



KubeCon



CloudNativeCon

Europe 2019

Building images efficiently and securely on Kubernetes with BuildKit

Akihiro Suda, NTT



KubeCon



CloudNativeCon

Europe 2019

Raise your hand if you have heard of BuildKit?





KubeCon



CloudNativeCon

Europe 2019

Raise your hand if you are already using BuildKit?





KubeCon



CloudNativeCon

Europe 2019

**Raise your hand if you are already running
BuildKit on Kubernetes?** 





KubeCon



CloudNativeCon

Europe 2019

Part 1

Introduction to BuildKit

BuildKit: next-generation docker build



KubeCon



CloudNativeCon

Europe 2019

- Concurrent multi-stage build
- Efficient caching
- Secure access to private assets
- Flexible syntax for build definition
- Does not require root privileges

BuildKit: next-generation docker build



KubeCon



CloudNativeCon

Europe 2019

- BuildKit is included in Docker since v18.06

```
$ export DOCKER_BUILDKIT=1  
$ docker build ...
```

- But this talk will focus on the standalone version of BuildKit (`buildkitd` & `buildctl`)
 - No dependency on Docker

LLB DAG



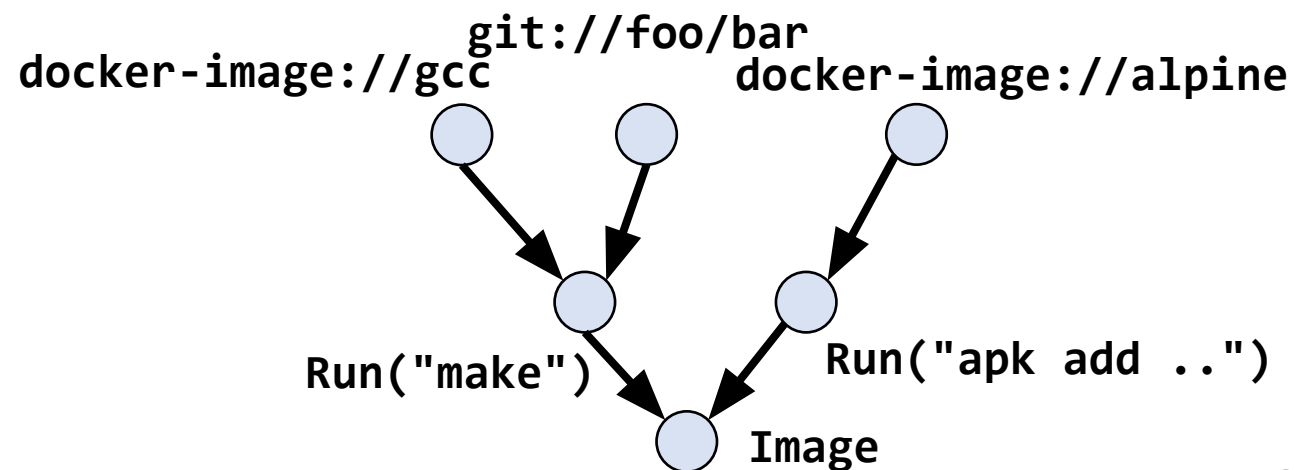
KubeCon



CloudNativeCon

Europe 2019

- LLB is to Dockerfile what LLVM IR is to C
- Typically compiled from Dockerfile
- Accurate dependency expression with DAG structure
 - Efficient caching
 - Concurrent execution



LLB DAG



KubeCon



CloudNativeCon

Europe 2019

```
FROM golang AS stage0
```

```
...
```

```
RUN go build -o /foo ...
```

```
FROM clang AS stage1
```

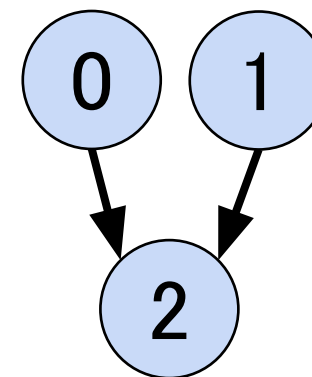
```
...
```

```
RUN clang -o /bar ...
```

```
FROM debian AS stage2
```

```
COPY --from=stage0 /foo /usr/local/bin/foo
```

```
COPY --from=stage1 /bar /usr/local/bin/bar
```

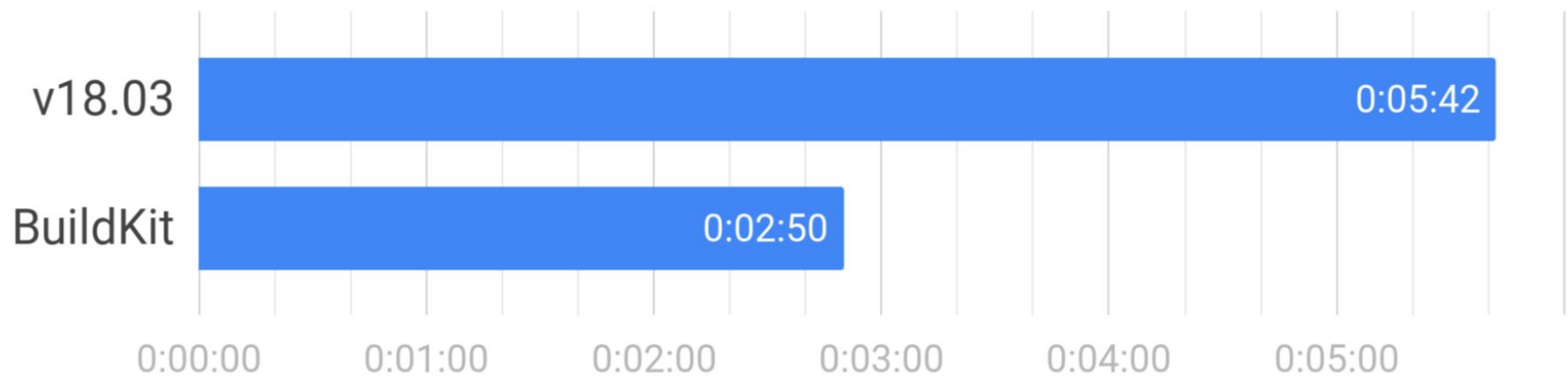


BuildKit

Performance example

Based on github.com/moby/moby Dockerfile, master branch. **Smaller** is better.

Time for full build from empty state



2.0x
faster

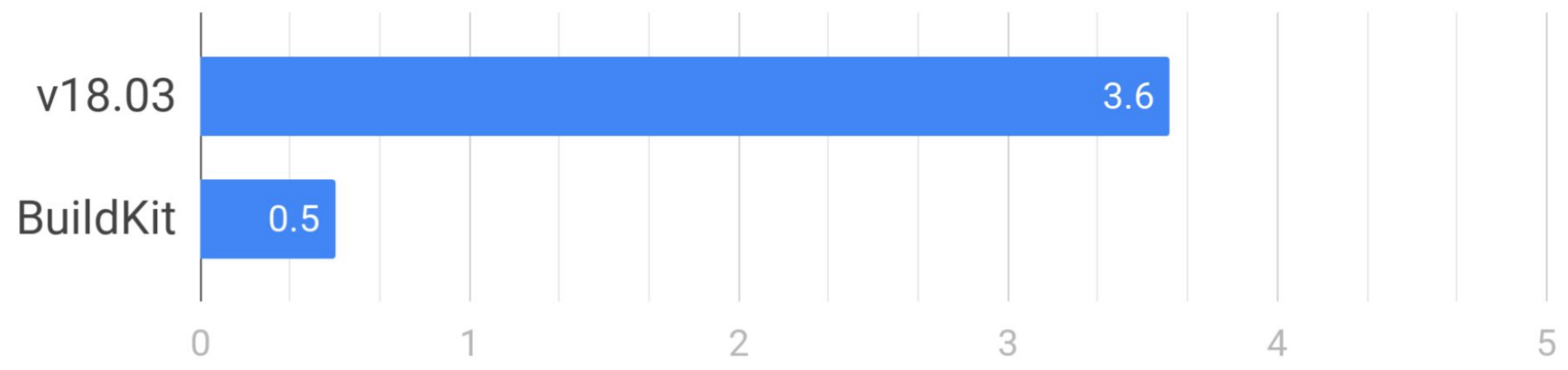
Measured on DO 4vcpu droplet

BuildKit

Performance example

Based on github.com/moby/moby Dockerfile, master branch. **Smaller** is better.

Repeated build with matching cache



7.2x
faster

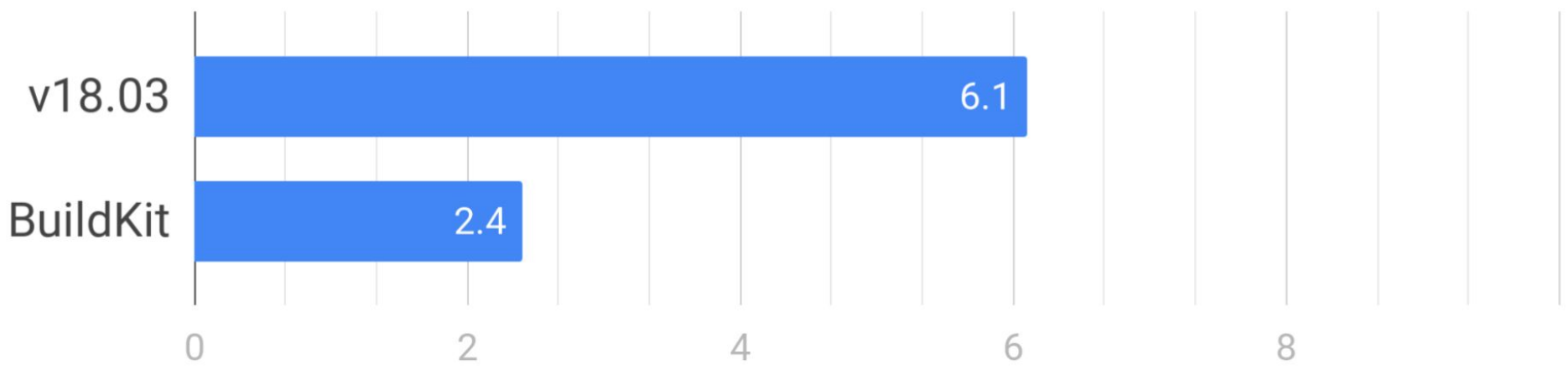
Measured on DO 4vcpu droplet

BuildKit

Performance example

Based on github.com/moby/moby Dockerfile, master branch. **Smaller** is better.

Repeated build with new source code



**2.5x
faster**

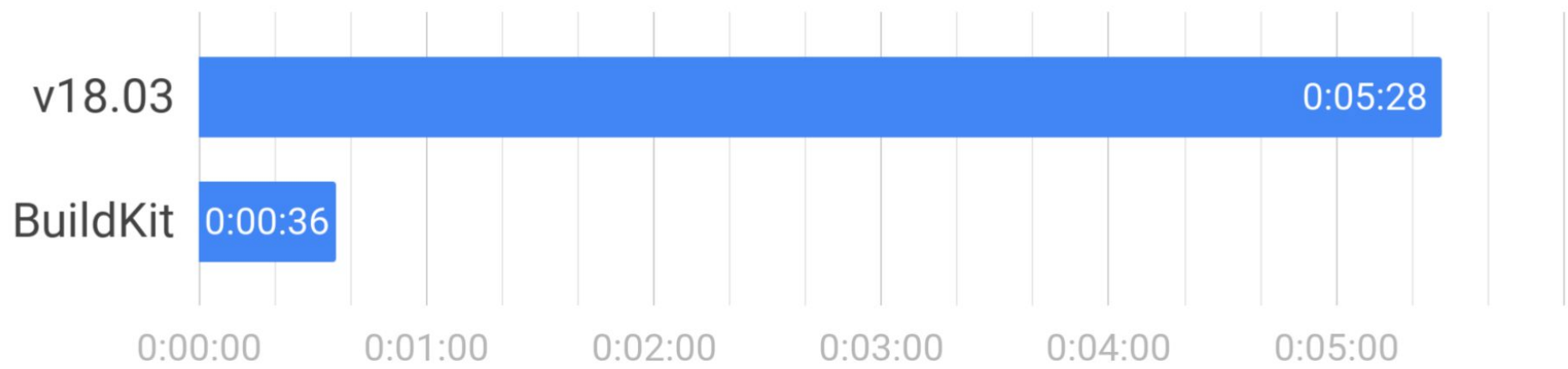
Measured on DO 4vcpu droplet

BuildKit

Performance example

Based on github.com/moby/moby Dockerfile, master branch. **Smaller** is better.

Fresh build with --cache-from from remote source



9.1x
faster

Measured on DO 4vcpu droplet

RUN --mount=type=cache



KubeCon



CloudNativeCon

Europe 2019

- Allows preserving caches of compilers and package managers

```
# syntax = docker/dockerfile:1.1-experimental
...
RUN --mount=type=cache,target=/root/.cache go build
...
```

RUN --mount=type=cache



KubeCon



CloudNativeCon

Europe 2019

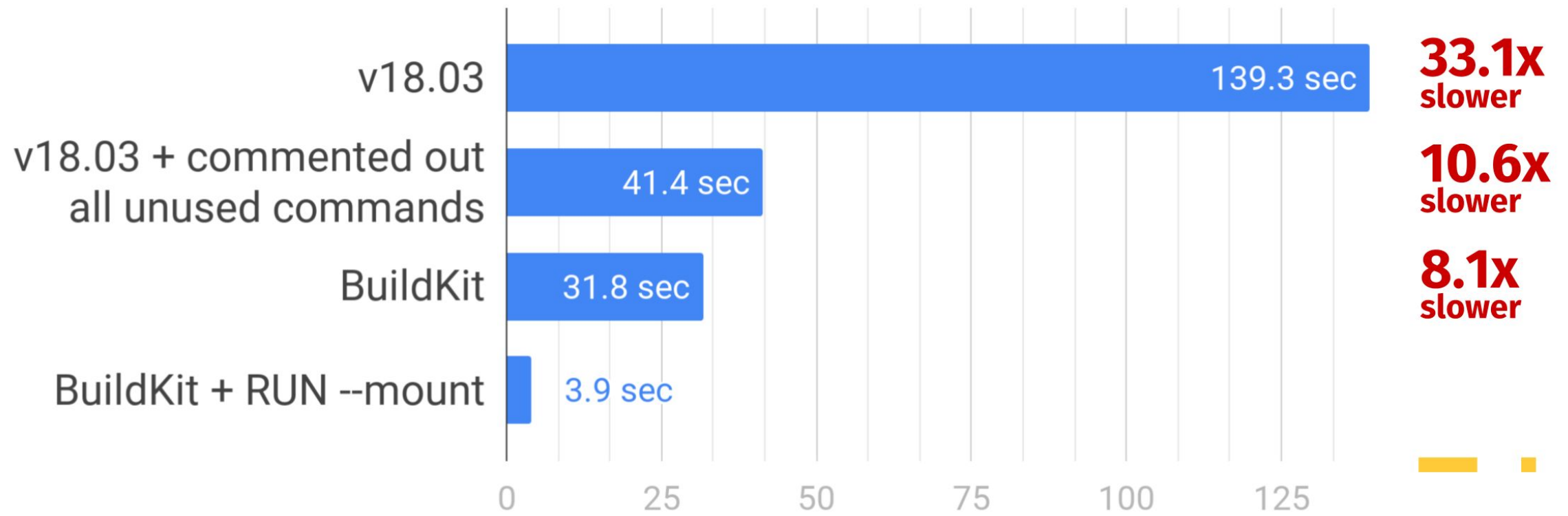
- Allows preserving caches of compilers and package managers

```
# syntax = docker/dockerfile:1.1-experimental
FROM ubuntu
RUN rm -f /etc/apt/apt.conf.d/docker-clean; \
    echo 'Binary::apt::APT::Keep-Downloaded-Packages "true";' > \
    /etc/apt/apt.conf.d/keep-cache
RUN \
    --mount=type=cache,target=/var/cache/apt \
    --mount=type=cache,target=/var/lib/apt \
    apt update && apt install -y gcc
```

Dockerfile syntax directive

Example: RUN --mount

moby/buildkit Dockerfile: time to binary rebuild after code change



Measured on DO 4vcpu droplet

RUN --mount=type=secret



KubeCon



CloudNativeCon

Europe 2019

- Allows accessing private assets without leaking credential in the image

```
# syntax = docker/dockerfile:1.1-experimental
...
RUN --mount=type=secret,id=aws,target=/root/.aws/credentials \
    aws s3 cp s3://... ..
```

```
$ buildctl build --secret id=aws,src=~/.aws/credentials ...
```

RUN --mount=type=secret



KubeCon



CloudNativeCon

Europe 2019

- Note: DON'T do this!

```
...  
COPY my_aws_credentials /root/.aws/credentials  
RUN aws s3 cp s3://... ..  
RUN rm -f /root/.aws/credentials  
...
```

RUN --mount=type=secret



KubeCon



CloudNativeCon

Europe 2019

- Note: DON'T do this either!

```
$ docker build \  
  --build-arg \  
  MY_AWS_CREDENTIALS=$(cat ~/.aws/credentials)
```

RUN --mount=type=ssh



KubeCon



CloudNativeCon

Europe 2019

- Akin to `--mount=type=secret` but specific to SSH
- Supports passphrase

```
# syntax = docker/dockerfile:1.1-experimental
...
RUN --mount=type=ssh git clone ssh://github.com/...
```

```
$ eval $(ssh-agent)
$ ssh-add ~/.ssh/id_rsa
(Enter your passphrase)
$ buildctl build --ssh default=$SSH_AUTH_SOCK ...
```


Non-Dockerfiles



KubeCon



CloudNativeCon

Europe 2019

- LLB can be also compiled from non-Dockerfiles
- Several languages are being proposed
 - Buildpacks
 - Mockfile
 - Gockerfile
- You can also create your own language

Buildpacks



KubeCon



CloudNativeCon

Europe 2019

- Ported from Heroku/CloudFoundry Buildpacks
- No support for Cloud Native Buildpacks yet

```
# syntax = tonistiigi/pack
---
applications:
- name: myapp
  memory: 128MB
  disk_quota: 256MB
  random-route: true
  buildpack: python_buildpack
  command: python hello.py
```

Mockerfile



KubeCon



CloudNativeCon

Europe 2019

- apt-get in highly declarative YAML

```
# syntax = r2d4/mocker
apiVersion: v1alpha1
images:
- name: demo
  from: ubuntu:16.04
  package:
    repo:
      - deb [arch=amd64] http://storage.googleapis.com/bazel-apt stable jdk1.8
    gpg:
      - https://bazel.build/bazel-release.pub.gpg
  install:
    - bazel
```

Gockerfile



KubeCon



CloudNativeCon

Europe 2019

- Really simple
- Specific to Golang

```
# syntax = po3rin/gocker  
repo: github.com/foo/bar  
path: ./cmd/baz  
version: v0.0.1
```



KubeCon



CloudNativeCon

Europe 2019

Part 2

Deploying BuildKit on Kubernetes

Why build images on Kube?



KubeCon



CloudNativeCon

Europe 2019

Two kinds of motivation:

1. CI/CD

2. Developer Experience

Why build images on Kube?

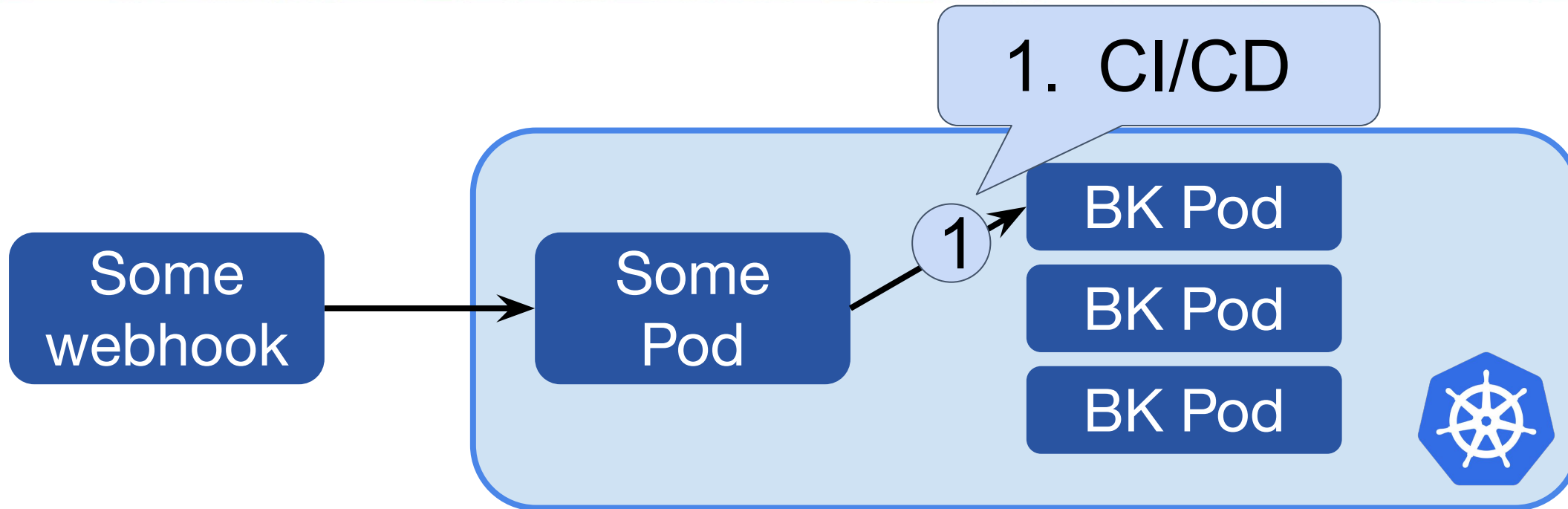


KubeCon



CloudNativeCon

Europe 2019



Why build images on Kube?

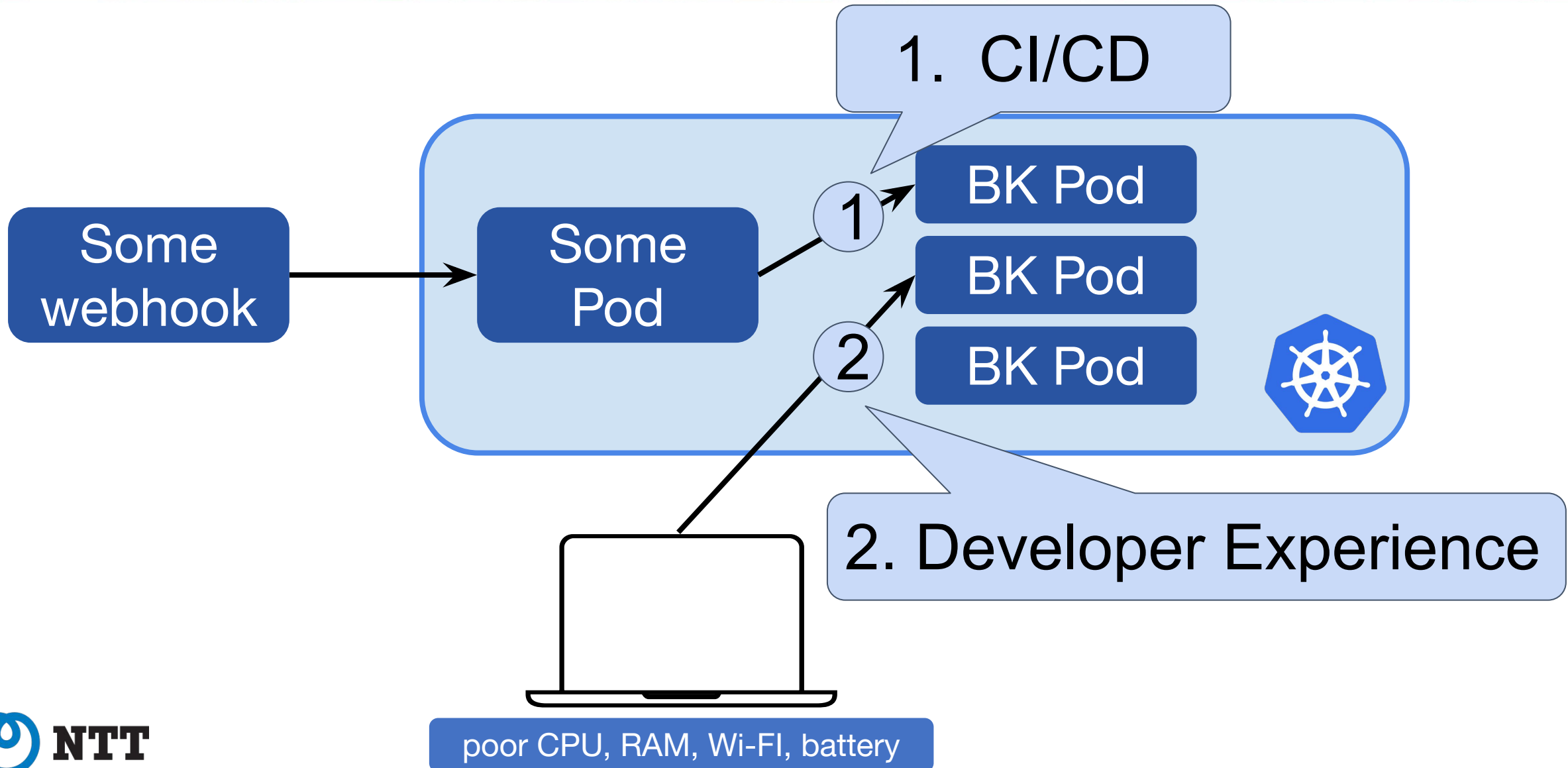


KubeCon



CloudNativeCon

Europe 2019



Issue with docker build on Kube



KubeCon



CloudNativeCon

Europe 2019

- The common pattern was to run `docker` Pod with `/var/run/docker.sock` `hostPath`
- Or run `docker:dind` Pod with `securityContext.privileged`
- Both are insecure



KubeCon



CloudNativeCon

Europe 2019

Part 2.1

Rootless mode

Rootless mode



KubeCon



CloudNativeCon

Europe 2019

- BuildKit can be executed as a non-root user
- No extra `securityContext` configuration needed
- Protect the host from potential BuildKit vulns



KubeCon



CloudNativeCon

Europe 2019

Demo

myth 1: requires securityContext.privileged



KubeCon



CloudNativeCon

Europe 2019

- Not true since BuildKit v0.4.0
- But you need to disable “Process Sandbox”:
launch `buildkitd` with
 - `--oci-worker-no-process-sandbox`
 - Disable unsharing PIDNS and mounting `/proc`

myth 1: requires securityContext.privileged



KubeCon



CloudNativeCon

Europe 2019

Process sandbox
(needs to be disabled)

Process sandbox

worker container (e.g. `RUN gcc ...`)

BuildKit daemon

Host

myth 1: requires securityContext.privileged

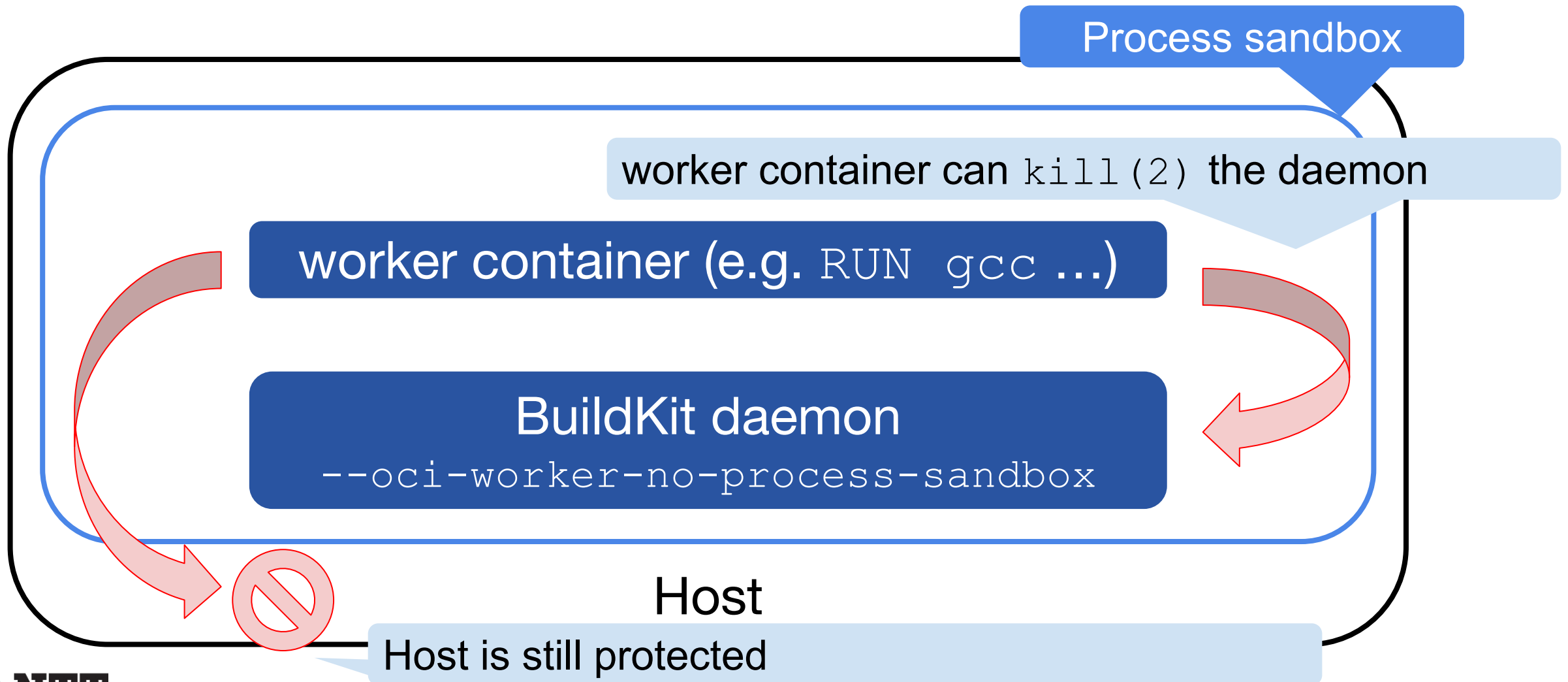


KubeCon



CloudNativeCon

Europe 2019



myth 1: requires securityContext.privileged



KubeCon



CloudNativeCon

Europe 2019

- To enable Process Sandbox,
`securityContext.procMount` needs to be set to
`Unmasked`
 - Requires Kubernetes v1.12+ with Docker v18.06+ /
containerd v1.2+ / CRI-O v1.12

myth 2: seccomp and AppArmor need to be disabled



KubeCon



CloudNativeCon

Europe 2019

myth 2: seccomp and AppArmor need to be disabled



KubeCon



CloudNativeCon

Europe 2019

- Not a myth :P
- seccomp (and AppArmor) is typically disabled by default on Kubernetes anyway
 - In Kubernetes world, seccomp is still in alpha status and AppArmor is in beta

myth 2: seccomp and AppArmor need to be disabled

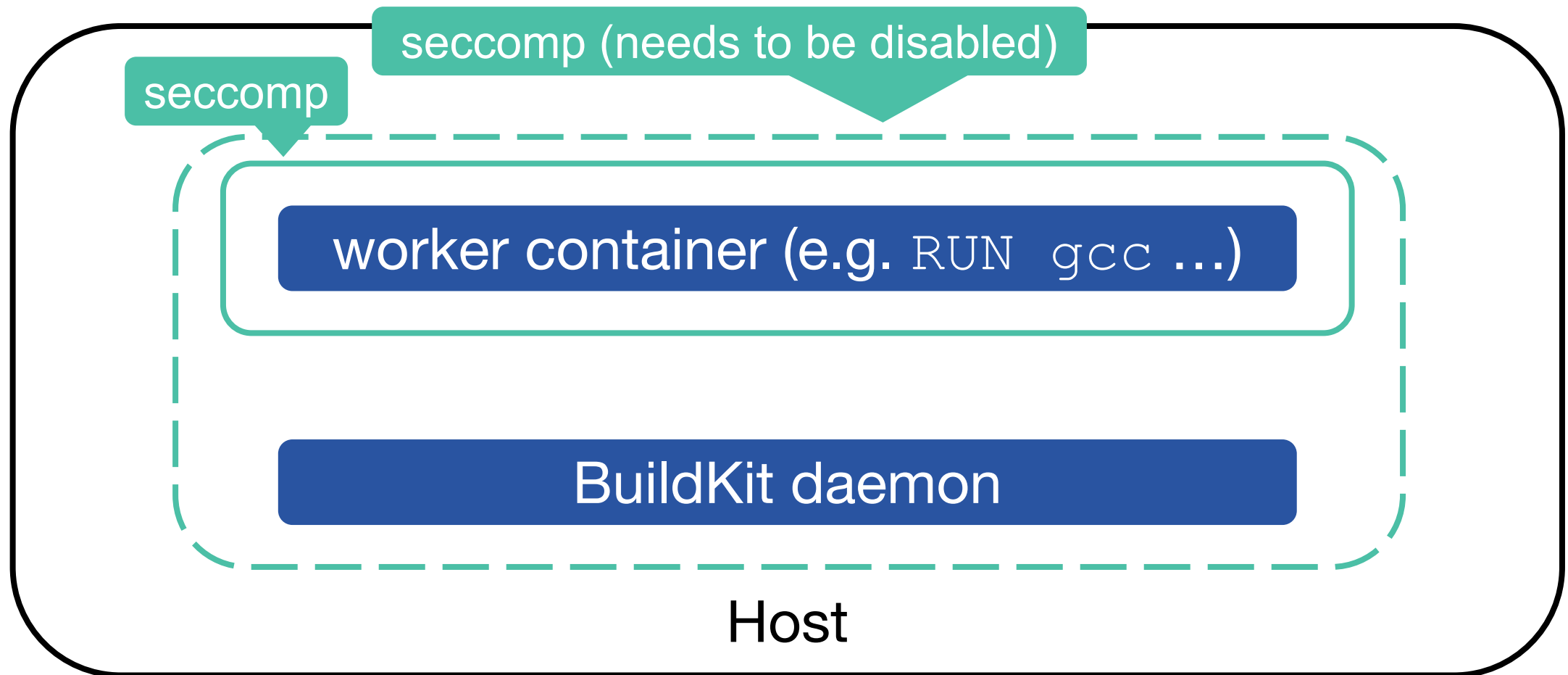


KubeCon



CloudNativeCon

Europe 2019



myth 2: seccomp and AppArmor need to be disabled



KubeCon



CloudNativeCon

Europe 2019

seccomp

worker containers are still protected with seccomp

worker container (e.g. RUN gcc ...)

BuildKit daemon

Host

Future work: gVisor integration?



KubeCon

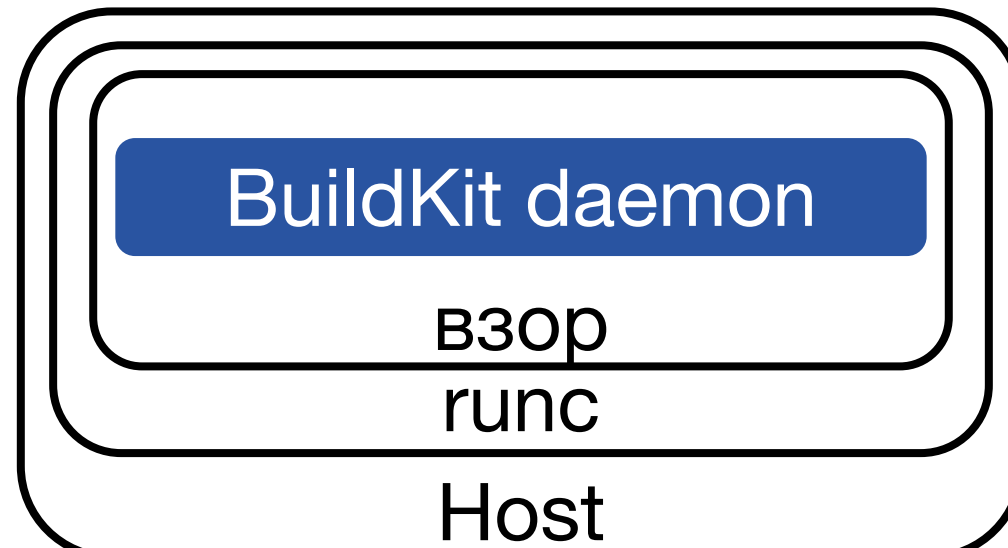


CloudNativeCon

Europe 2019

- gVisor: Yet another Linux kernel implementation in userspace
- B3op (*vzor*): gVisor-based sandbox for runc containers

<https://github.com/tonistiigi/vzor>



Future work: gVisor integration?



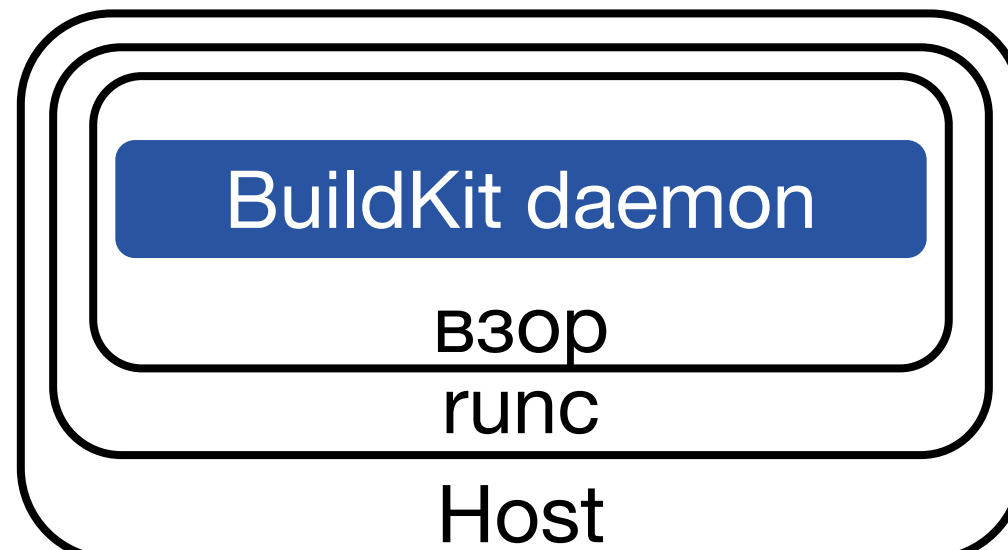
KubeCon



CloudNativeCon

Europe 2019

- No need to disable seccomp/AppArmor for runc
- Can also mitigate kernel vulns



Future work: gVisor integration?



KubeCon



CloudNativeCon

Europe 2019

- Currently BuildKit fails with `EINVAL` due to syscall incompatibility
- Or User-Mode Linux?
 - Full Linux compatibility
 - 20 yo, still alive :)

Rootless BuildKit vs Kaniko



KubeCon



CloudNativeCon

Europe 2019

- Kaniko runs as the root but “unprivileged”
 - No need to disable seccomp and AppArmor
- Kaniko might be able to mitigate some vuln that Rootless BuildKit cannot mitigate - and vice versa
 - Rootless BuildKit might be weak against kernel vulns
 - Kaniko might be weak against runc vulns



KubeCon



CloudNativeCon

Europe 2019

Part 2.2

Deployment strategy

Deployment strategy



KubeCon



CloudNativeCon

Europe 2019

Deployment?

DaemonSet?



StatefulSet?

Job?

Deployment strategy



KubeCon



CloudNativeCon

Europe 2019

- Deployment
 - Most typical deployment
- DaemonSet
 - **Better Pod placement**
 - **But unlikely to hit daemon-local cache if you have a bunch of replicas**
 - **So might be not always optimal for large clusters w/ complex Dockerfiles**

Deployment strategy



KubeCon



CloudNativeCon

Europe 2019

- StatefulSet
 - **Consistent Pod names**
 - **Good for Consistent Hashing (discussed later)**

Deployment strategy



KubeCon



CloudNativeCon

Europe 2019

- Job (“Daemonless”)
 - `buildctl` and ephemeral `buildkitd` in a single ephemeral Pod
 - No need to manage the life cycles of the daemons
 - Needs PR: [moby/buildkit#979](https://github.com/moby/buildkit/pull/979)
 - or github.com/genuinetools/img (lacks some upstream features)

How to connect to BuildKit?



KubeCon



CloudNativeCon

Europe 2019

- BuildKit daemon can listen on TCP (with TLS)
- The entire operation (build & push) just needs a single gRPC connection
- So you can create `Kubernetes Service` for connecting to `BuildKit Deployment / DaemonSet / StatefulSet`

How to connect to BuildKit?

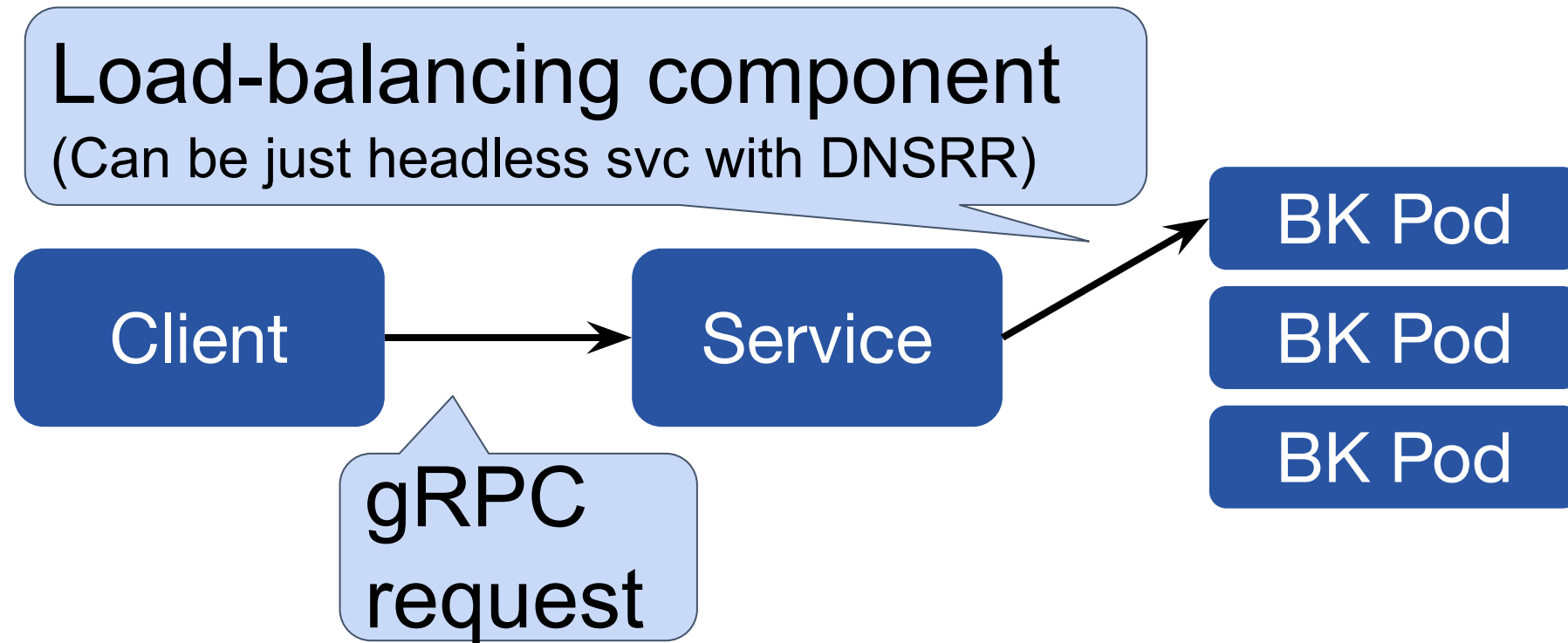


KubeCon



CloudNativeCon

Europe 2019



How to connect to BuildKit?



KubeCon



CloudNativeCon

Europe 2019

- But you don't need to necessarily create `Service`
- `buildctl` CLI can directly connect to a daemon in a `Pod` **without** `Service`
 - Internally invokes `kubectl exec`

How to connect to BuildKit?



KubeCon



CloudNativeCon

Europe 2019

```
$ kubectl run \  
  --generator=run-pod/v1 \  
  --image=moby/buildkit:master-rootless \  
  bk -- --oci-worker-no-process-sandbox
```

```
$ export BUILDKIT_HOST=kube-pod://bk  
$ buildctl build ...
```

Coming soon: docker buildx for Kube



KubeCon



CloudNativeCon

Europe 2019

- `docker buildx` is the next generation CLI for integrating BuildKit to Docker
 - Supports building multi-arch image with remote ARM machines
 - “Bake”: compose-like build
- `docker buildx` will support connecting to BuildKit on Kubernetes in the same UX



KubeCon



CloudNativeCon

Europe 2019

Part 2.3

Caching

Remote cache



KubeCon



CloudNativeCon

Europe 2019

- Cache can be shared via either registry or shared FS
- Similar to classic `docker build --cache-from` but more chance of hitting cache
- For building non-container artifacts (it's a valid use-case), FS cache might be useful

Remote cache

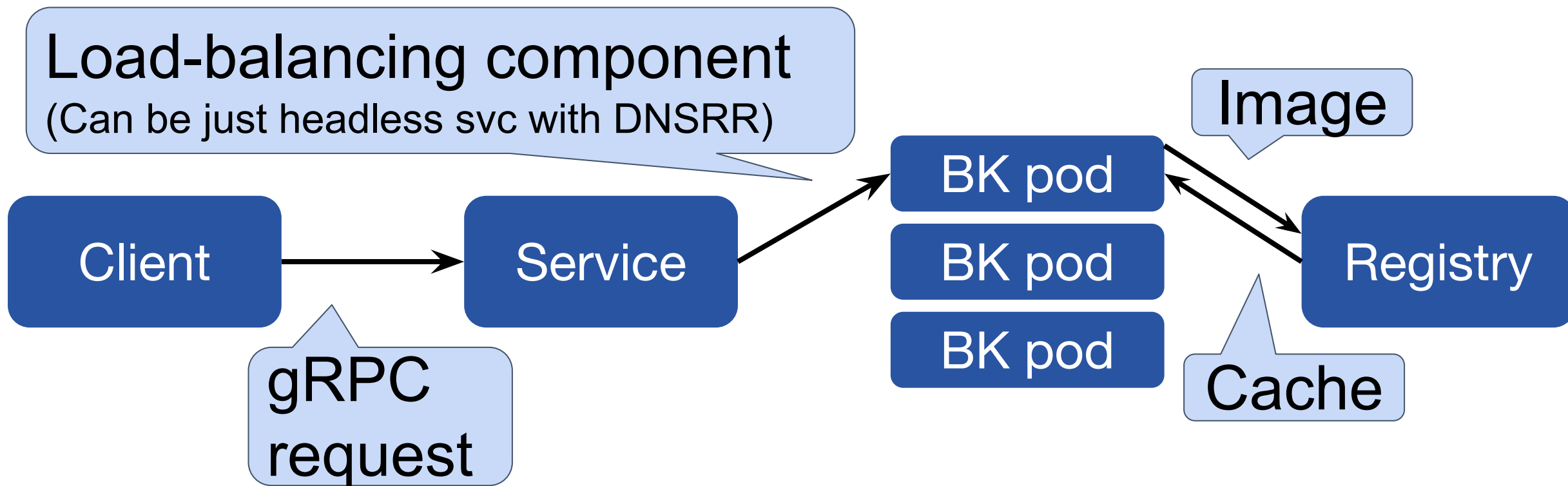


KubeCon



CloudNativeCon

Europe 2019



Remote cache



KubeCon



CloudNativeCon

Europe 2019

- Remote cache might be slow compared to the daemon-local cache
- Example from Part 1 slides:
 - No cache: 2m50s
 - Remote cache: : 36s
 - Daemon-local cache: 0.5s

Consistent hashing



KubeCon



CloudNativeCon

Europe 2019

- Consistent hashing allows sticking a build request to a specific Pod in StatefulSet
- So the build request can always hit the daemon-local cache in the Pod

Consistent hashing

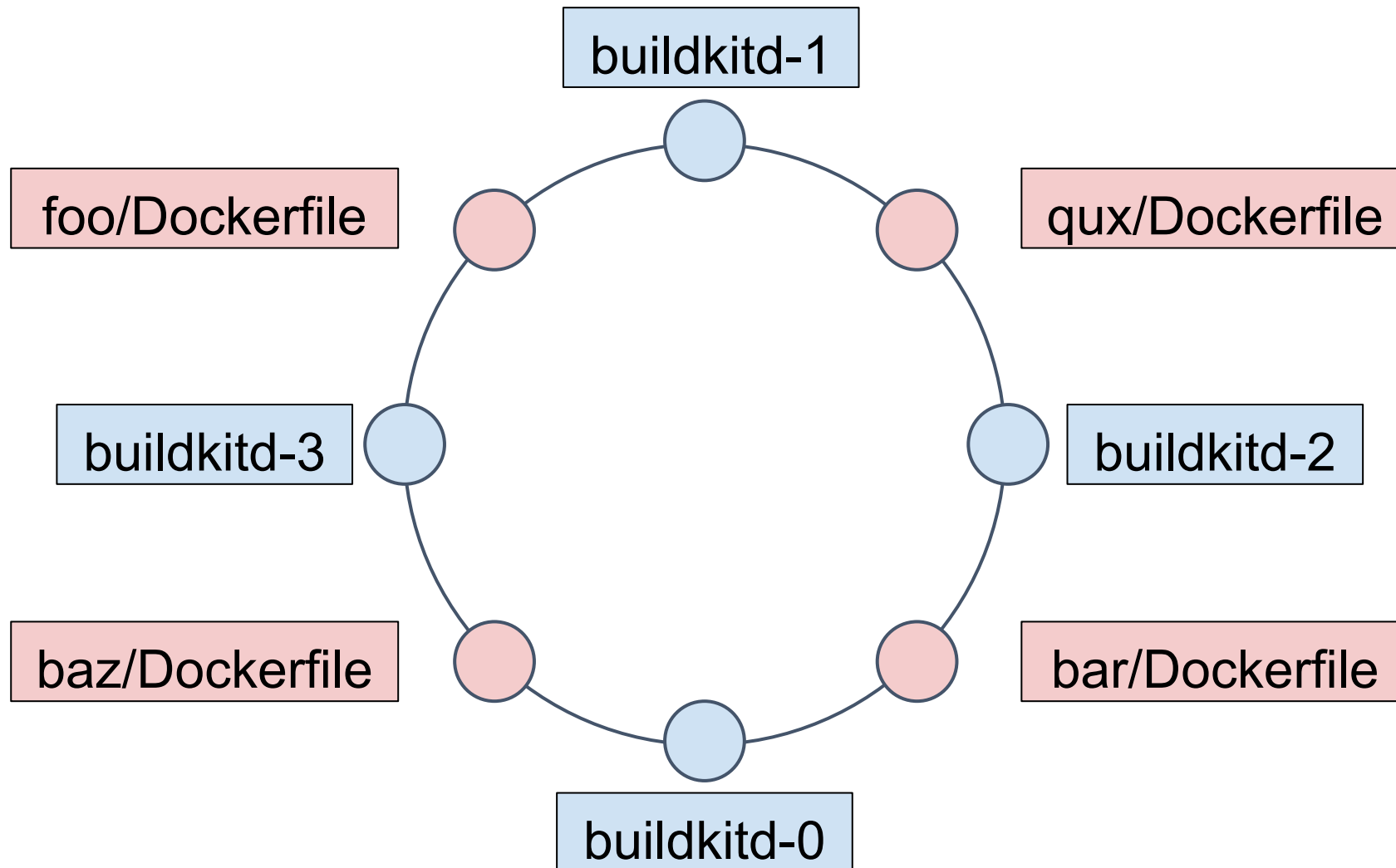


KubeCon



CloudNativeCon

Europe 2019



Consistent hashing

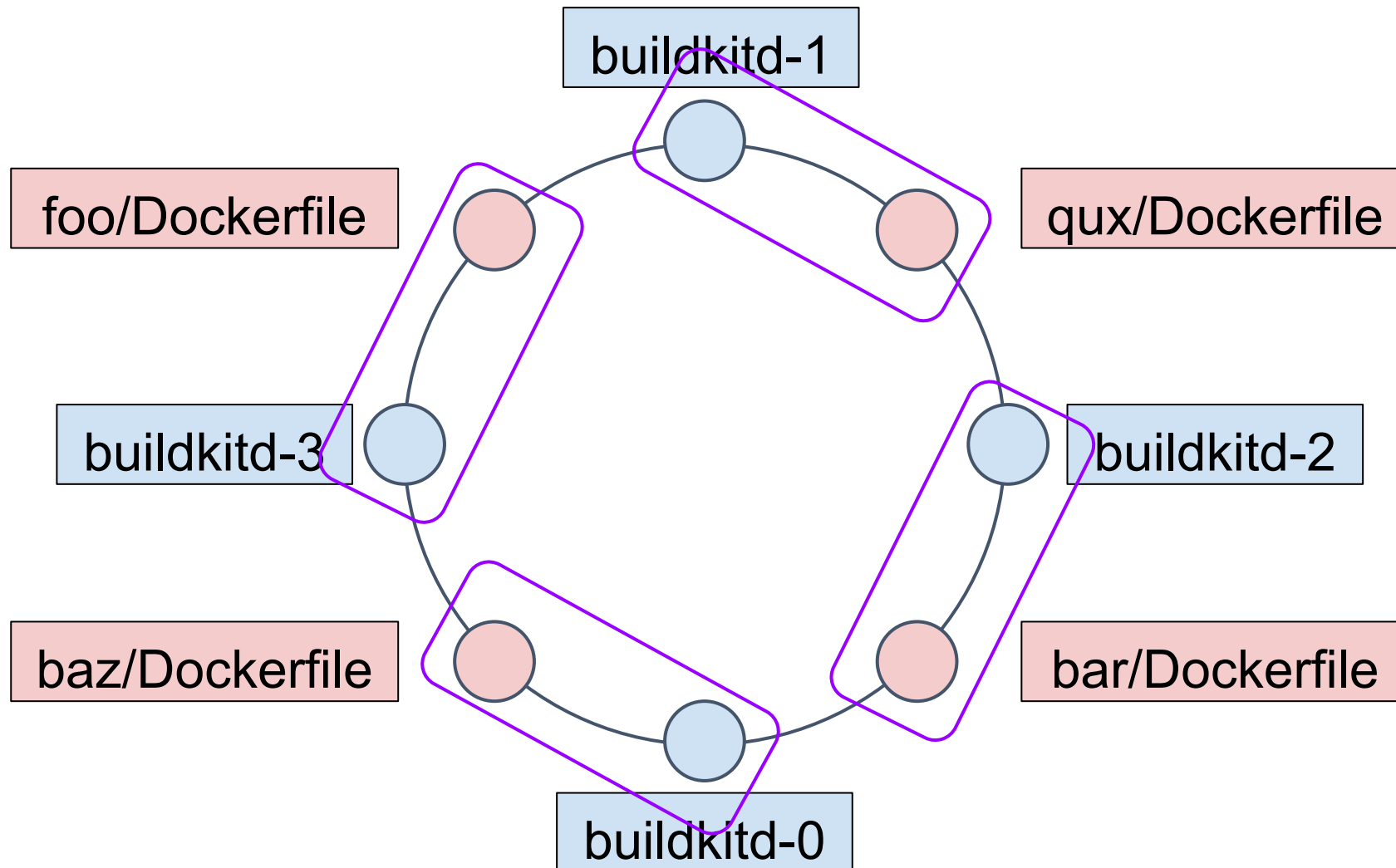


KubeCon



CloudNativeCon

Europe 2019



Consistent hashing

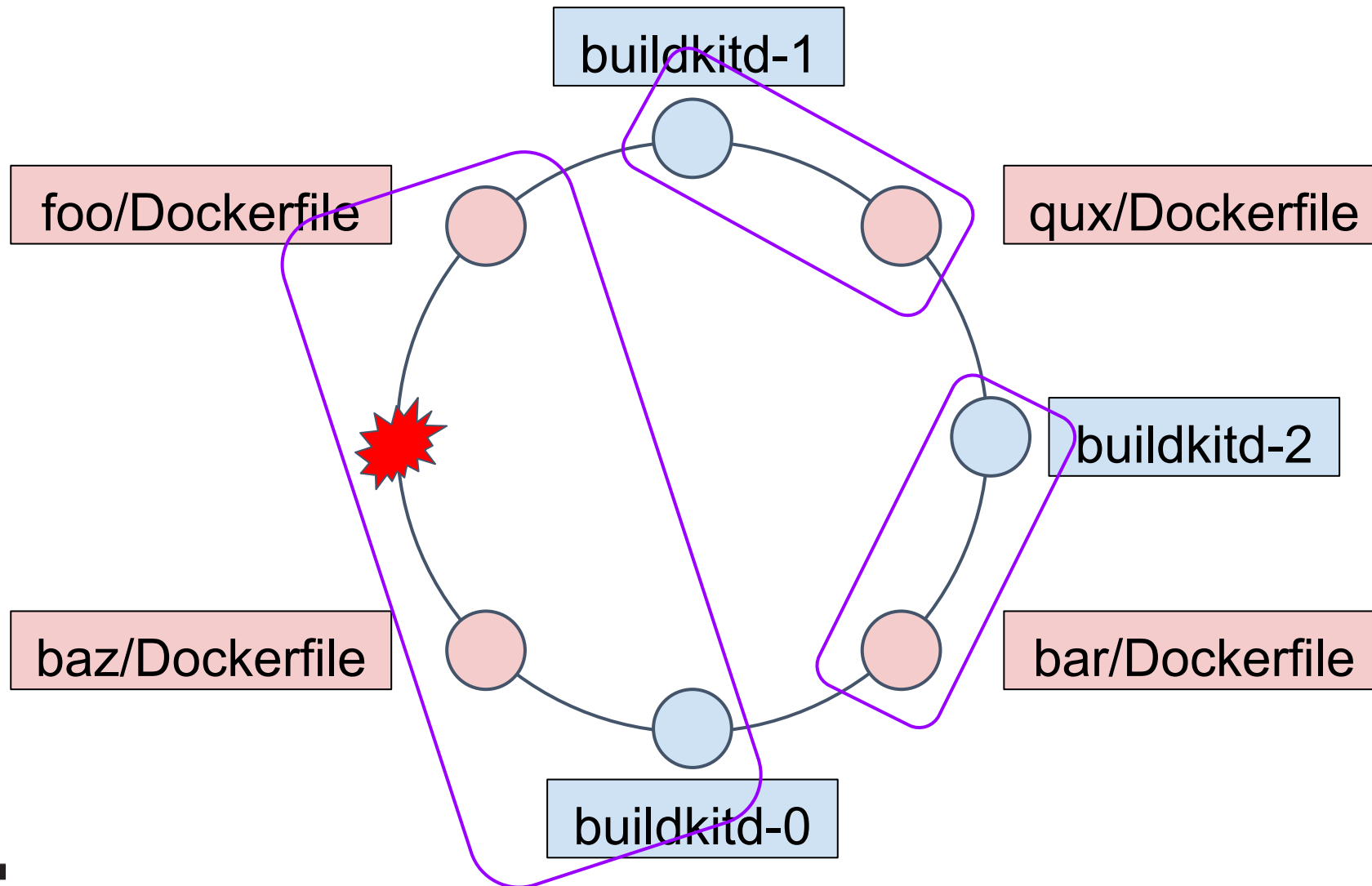


KubeCon



CloudNativeCon

Europe 2019



Consistent hashing

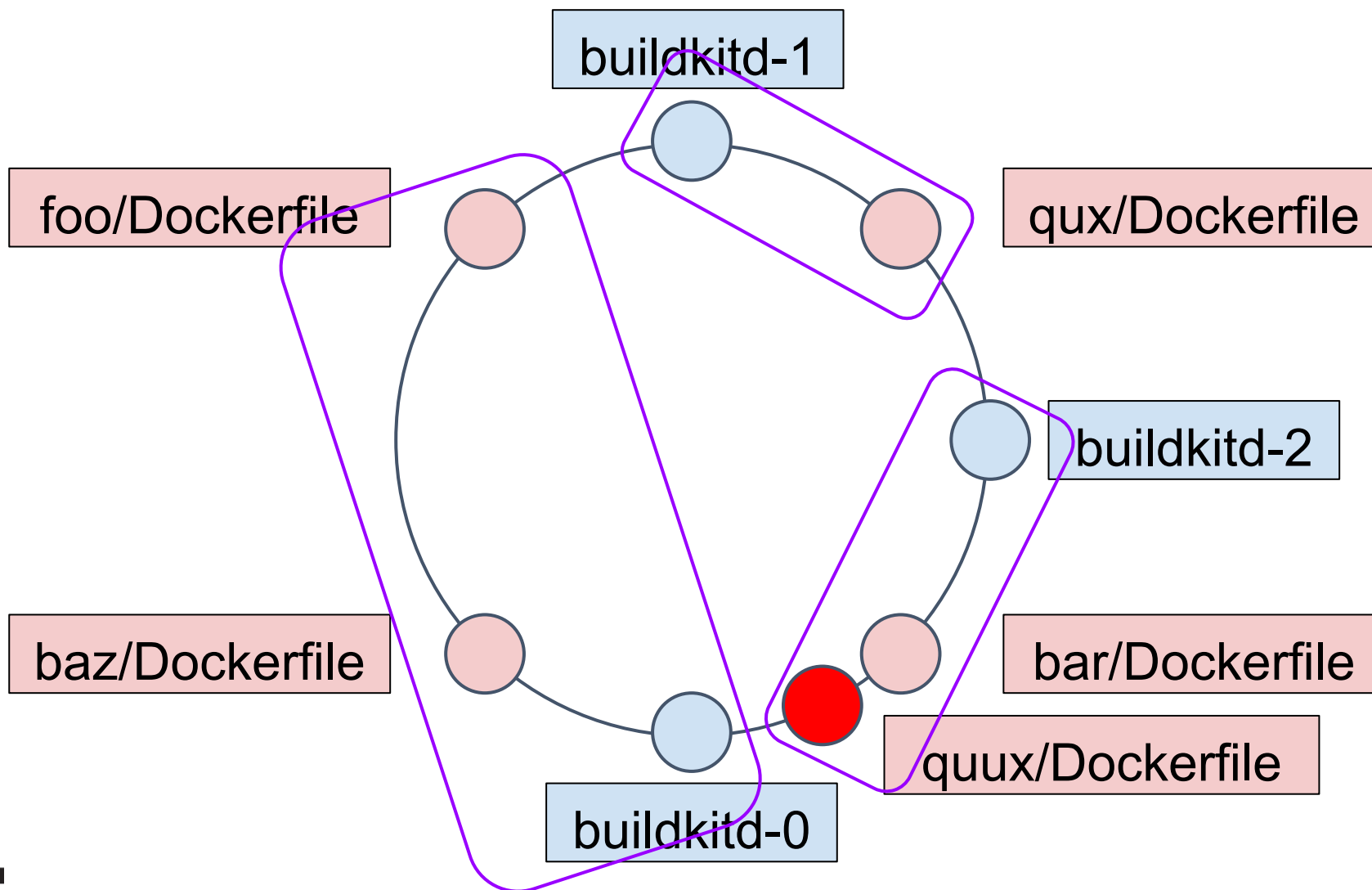


KubeCon



CloudNativeCon

Europe 2019



Consistent hashing



KubeCon



CloudNativeCon

Europe 2019

- Caveats:
 - High I/O overhead on specific set of nodes
 - Some nodes might not be used at all
- **See** `examples/kube-consistent-hashing` **in the** `moby/buildkit` **repo**

Remote cache vs Consistent hashing?



KubeCon



CloudNativeCon

Europe 2019

- If your cache registry is fast enough for your Dockerfiles, remote cache w/ load-balancing might be better
- If you don't like transferring cache, consistent hashing might be better



KubeCon



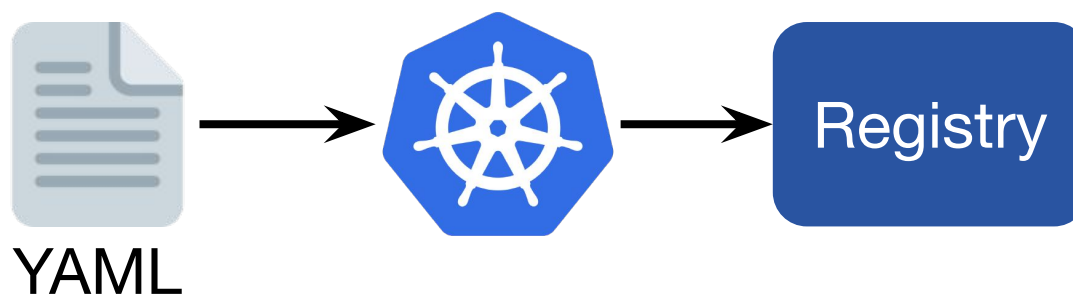
CloudNativeCon

Europe 2019

Part 2.4

CRD

YAMLIFY ALL THE THINGS



Container Builder Interface (CBI)



- The first common build CRD
- Supports Docker, BuildKit, Buildah, kaniko, img, Google Cloud Container Builder, Azure Container Registry Build, and OpenShift S2I
- Complex design with a bunch of microservices
- Now being deprecated

Container Builder Interface (CBI)

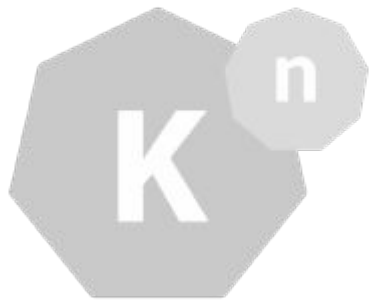


- Simpler than CBI and easily extensible
- The build component (not entire Knative) might be going to be deprecated in favor of Tekton



- Spun out from Knative
- Much more simple and extensible

Container Builder Interface (CBI)



- Simpler than CBI and easily extensible
- The build component (not entire Knative) might be going to be deprecated in favor of Tekton



- Spun out from Knative
- Much more simple and extensible

Tekton



KubeCon



CloudNativeCon

Europe 2019

```
apiVersion: tekton.dev/v1alpha1
```

```
kind: TaskRun
```

```
metadata:
```

```
  name: foobar
```

```
spec:
```

```
  taskRef:
```

```
    name: buildkit
```

The interface is same as other image builders (Buildah, Kaniko, and Makisu)

```
...
```

inputs:

resources:

- name: source

resourceSpec:

type: git

params:

- name: url

value: `git@github.com:foo/bar.git`

SSH credential is loaded from the `Secret` associated with the `ServiceAccount`

outputs:

resources:

- name: image

resourceSpec:

type: image

params:

- name: url

value: **registry.example.com/foo/bar**

Registry credential is loaded from the Secret associated with the ServiceAccount

Wrap-up



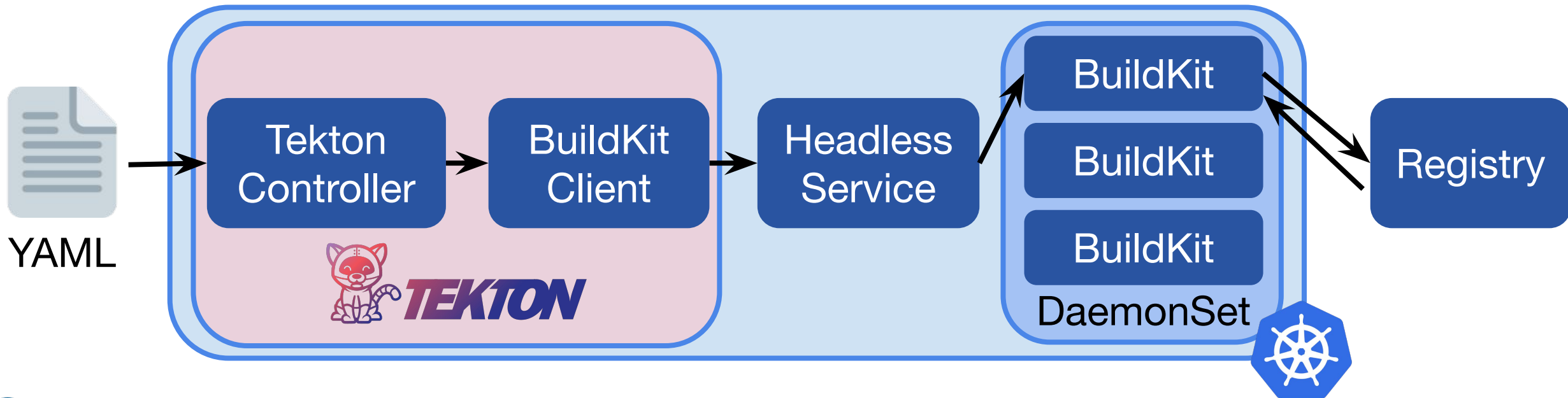
KubeCon



CloudNativeCon

Europe 2019

- BuildKit is fast and secure
- Several deployment plans, w/ and w/o daemon
- Example:





KubeCon



CloudNativeCon

Europe 2019

Join us: <https://github.com/moby/buildkit>