

Toro, a Dedicated Kernel for Microservices

Matias E. Vara Larsen Silicon-Gears Cesar Bernardini Barracuda Networks

What is a microservice?



Service Instance per Virtual Machine (VM)











Ingredients for ToroKernel

- User application within the kernel
- Cooperative Scheduler
- Dedicated Resources
- Single Thread Event Loop Networking

User application within the Kernel

- In general purpose OS, the user application executes as a process in the less privileged mode, e.g., ring3 in x86.
- The communication from the user application to the kernel relies on syscalls
- Context switch needs to switch from kernel mode to user mode, e.g., interruption and scheduling

User application within the Kernel

- In general purpose OS, the user application executes as a process in the less privileged mode, e.g., ring3 in x86.
- The communication from the user application to the kernel relies on syscalls
- Context switch needs to switch from kernel mode to user mode, e.g., interruption and scheduling



User application within the Kernel Proposal

- 1. Run kernel and user application in the most privileged level.
- 2. Rely on the hypervisor to isolate the context of each VM
- 3. Use a flat memory model that is shared by the kernel and the user application
- 4. Use only threads
- 5. Provide a simple kernel API dedicated to microservices







Ingredients for ToroKernel

- User application within the kernel
- Cooperative Scheduler
- Dedicated Resources
- Single Thread Event Loop Networking

Cooperative Scheduler

- In a General Purpose OS, the scheduler is in charge to distribute the CPU time for each process
- When a task has consumed its time, the scheduler switches to the next ready task
- The scheduler relies on a timer

Cooperative Scheduler

- In a General Purpose OS, the scheduler is in charge to distribute the CPU time for each process
- When a task has consumed its time, the scheduler switches to the next ready task
- The scheduler relies on a timer

Can we do it better for a single purpose kernel?

Cooperative Scheduler Proposal

1. Cooperative scheduler, i.e., each thread decides when to yield the CPU

2. Simple scheduler, i.e., the scheduler chooses the first thread in ready state

- 2. One scheduler per core
- 3. Remote creation of threads by relying on a lock-free algorithm







Core 2

Core 1



Ingredients for ToroKernel

- User application within the kernel
- Cooperative Scheduler
- Dedicated Resources
- Single Thread Event Loop Networking

- In a multicore system, the problematic resource is the shared memory. The use of shared memory causes:
 - Overhead in the memory bus
 - Overhead in the cache to keep it coherent
 - Overhead to guaranty mutual exclusion, e.g., use of spin-locks

Dedicated Resources Proposal

- Toro improves memory access by keeping the resources locals:
 - The memory is dedicated per core
 - The kernel data structures are dedicated per core
 - The access to kernel data structures is lock free









Network Stack

Network Card

Memory 2

Core 2

Messages, shared memory or any mechanism to communicate between threads







In Toro, microservices are first-class objects

Ingredients for ToroKernel

- User application within the kernel
- Cooperative Scheduler
- Dedicated Resources
- Single Thread Event Loop Networking

- Toro networking is based on the single thread event loop model [1], i.e., one thread per microservice
- The kernel provides a dedicated API to create microservices
- The kernel implements the microservice and improves the CPU usage[2]

[1] Node JS Architecture – Single Threaded Event Loop[2] Reducing CPU usage of a Toro Appliance, FOSDEM'18











HelloWorld Microservice Example

- We implement a simple microservice that responds "Hello World"
- We implement it by using three approaches: Docker, Ubuntu guest (KVM) and Toro guest (KVM)
- We compare these approaches in term of:
 - Deploying Time
 - Bootstrap Time
 - Image Size
 - CPU usage
 - Time per Request -

- Microservice footprint

HelloWorld Microservice Example setup

Docker	General OS/KVM	Toro/KVM
 Cpu limit Memory limit	 Cpu limit Memory limit	 CPU limit Memory limit
NGINXUWSGI (4 processes)Flask	NGINXUWSGI (4 processes)Flask	Toro WWW server

Deploying Time

• Time required to build an image within the microservice



Bootstrap Time

• Time to boot and to answer the first request



Image Size

• Size of the image that contains the microservice and its dependencies.



CPU Usage

TORO CPU Usage



End-User Delay

 Benchmarking with ab and measuring the Time per Request (mean) [ms]

		Number of Concurrent Request		
Approach	CPUs	200	500	1000
Docker	4	139.980 ms	333.937 ms	801.422
Ubuntu/KVM	4	94.507 ms	238.149 ms	560.513 ms
TORO/KVM	1	120.065 ms	301.736 ms	596.792 ms

Take away lessons

- Minimal image size (< 4MB)
 - NGINX docker image is 15-times the size of a Toro image
- Continuous Integration: 1 sec to re-deploy a microservice
 - Deploy an OS w/similar configuration takes 300 sec, with docker ~50 sec
- Time per Requests
 - Comparable level with cutting edge technology (NGINX)
- CPU Usage
 - Comparable with Docker
 - Toro is 100% isolated from the host OS, Docker is not.

Summary

- Toro is a kernel dedicated to run microservices
- Toro provides a dedicated API to specify microservices
- Toro design is improved in four main points:
 - Booting time and building time
 - communication to kernel
 - memory access
 - networking

Future Work

- Ease tooling to develop, test and debug microservices
- Investigate new use-cases
- Investigate the porting of applications
- Investigate new ideas to improve the network stack for microservices, e.g., improve socket scheduling for http, resource allocation algorithm



QA

- http://www.torokernel.io
- torokernel@gmail.com
- Twitter @torokernel
- Torokernel wiki at github
 - My first Three examples with Toro
- Test Toro in 5 minutes (or less...)
 - torokernel-docker-qemu-webservices at Github



Thanks!

