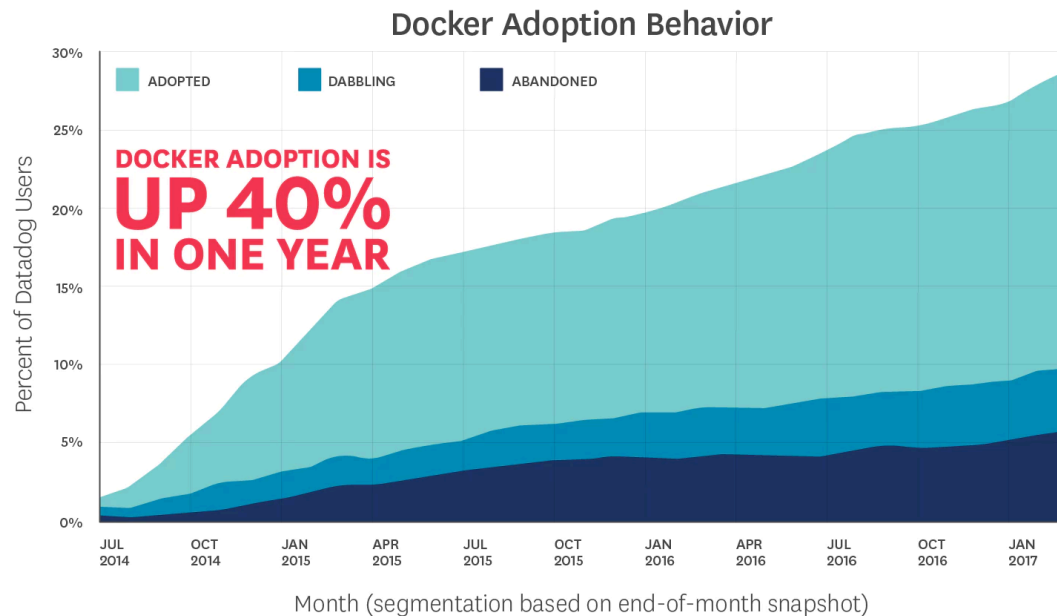
UNIVERSITY OF  
TEXAS  
ARLINGTON

# Containers

- A lightweight alternative to virtual machines for application packaging
  - Ease of management
  - Fast and cross-platform deployment
  - On-demand elasticity
  - **Low overhead, high efficiency and density**

# The Rise of Containers

- Increasingly and widely adopted in datacenters
  - Google Search launches about **7,000** containers every second



Source: Datadog



## Service Providers



## Dev Tools



## Official Repositories



## Operating Systems



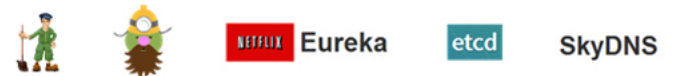
## Configuration Management



## Big Data



## Service Discovery



## Orchestration



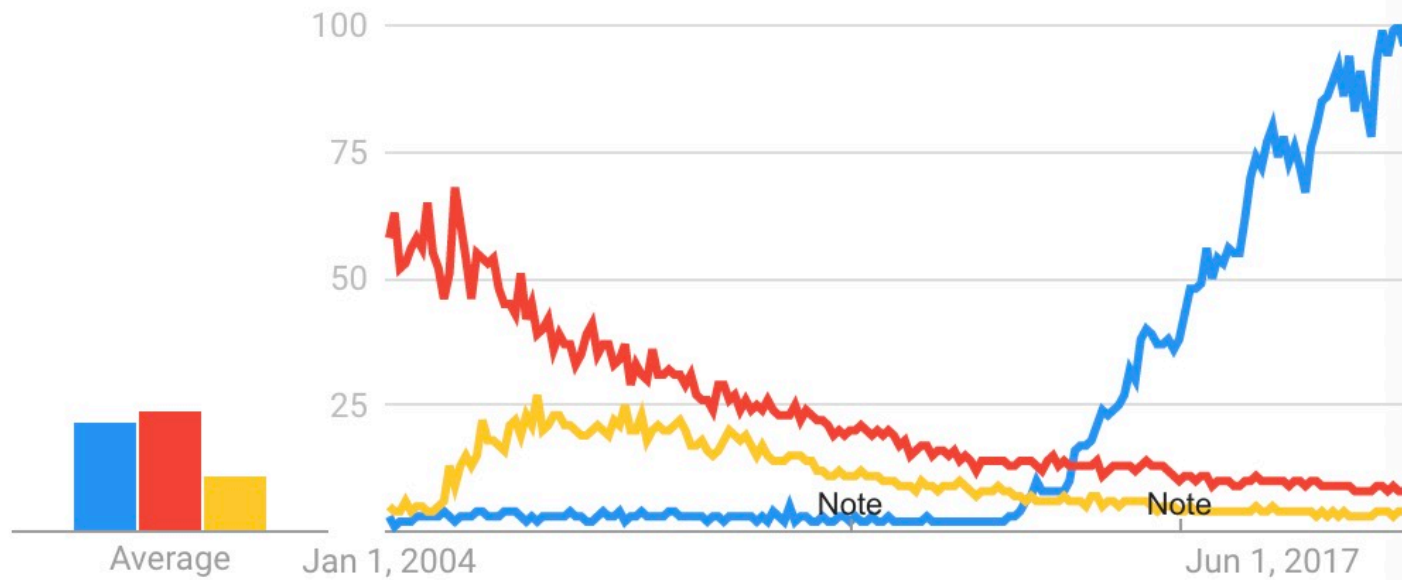
## System Integrators



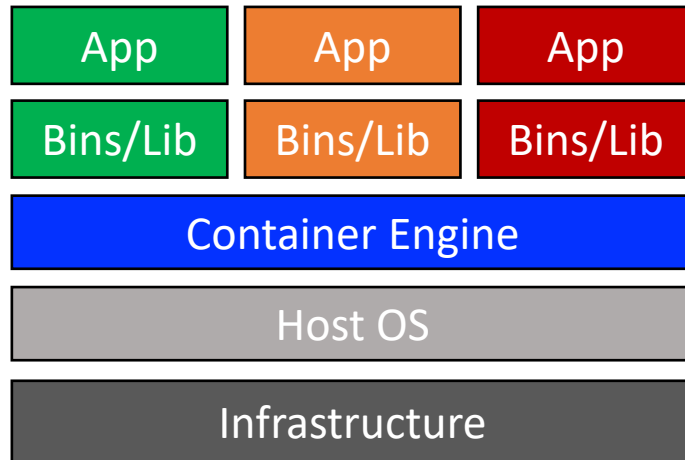
## Interest over time

Google Trends

● docker ● KVM ● Xen

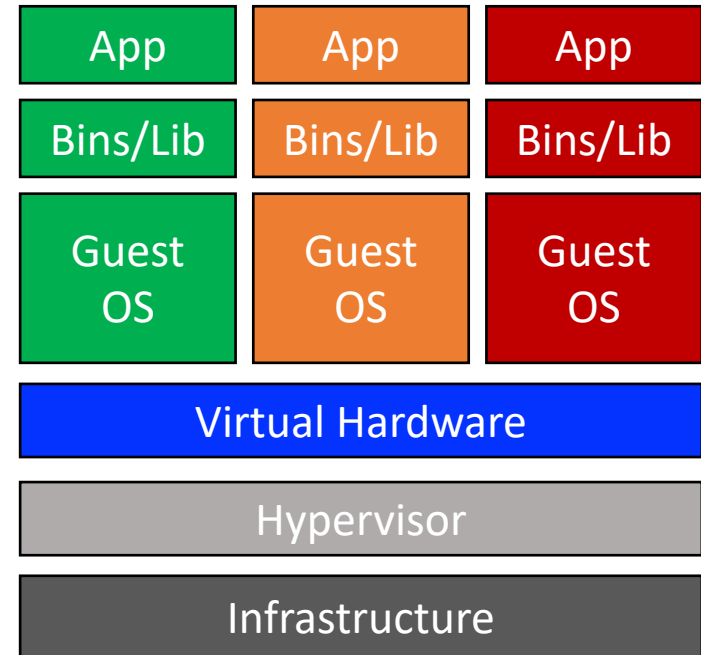


# Container vs. Virtual Machine



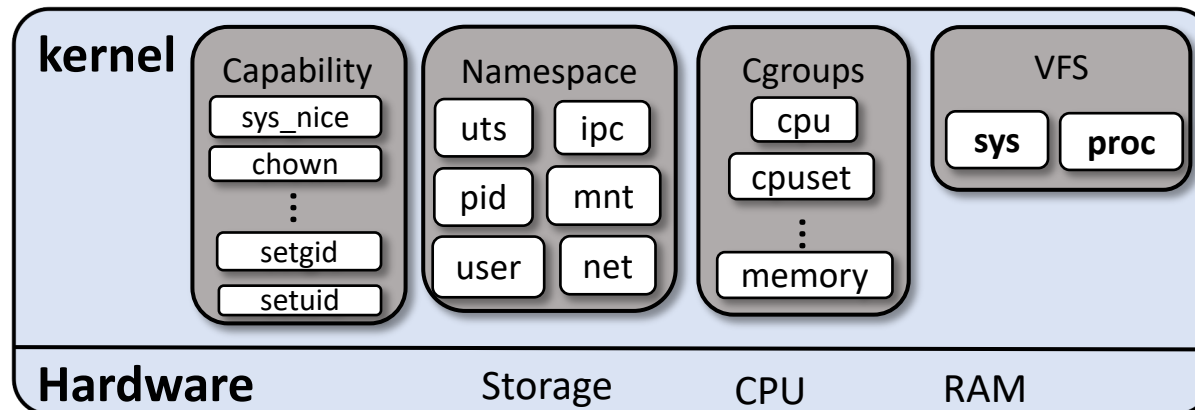
Container

**Shallower virtualization stack,  
less overhead while still  
providing application isolation**



Virtual machine

# Docker Containers



**Docker containers share the same OS kernel and are isolated through namespaces and control groups (cgroups)**

# Preliminaries - cgroups

- CPU isolation
  - **Mask** – specify on which CPUs a container can run
  - **Limit** – the maximum amount of CPU allocation
  - **Share** – the relative weight in CPU allocation
- Memory isolation
  - **Hard limit** – the upper bound on memory allocation, beyond which swapping is triggered
  - **Soft limit** – If system has free memory, containers can use more memory than soft limit. But if there is memory shortage, containers that are beyond their soft limit first reclaim memory



While weaker isolation between containers helps attain near-native performance, it creates an important ***semantic gap***



## Application

The absence of the virtual hardware abstraction allows containers to *observe* the *total available* resources on a host



Containers are only limited to *access* a *subset* of system resources due to resource capping or multiplexing

## System admin

Many applications behave differently  
when containerized due to the gap  
between *resource availability* and  
*constraints*

# Resource Availability vs. Constraints

kubernetes

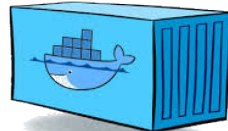


Create a container with 2 cores

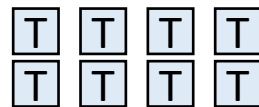
Create 3 containers with 2 cores



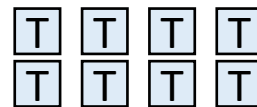
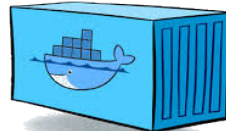
cpu\_share=2048



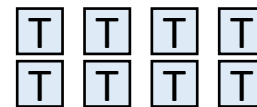
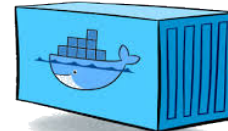
I own 8 cores



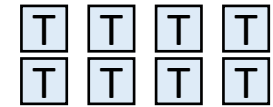
cpu\_share=2048



cpu\_share=2048



cpu\_share=2048



**Oversubscribed !!**

# Resource Availability vs. Constraints

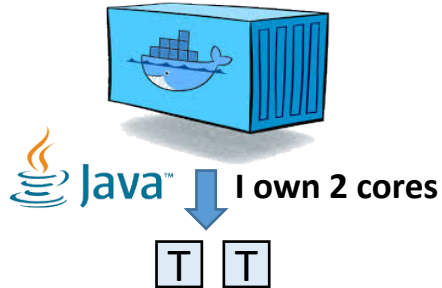
kubernetes



Create a container with 2 cores



cpu\_share=2048

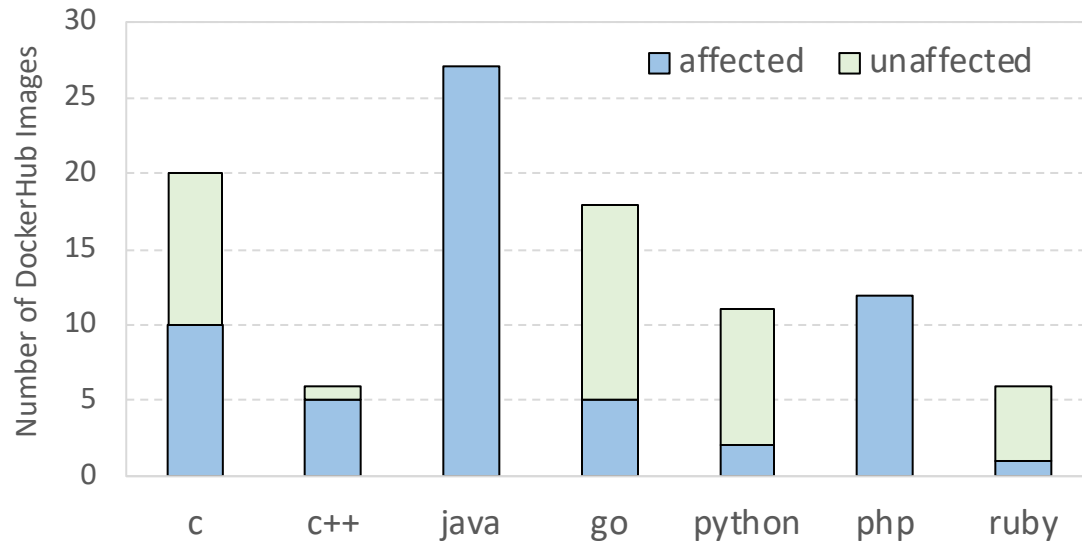


**Underutilized !!**

# Resource Discovery in Native Linux

- Containerized applications detect resource availability through the OS kernel
  - Many runtime systems use the `sysconf` in the GNU C library (glibc) to probe system resource
  - `sysconf` queries `sysfs` to determine the number of online CPUs and available memory size

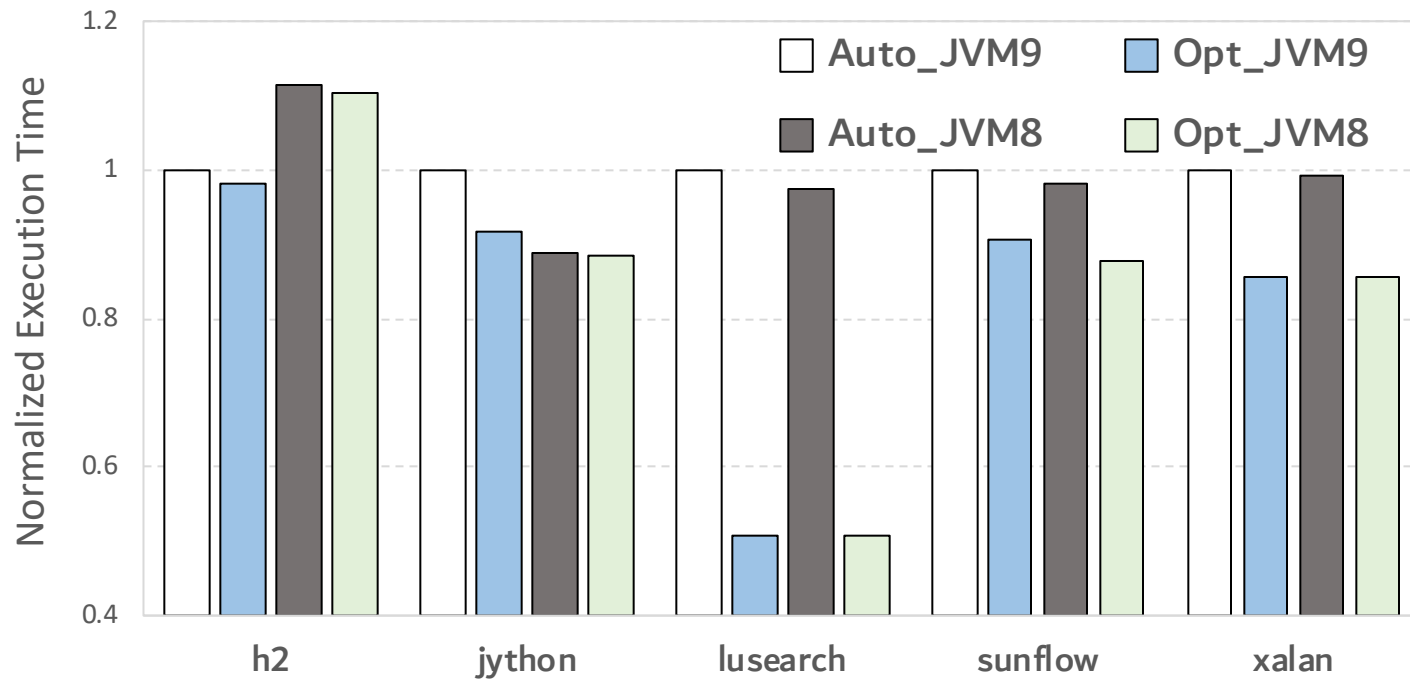
# The Prevalence of the Issue



Analysis of the top 100 application images on DockerHub.

More than **60%** of the top 100 apps on DockerHub automatically configure app based on detected resource availability and are **affected** by the semantic gap

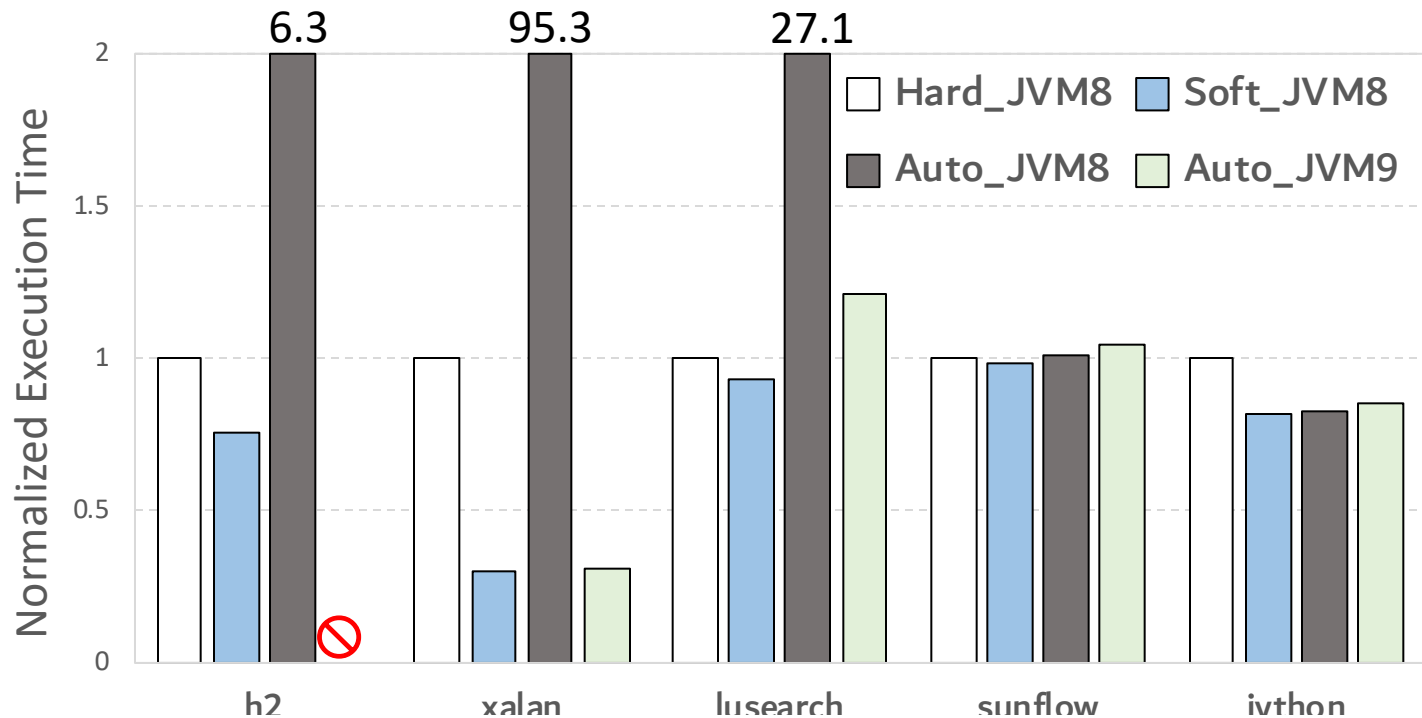
# CPU Oversubscription



Oversubscribing CPU leads to  
considerable performance degradation



# Memory Oversubscription



Oversubscribing memory can lead to program crashes or more severe degradation

# Problem Statement

- *Objective*

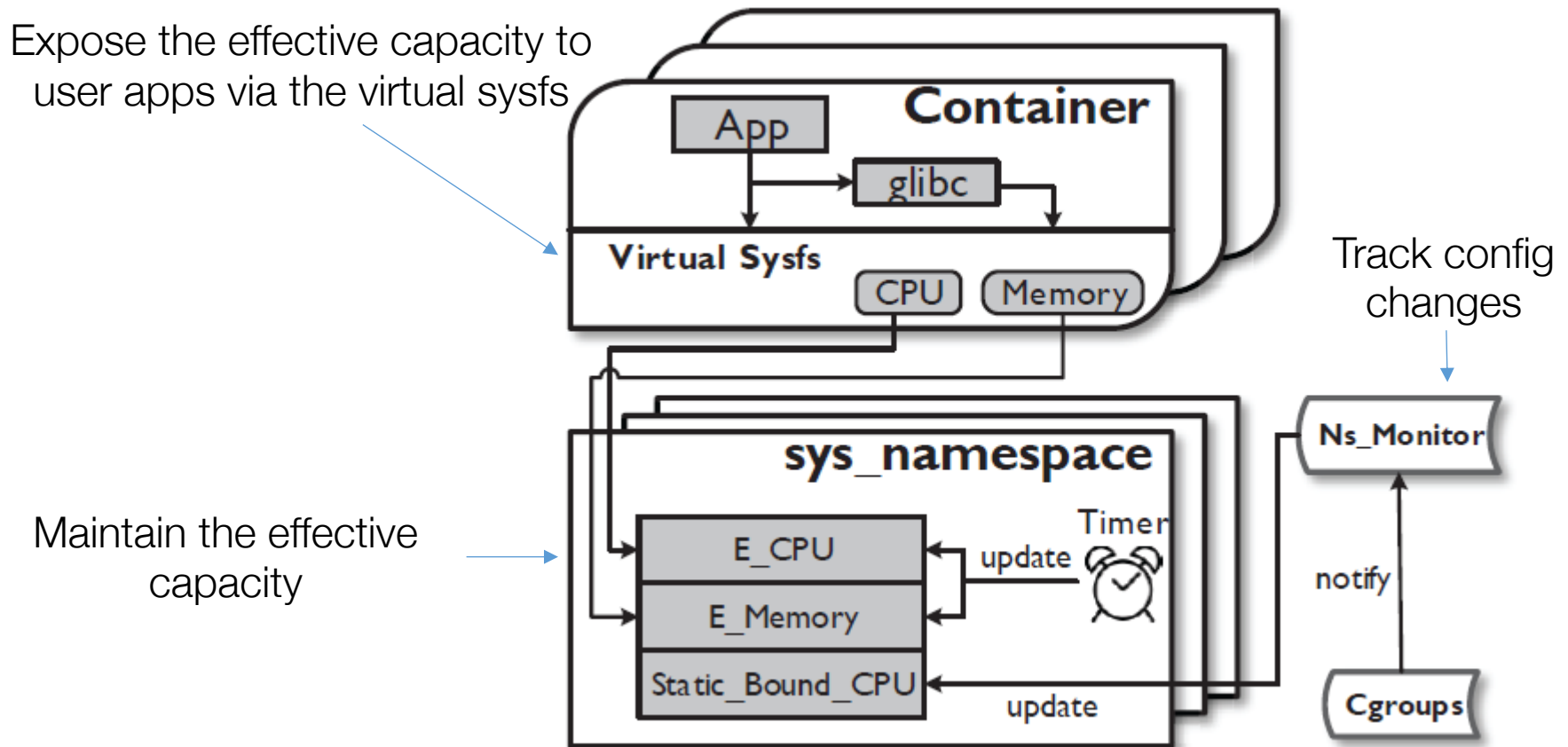
- Bridge the semantic gap between application-observed resource availability and the actual resource allocation
- Build elastic apps that efficiently utilize variable resources

- *Challenges*

- How to determine a container's real-time, actual resource allocation given various resource management schemes, e.g., limits, shares and resource affinity?
- How to develop a general method to export such information to apps?
- How can apps be made elastic, leveraging variable resource allocation?

# Adaptive Per-container Resource View

System architecture that enables the resource view :  
a new **sys\_namespace**, a **virtual sysfs**, and a system-wide daemon **Ns\_Monitor**.



# Effective CPU– Static Bounds

- Lower bound

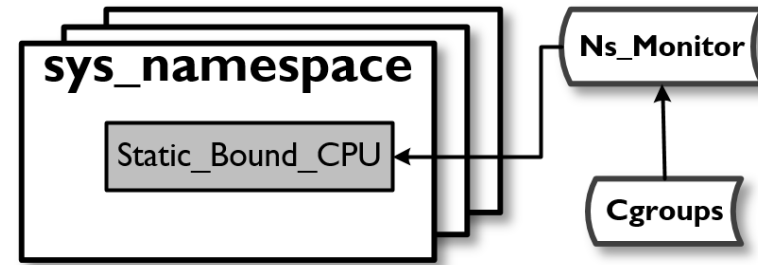
- Reservation, guaranteed CPU allocation

- $\min(\text{limit}, |M_i|, \left\lceil \frac{w_i}{\sum w_j} \cdot |P| \right\rceil)$

- Upper bound

- The maximum CPU usage

- $\min(\text{limit}, |M_i|)$

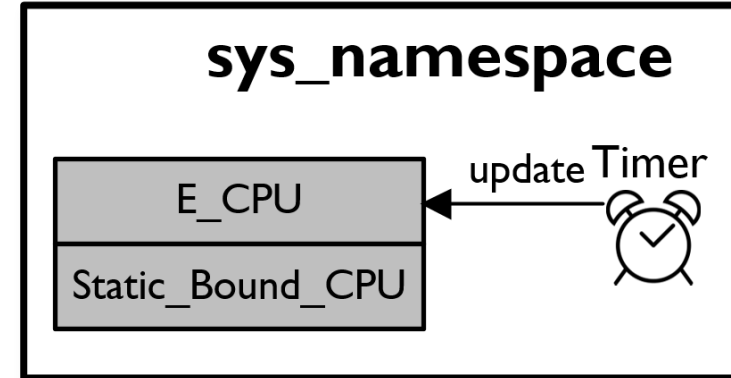


The lower and upper bounds do not change  
unless administrator reconfigures the container

# Effective CPU - Dynamic CPU Adjustment

- Periodic CPU adjustment

- [lower bound, upper bound]
- Effective CPU = # of CPU a container can fully utilize
- Expressed as a discrete *CPU count*



```
If  $p_{\text{slack}} > 0$  then
  if  $\frac{u_i}{E_{\text{CPU}_i} \cdot t} > 95\%$  and  $E_{\text{CPU}_i} < \text{UPPER\_CPU}_i$  then
     $E_{\text{CPU}_i} + +$ 
  end if
Else
  if  $E_{\text{CPU}_i} > \text{LOWER\_CPU}_i$  then
     $E_{\text{CPU}_i} - -$ 
  end if
End if
```

# Effective Memory

- Memory adjustment

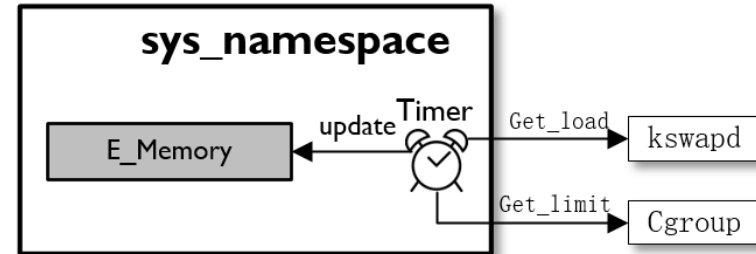
- [soft\_limit, hard limit]

- Gradually increase memory by 10% if

1. Container memory usage > 90%

2. The projected system-wide memory increase does not bring system-wide free memory below high\_watermark

- $\Delta_{system} = \frac{\text{change in system free mem}}{\text{change in container mem usage}} \times \Delta_{mem}$

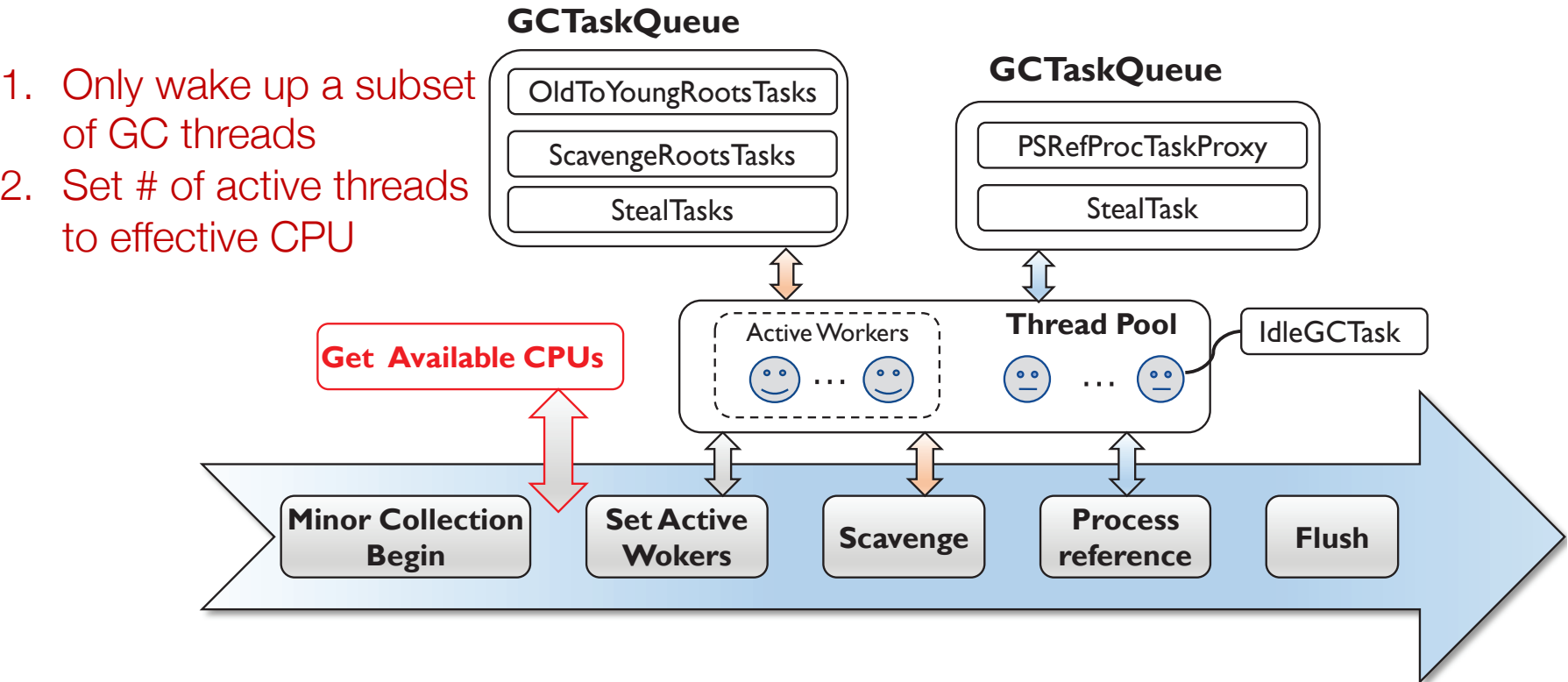


# Building Elastic Applications

- *Method*

- Instruct applications to check resource availability in per-container virtual `sysfs` and adjust program-level resource management accordingly
- CPU elasticity - adjust **thread-level parallelism**
- Memory elasticity – adjust **memory bounds** in runtime systems

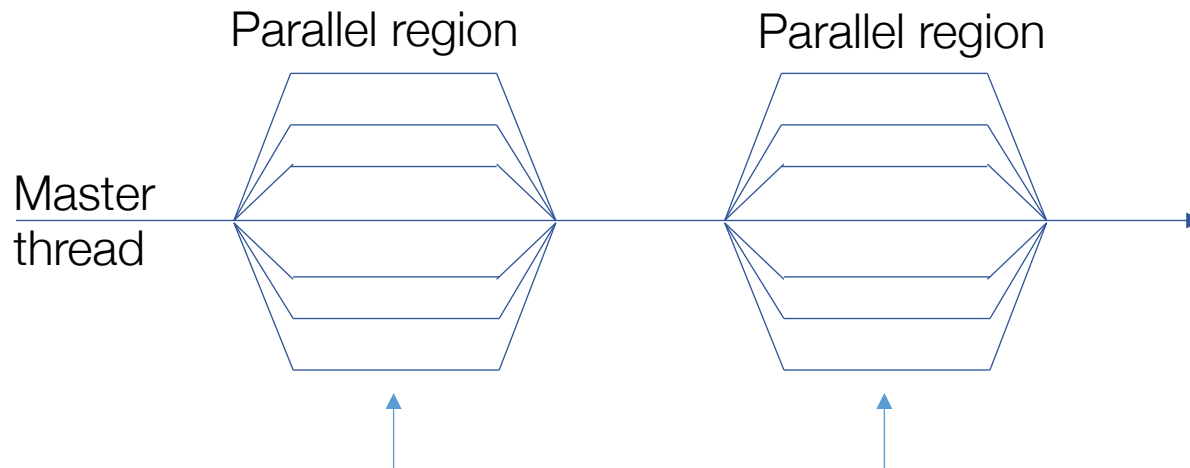
# Dynamic Parallelism – HotSpot JVM



The process of a Parallel GC



# Dynamic Parallelism – OpenMP



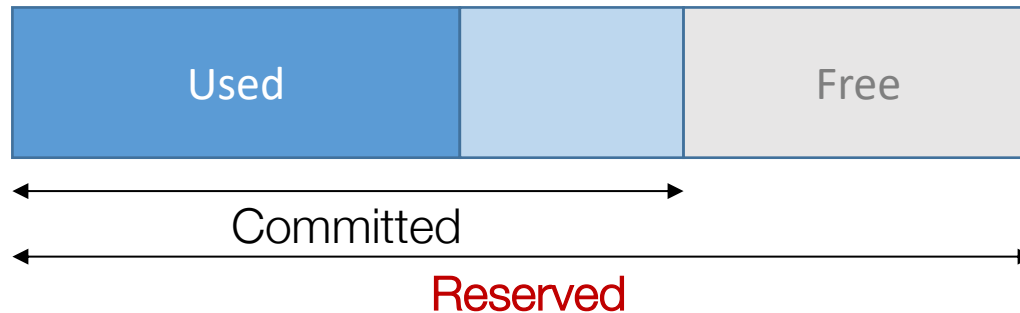
$\text{gomp\_dynamic\_max} : n\_thread = n\_onln - loadavg$



$\text{gomp\_dynamic\_max} : n\_thread = E\_CPU$

# Elastic Memory – HotSpot Heap Management

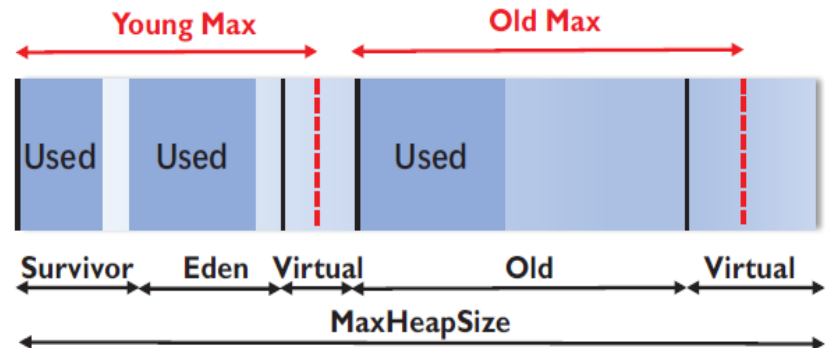
- The structure of JVM heap



- Adaptive heap sizing algorithm
  - Dynamically change the committed size
  - Expand or shrink the heap based on feedbacks of GCs.

# Elastic Memory – cont'd

- Scale up
  - Set reserved memory to effective memory and rely on the sizing algorithm to expand committed space
- Scale down
  - Trigger GC to free memory if: 1) new memory size is smaller than used and 2) adjustment is needed to maintain the ratio of young and old generations



# Evaluation

## Experimental Settings

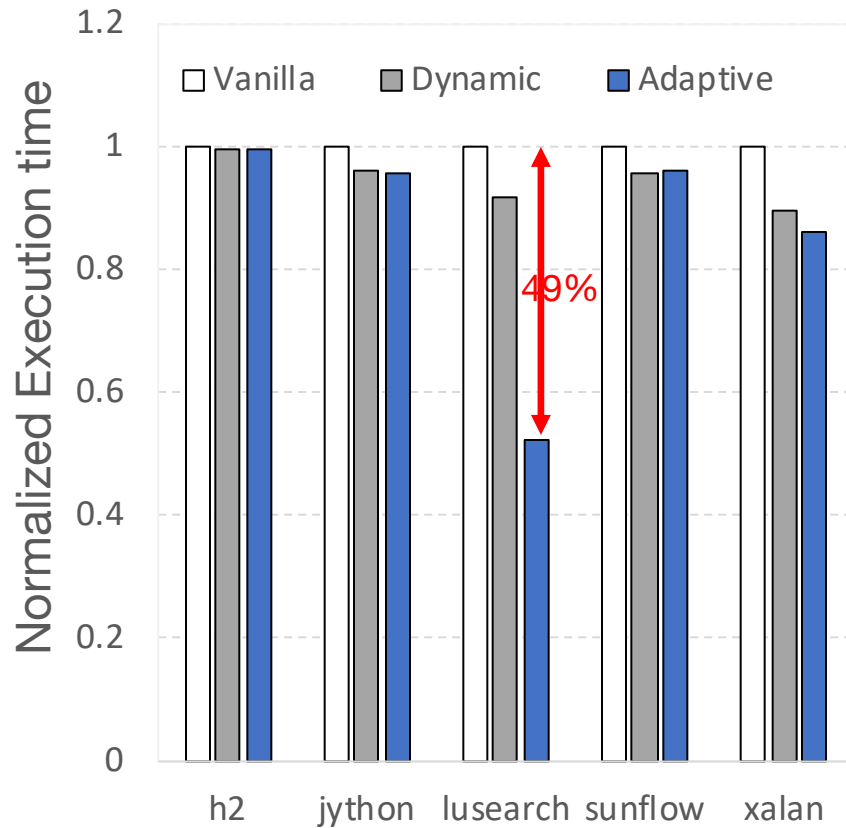
Hardware	
Server	PowerEdge R730
CPU	dual 10-core Intel Xeon 2.30 GHz processors
RAM	128GB
Storage	1TB SATA hard drive
Software	
System	CentOS 7 64bit
Kernel	4.12.3
Docker	Version 17.06.1
Benchmarks	Dacapo, PECjvm2008, HiBench, NAS parallel benchmarks

# Methodology

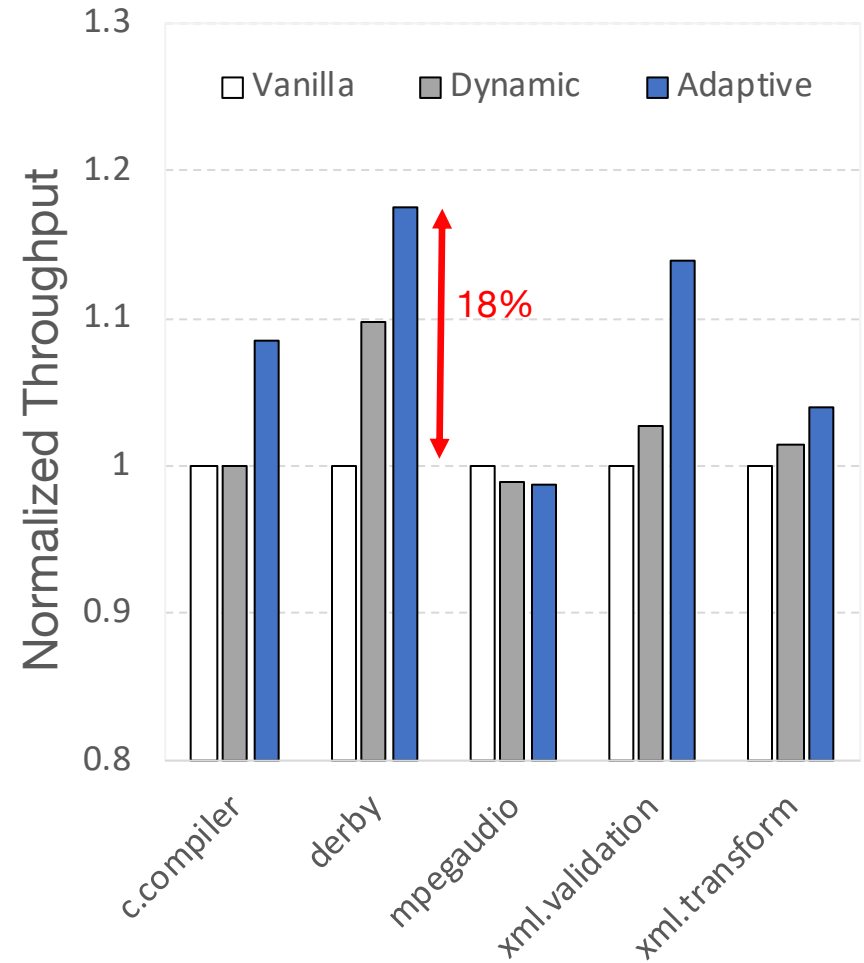
- Study various workload consolidation scenarios to evaluate
  - The (performance) benefits of per-container resource view
  - How well do applications adapt to variable resource availability
- Evaluated approaches
  - **Vanilla**: original JVM 8 and OpenMP without container awareness
  - **Dynamic**: JVM 8 and OpenMP with dynamic threads enabled and JVM 8 with adaptive heap sizing
  - **JVM 9 and JVM 10**: The new versions of JVMs that add container awareness. JVM 9 detects static CPU limit and affinity; JVM 10 detects static CPU share
  - **Adaptive**: Our proposed elastic apps based on effective CPU and memory

# Improvement on Application Performance

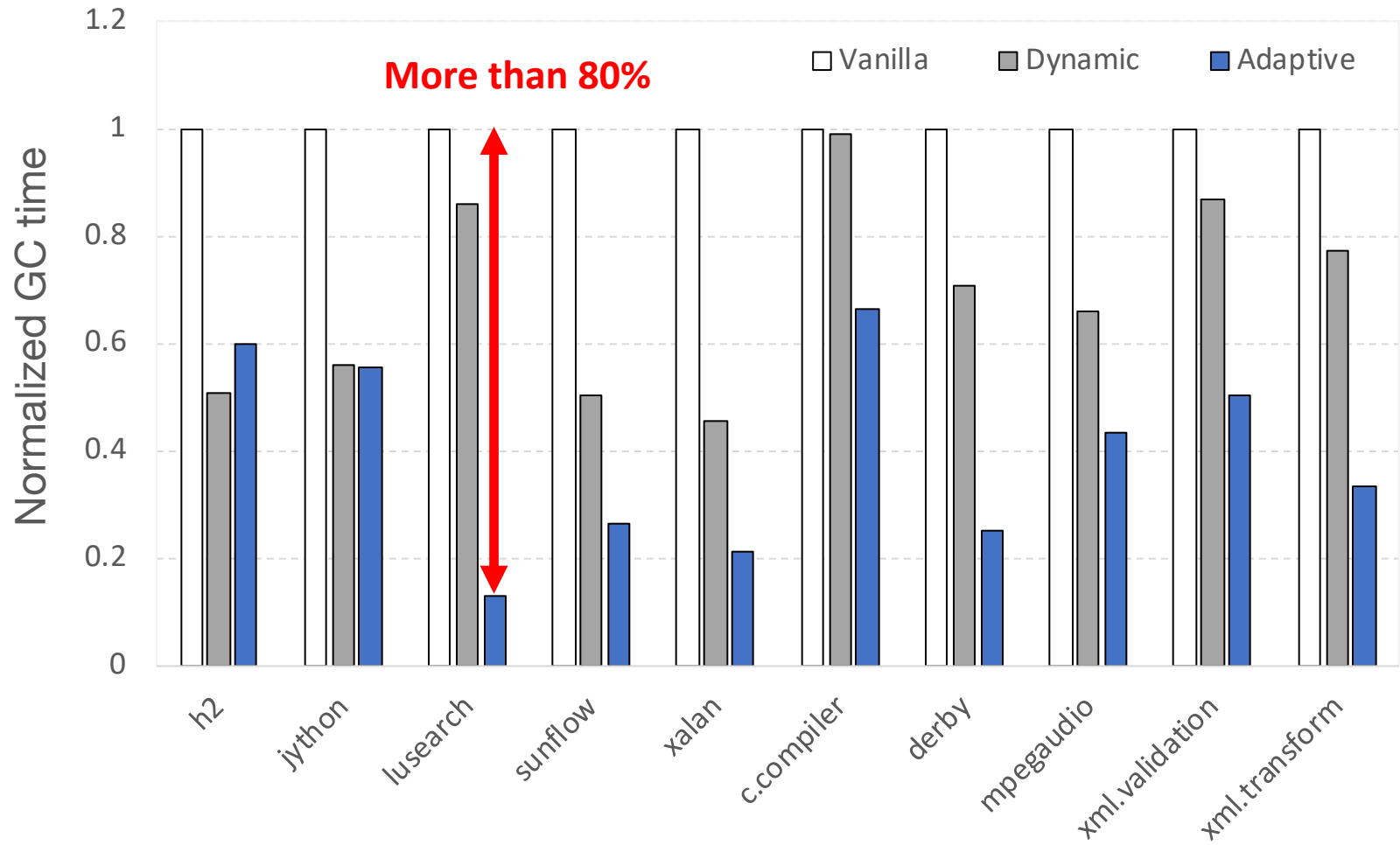
Lower is better



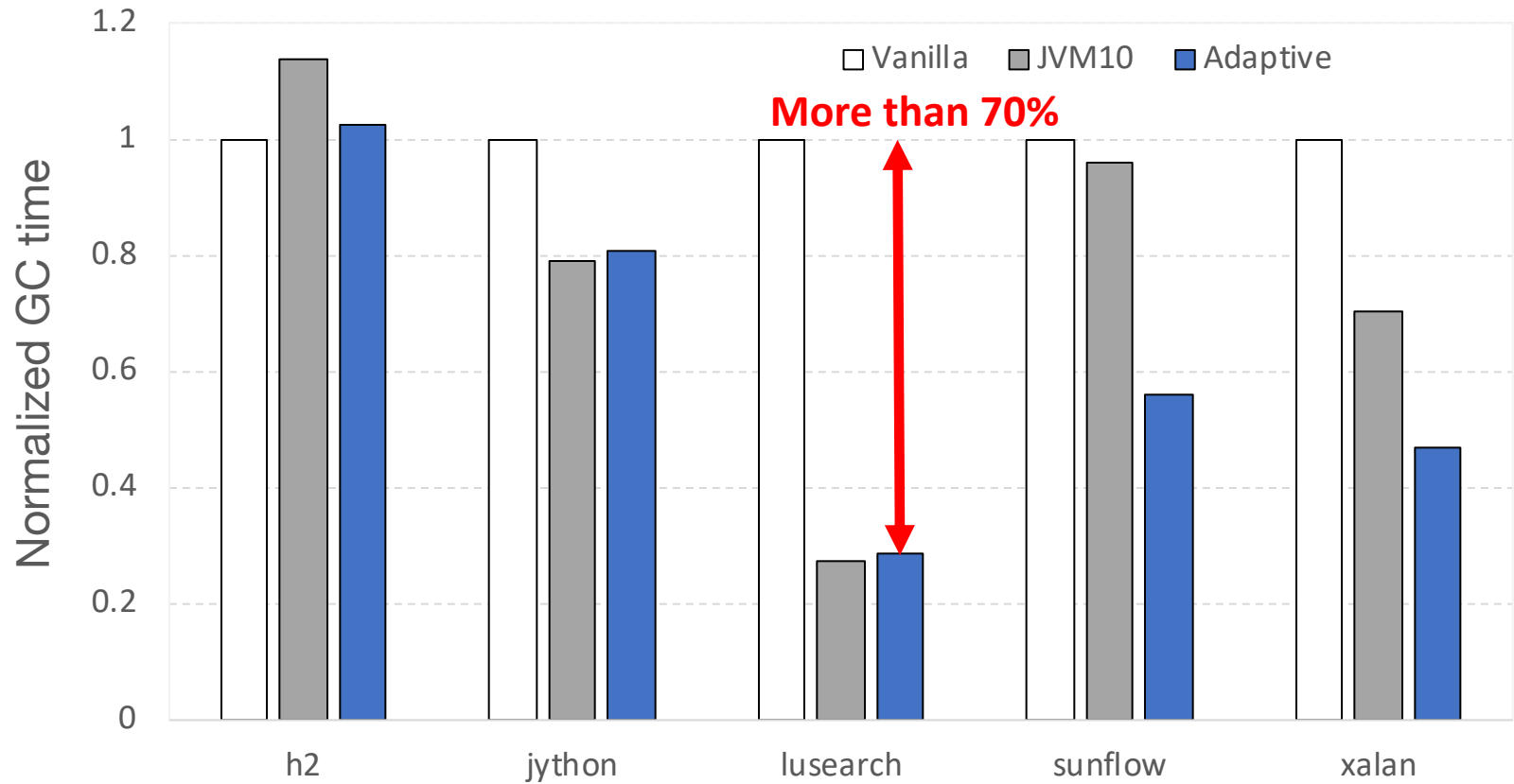
Higher is better



# Improvement on GC time

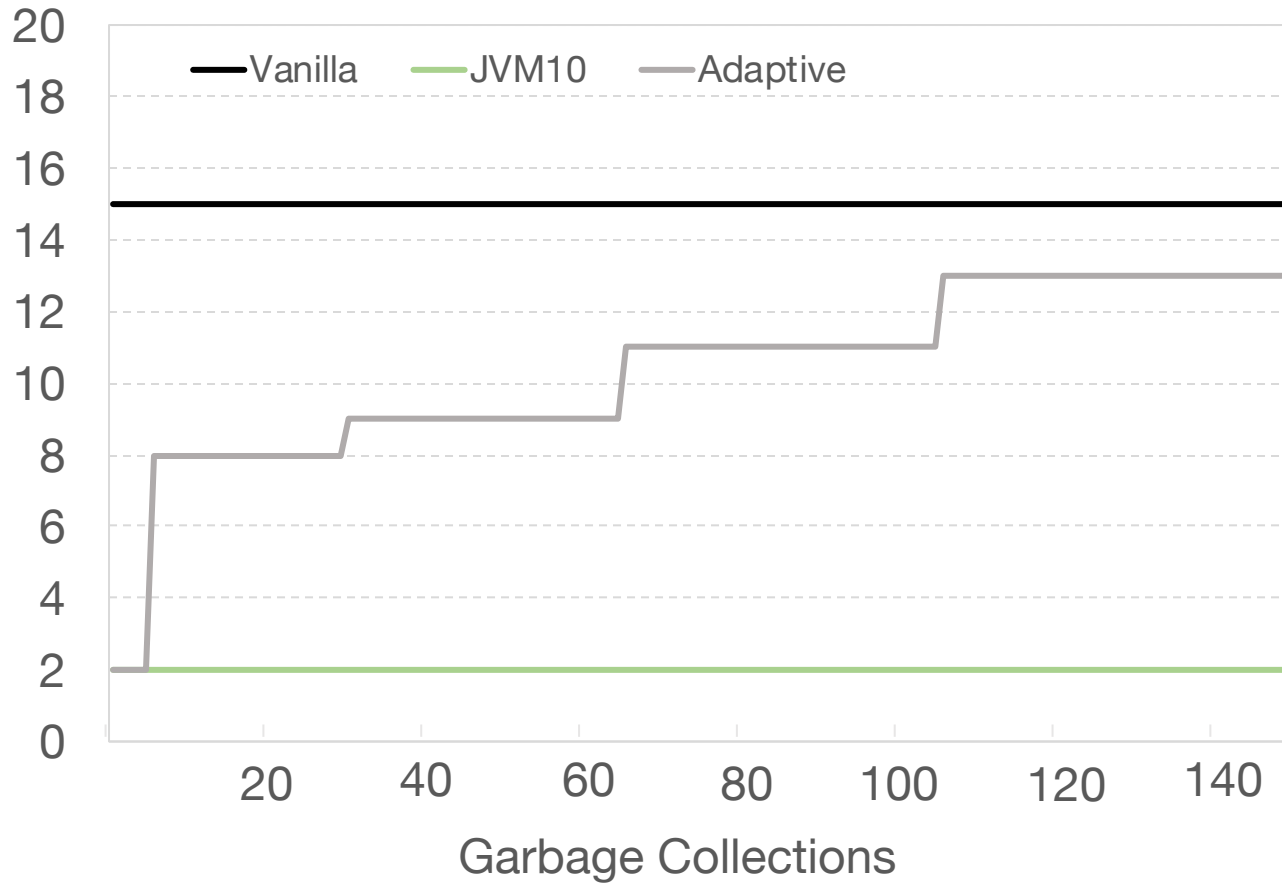


# Improvement with varying CPU availability



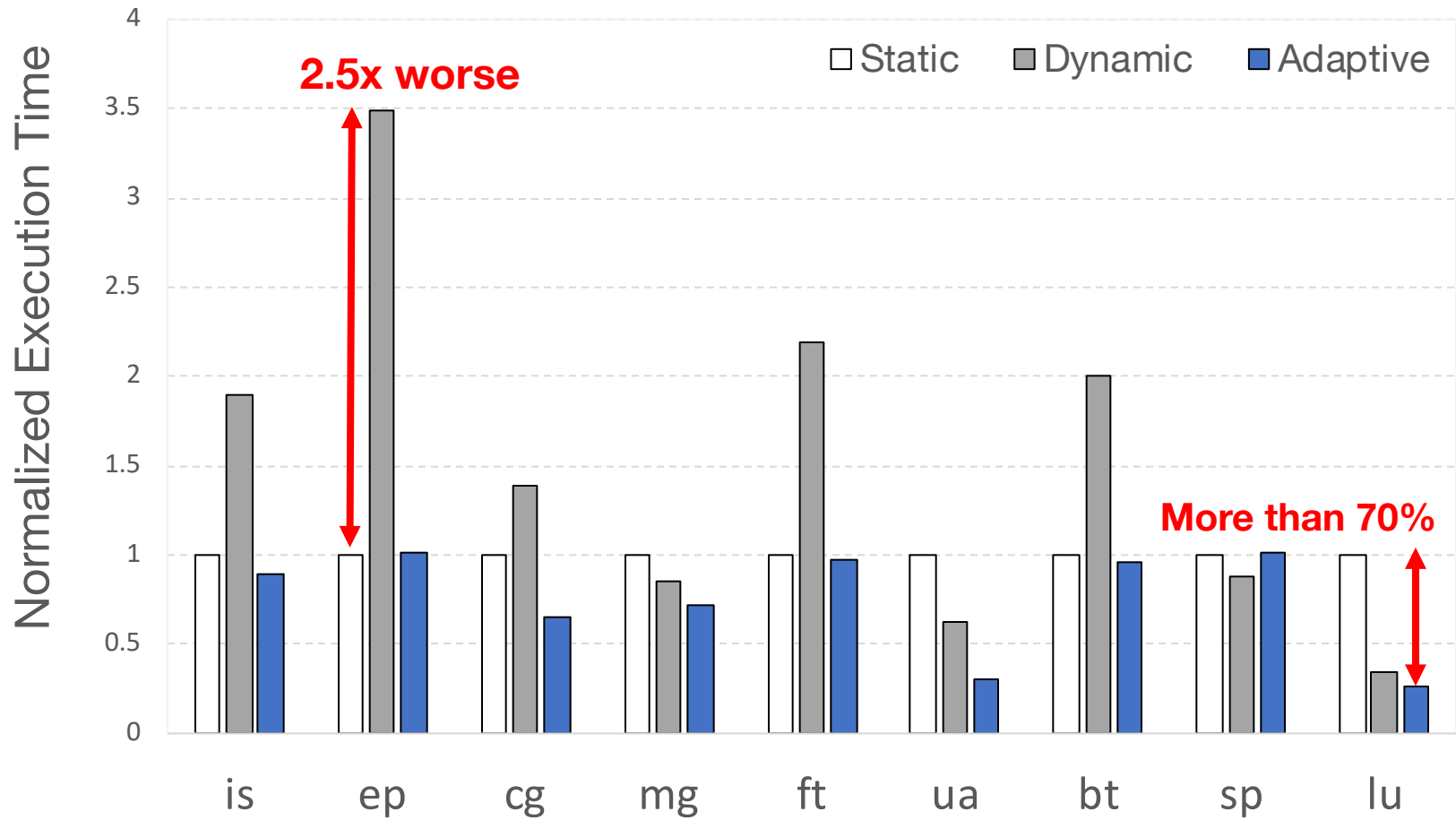


# Results on Different Versions of JVM

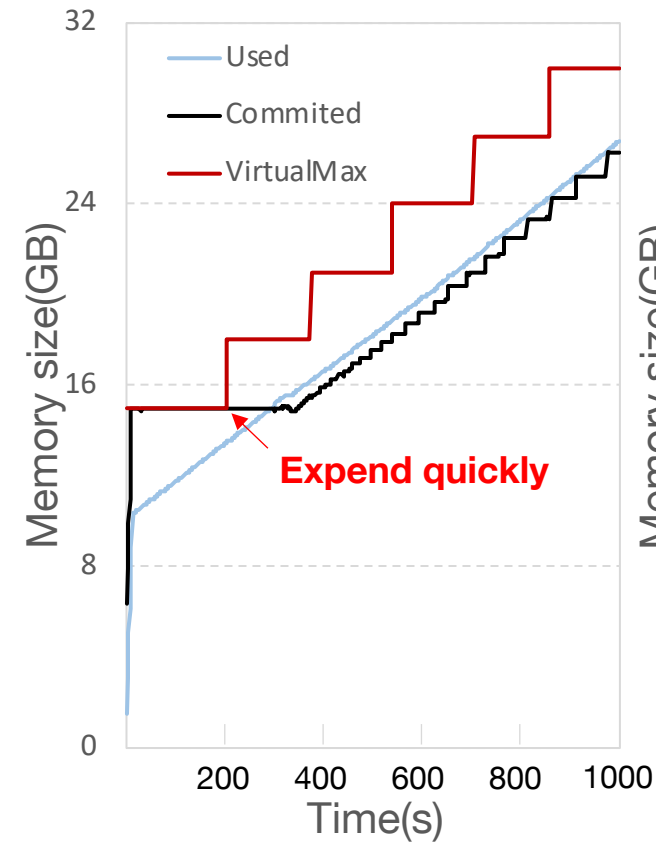


Number of GC threads in sunflow

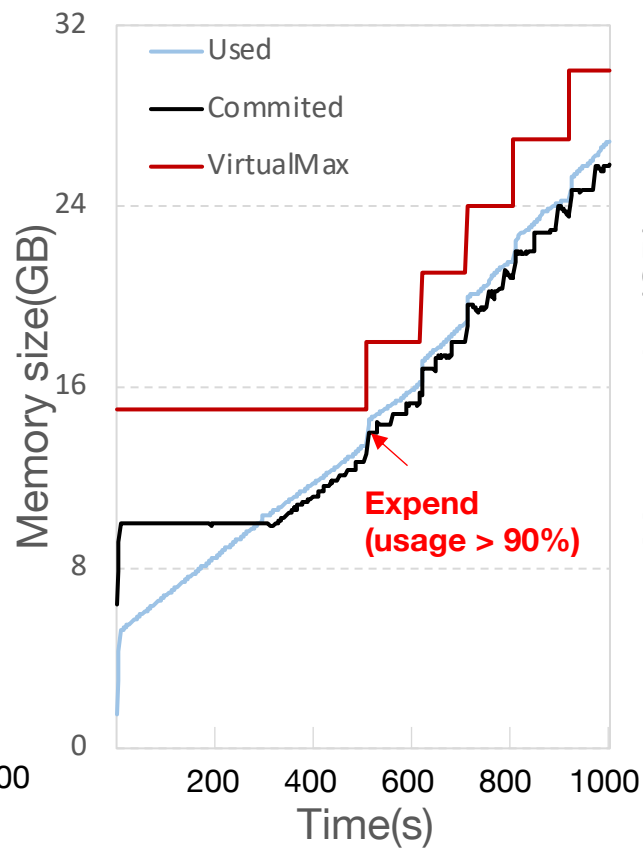
# Improvement on OpenMP



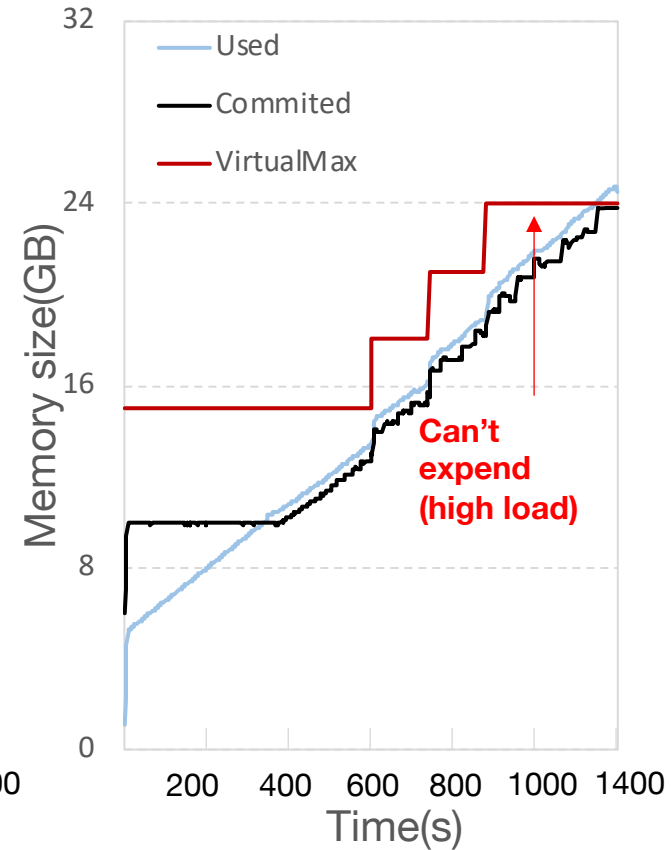
# Results on Elastic Heap



Vanilla JVM(Single)



Elastic JVM(Single)



Elastic JVM(Multiple)

# Conclusions and further research

## Motivation:

We identified an important semantic between *resource availability* and *constraints that can lead to* containerized applications performing inefficiently.

## Solution:

We created new sys\_namespace to propose per-container view of resources and developed a virtual sysfs to expose resources view to user space application.

## Results:

Experiment results showed consistent performance improvement for Java and OpenMP applications.

## Limitations :

Our system is still inconvenient for applications or runtime systems without mechanisms to alter thread-level parallelism or memory allocation at runtime.

Thanks!

*Questions?*