

Measurement and Attestation Schemes for Container Sandboxes

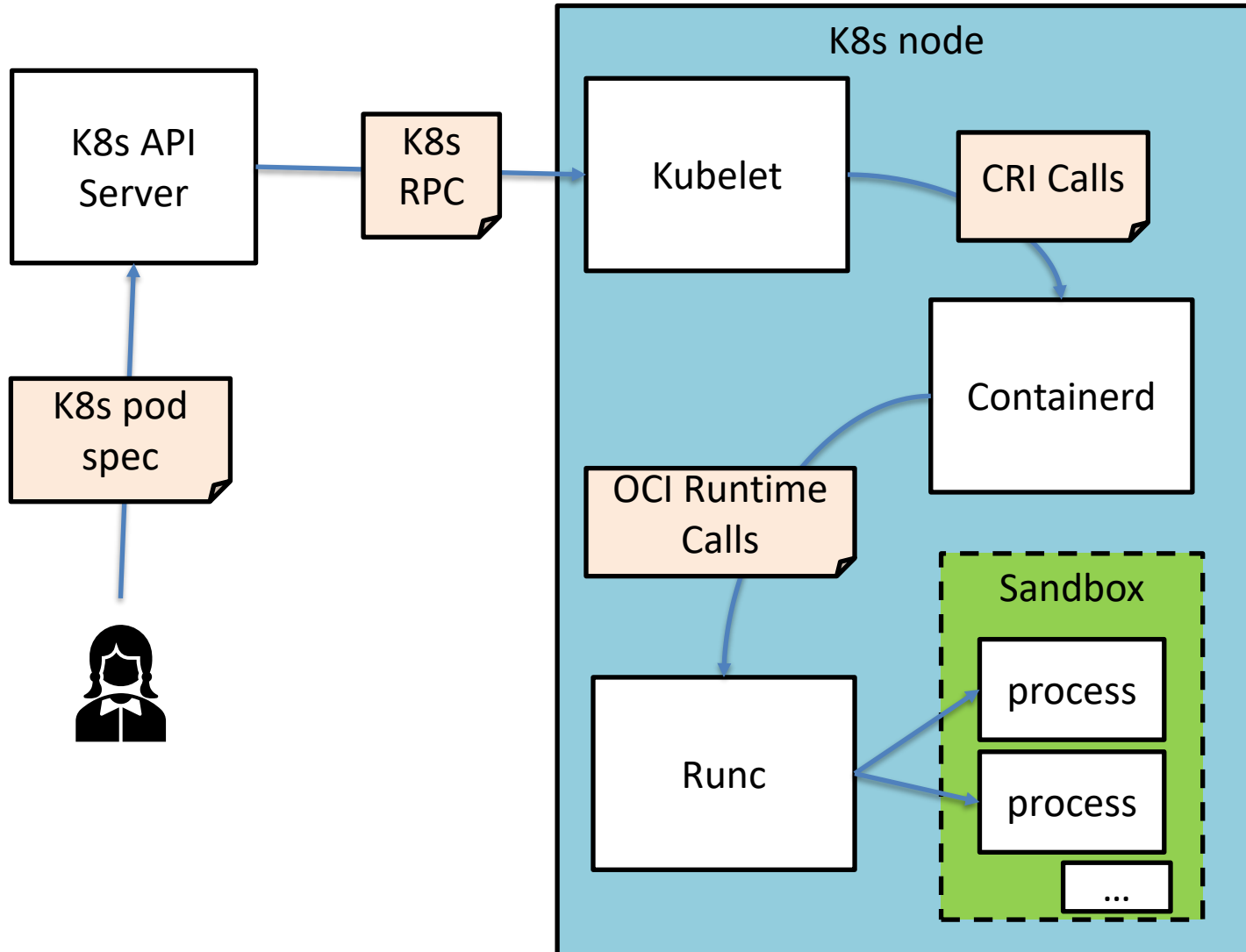
Magnus Kulke, swe @Azure Core Linux

FOSDEM 25, Attestation Devroom

Context

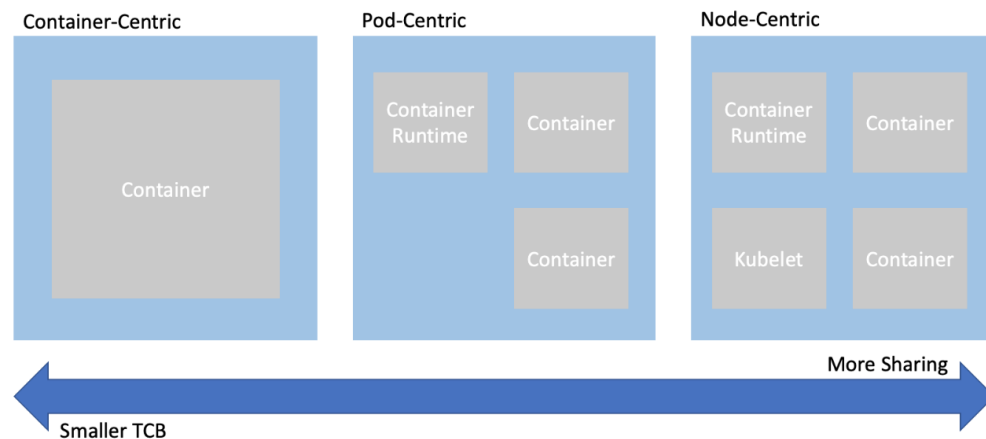
- Confidential Containers (CoCo)
 - CNCF project
 - vendor neutral
 - Facilitate confidential computing in the container ecosystem
- Confidential Computing is (mostly) a VM technology
 - Containers usually do not run in VMs
 - Kata Containers adds virtualization as isolation layer

What is a typical container launch?

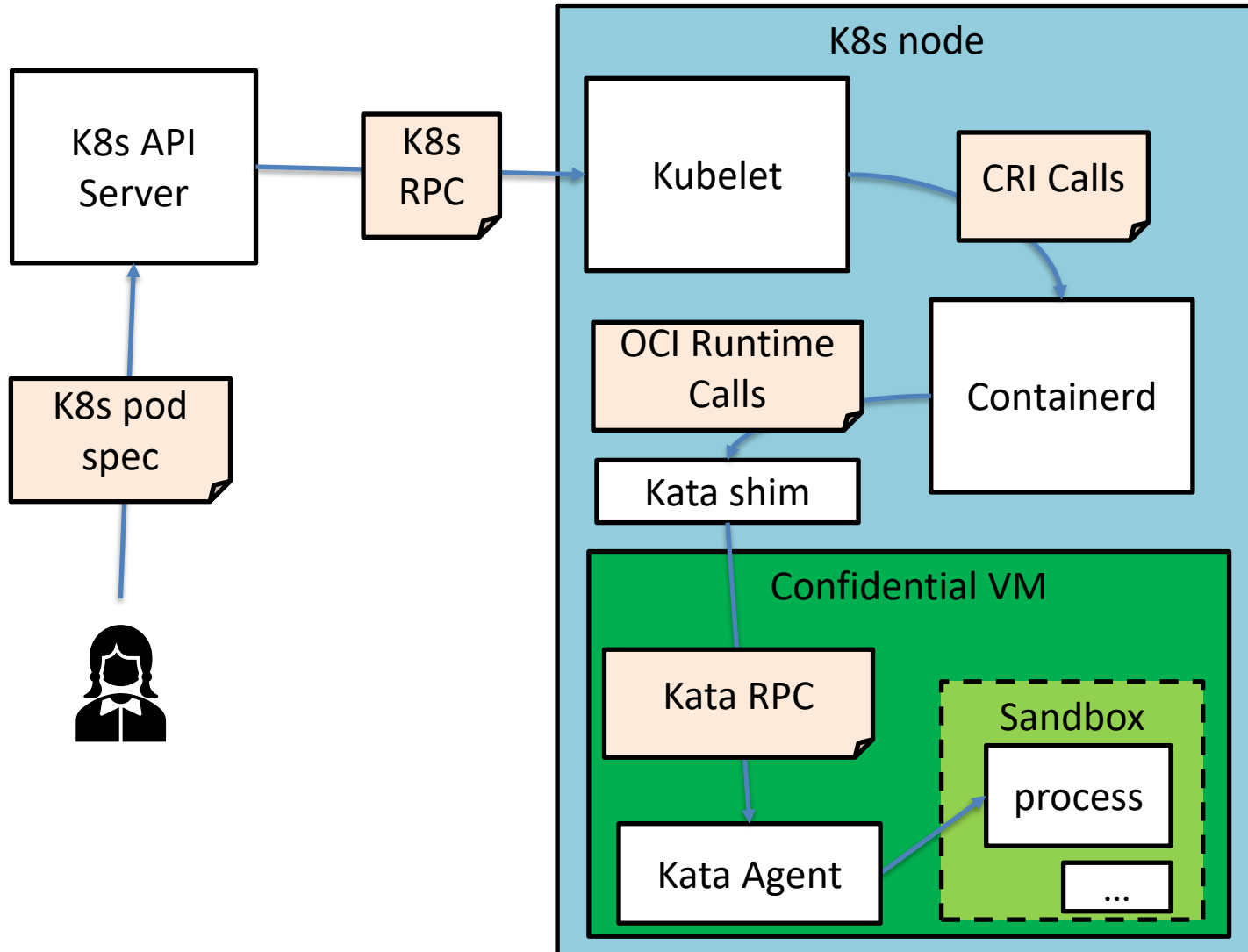


Sandbox?

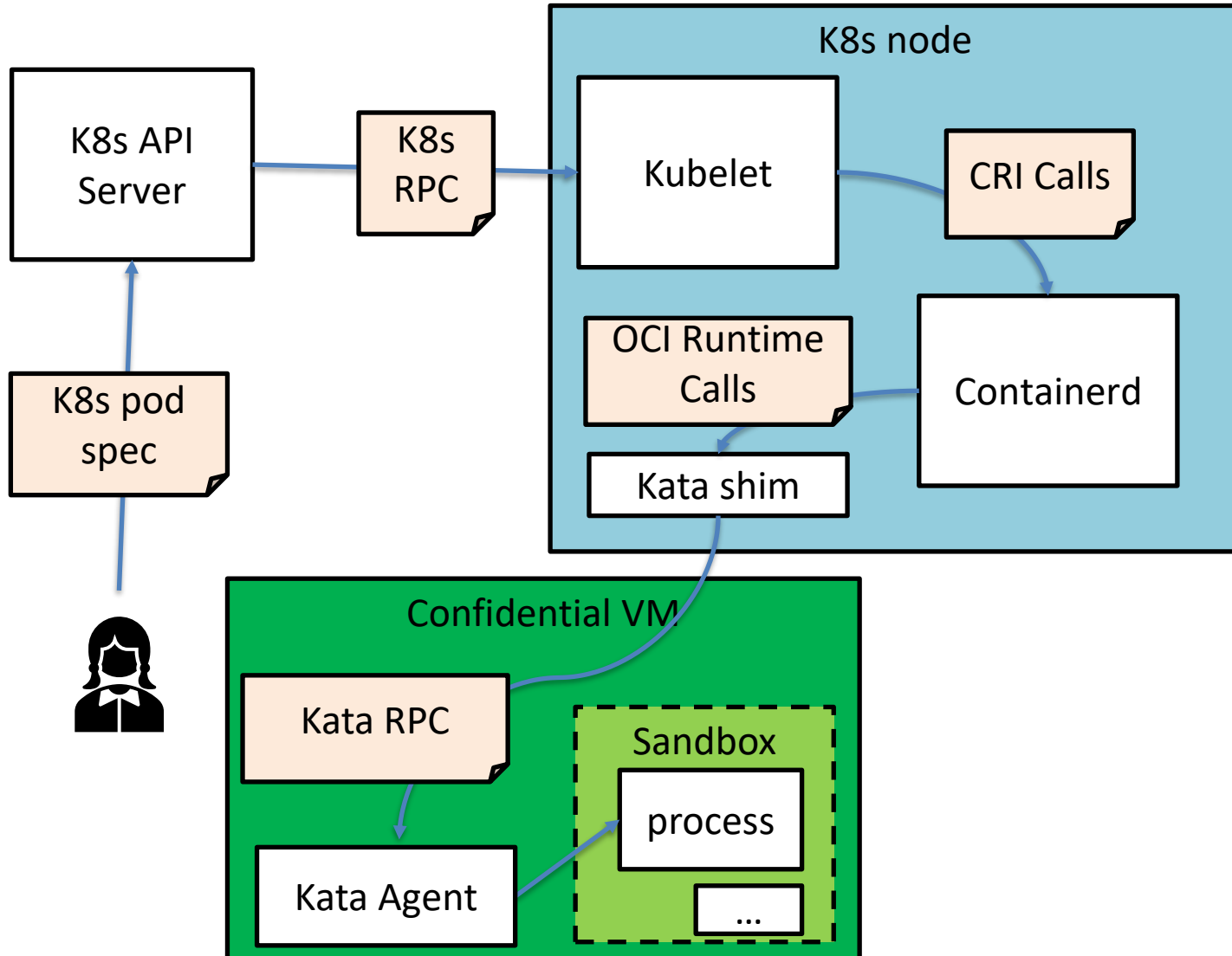
- Pod: Kubernetes deployment „atom“.
- Set of collocated processes (containers) that share namespaces and resources
- Good abstraction to introduce confidentiality boundaries



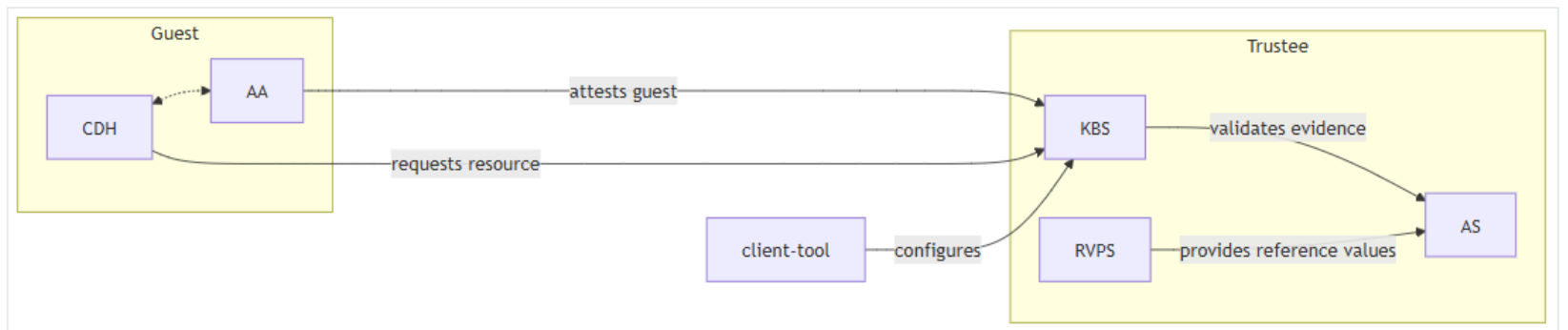
Confidential container launch



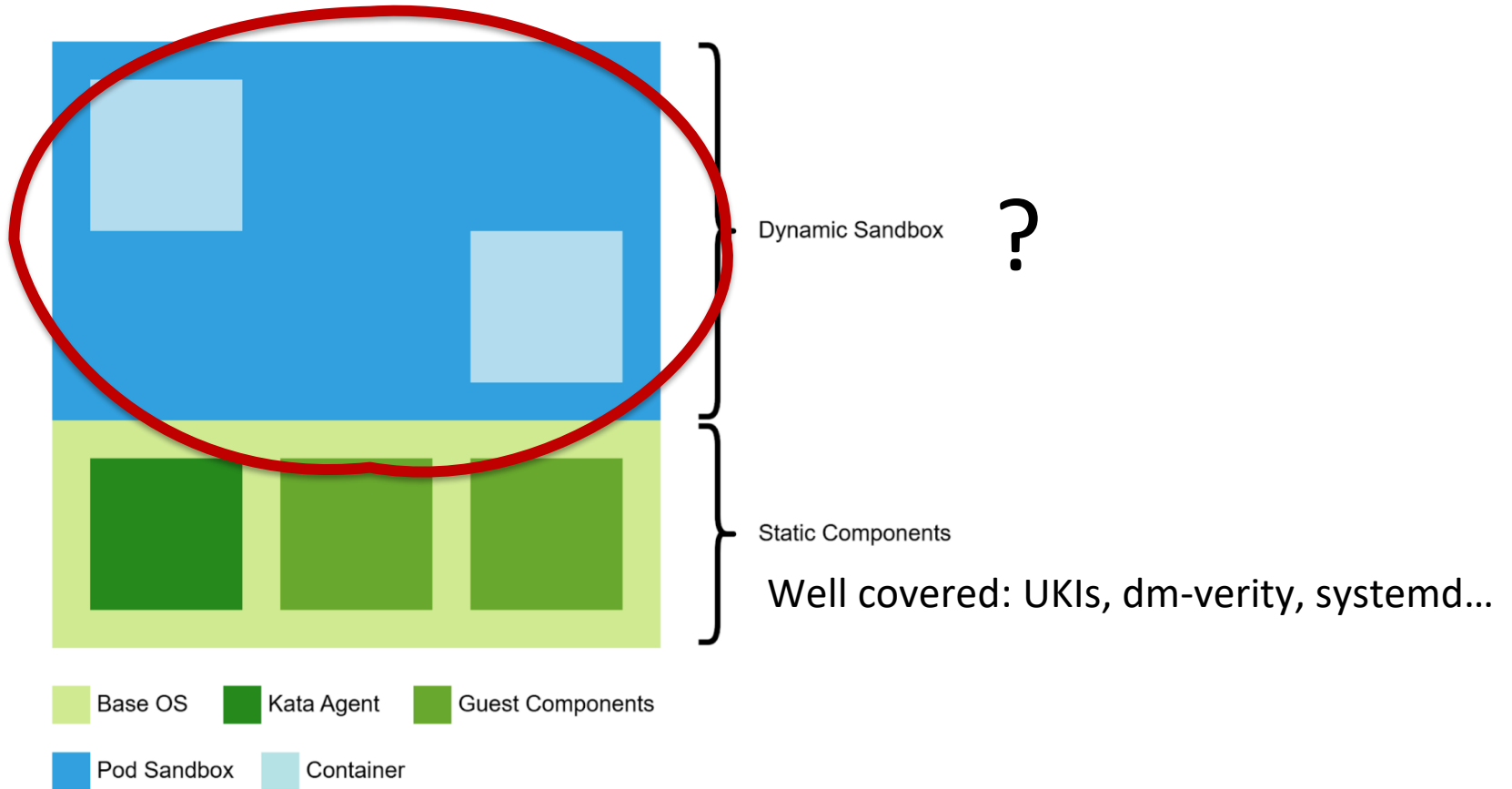
... or with a remote CVM



Attestation Architecture



Static and dynamic components



OCI images are content addressable

```
wslhost.exe
$ oras manifest fetch ghcr.io/mkulke/nginx-encrypted@sha256:5a81641ff9363a63c3f0a1417d29b527ff6e155206a720239360cc6c0722696e > manifest.json
$ jq '.layers[0].digest' < manifest.json
"sha256:d5e2d29403b03b4e74953d6bab263777d753316e6c32ff7d9fe4efa4eaba9e53"
$

$ sha256sum -b manifest.json
5a81641ff9363a63c3f0a1417d29b527ff6e155206a720239360cc6c0722696e *manifest.json
$
```

[1] 0: bash* "magnuskulke@DESKTOP-2" 19:45 30-Jan-25

Sandbox with imperative control

```
wslhost.exe
t/agent_log_rpc.txt t/nginx.yaml buffers
5 ttRPC server started
4 rpc call from shim to agent: "create_sandbox"
3 rpc call from shim to agent: "get_guest_details"
2 rpc call from shim to agent: "copy_file"
1 rpc call from shim to agent: "create_container"
0 rpc call from shim to agent: "start_container"
1 rpc call from shim to agent: "wait_process"
2 rpc call from shim to agent: "copy_file"
3 ...
4 rpc call from shim to agent: "create_container"
5 rpc call from shim to agent: "start_container"
6 rpc call from shim to agent: "stats_container"
7 ...
~
~
~
~
~
~
~
~
~
~
8 apiVersion: apps/v1
7 kind: Deployment
6 metadata:
5   name: nginx
4   namespace: default
3 spec:
2   selector:
1     matchLabels:
0       app: nginx
1   replicas: 1
2   template:
3     metadata:
4       labels:
5         app: nginx
6     spec:
7       runtimeClassName: kata-remote
8       containers:
9         - name: nginx
10           image: bitnami/nginx:1.14
11           ports:
12             - containerPort: 80
13           imagePullPolicy: Always
~
tmp/agent_log_rpc.txt 75W 46% ln :6/13 :46 N... ./tmp/nginx.yaml 40% ln :9/22 :14
```

Attesting a container environment

Objective:

Ensure that only intended operations are executed within the sandbox (before releasing a secret)

Requires:

A comprehensive measurement of the “container workload”

Challenges in dynamic environments

- Dynamic Nature of Pods:
 - Pods can have containers created, deleted, or updated imperatively.
 - Dynamisms make it challenging to guarantee integrity.
- Kubernetes Control Plane:
 - Can and will adjust a user's pod spec
 - Examples: env variables, admission controllers

Options

- Lock k8s control plane:
 - Allow only “trusted” (predictable) operations
 - K8s api surface is huge, increasing constantly
 - Requires lots of glue code and ceremony
- Effort underway: “split-api”

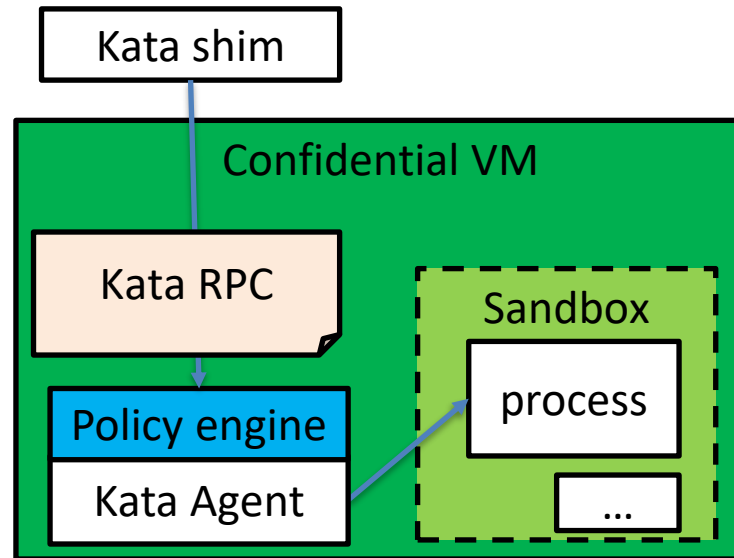
Options

- Keep a log
 - (Somewhat) like linux' IMA
 - Record Kata RPC + payloads into replayable log
 - Not all TEEs provide registers that can be extended at runtime
 - Some payloads are not predictable, b/c controlled by the env
 - Verification is not trivial

Options

- Policy in the TEE
 - Describe invariants (image digest)
 - Allow “acceptable” dynamism (env: SERVICE_*)
 - Reject Kata RPCs by default
 - cherrypick what’s required
- Currently implemented in Kata-Agent
 - Engine based on Rego (popular in container-land)
 - genpolicy tool to generate policy from a pod spec

Plugging policy eval into the workflow



Example Policy

```
1 package agent_policy
2
3 import future.keywords.in
4 import future.keywords.if
5 import future.keywords.every
6
7 default CopyFileRequest := true
8 default DestroySandboxRequest := true
9 default CreateSandboxRequest := true
10 default GuestDetailsRequest := true
11 default ReadStreamRequest := true
12 default RemoveContainerRequest := true
13 default SignalProcessRequest := true
14 default StartContainerRequest := true
15 default StatsContainerRequest := true
16 default WaitProcessRequest := true
17
18 default CreateContainerRequest := false
19 default ExecProcessRequest := false
20
21 CreateContainerRequest if {
22     every storage in input.storages {
23         some allowed_image in policy_data.allowed_images
24         storage.source == allowed_image
25     }
26 }
27
28 ExecProcessRequest if {
29     input_command = concat(" ", input.process.Args)
30     some allowed_command in policy_data.allowed_commands
31     input_command == allowed_command
32 }
33
34 policy_data := {
35     "allowed_commands": [
36         "whoami",
37         "false",
38         "curl -s http://127.0.0.1:8006/aa/token?token_type=kbs",
39     ],
40     "allowed_images": [
41         "pause",
42         "docker.io/library/nginx@sha256:e56797eab4a5300158cc015296229e1",
43     ],
44 }
```

How to provide a policy to the TEE?

- Policy is specific per workload
- CVM images are generic
- Provide it as measured configuration at launch
- Link it to the TEE HW evidence
 - Put hash in HOSTDATA (SEV-SNP), MRCONFIGID (TDX), part of signed HW evidence (verify in TEE)
 - Extend runtime registers (vTPM)

Init-Data Specification

- Measured configuration for CoCo
- TOML dict of path/file content
- Currently being implemented
- Available for some TEEs
- Embed into Pod spec as annotation

```
1 | vim init-data.toml
2 | INIT_DATA_B64="$(cat "init-data.toml" | base64 -w0)"
3 | cat nginx-cc.yaml | jq \
4 |   --arg initdata "$INIT_DATA_B64" \
5 |   '.spec.template.metadata.annotations = { "io.katacontainers.config.r
6 | | kubectl apply -f -
```

Initdata Example

```
1 algorithm = "sha256"
2 version = "0.1.0"
3
4 [data]
5 "aa.toml" = '''
6 [token_configs]
7 [token_configs.kbs]
8 url = 'http://my-as:8080'
9 cert = """
10 -----BEGIN CERTIFICATE-----
11 MIIDEjCCAFqgAwIBAgIUZYcKIJD3QB/LG0FnacDyR1KhoikwDQYJKoZIhvcNAQEL
12 ...
13 4La0LJGguzEN7y9P59TS4b3E9xFyTg==
14 -----END CERTIFICATE-----
15 """
16 ...
17
18 "cdh.toml" = '''
19 socket = 'unix:///run/confidential-containers/cdh.sock'
20 credentials = []
21 ...
22 ...
23
24 "policy.rego" = '''
25 package agent_policy
26 ...
```

Challenges

- Policy is stateless, declarative
 - Kata RPC is imperative
 - What about more complex orchestration?
 - launch container x first (init container) then container y
- Ongoing effort: stateful policies

Challenges

- Practical problems
 - Size of policies?
 - Policies can be quite large
 - Pod annotation has limits
 - Compression, splitting, bundling a library
 - User experience is subpar
 - Rego is modelled after Datalog
 - Unusual paradigms
 - Not trivial to write large policies

Challenges

- Conceptual problems, maybe?
 - Have to track kata's RPC interface closely
 - New exploit vectors can be introduced inadvertently
 - Kata is not just for CoCo use case
 - Need to keep tabs on API changes in semantics and implementation

Challenges

- Runtime measurements
 - (very) long running workloads in TEEs
 - Examples: LLM inference, training tasks
 - Continuous measurement to catch drift
 - Not all TEEs have PCRs/RTMRs
 - Can be retrofitted via privilege levels + paravisor/SVSMs.
- “Composite” TEEs
 - Confidential GPUs + Confidential CPUs
 - Potentially more, e.g. accelerated NICs
 - Attest individually? Chained?

Recap

- Attestation for container sandboxes is tricky due to inherent dynamic nature.
- “Offloading” verification to a policy is a viable mitigation
- Few challenges remain, most seem manageable
- But policy is maybe not fully adequate

thx!

References

- [Confidential Containers](#)
- [Kata Containers - Open Source Container Runtime](#)
- [Policing a Sandbox | Microsoft Community Hub](#)
- [CoCo Initdata spec](#)