

**RESILIENCE  
REALIZED**



KubeCon



CloudNativeCon

North America 2021



KubeCon



CloudNativeCon

North America 2021

RESILIENCE  
REALIZED

# containerd introduction and deep dive

*Mike Brown (IBM)*

*Phil Estes (AWS)*

*Derek McGowan (Apple)*

*Maksym Pavlenko (Apple)*

# Status Update

The slide has a yellow gradient background with a pixelated texture. At the top right are the KubeCon and CloudNativeCon logos with the text "North America 2019". The main title "Containererd Mini-Summit" is in large bold black font. Below it is a subtitle with names and companies: "Phil Estes, IBM; Lantao Liu, Google; Derek McGowan, Docker; & Yu-Ju Hong, Google". The bottom features a decorative illustration of palm trees, a beach umbrella, surfboards, and a ship's wheel against a sunset background. A video control bar at the bottom includes icons for play, volume, and chapters, along with the timestamp "0:01 / 1:33:00".

**Containererd Mini-Summit**

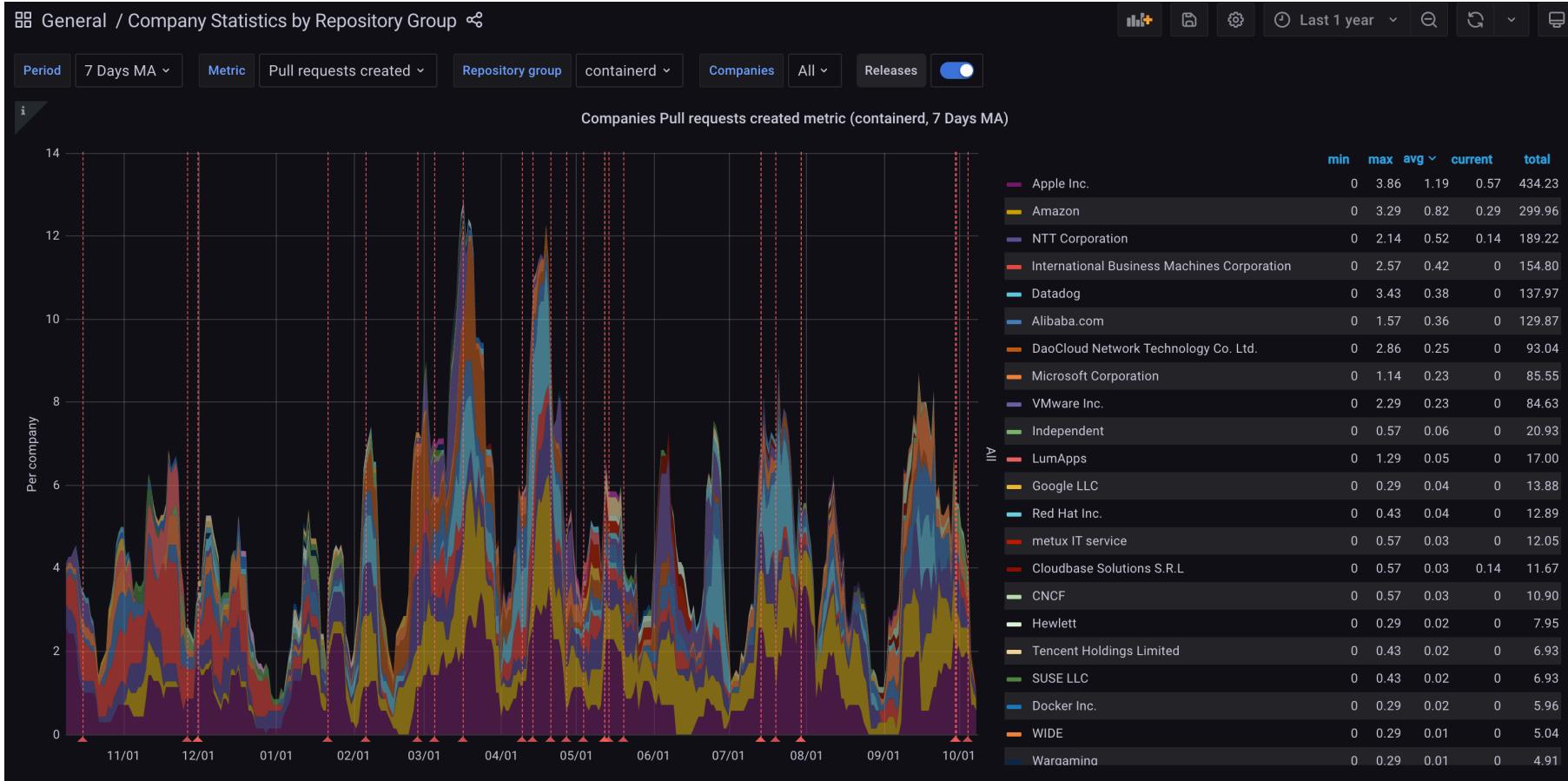
*Phil Estes, IBM; Lantao Liu, Google; Derek McGowan, Docker; & Yu-Ju Hong, Google*

- **2 years** since our last in-person update (KubeCon San Diego)
- Community has done **3 virtual** KubeCon/CloudNativeCon updates

# In the Last Two Years...

- **RELEASES**
  - 3 major release lines (1.4.x, and 1.5.x; 1.6.0-beta.0 **yesterday**)
  - 29 total releases across 3 supported release branches!
  - Releases **fully automated** through GitHub Actions (full CI migration)
  - **Windows** official release packaging; ARM64 now supported/tested
  - Security **disclosure/release process** implemented (and tested with 3 CVEs)
- **COMMUNITY**
  - Governance streamlined; **security advisor** role added
  - Many new reviewers from multiple employers
  - Over 370 unique code contributors; 10% of those submitted 10 PRs or more

# Community



- **67 unique companies** represented by contributors
- **Commitors:** 13 members from 9 companies
- **Reviewers:** 14 members from 10 companies

# Project Growth

- **NON-CORE SUBPROJECTS**
  - **stargz-snapshotter** - remote lazy-loading snapshotter
  - **imgcrypt** - container layer encryption support
  - **nerdctl** - more complete Docker client replacement using containerd. Offers packaged support for stargz, rootless, encryption, etc.
- **WHAT'S NEW**
  - CRIv1, more metrics, OpenTelemetry (initial support), shim plugins
  - Confidential Computing: Inclavare project using containerd
  - Upcoming: Sandbox API, Shim runtime/task/API rework, and more!

# Parts of containerd

## Client

- ctr
- nerdctl
- Go library

## containerd daemon

- API server
- CRI plugin
- Resource managers
  - Data storage
  - Garbage Collection
  - Shim Management

## containerd shim

- Per container or pod instance
- Manages running processes

# Architecture

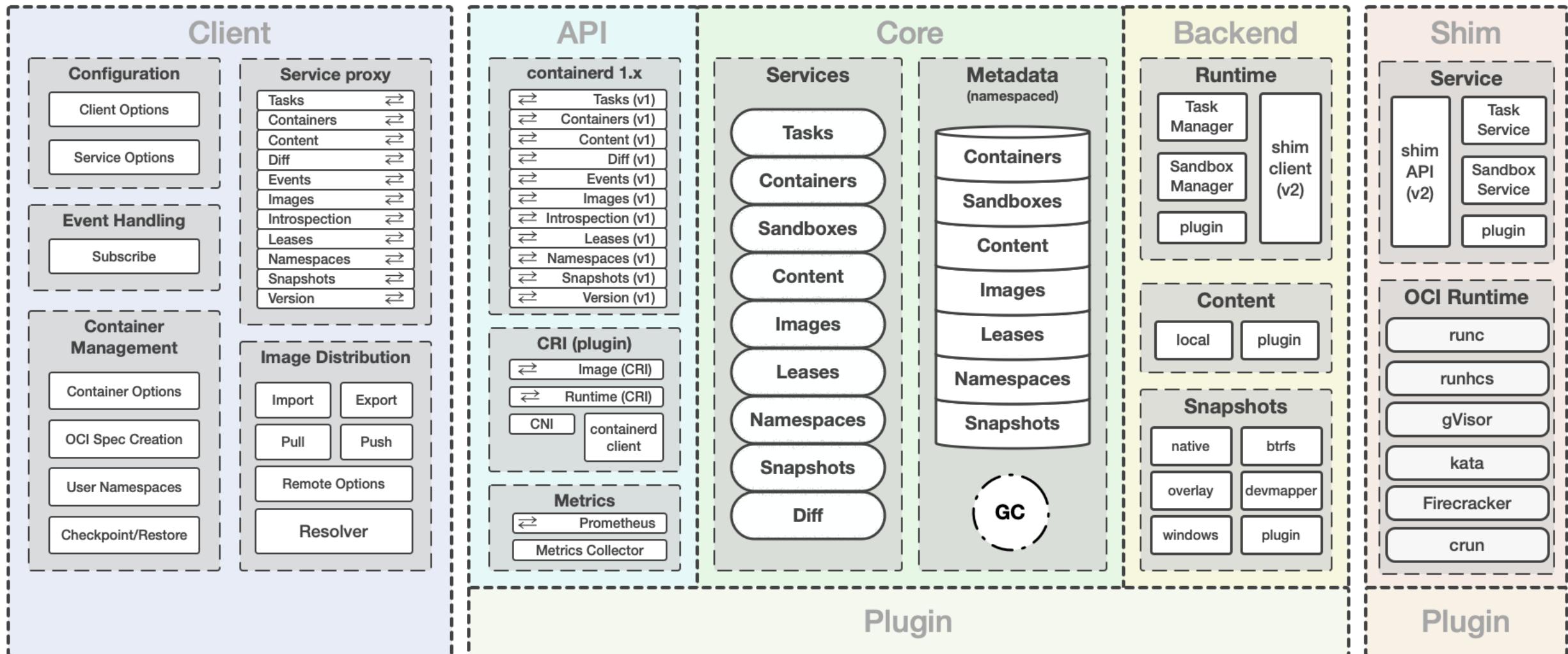


KubeCon



CloudNativeCon

North America 2021



# Architecture (Client)



KubeCon

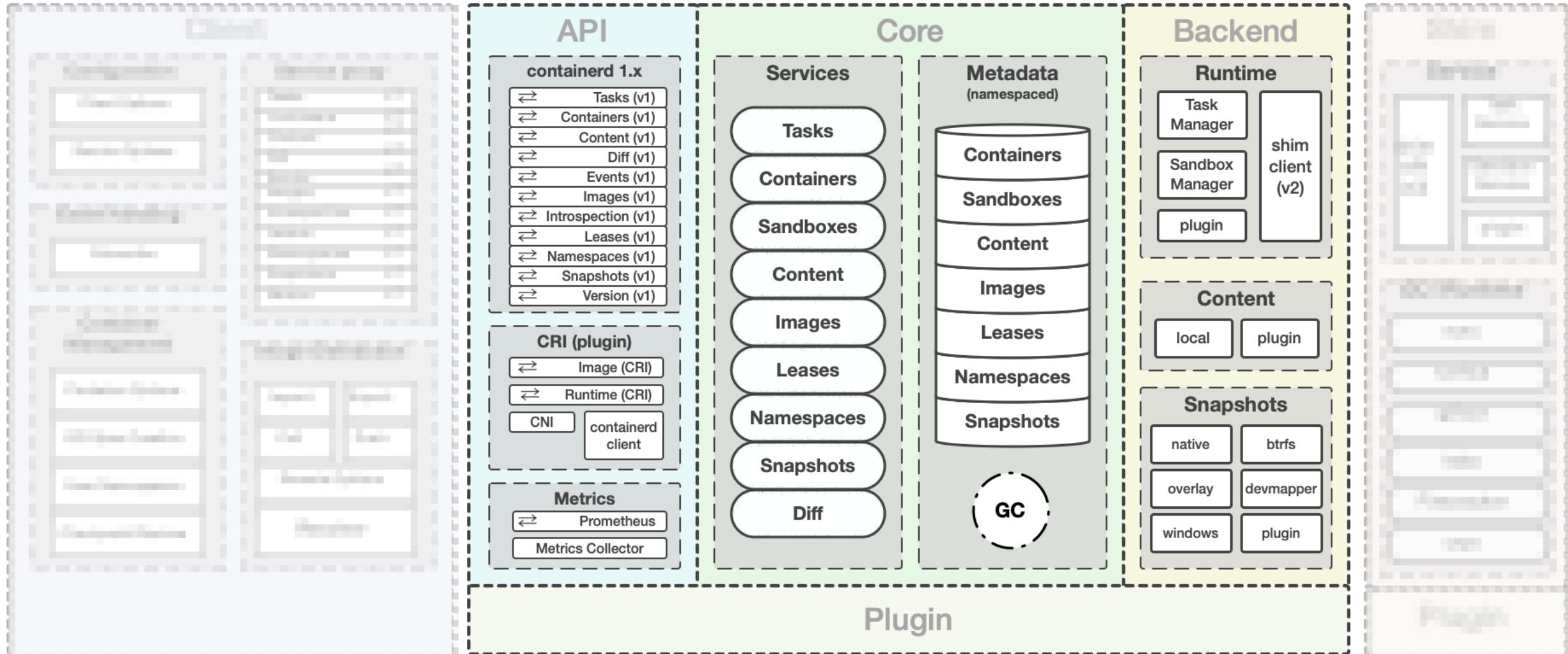


CloudNativeCon

North America 2021



# Architecture (Core)



# Architecture (Runtime/Shim)



# Pull Flow

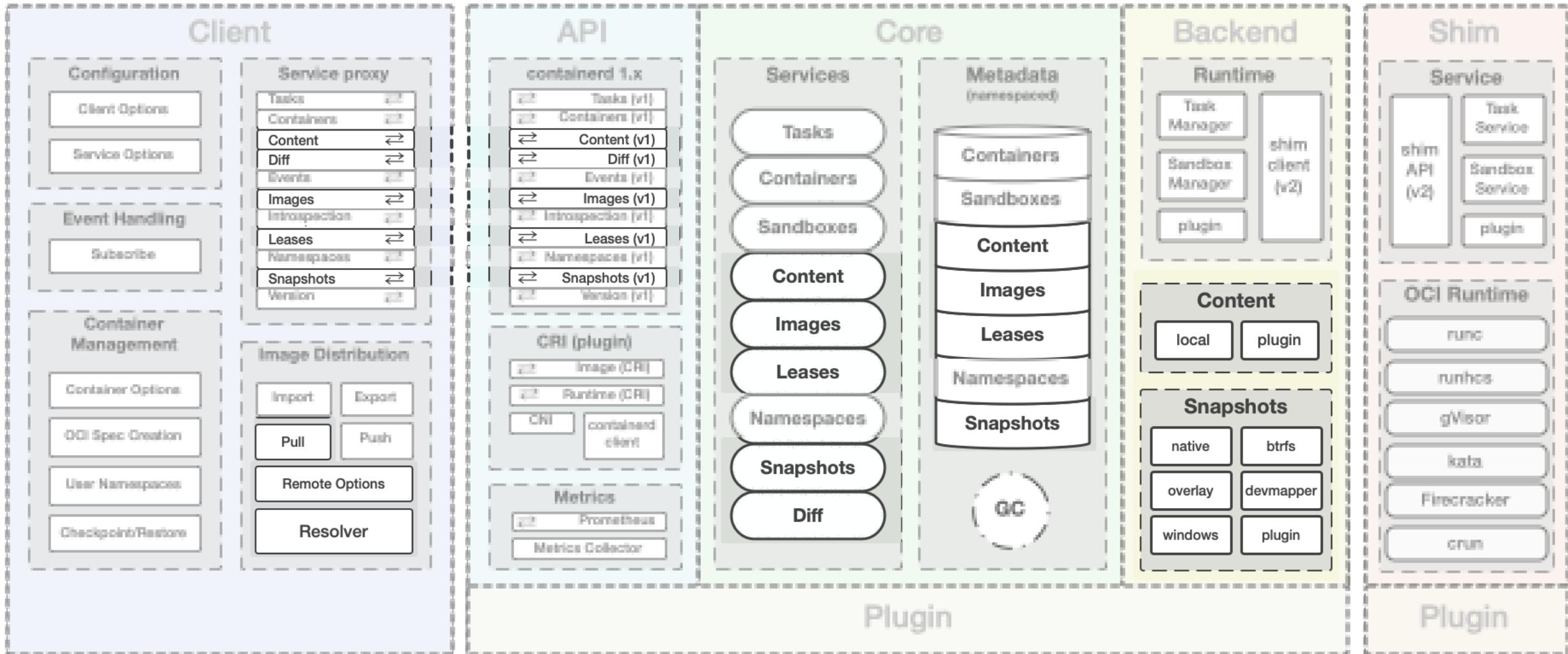


KubeCon

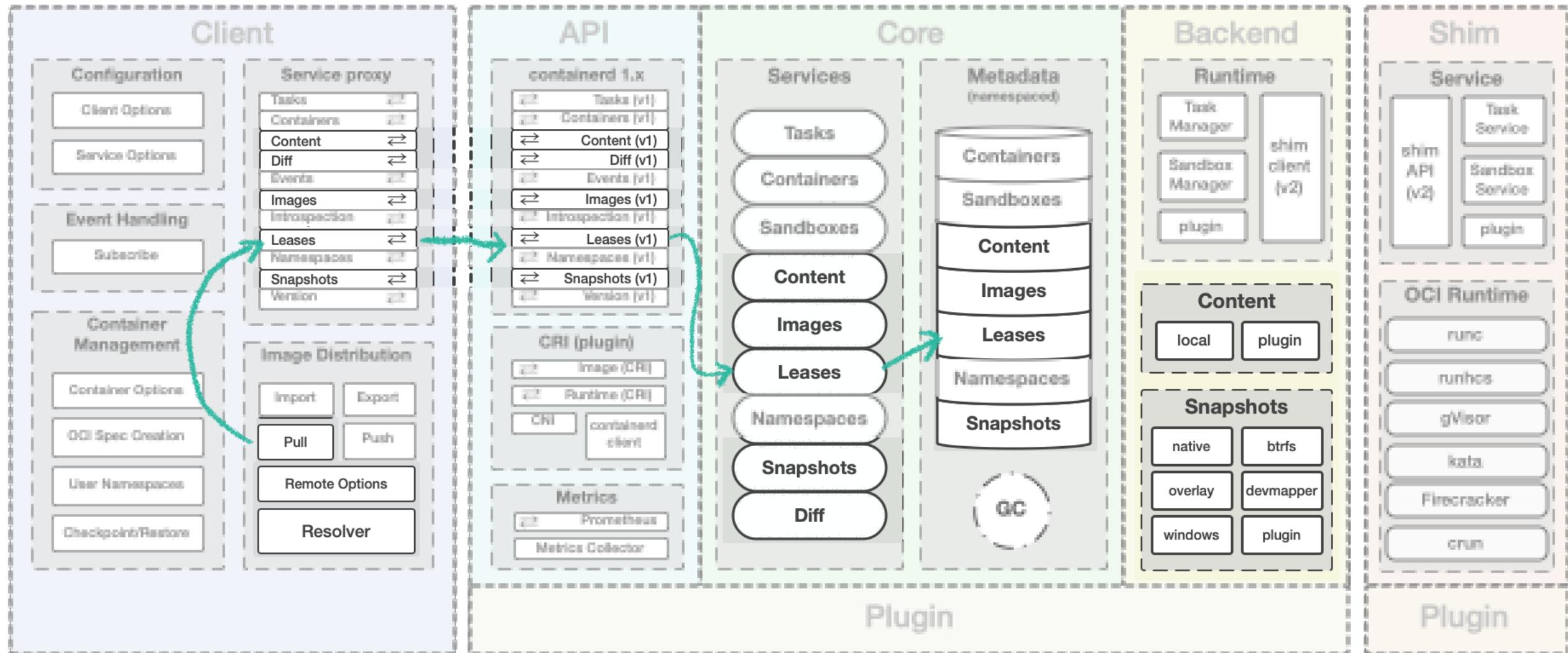


CloudNativeCon

North America 2021



# Pull Flow



# Pull Flow

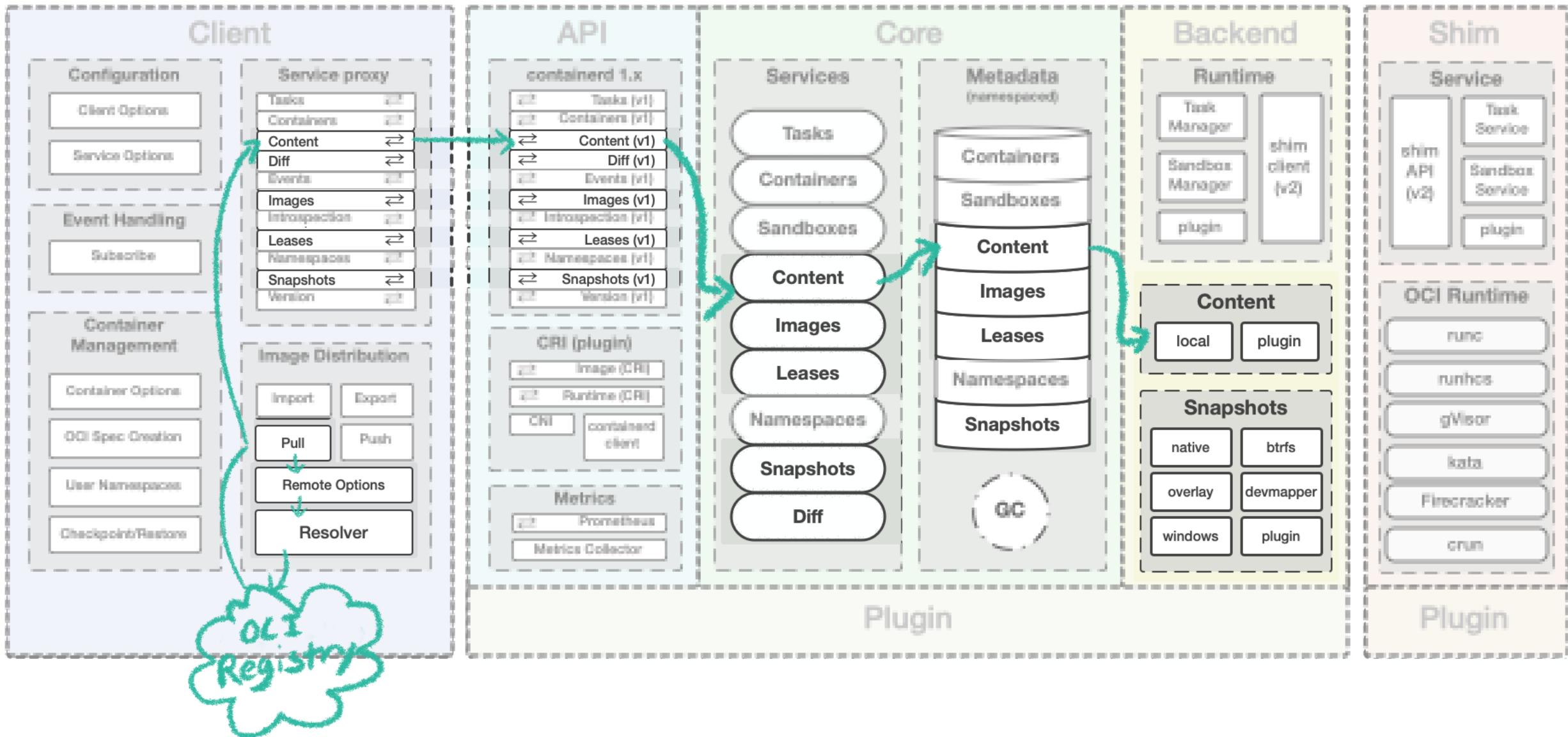


KubeCon



CloudNativeCon

North America 2021



# Pull Flow

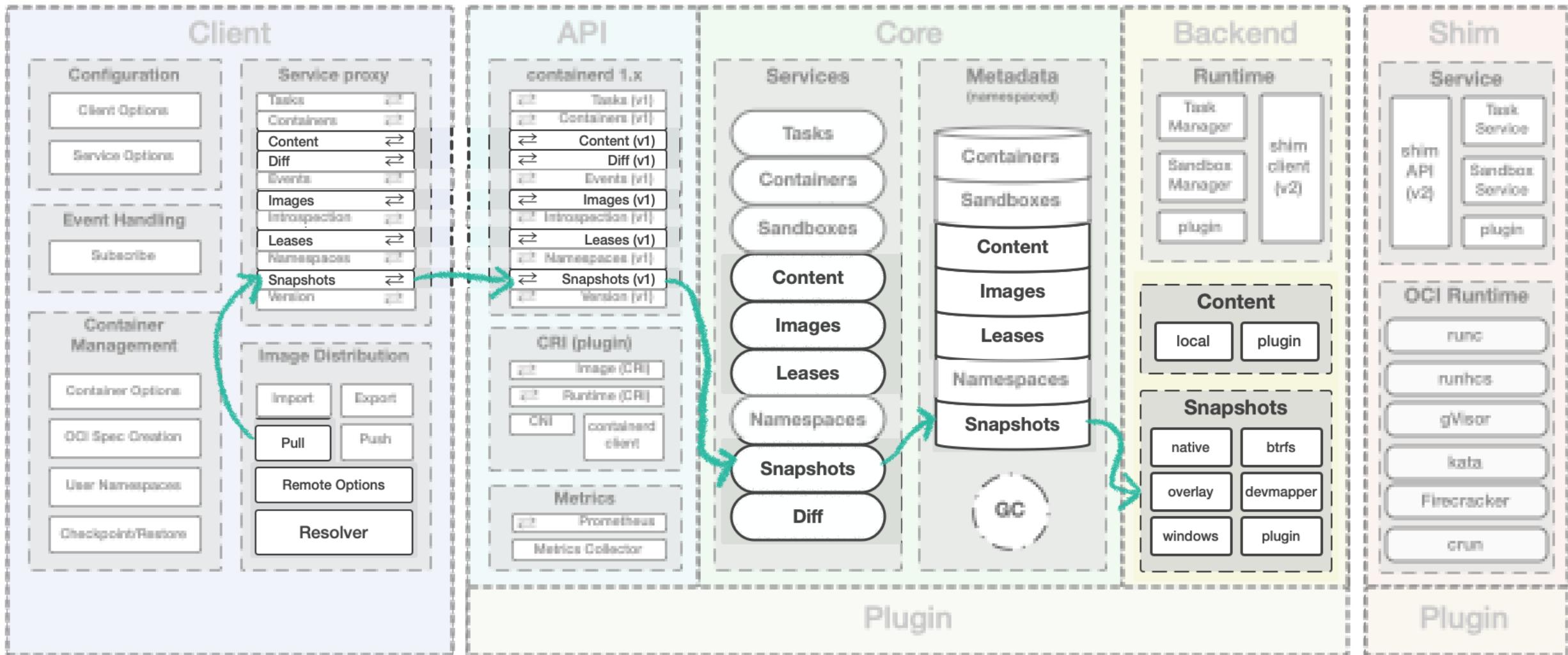


KubeCon



CloudNativeCon

North America 2021



# Pull Flow

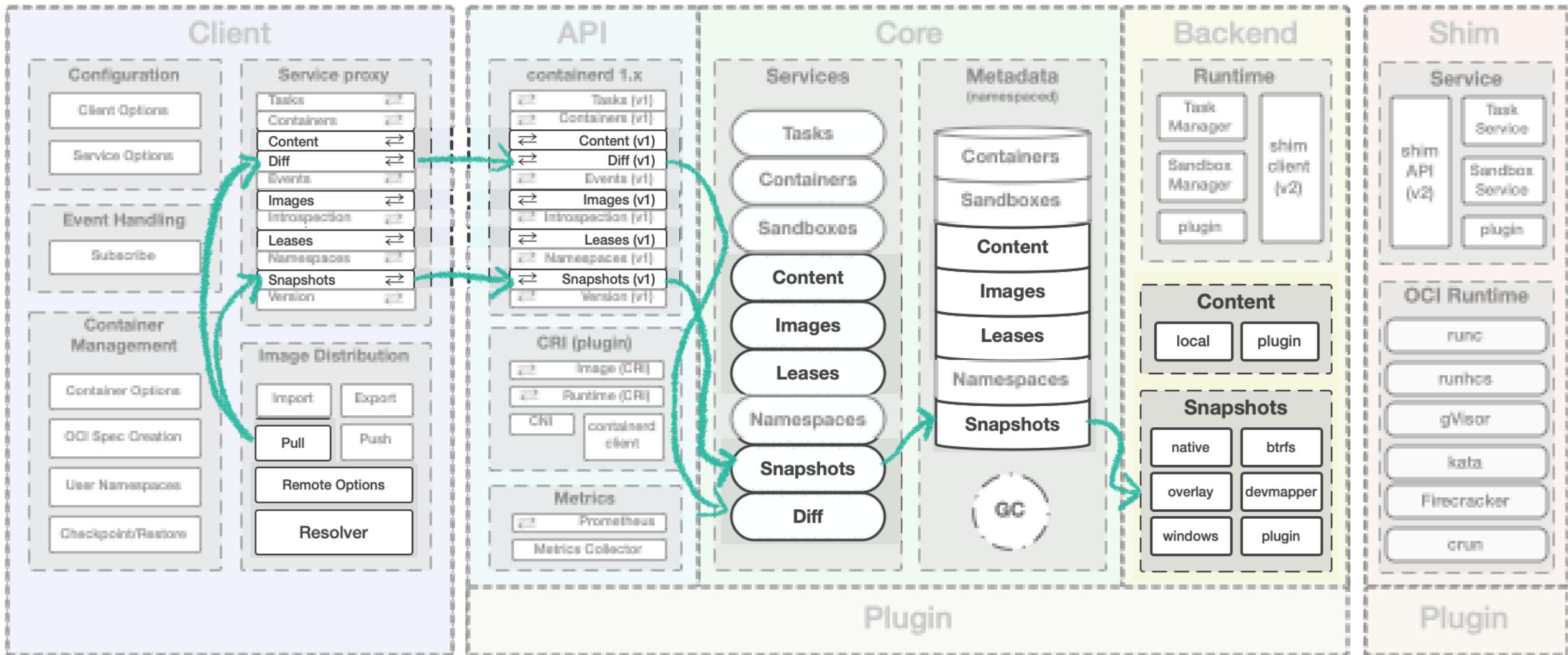


KubeCon



CloudNativeCon

North America 2021



# Pull Flow

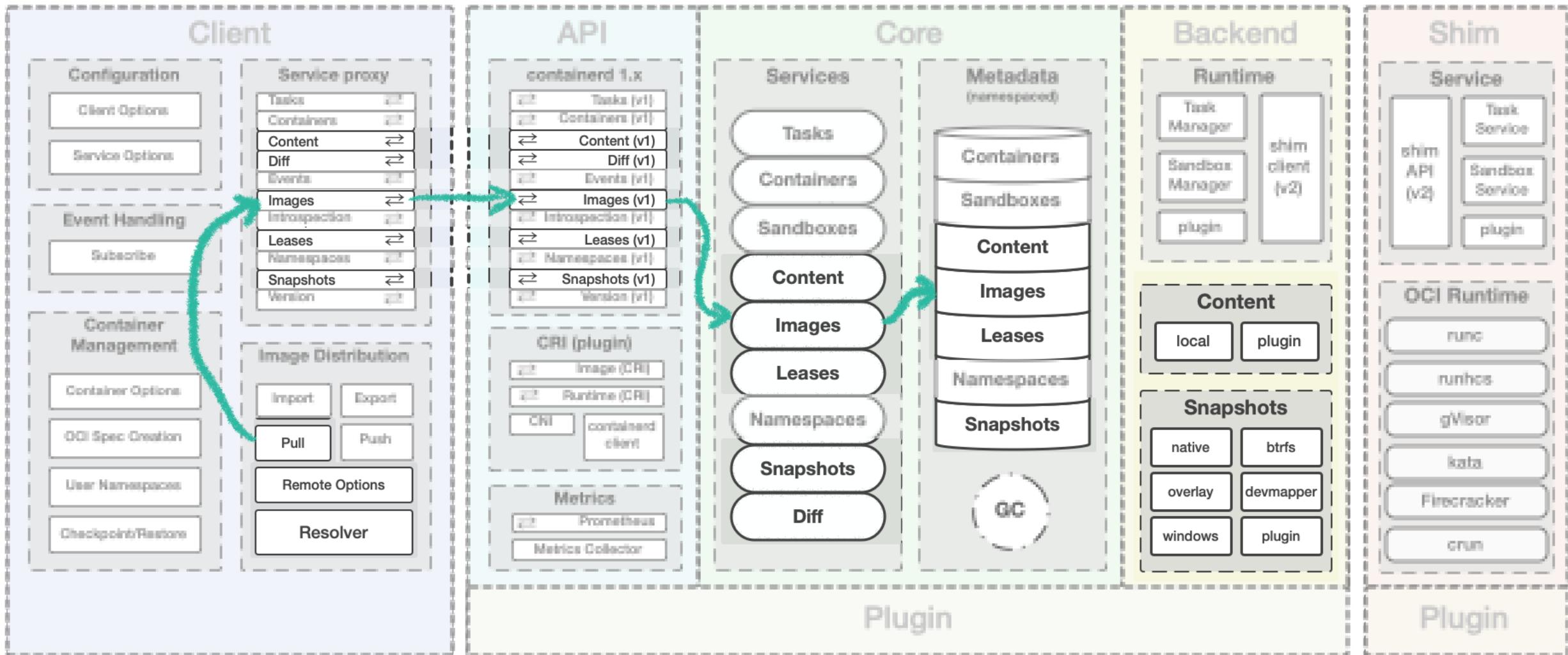


KubeCon



CloudNativeCon

North America 2021



# Run Flow

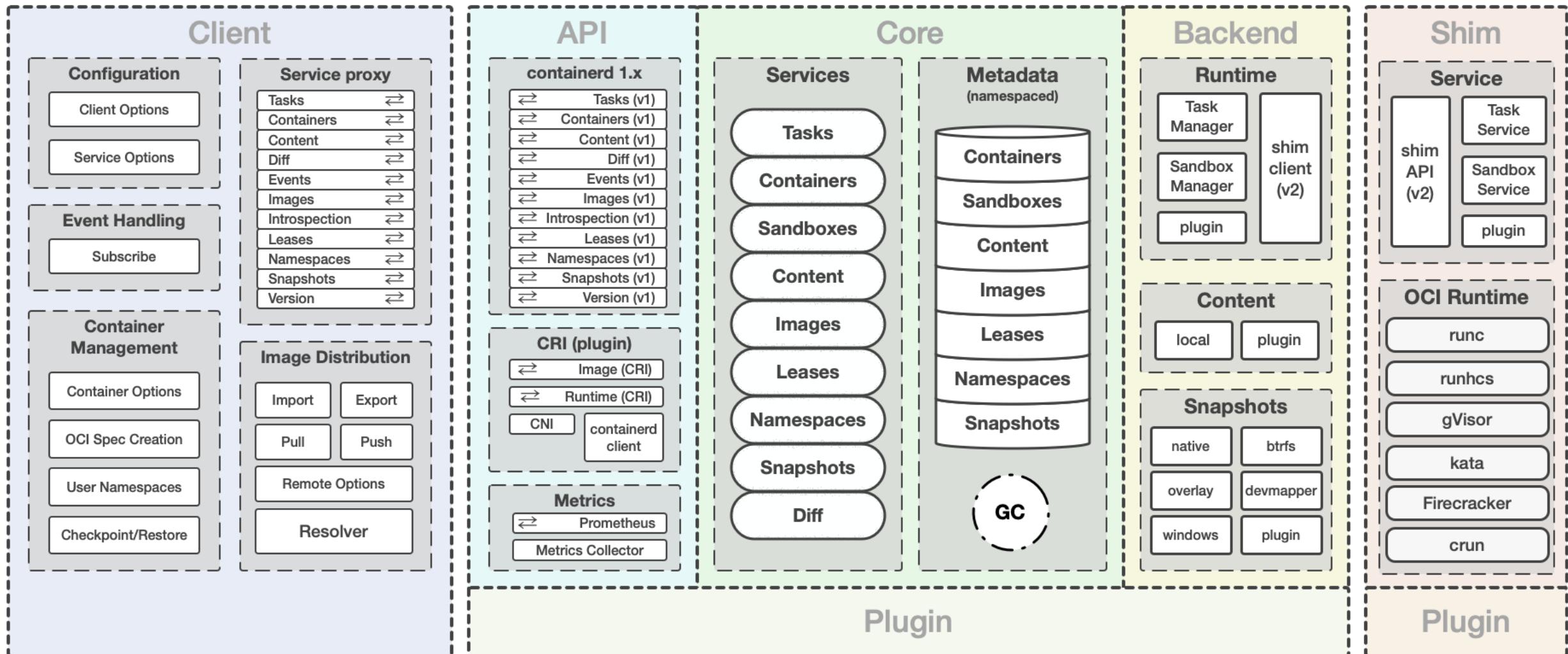


KubeCon



CloudNativeCon

North America 2021



# Run Flow

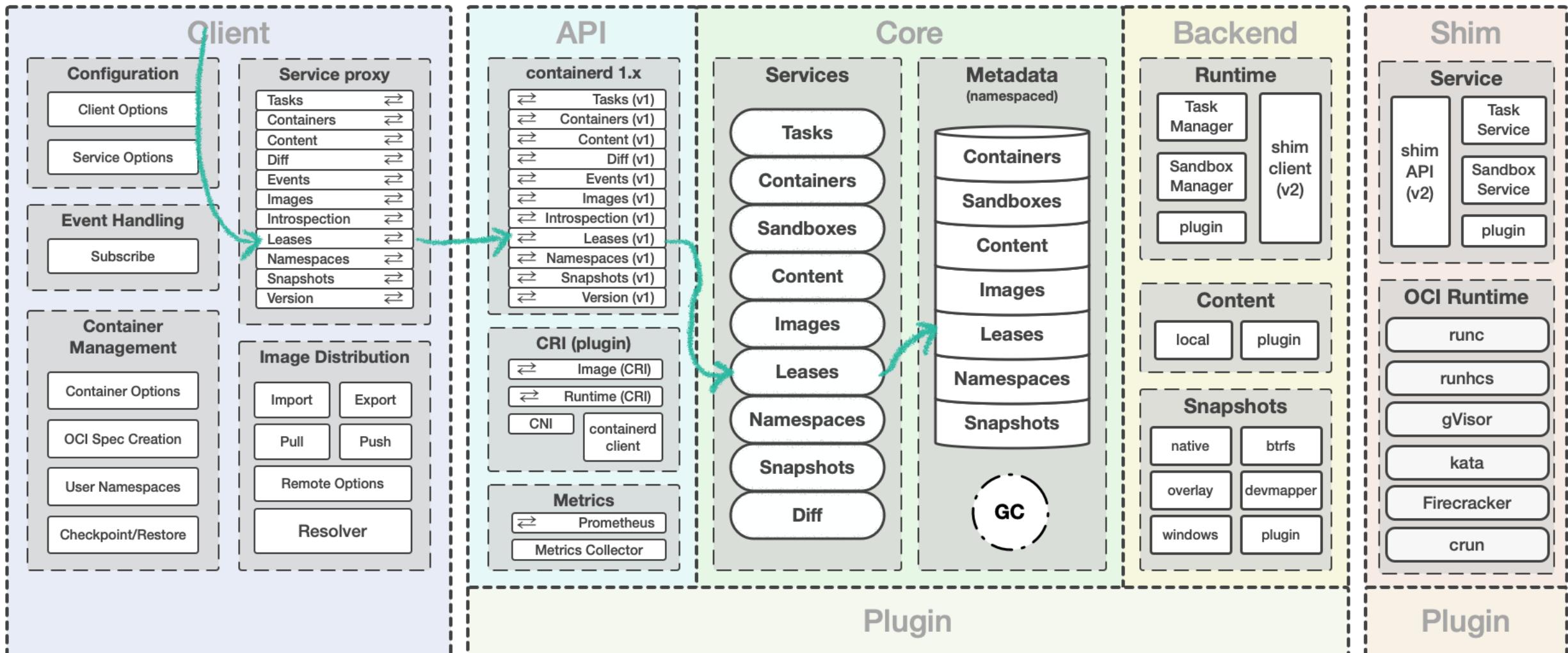


KubeCon



CloudNativeCon

North America 2021



# Run Flow

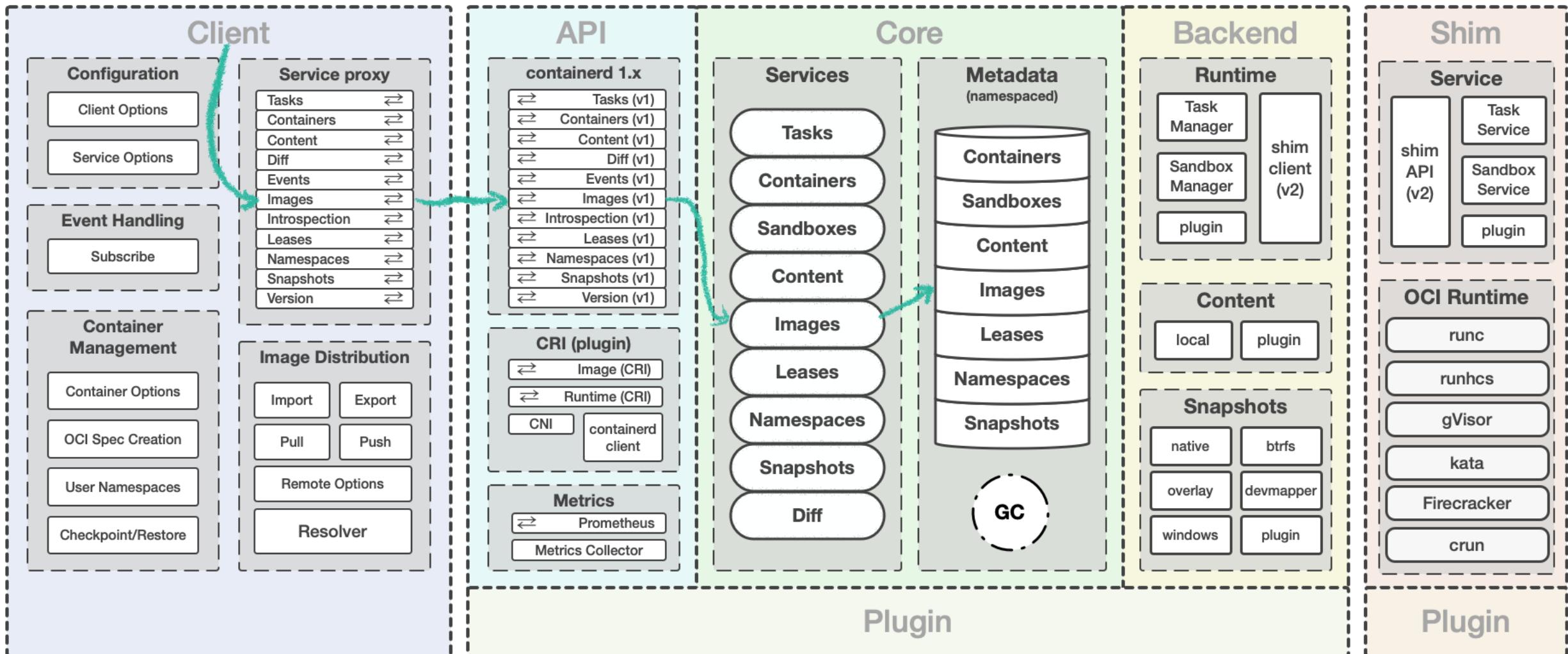


KubeCon



CloudNativeCon

North America 2021



# Run Flow

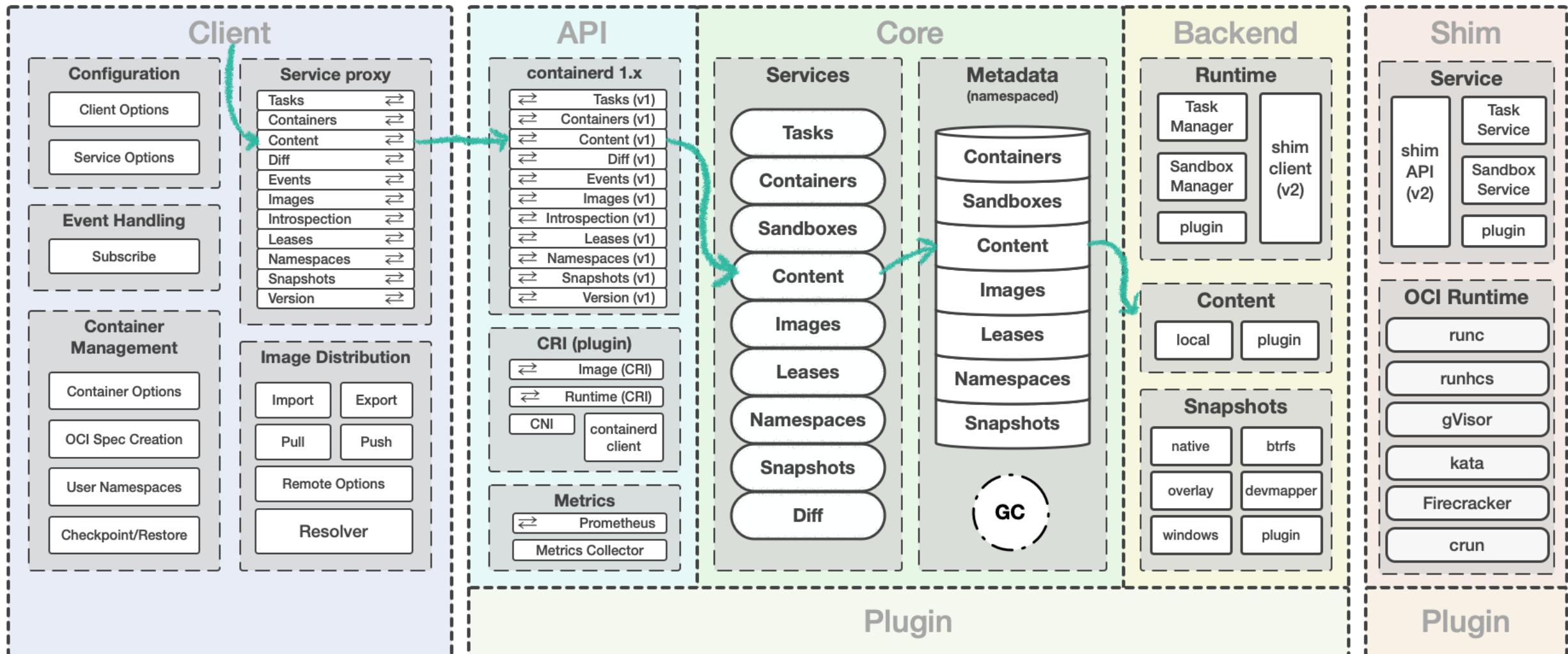


KubeCon



CloudNativeCon

North America 2021



# Run Flow

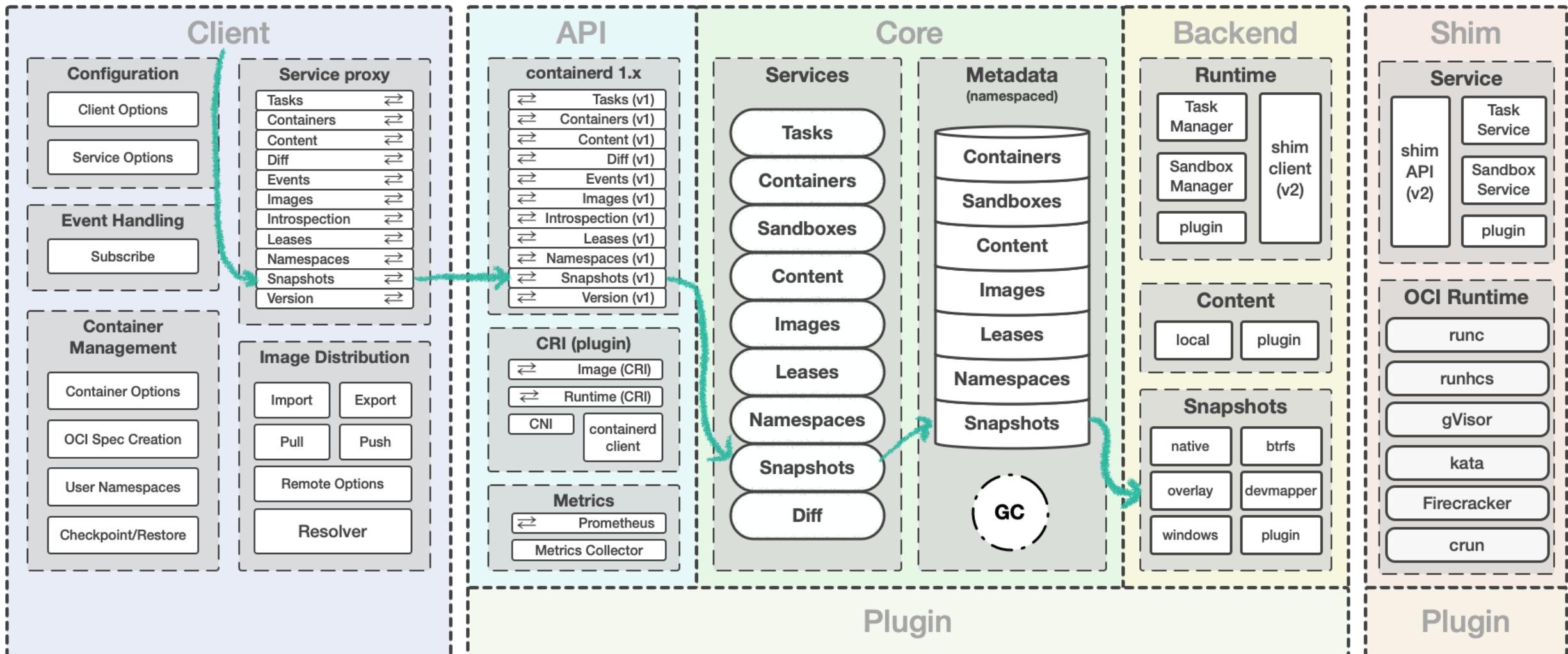


KubeCon



CloudNativeCon

North America 2021



# Run Flow

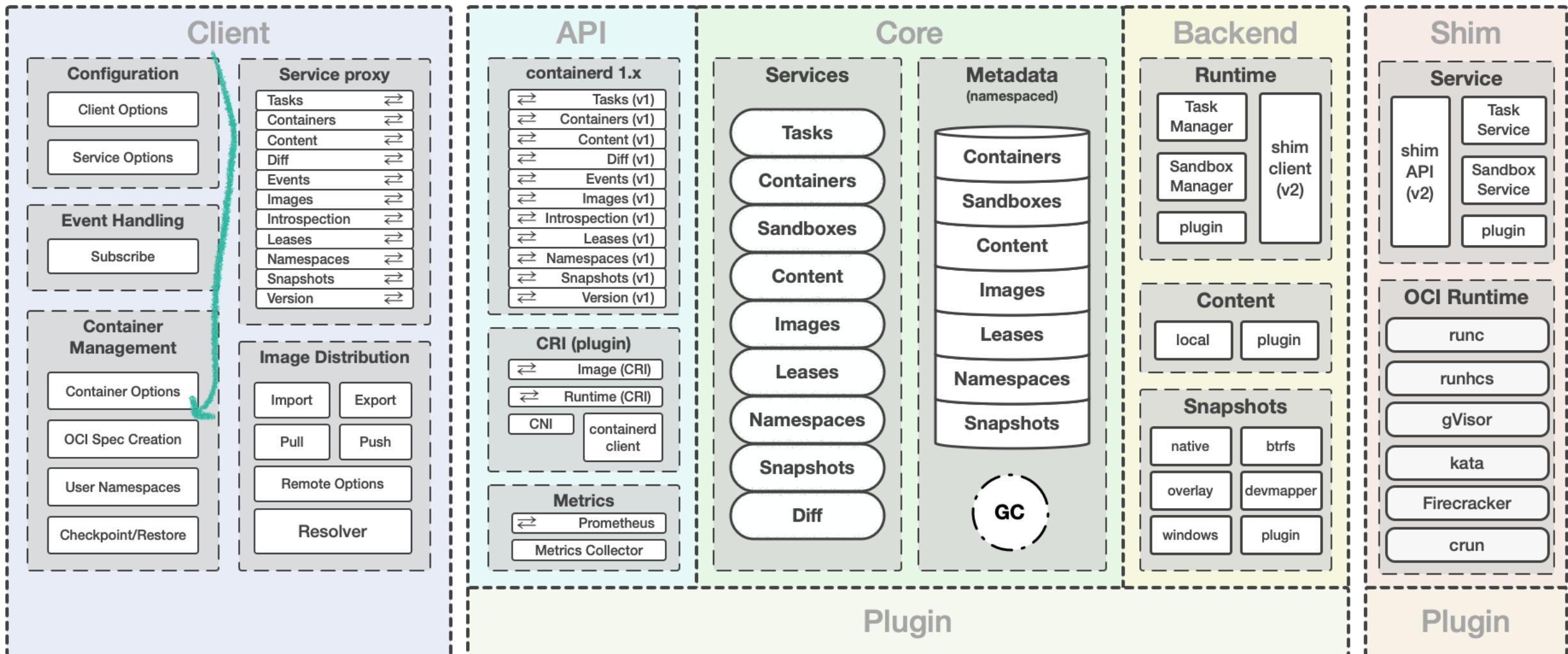


KubeCon



CloudNativeCon

North America 2021



# Run Flow

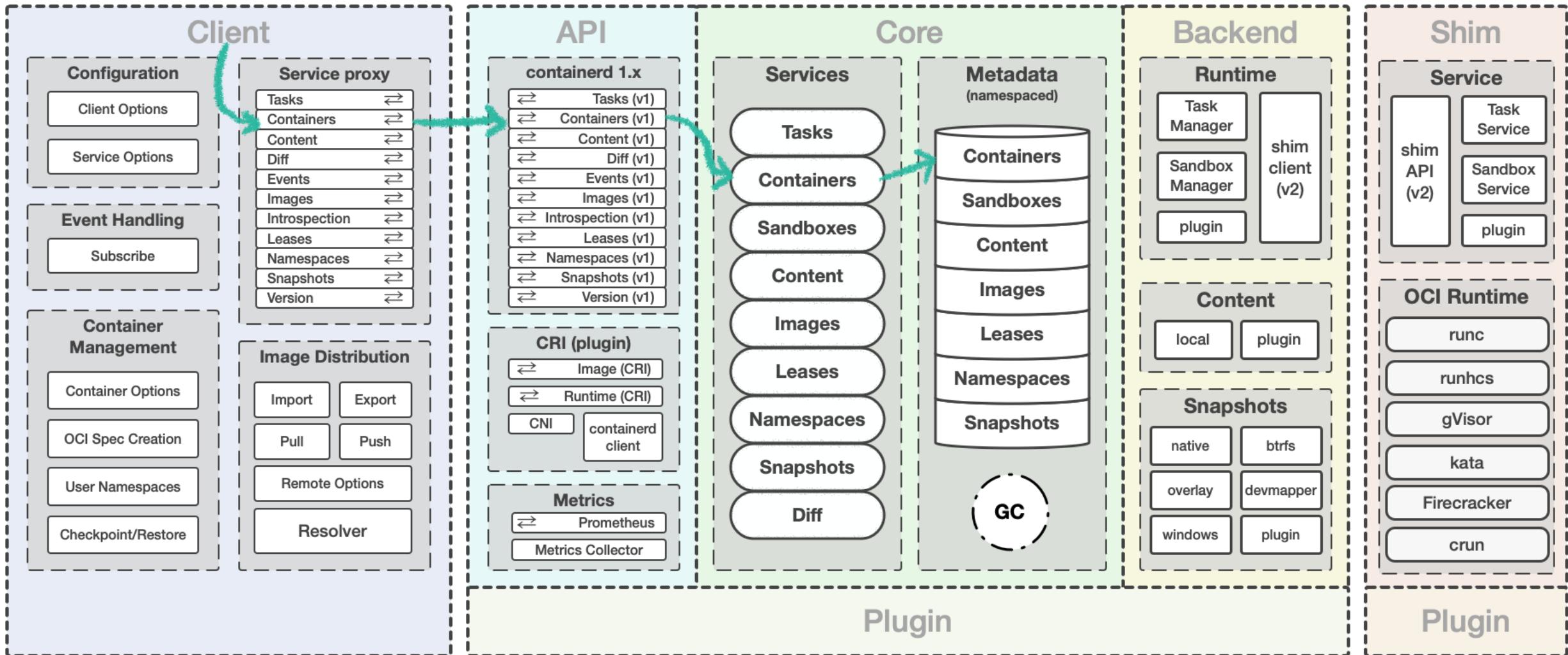


KubeCon



CloudNativeCon

North America 2021



# Run Flow

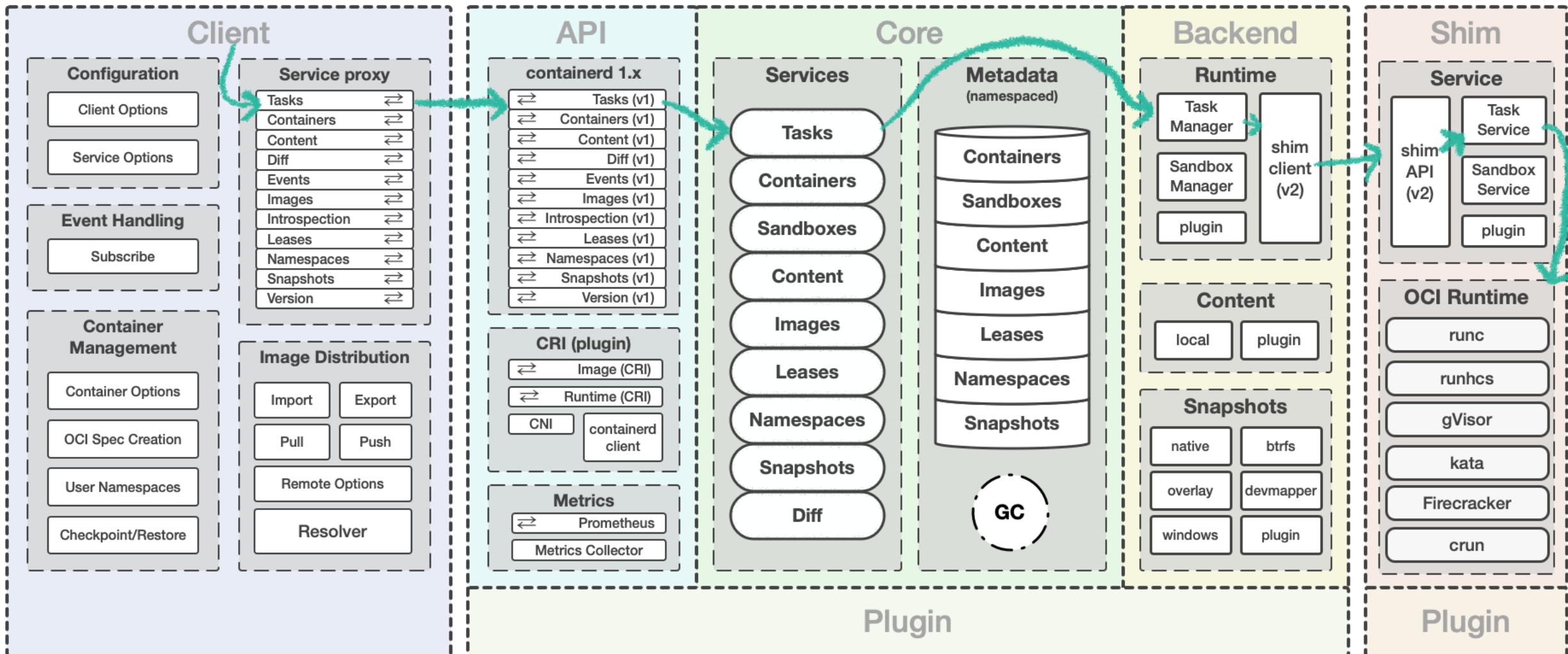


KubeCon



CloudNativeCon

North America 2021



# Run Flow

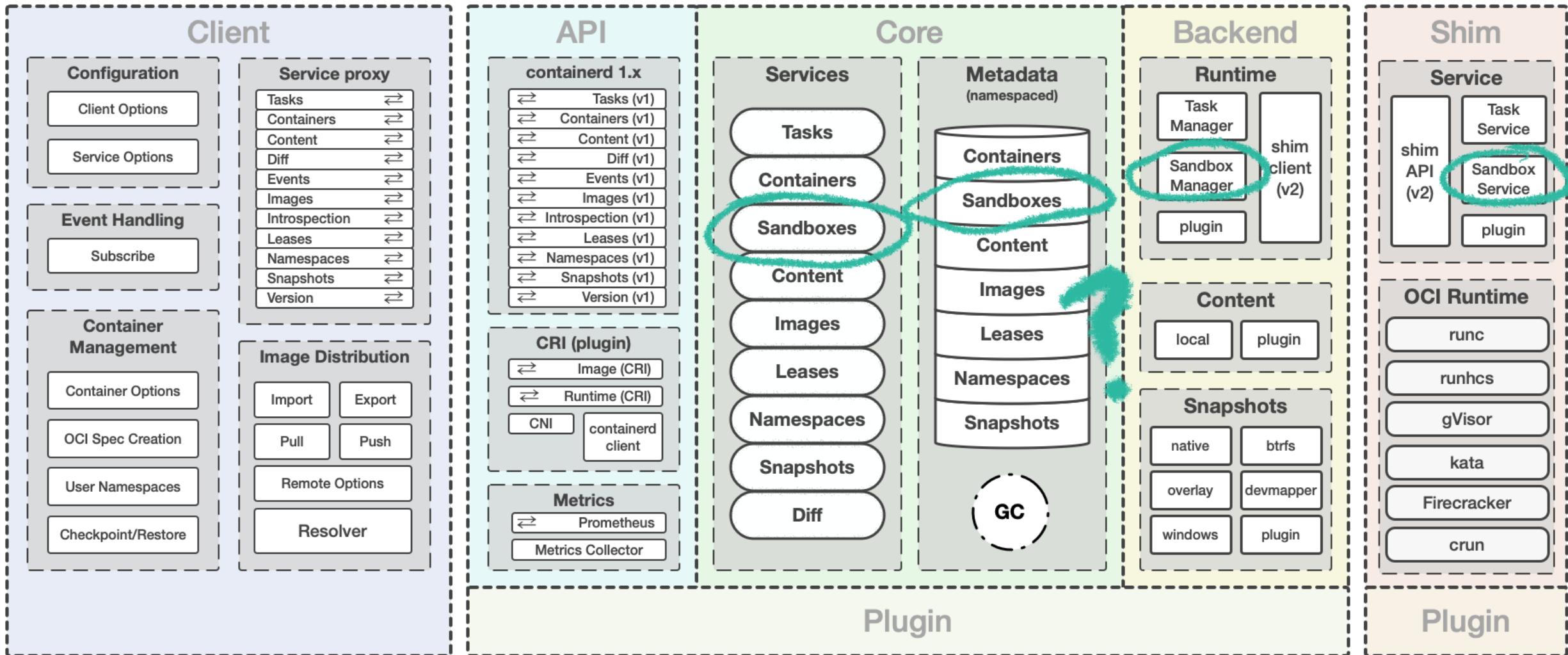


KubeCon



CloudNativeCon

North America 2021



- CRI arch refresher
  - Kubelet registers with the api-server to provide node level resource objects.
  - Kubelet uses the Container Runtime Interface (CRI) to communicate with a configured container runtime integration.
  - The container runtime services CRI requests to manage pods, containers, and container images.
  - Pod Spec may select which runtime engine the container runtime should use. (E.g. runc, crun, kata containers, gVisor, Firecracker, ...)

- Kubelet has deprecated the internal docker shim, which was partially based on CRI - work underway for an external CRI docker shim
- Changes underway
  - Test bucket migration away from internal docker shim (e2e, e2e-node, ... and critest(cri-tools)) - we need help here.  
**Shout out to the SIG-Node team for all their help.**
  - Simultaneous support for v1 and v1alpha CRI APIs (Kubelet and container runtimes with restart support for either Kubelet or the container runtime and through upgrades)

- Confidential Computing (CC) / Multi-Tenant Use Cases
  - Kubelet's Ensure Pod Container Images Process - ensures that a pod's container images are available on the node via CRI image services.
    - Pod pull image policies include "never pull", "pull if not present", and "always pull"
    - Pull image is performed based on policy, image status on the node, and if the registry requires authentication the POD's image pull secrets are used.
    - [KEP #1608](#) Ensure Secret Pulled Images (phase 1) is scheduled for k8s v1.23 and will ensure that a pod using a container image should have access to the image.
  - Kubernetes has a per runtime handler which is pod specified/scoped for choosing runtime engine config (kata, gVisor, runc-config-a, runc-config-b, etc..)
    - We also need that runtime handler to be passed over CRI image services such that image metadata, container snapshots, and storage can be scoped similarly

- Faster cycle times for PODs
  - Probes
    - Kubernetes uses liveness, readiness and startup probes to help in the managing of containers.
    - These probes operate at the per/second(s) level
    - Certain workloads need us to move from per/second(s) probes to fine grained per/millisecond(s) probes KEP ([wip](#)) code ([wip](#))
  - Kubelet currently uses a lazy update cycle process for determining state for PODs, Containers, and Container Images
    - Discussions are underway for introducing a KEP to optionally add a subscription service for POD/Container status changes and possibly also for container image/storage updates

# Upcoming containerd 1.6+

- Big focus on shims
  - More runtime features
  - Better microVM support
  - New ways to manage images
- Revisit runtime for new use cases
  - How to support new use cases?
  - . . . and keep things flexible
  - . . . and keep things backward compatible

# New shim features

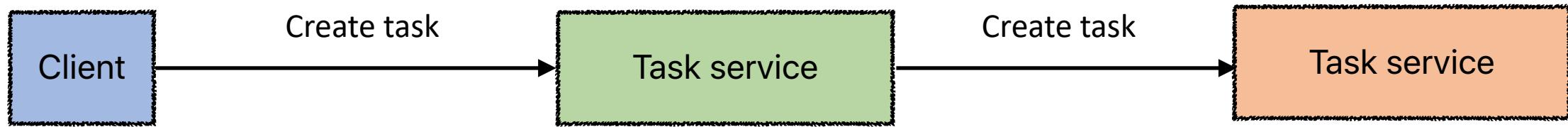
- Sandbox API
  - Better microVM lifecycle/runtime support
- Confidential containers
  - Hide container image contents from the host
- Port forwarding
- New custom APIs?

# New shim runtime

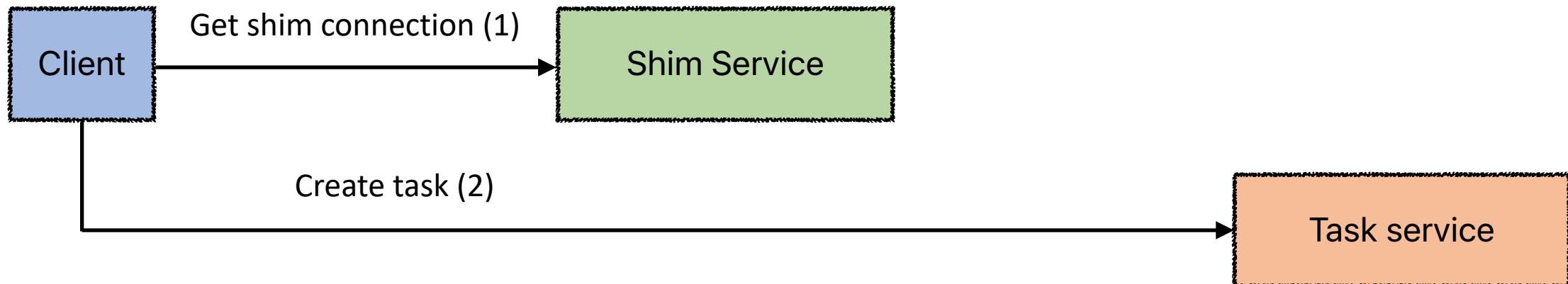
- Plugin support
  - Extend shim with new plugins
- New Shim service
  - Manage shims independently from Tasks
- New client API to manage Tasks
  - Use new shim API, but keep backward compatibility
  - Support new shim features from client

# Shim service call flow

1.5



1.6+



# Shim service

- Manage shim lifecycle
  - Start
  - Stop
  - Clean
- Foundation for higher level services
  - Tasks
  - Sandbox
- Call flow defined by Client

```
type ShimService interface {
    StartShim(...) // call `shim --start`
    DeleteShim(...) // call `shim --delete`
    GetConn(shimId string) // get connection
}
```

# Sandbox API

- Manage microVMs
  - Launch VM before Container
  - Resize VM instance
  - Suspend/resume
- Needs shim service as base
  - Different shim lifecycle
  - Driven by containerd client

# Thank you!