



Multi-Level Tracing of Containerized Application Orchestrator

Benjamin Saint-Cyr

Polytechnique Montréal
DORSAL

Introduction

Cloud Computing

- 1 aggregation and distribution of computing resources in a homogeneous way to clients
- 2 Allows customers to scale up operations instantly
- 3 Also facilitate the restitution of resources as needed
- 4 Flexible billing based on usage

Introduction

Containers

- 1 Containers have emerged as a more efficient method for resource allocation and isolation
- 2 Share the same Kernel
- 3 Containers simplify packaging and distribution services
- 4 Provide a consistent runtime environment in both development and production



Introduction

- 1 Orchestrators are an easy way to deploy containers to the cloud
- 2 They enable the deployment of services using descriptive configurations
- 3 Orchestrators facilitate achieving and maintaining the desired state
- 4 Also offers utilities for automatic scaling, implementing rollout strategies, and managing configurations and secrets



Introduction

- ① Performance is an important factor
- ② Reduced latency enhances the quality of service provided
- ③ Quicker boot times improves elastic scaling
- ④ More efficient resource utilization directly leads to reduced hosting costs



Research Objective

To observe and analyze the behaviour of a distributed system of orchestrated containers, with the aim of facilitating the diagnosis of performance issues



Methodology

Performance Aspects

- 1 Quicker boot times enhance elasticity in scaling
- 2 Better utilization of resources improves Quality of Service (QoS) and reduces costs



Methodology

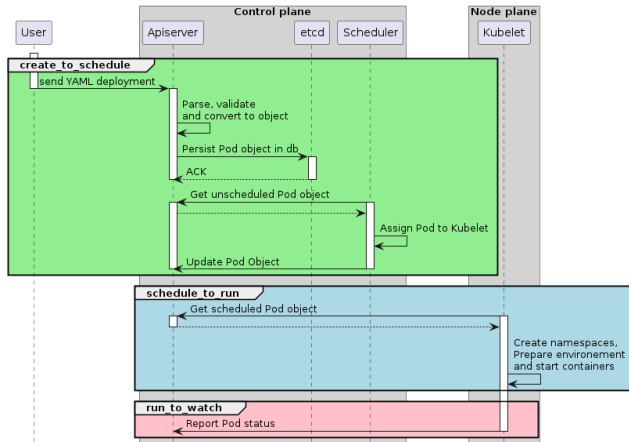


Figure: Components and phases of a pod start up

Methodology

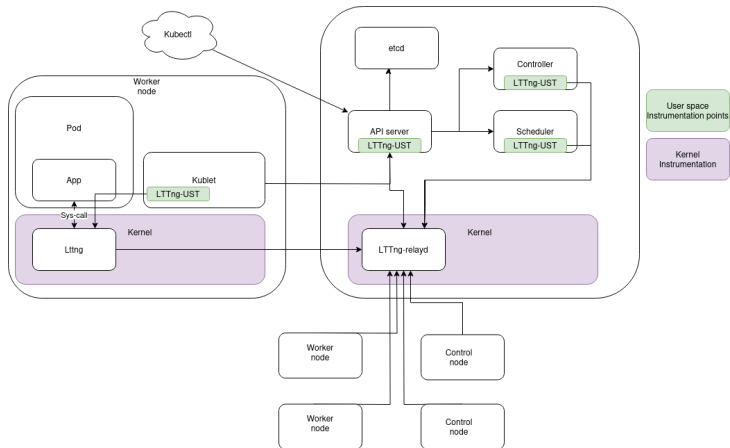


Figure: Tracing Architecture



Methodology

Common Issue During Execution : Managing CPU Limits

- ① Limits are specified in cores, but enforced as quotas
- ② A container can use more cores than allocated

$$Quota = Cores\ Limit \times Period$$

$$Throttled\ Time = Period \times (Cores\ Used - Cores\ Limit)$$



Methodology

Expectation

Running	Running
Running	Running
Idle	Idle
Idle	Idle
Idle	Idle
Idle	Idle
Idle	Idle
Idle	Idle

Reality

Running	PREEMPTED	Running	PREEMPTED
Running	PREEMPTED	Running	PREEMPTED
Running	PREEMPTED	Running	PREEMPTED
Running	PREEMPTED	Running	PREEMPTED
Running	PREEMPTED	Running	PREEMPTED
Running	PREEMPTED	Running	PREEMPTED
Running	PREEMPTED	Running	PREEMPTED
Running	PREEMPTED	Running	PREEMPTED



Methodology



The app needed 100ms
The app was throttled 120ms

Figure: Example of Latency Introduced by CPU Limit



Algorithm to characterize CPU usage per pod

Algorithm 1 CPU usage calculation per Cgroup

```

1: Define the function normalize(prevTime, time, cpuTime) as  $\frac{cpuTime}{time - prevTime} \times 100$ 
2: procedure ANALYZECPUUSAGE(fetchParameters, stateSystem)
3:   Initialize totalCpu  $\leftarrow 0$ 
4:   Initialize prevTime  $\leftarrow$  Get the start time from fetchParameters
5:   Initialize xAxis  $\leftarrow$  Initialize xAxis based on the time range and resolution
6:   Initialize map threadUsage  $\leftarrow$  Initialize the cpu usage of each thread based on stateSystem
7:   Initialize map CgroupUsage  $\leftarrow$  Initialize the cpu usage of each cgroup based on stateSystem
8:   for  $i \leftarrow 0$  to xAxis.length - 1 do
9:     time  $\leftarrow$  xAxis[i]
10:    for all  $(threadName, cpuTime) \in$  fetchThreadsCpuUsage(time) do
11:      totalCpu  $\leftarrow$  totalCpu + cpuTime
12:      normalizedThreadCpuUsage  $\leftarrow$  normalize(prevTime, time, cpuTime)
13:      threadUsage[threadName][i]  $\leftarrow$  normalizedThreadCpuUsage
14:      if TidToCgroup.contains(threadName) then
15:        key  $\leftarrow$  TidToCgroup.get(threadName)
16:        CgroupUsage[key][i]  $\leftarrow$  CgroupUsage[key][i] +
           normalizedThreadCpuUsage
17:      end if
18:    end for
19:    prevTime  $\leftarrow$  time
20:  end for
21:  return CgroupUsage, threadUsage
22: end procedure

```



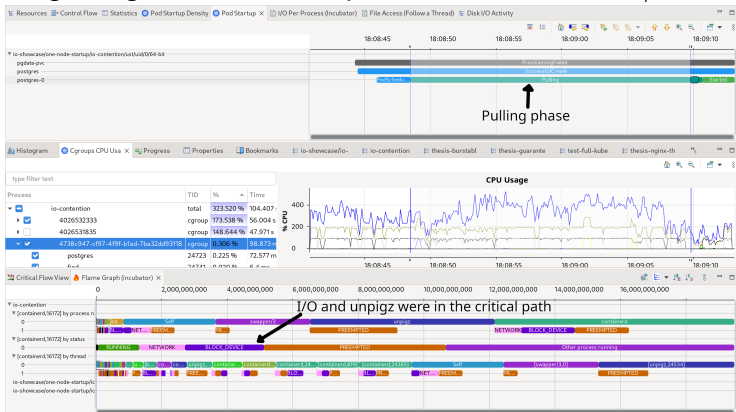
Results

Pod Lifecycle



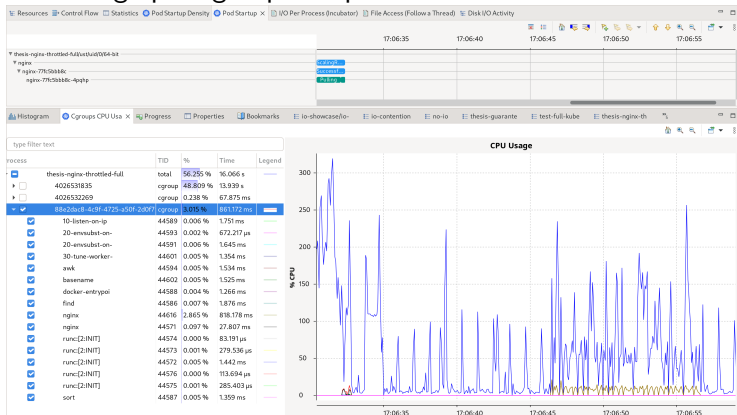
Results

Diagnosing Slow Pod Startup Issues Attributable to I/O Contention



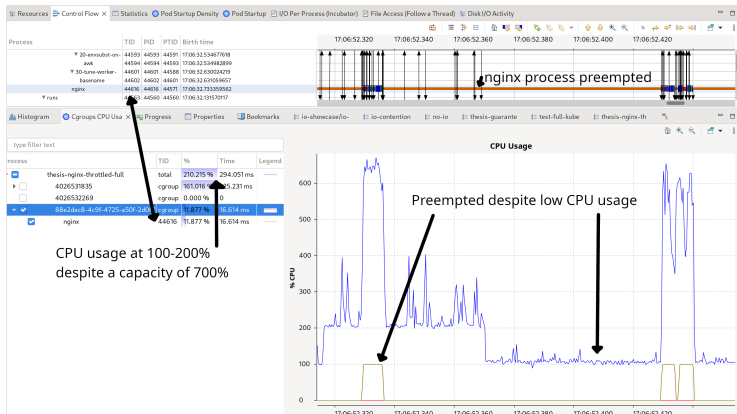
Results

CPU usage per cgroup and pod



Results

Diagnose latency due to limits



Conclusion

- 1 Demonstrates how tracing aids in observing the internal states of the orchestrator
- 2 Diagnosing issues at both the orchestrator and kernel levels
- 3 Highlights the utility of multi-level tracing for diagnostic purposes



Conclusion

Future work

- 1 Reuse methodology for other resources: Network, memory, disk, etc.
- 2 Use the trace data to characterize workload and improve scheduling



Conclusion

Thank you!
Questions?

