



---

# Serverless Confidential Containers: Challenges and Opportunities

---

**Carlos Segarra**

(w/ Tobin Feldman-Fitzthum and Daniele Buono)

**Large-Scale Data & Systems (LSDS) Group - Imperial College London**

**Visiting IBM TJ Watson (Sep'23 – Nov'23)**

<https://carlosgarra.com>

[<cs1620@ic.ac.uk>](mailto:cs1620@ic.ac.uk)



# Agenda

---

## **1. Introduction to Confidential Serverless**

- Characterising serverless functions: Cold/Warm starts and burstiness
- Problems with existing serverless offerings

## **2. Background (Will skip most of it):**

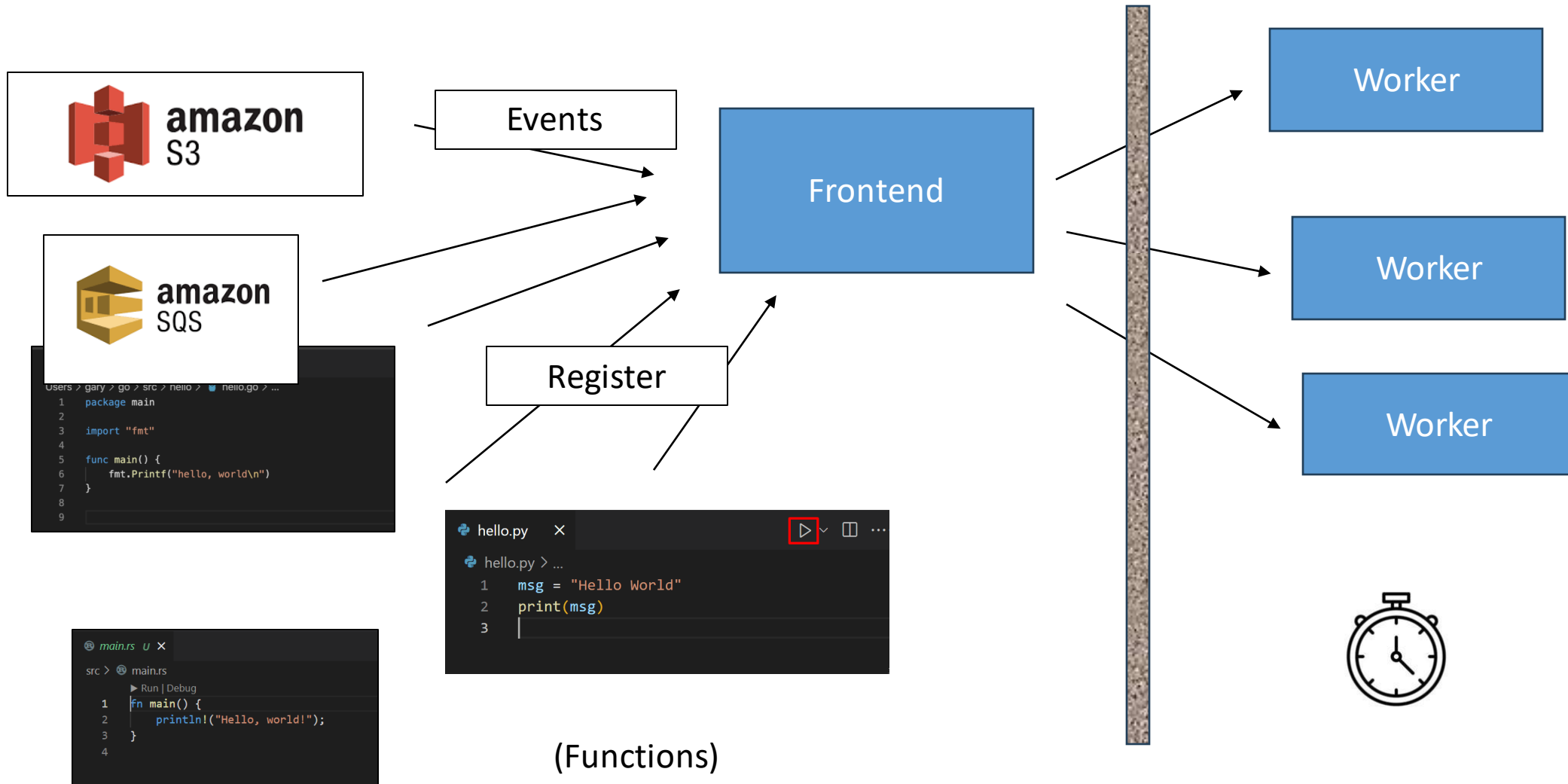
- Design space for confidential serverless
- Kata and Confidential Containers

## **3. PoC: Knative on Confidential Containers (Will skip most of it)**

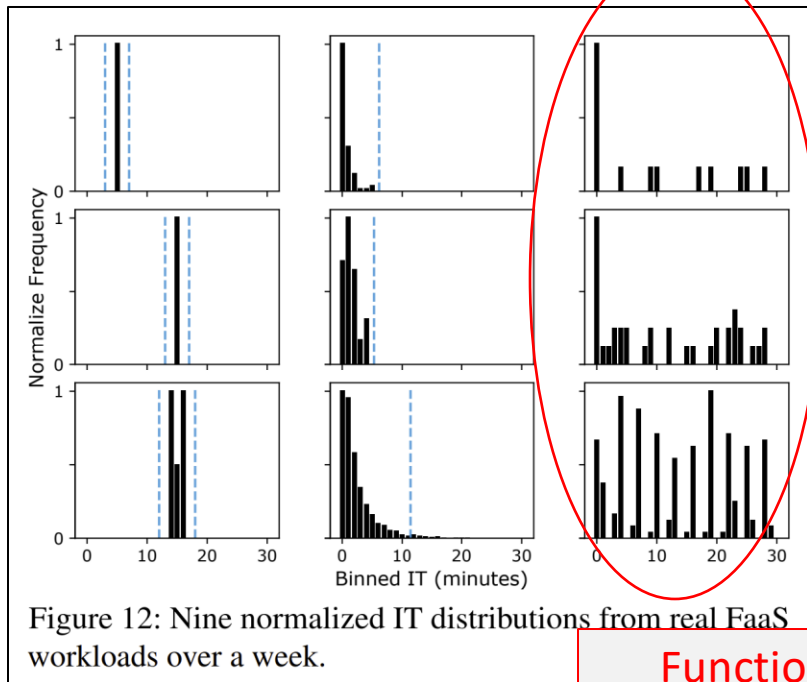
## **4. Evaluation**

- Cold-Starts
- Warm-Starts
- Instantiation Throughput

# Introduction: Serverless Functions

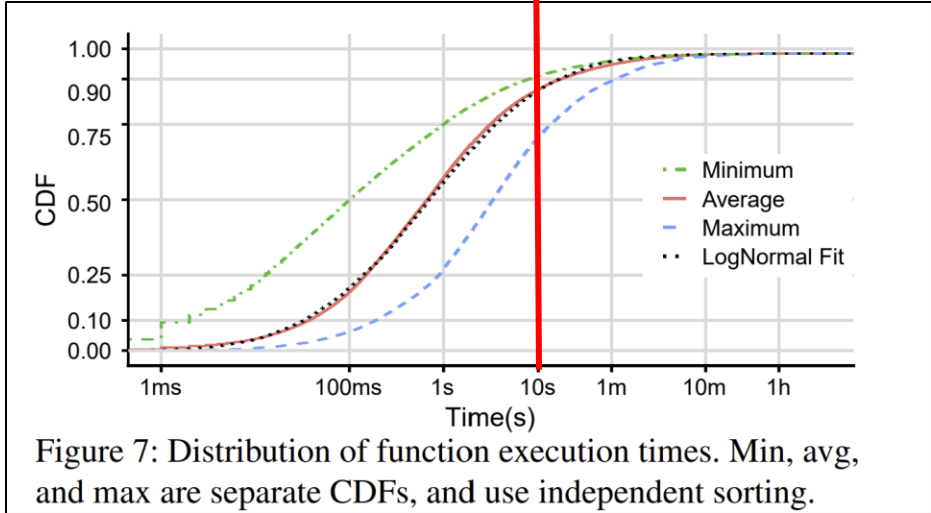


# Introduction: Characterizing Serverless Functions



Functions are bursty!

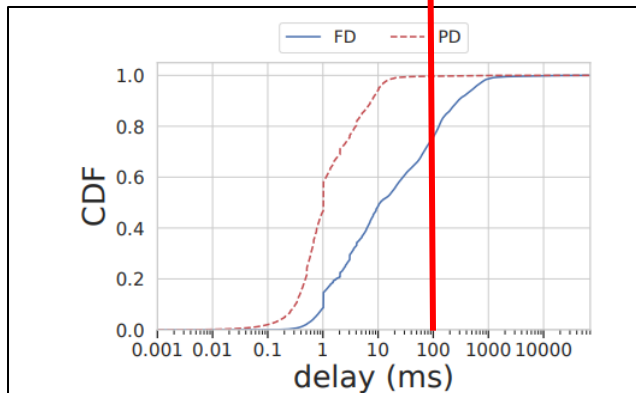
Functions are short-lived!  
(90% shorter than 10s)



[ATC'20] Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider

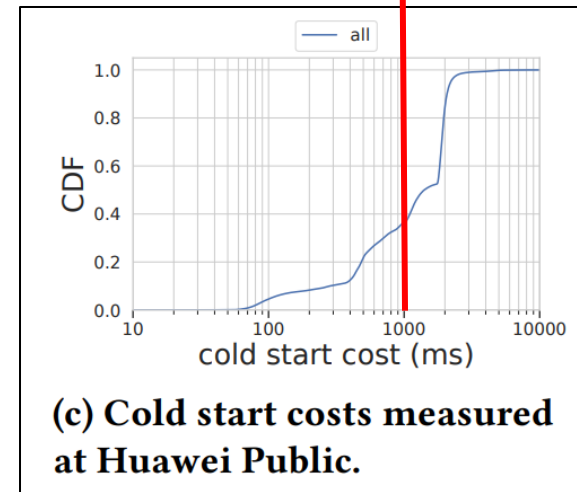
# Introduction: Characterizing Serverless Functions

Yet 80% of functions execute in < 100 ms !



**(a) Function Delay, Platform Delay in Huawei Private.**

Only 40% of functions take less than 1s to start



**(c) Cold start costs measured at Huawei Public.**

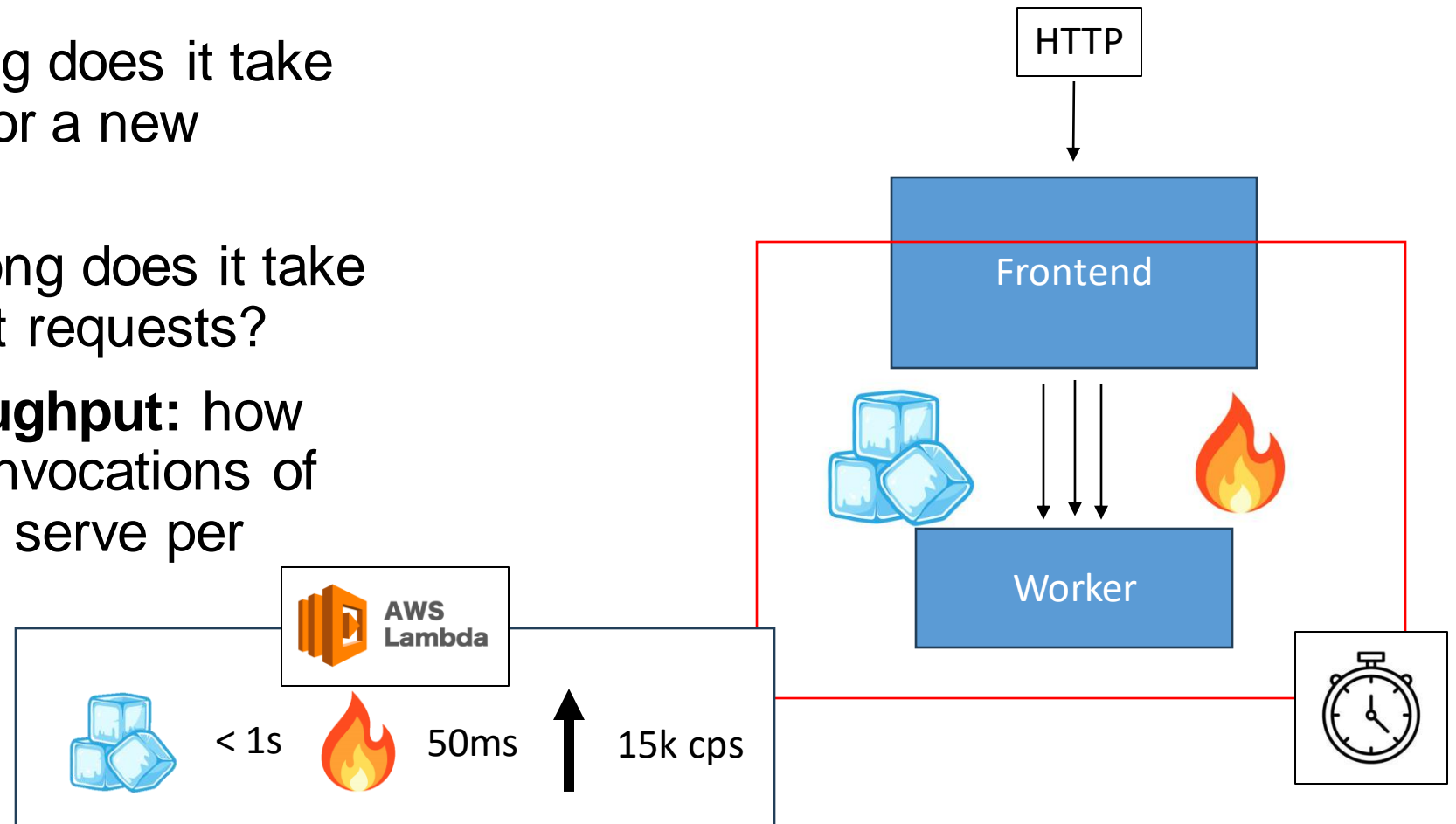
[SoCC'23] How Does It Function? Characterizing Long-term Trends in Production Serverless Workloads

# Introduction: Problems in Serverless

**Cold-Start:** how long does it take to serve a request for a new function?

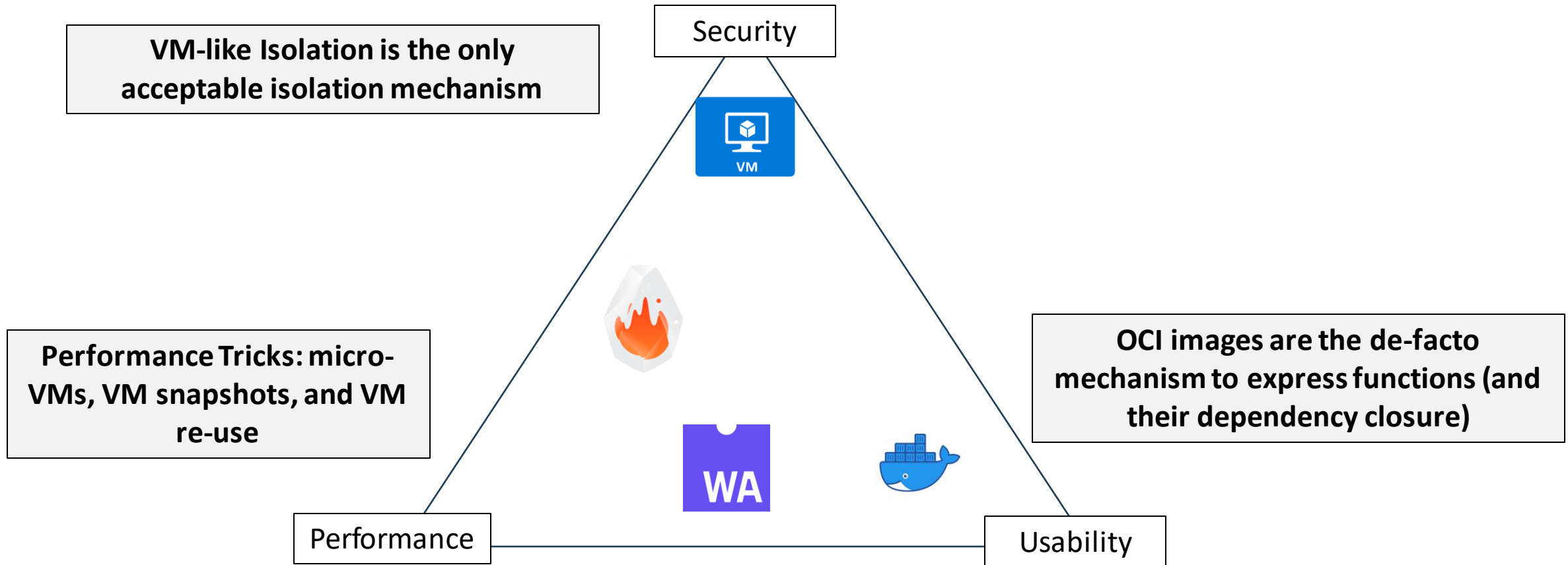
**Warm-Start:** how long does it take to serve subsequent requests?

**Instantiation Throughput:** how many (concurrent) invocations of this function can we serve per second?



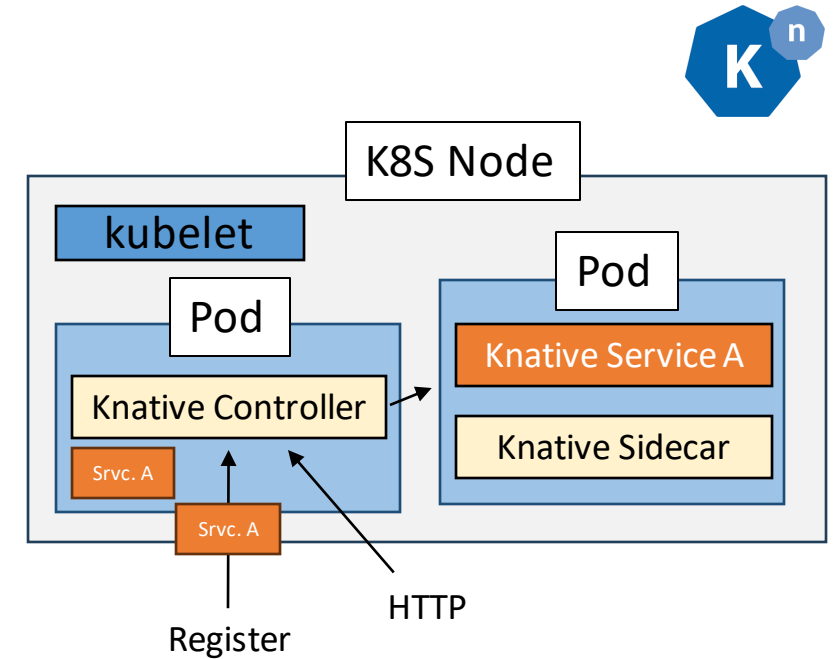
[ATC'23] On-demand Container Loading in AWS Lambda

# Introduction: Inter-Function Isolation in Serverless



# Introduction: More Problems in Serverless!

Inter-function isolation is fine, but not enough!



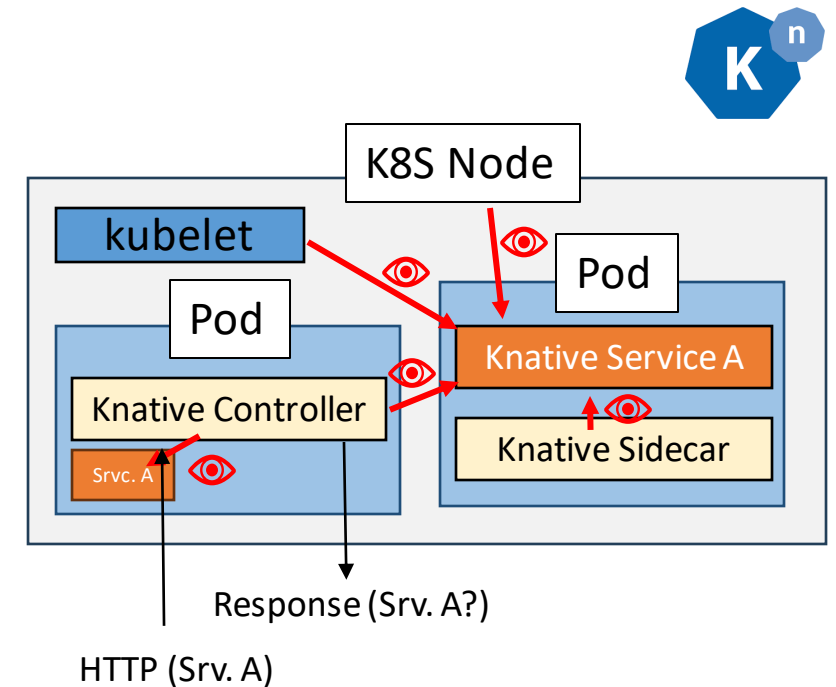


# Introduction: More Problems in Serverless!

**Inter-function isolation is fine, but not enough!**

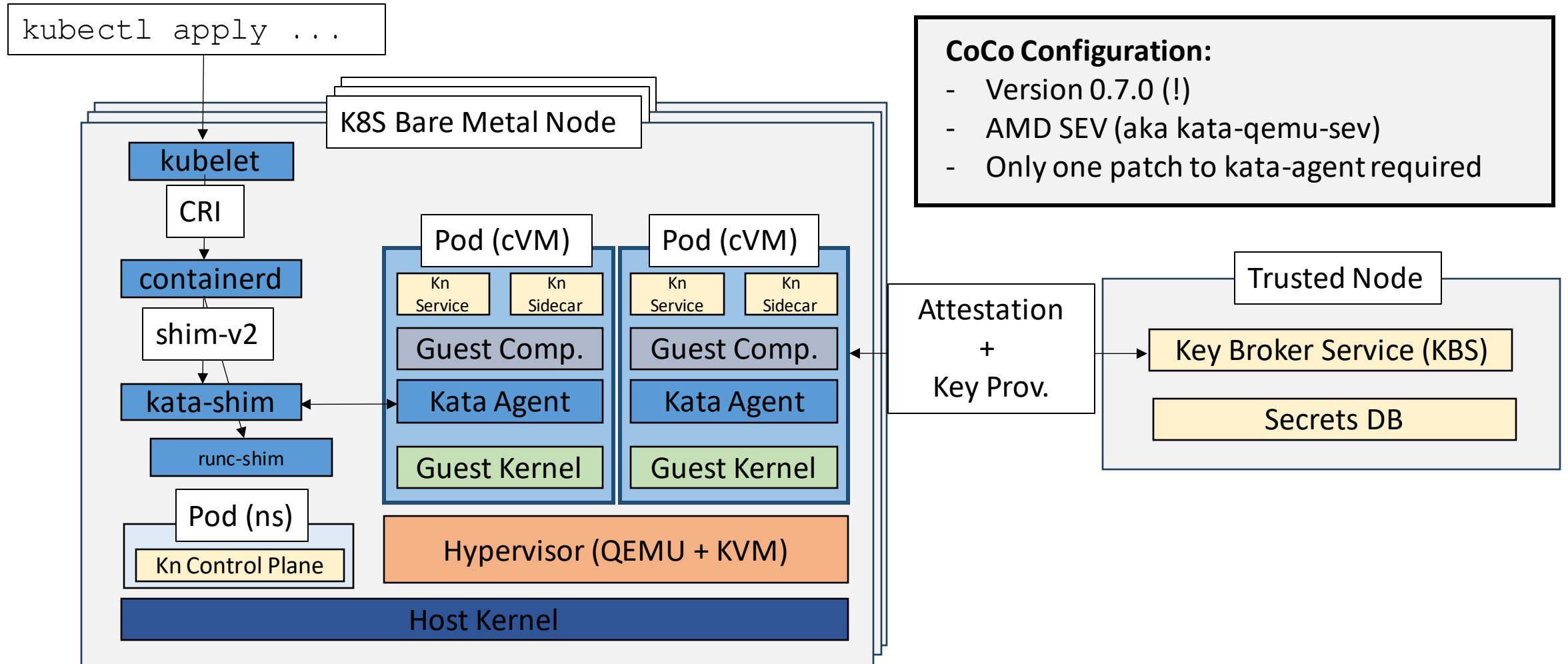
We need isolation from the host environment to guarantee...

- Data Confidentiality
- Code Confidentiality
- Execution Integrity

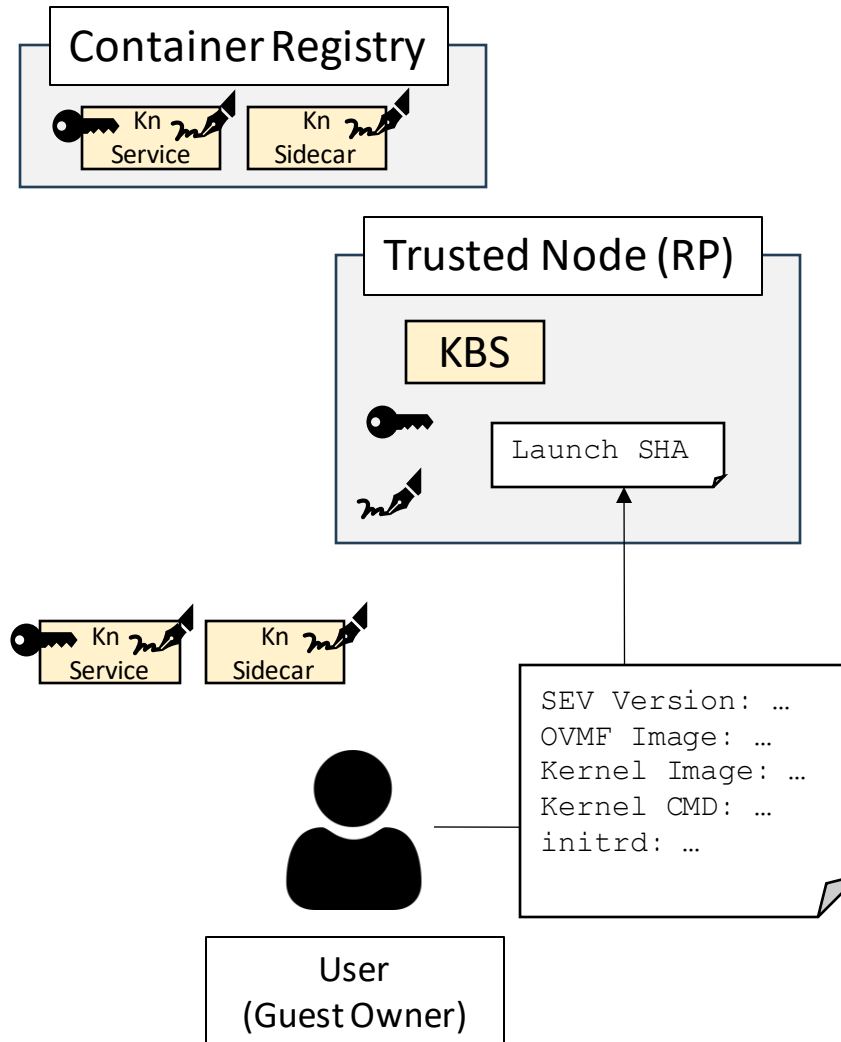
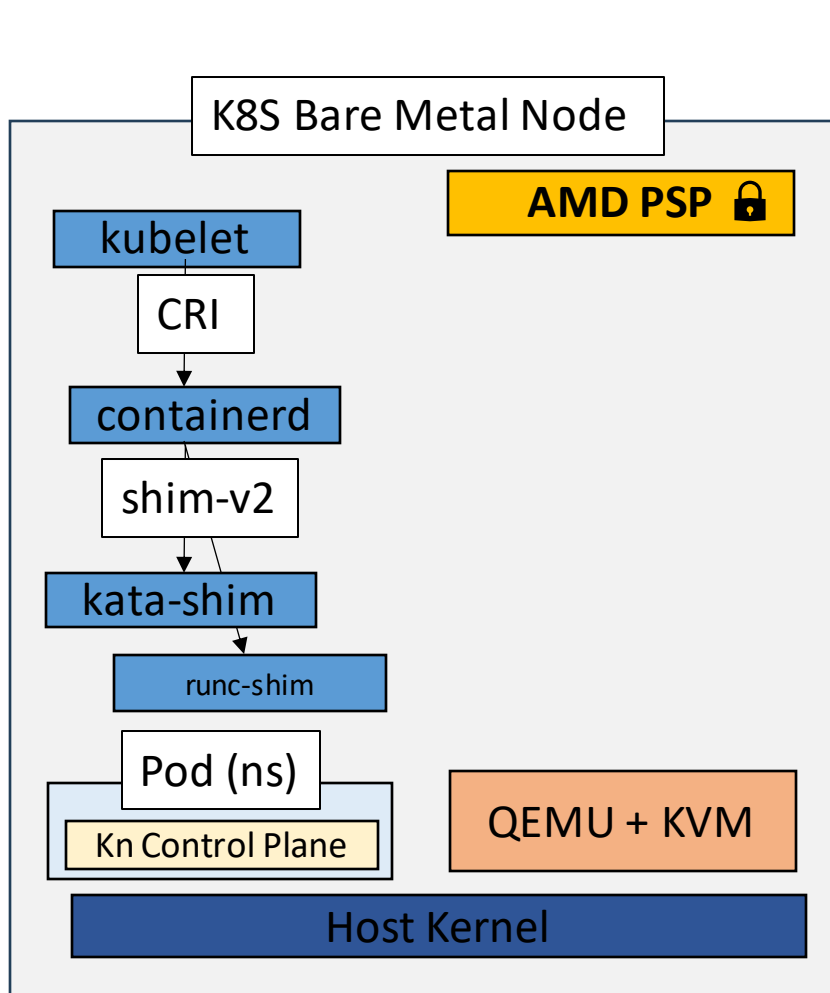


**Confidential Computing**

# PoC: Knative on Confidential Containers



# PoC: Attestation of Knative on CoCo (AMD SEV)

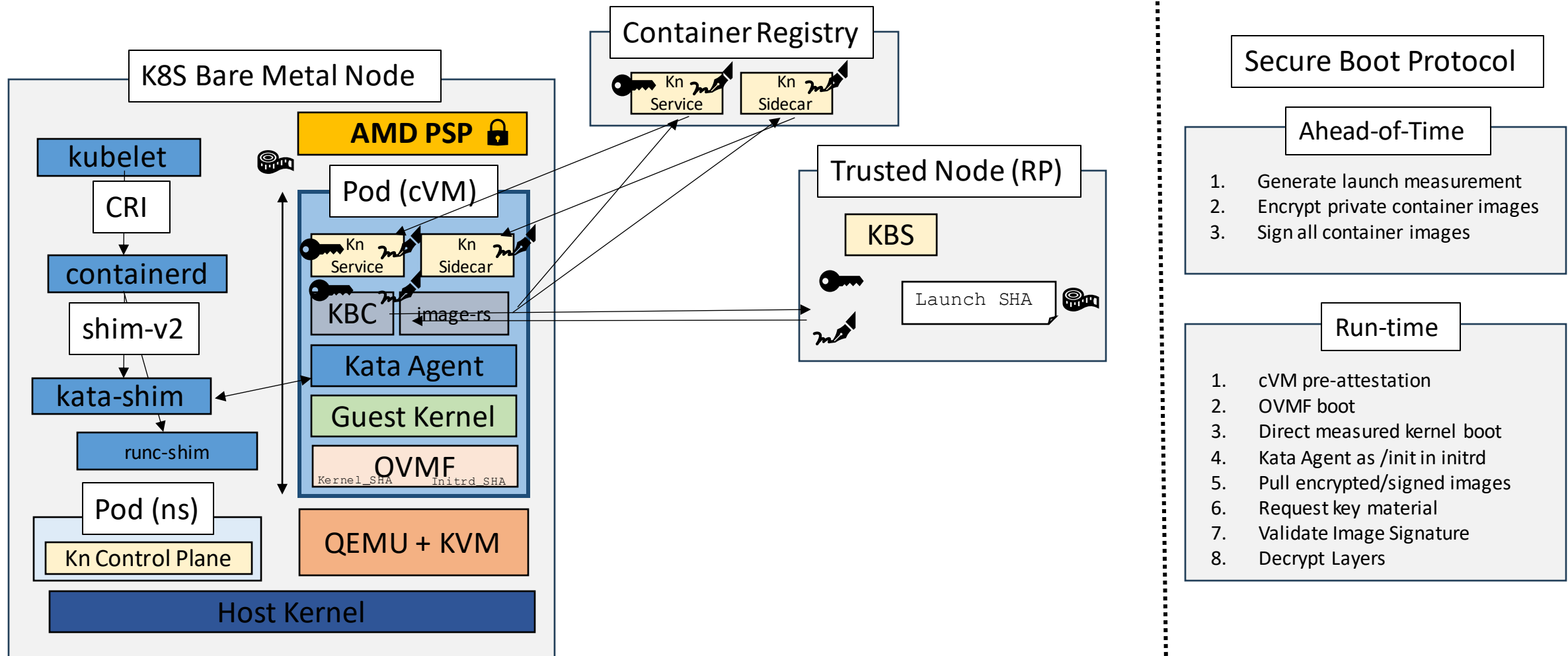


## Secure Boot Protocol

### Ahead-of-Time




1. Generate launch measurement
2. Encrypt private container images
3. Sign all container images

# PoC: Attestation of Knative on CoCo (AMD SEV)

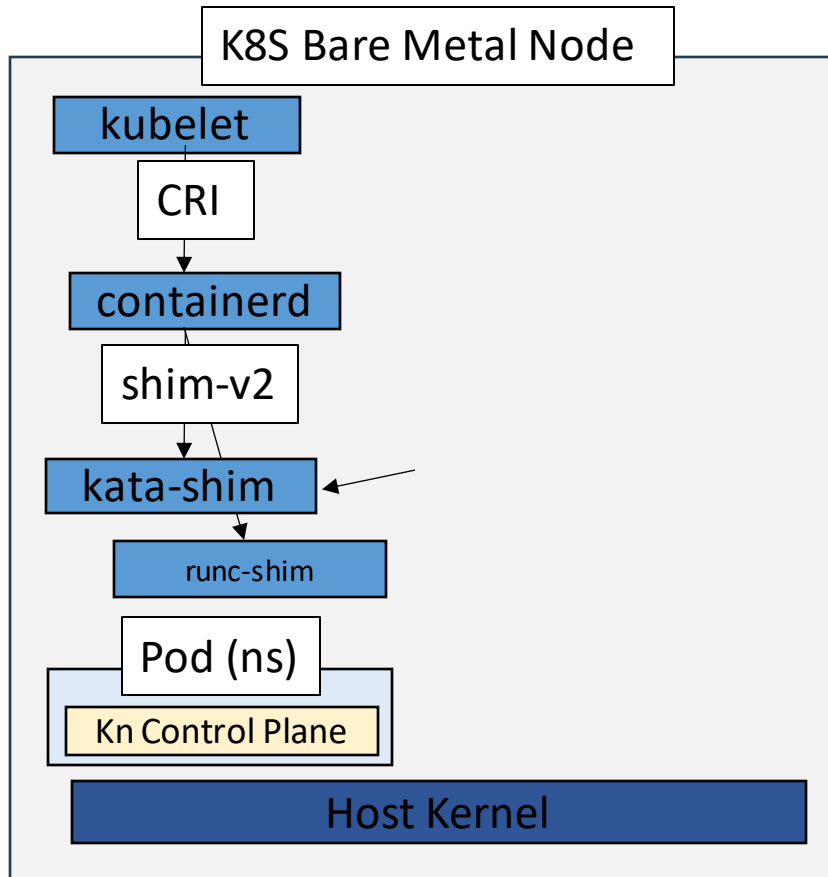


# Evaluation

We want to evaluate the feasibility of our PoC according to the three key metrics we identified for serverless:

	1. Cold Start Times	2. Warm Start Times	3. Instantiation Throughput
 <b>K8S</b> <b>RUNC</b>	6s	1s	1 fps
 <b>kata</b> containers	7s	2s	0.5 fps
 <b>CONFIDENTIAL</b> <b>CONTAINERS</b>	??	??	??

# Evaluation: Baselines

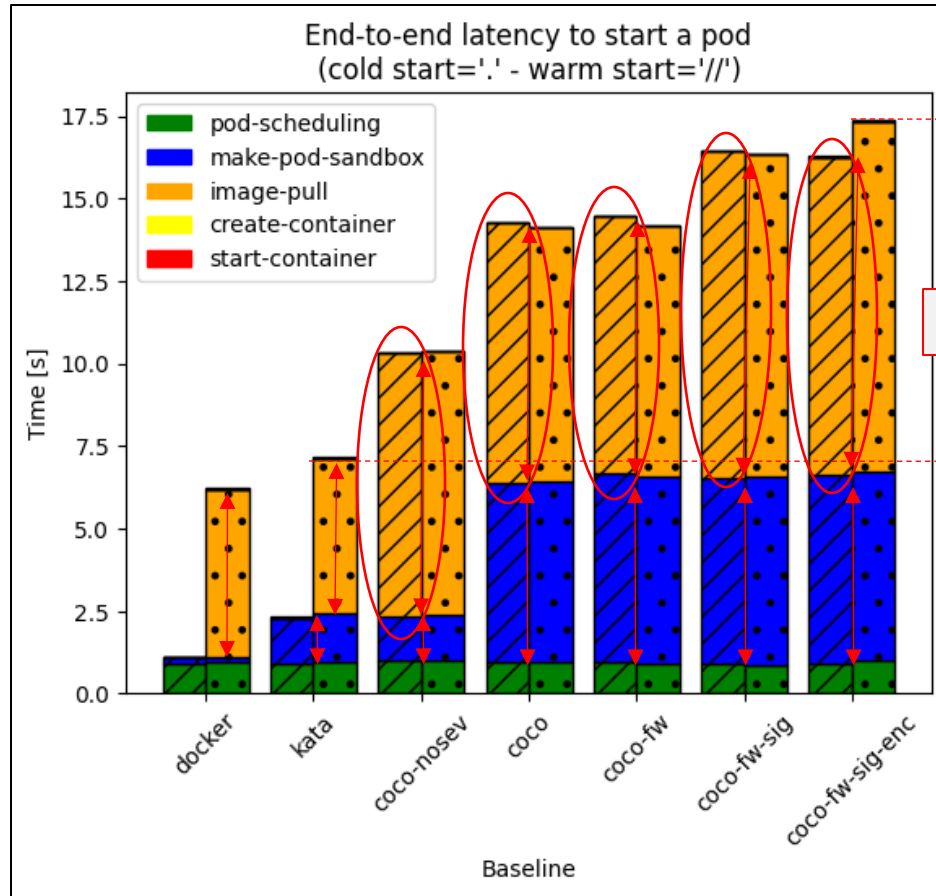


0. **docker (i.e. runc):** no VMs
1. **kata:** VMs
2. **coco-nosev:** + pull in guest
3. **coco-nosev-ovmf:** + OVMF
4. **coco:** + SEV
5. **coco-fw:** + HW att
6. **coco-fw-sig:** + image signature
7. **coco-fw-sig-enc:** + image enc.

Knative Service is a simple "Hello World" in Python

Code, patches, and evaluation scripts are available and  
(hopefully) reproducible:  
<https://github.com/csegarragonz/coco-serverless>

# Evaluation: Cold/Warm Starts



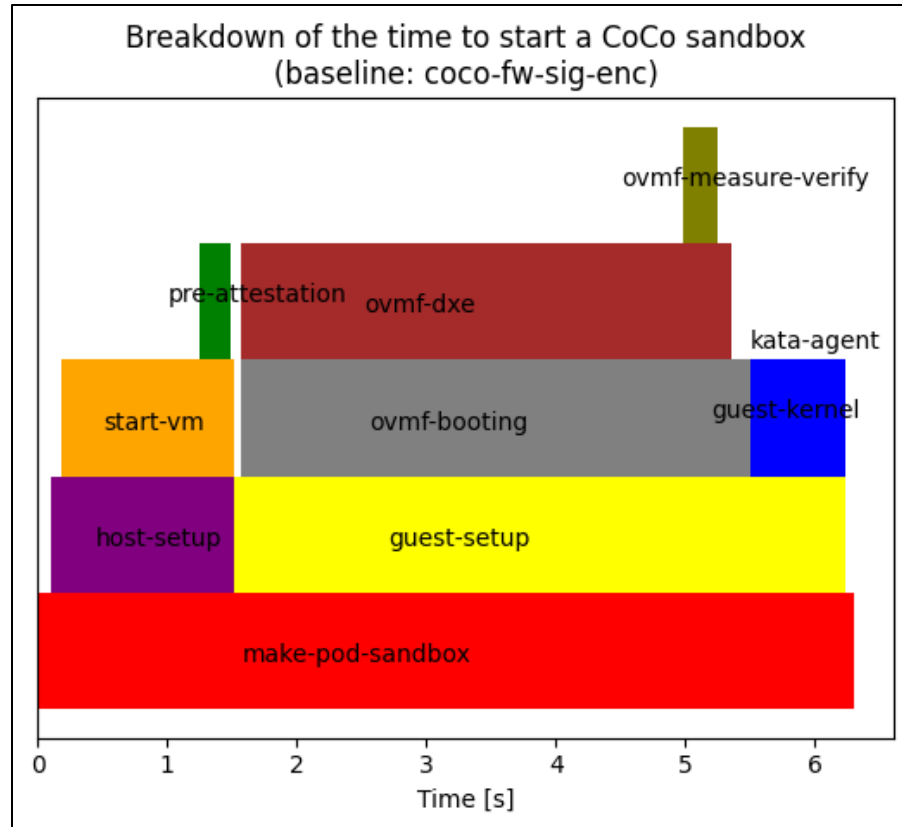
## Observations:

1. Why is VM start-up 4x slower with SEV?
2. Why is image pulling 2-3x slower w.r.t docker?
3. Why are there no warm starts?

Additional 10 s!

Let us address these questions one by one

# Evaluation: VM Start-Up in detail



Q1: Why is VM start-up 3x slower with SEV?

A: It seems that we spend a lot of time in OVMF...

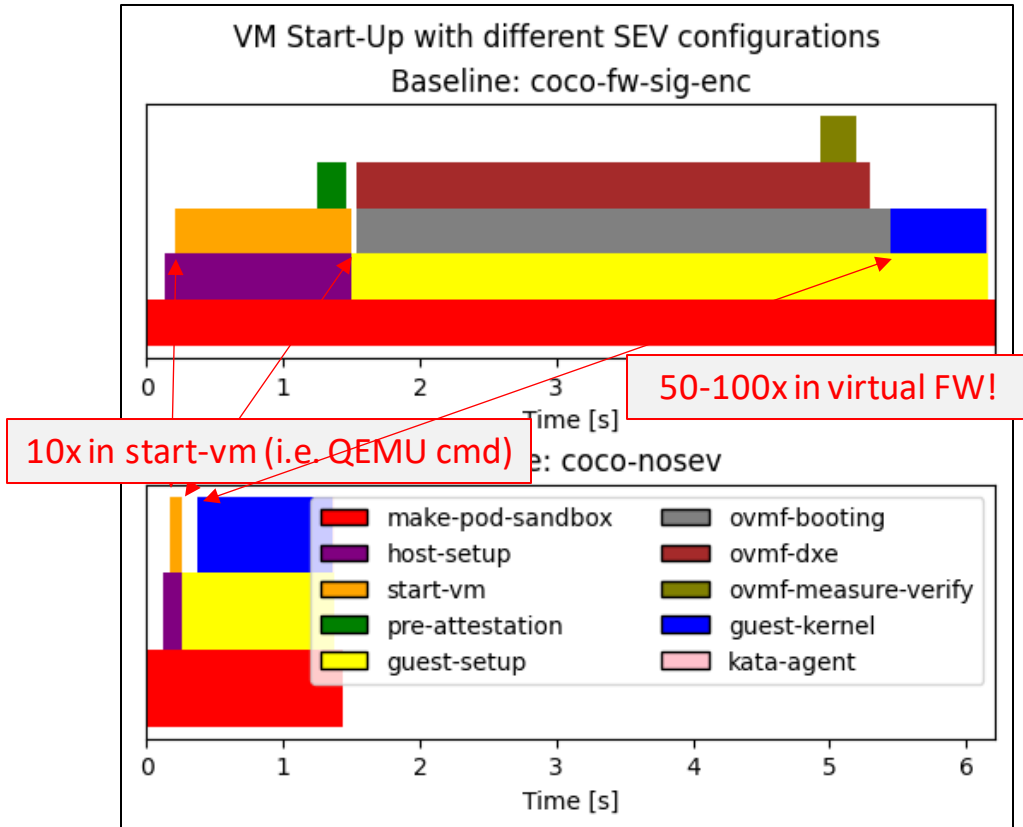


# Evaluation: VM Start-Up in detail

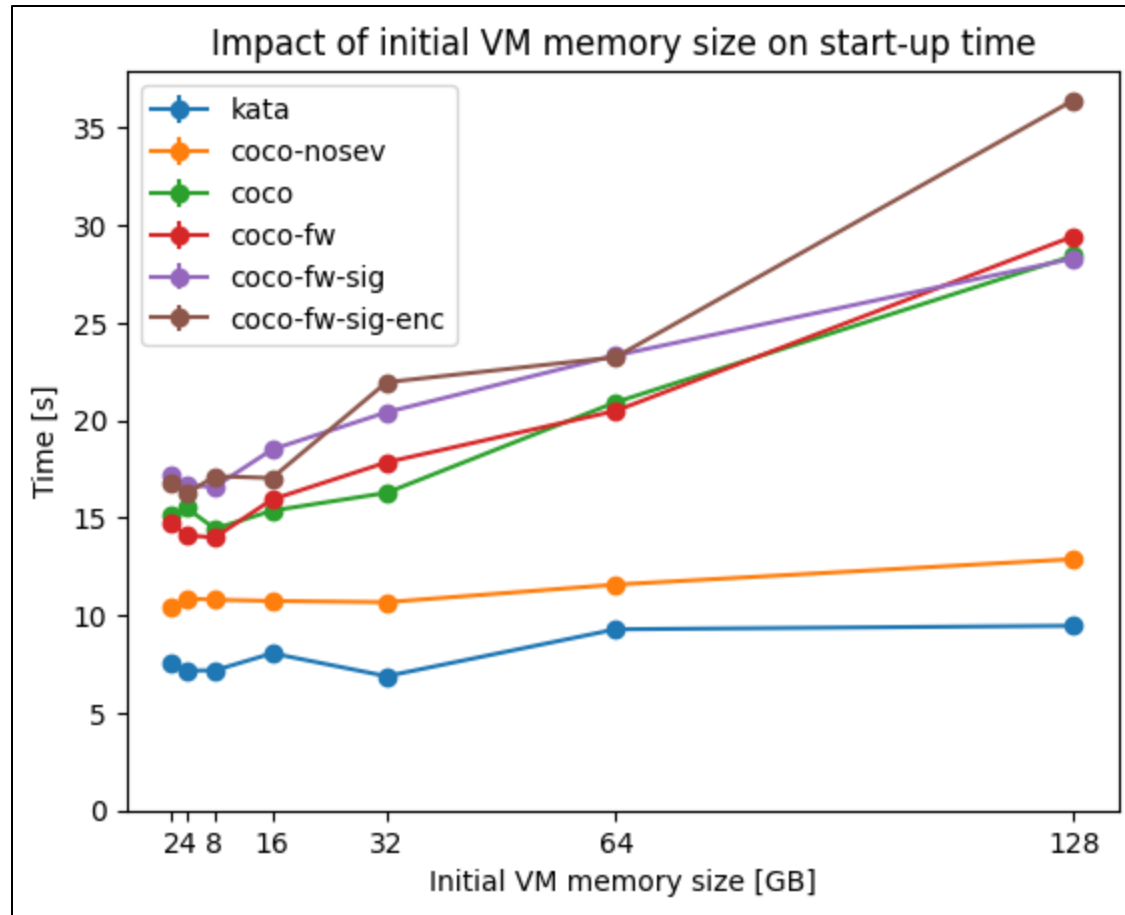
Q1: Why is VM start-up 3x slower with SEV?

A: It seems that we spend a lot of time in OVMF...

A: Compared to a non-SEV VM (w/ SeaBIOS) we spend:



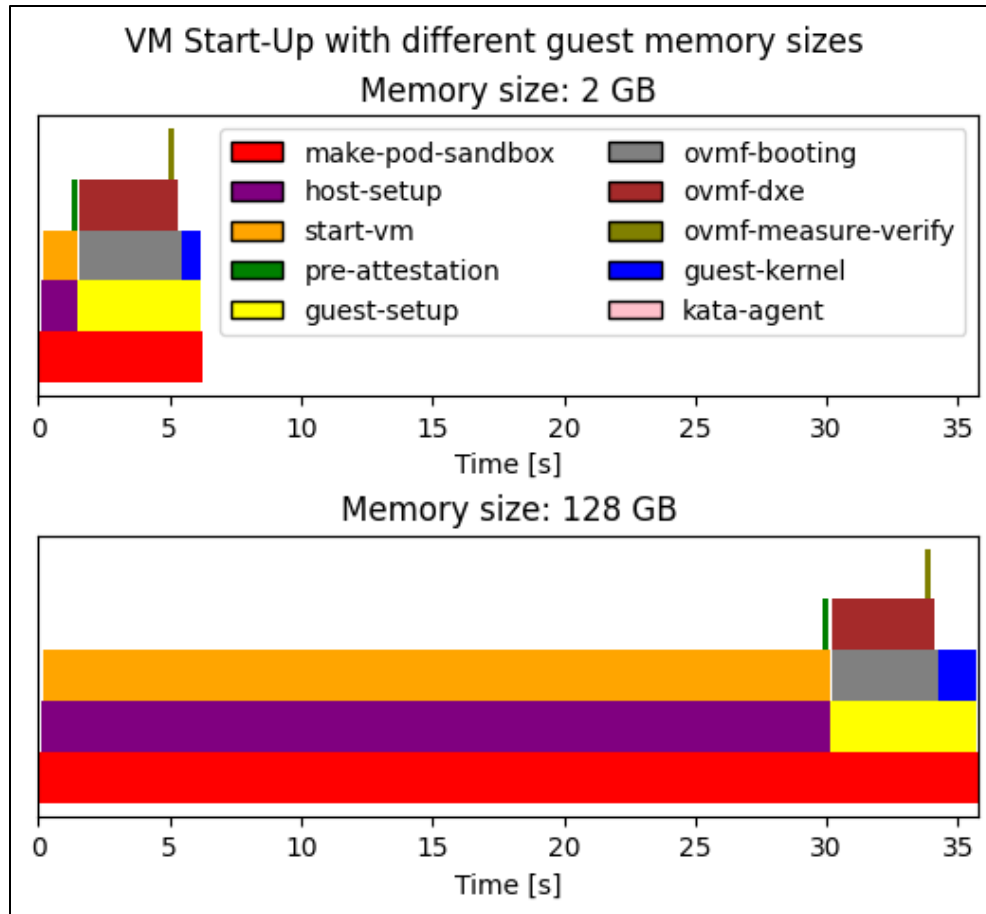
# Evaluation: VM Start-Up in detail (ctd)



**Q1: Why is VM start-up 3x slower with SEV?**

**A: During the start-vm phase, QEMU provisions all the memory pages assigned to the guest**

# Evaluation: VM Start-Up in detail (ctd)



**Q1: Why is VM start-up 3x slower with SEV?**

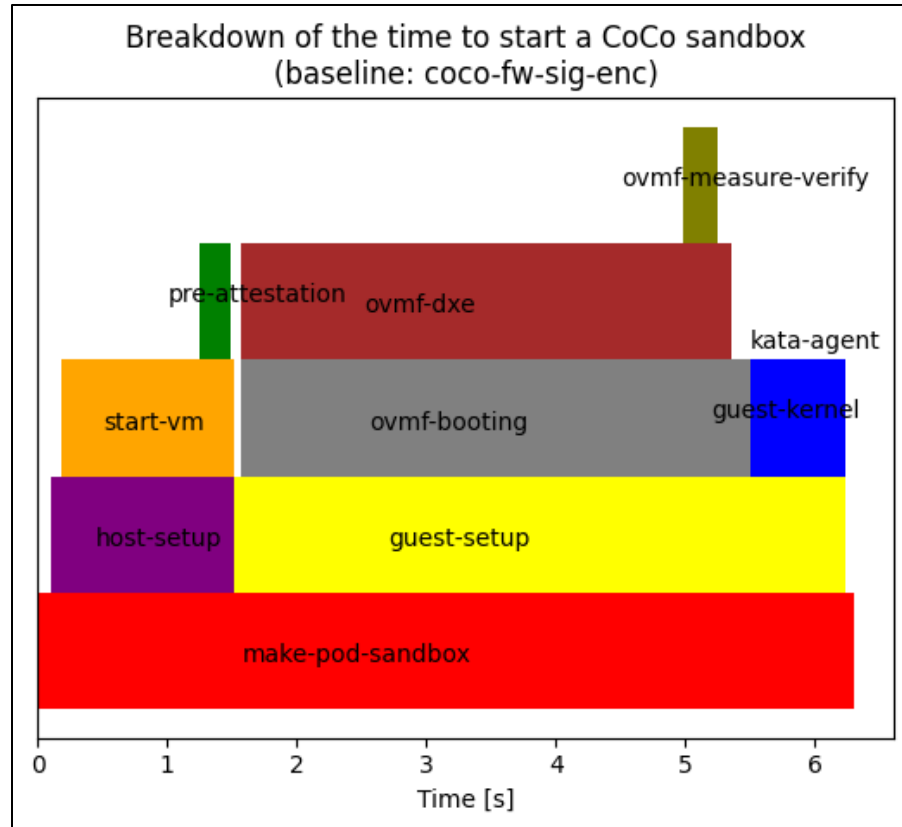
**A: During the start-vm phase, the PSP provisions all the memory pages assigned to the guest**

**Suggested Solution:**

- Can we assign memory pages lazily, off the hot-path?

**Serverless CoCo Task 1: Optimize cVM provisioning**

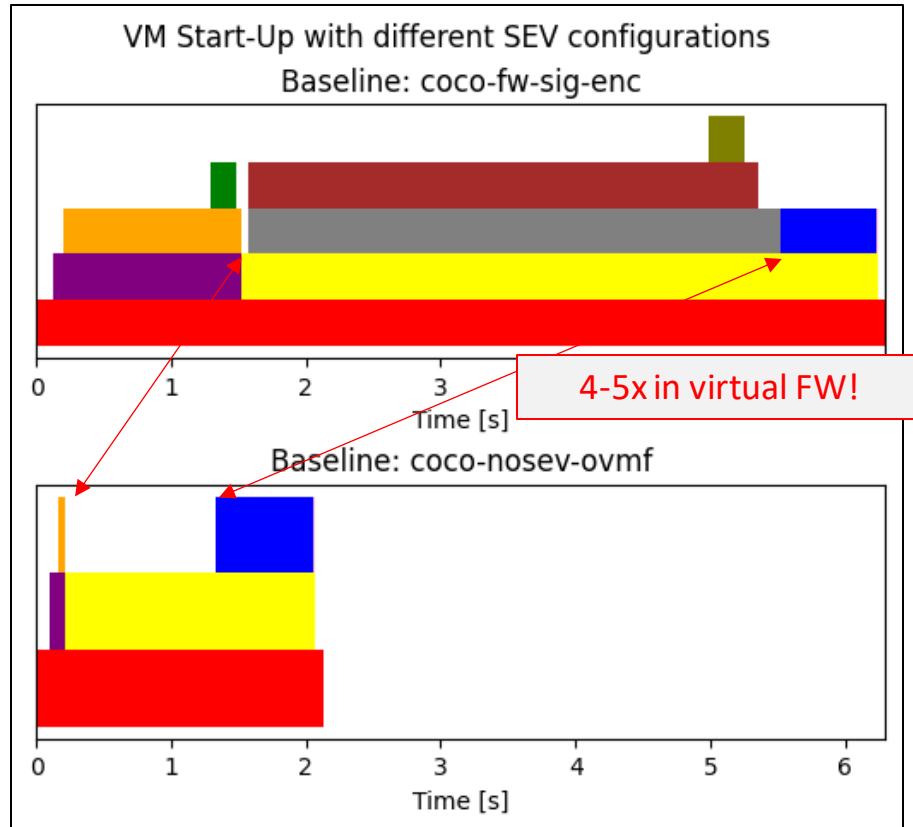
# Evaluation: VM Start-Up in detail



Q1: Why is VM start-up 3x slower with SEV?

A: It seems that we spend a lot of time in OVMF...

# Evaluation: VM Start-Up in detail



**Q1: Why is VM start-up 3x slower with SEV?**

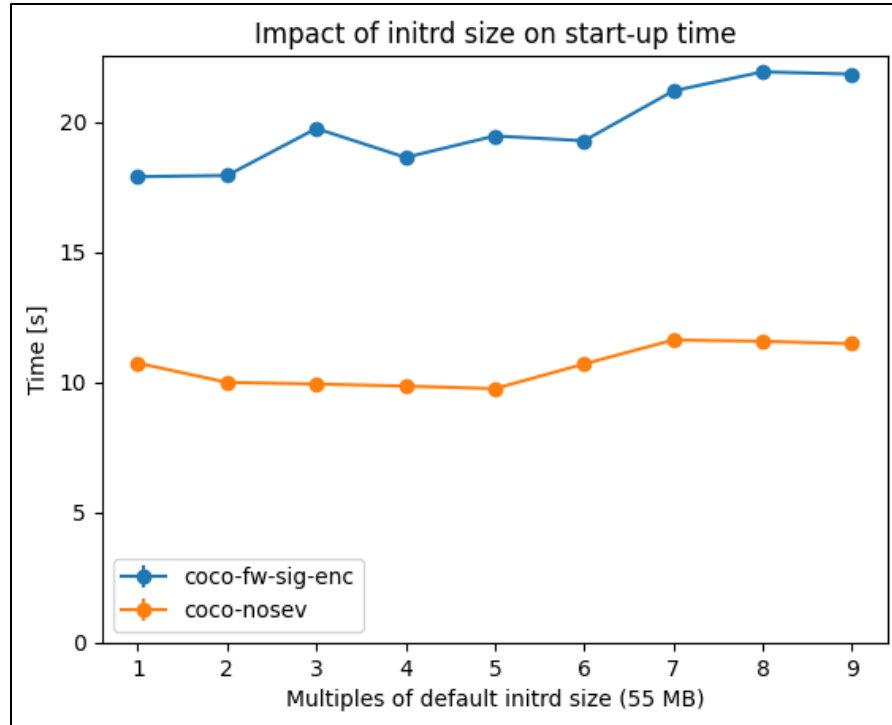
**A: It seems that we spend a lot of time in OVMF...**

**A: Compared to a non-SEV VM (w/ OVMF) we spend:**

**Q: What is the difference between SEV/non-SEV OVMF?**

**A: For SEV, we measure and verify kernel/initrd/cmdline**

# Evaluation: VM Start-Up in detail



**Q1: Why is VM start-up 3x slower with SEV?**

**A: It seems that we spend a lot of time in OVMF...**

**A: Compared to a non-SEV VM (w/ OVMF) we spend:**

**Q: What is the difference between SEV/non-SEV OVMF?**

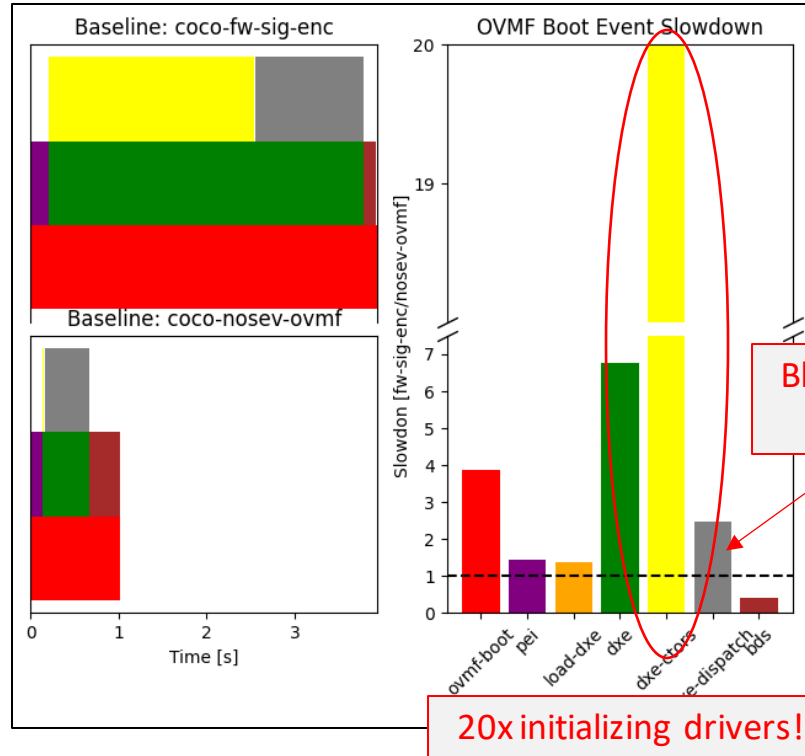
**A: For SEV, we measure and verify kernel/initrd/cmdline**

# Evaluation: VM Start-Up in detail

Q1: Why is VM start-up 3x slower with SEV?

A: It seems that we spend a lot of time in OVMF...

A: Compared to a non-SEV VM (w/ OVMF) we spend:

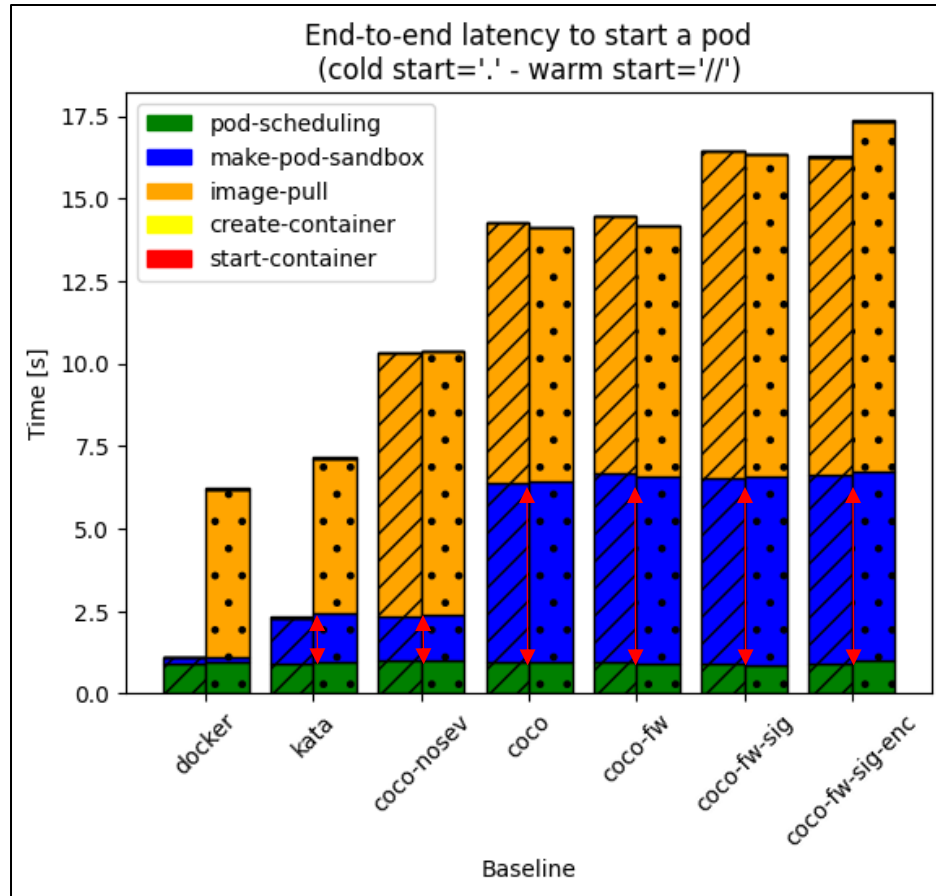


Blob measurement/verification happens here

This behaviour is unexpected

20x initializing drivers!

# Evaluation: Cold/Warm Starts



## Observations:

1. Why is VM start-up 4x slower with SEV?
2. Why is image pulling 2-3x slower w.r.t docker?
3. Why are there no warm starts?

**Problem:** Provisioning guest memory pages introduces 1-2 extra seconds (for 2GB of memory)

**Solution:** Hot-Plug guest memory pages (or provision off the hot path)

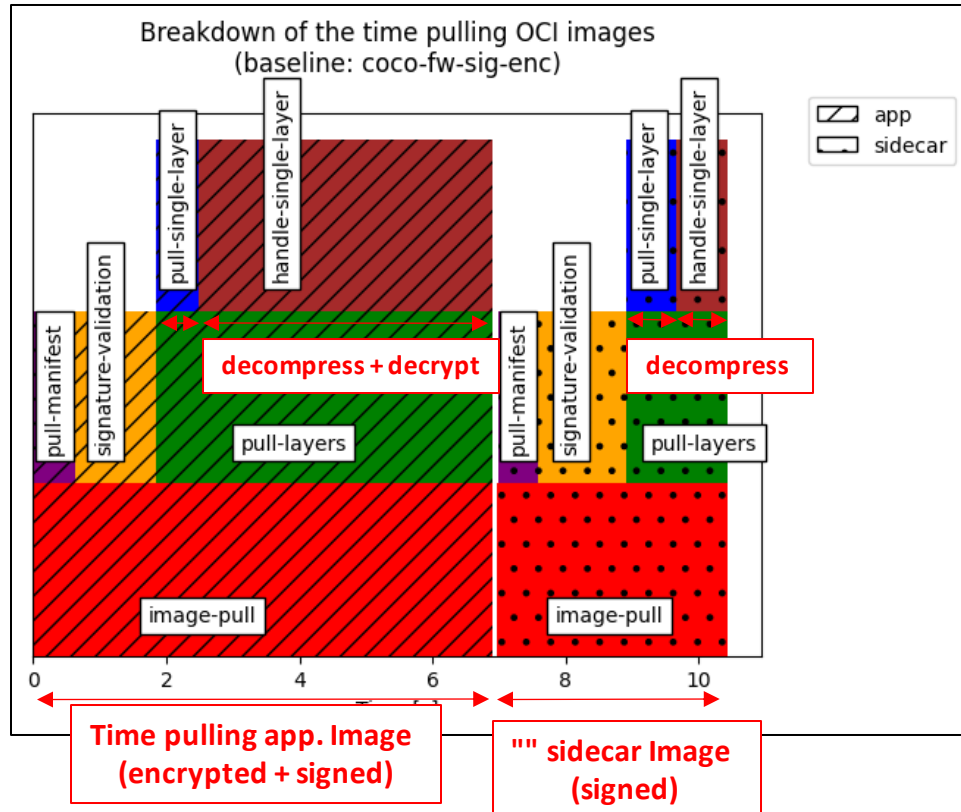
**Problem:** OVMF DXE driver initialization introduces 3-4 extra seconds

**Solution:** Not clear! Any ideas?



# Evaluation: VM Start-Up in detail

Q2: Why is image-pulling 2x slower w.r.t Docker?



A: containerd's PullImage becomes blocking!

A(ctd): Decrypting image layers is the bottleneck!

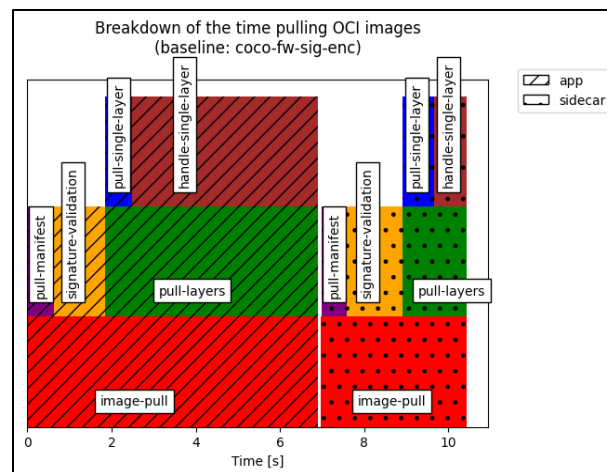
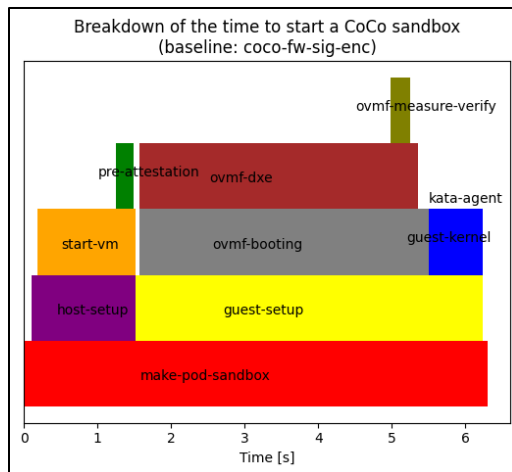
Suggested Solutions:

- Parallelize image pulling within pod
- HW-accelerated image decryption

Serverless CoCo Task 2: Optimize Image Pulling Time

# Evaluation: VM Start-Up in detail

Q3: Why are there no warm starts?



**Serverless CoCo Task 3: Design Secure CoCo sandbox re-use strategies**

**A: SEV guests are cryptographically bound to one "guest owner"**

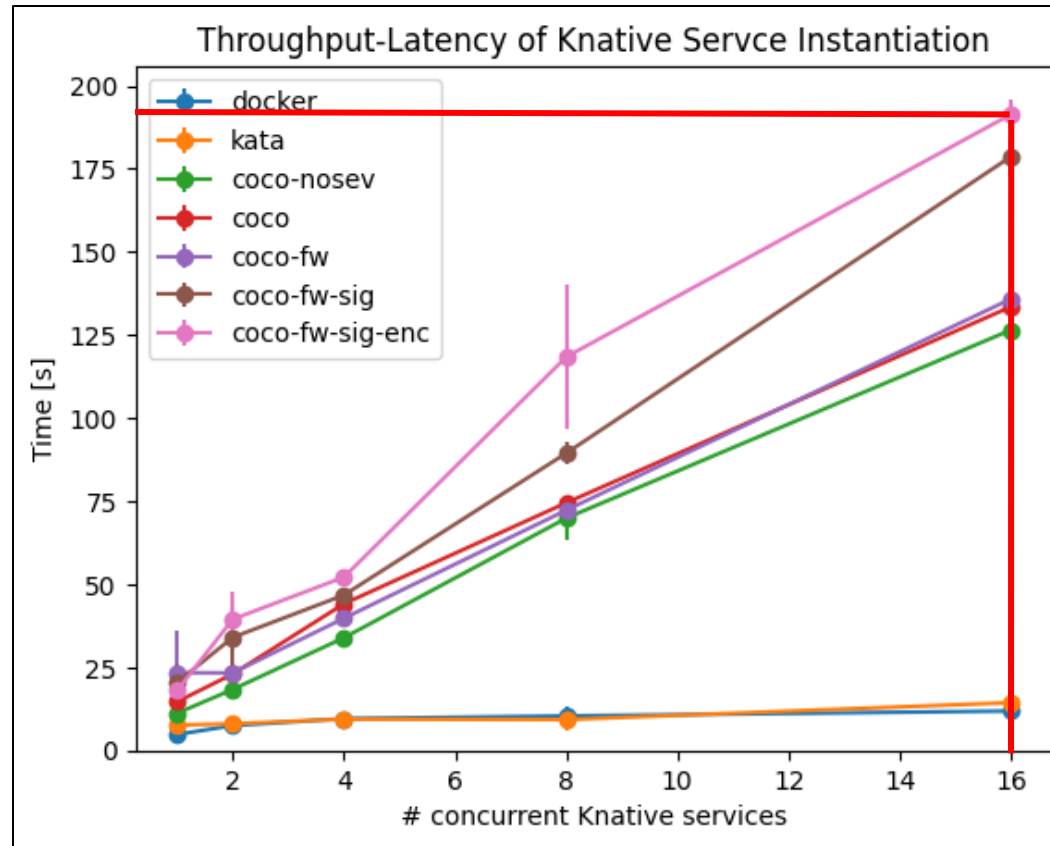
**A: Cannot rely on the host to mount container images**

**A: Cannot easily share (or lazy load) encrypted image layers**

## Suggested Solutions:

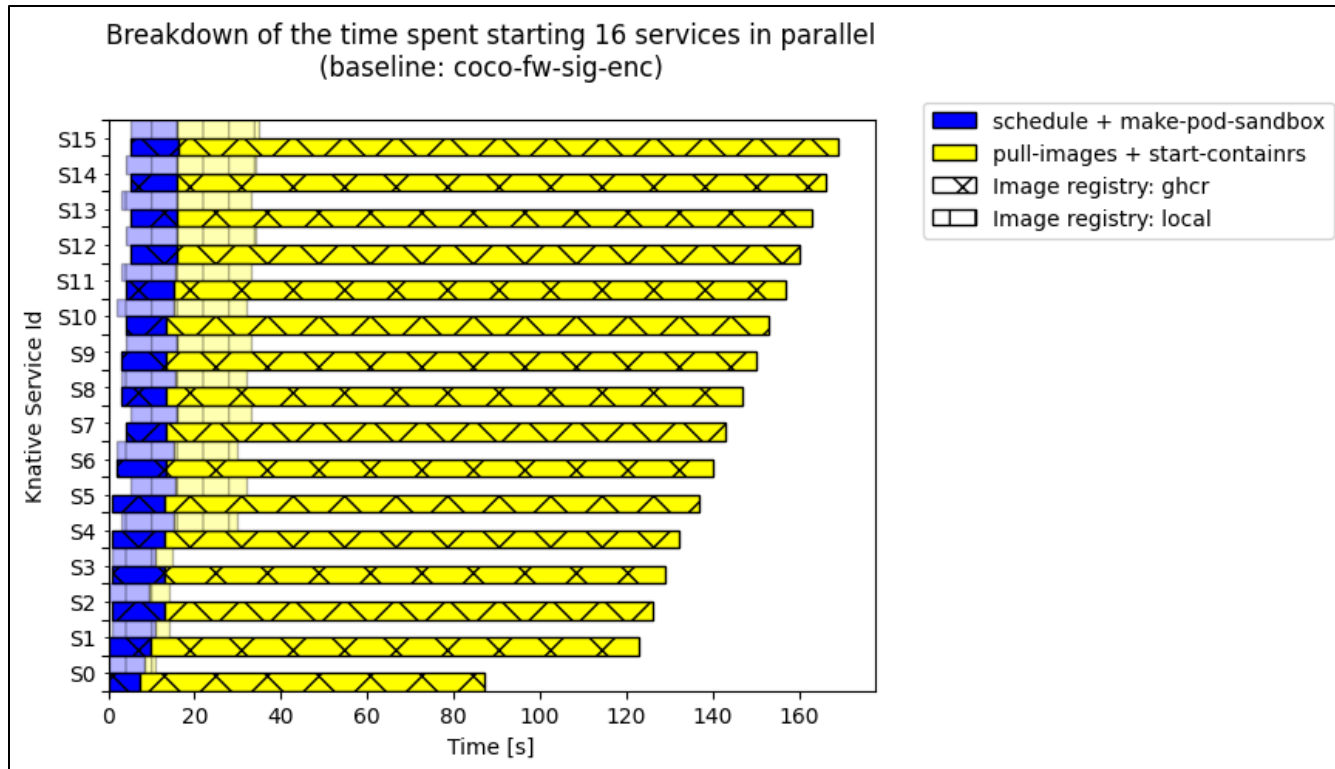
- Use the KBS as trusted relying-party in VM pre-warm
- Freeze the Kata Agent until pre-warmed VM is assigned
- Encrypted block-based lazy image loading (Nydus)
- Label image layers as encrypted or not

# Evaluation: Instantiation Throughput



Starting 16 concurrent functions takes > 3' !!

# Evaluation: Instantiation Throughput (ctd.)



**Q: Why Starting 16 concurrent functions takes > 3'?**

**A: We are being throttled by the registry!**




## Suggested Solutions:

- Node-local (and/or cluster-local) layer/block cache
- Actual image re-use (when possible!)

**Serverless CoCo Task 4: Improve Scalability of CoCo sandbox provisioning**

# Evaluation

We want to evaluate the feasibility of our PoC according to the three key metrics we identified for serverless:

	1. Cold Start Times	2. Warm Start Times	3. Instantiation Throughput
	6s	1s	1 fps
	7s	2s	0.5 fps
	17.5 s	17.5 s	~ 0.1 cps

# FYP CoCo: Summary

Slowdown



1. Cold Start Times

2.5x

## cVM Start-Up Overhead

- 1) Guest Memory Pages Provisioning
- 2) OVMF DXE Initialization

## Guest-Side Image Pulling

- 1) Serial (per-ctr) pulling
- 2) Image Layer Decryption

2. Warm Start Times

8.75x

## No CoCo pre-warming

- 1) SEV Guests <-> owner

## No Image Re-Use

- 1) Cannot mount images from host
- 2) Cannot share images between tenants
- 3) Cannot lazy load images

3. Instantiation Throughput

5x

## Registry Throttling

- 1) If all CoCo's pull from the guest, cannot scale w/out pass-through cache
- 2) Will benefit from improvements in warm starts



---

# Serverless Confidential Containers: Challenges and Opportunities

---

**Carlos Segarra**

(w/ Tobin Feldman-Fitzthum and Daniele Buono)

**Large-Scale Data & Systems (LSDS) Group - Imperial College London**

**Visiting IBM TJ Watson (Sep'23 – Nov'23)**

<https://carlosgarra.com>

[<cs1620@ic.ac.uk>](mailto:cs1620@ic.ac.uk)

