# Introduction

Alcor is a cloud native SDN platform that aims to provides high availability, high performance, and large-scale virtual networking control plane and management plane at a high resource provisioning rate.

| P | A | S | E |
|---|---|---|---|
| **Performance** | **Availability** | **Scalability** | **Extensibility** |
| • Throughput-optimal design to allow batched provisioning of network resources<br>• Fast provisioning path to support time-critical applications such as serviceless | • Always-on control plane without a single point of failure<br>• Cross-AZ resilience for services and data<br>• Fault-tolerant design with multiple resource provisioning paths | • Management of large numbers of network resources<br>• Scale to half a million hosts and tens of millions network ports | • Unified network management of both VMs and containers<br>• Plug-able model to support various implementations of data plane |

# Architecture Overview

## Alcor Management plane

- Expose REST APIs to clients
- Provide partition management, end-user monitoring, and billing & quota management

## Alcor Control Plane

- Offer multiple network management services including VPC, ELB and EIP
- Support multiple provisioning path including fast path, normal path and rescue path
- Drive network configurations to on-host Alcor agents in scale
- Alcor agents program data plane to establish network connection between containers and VMs in the same VPC

# Cloud-Native Control Plane

**Powered by Kubernetes**
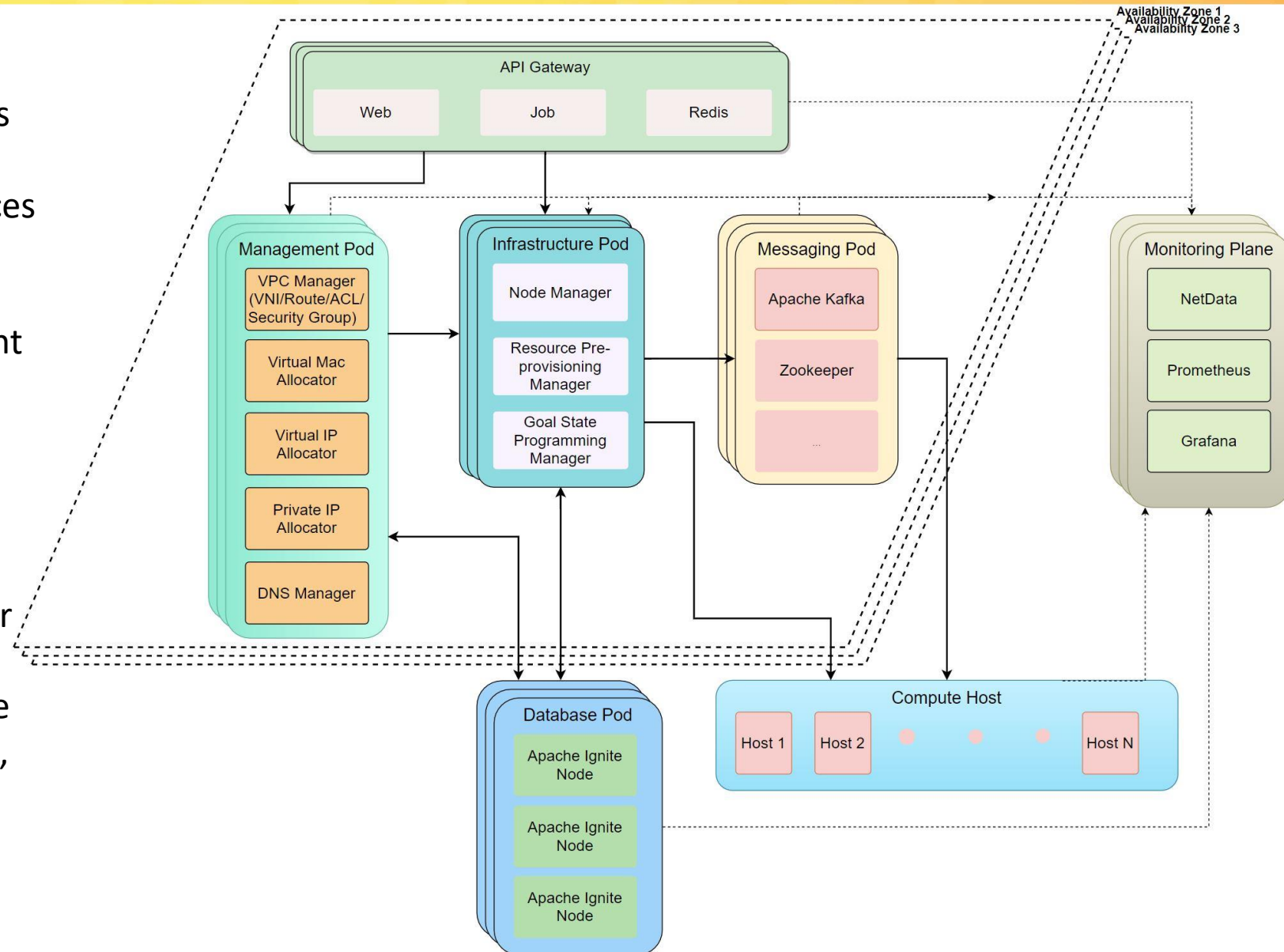
- Each controller instance is an Kubernetes application
- Each application contains multiple services
    - Customer resource management service
    - Infrastructure resource management service
    - Database service
    - Messaging service

**Distributed Micro-Services Architecture**

- Use Istio to secure, connect, and monitor control plane micro-services
- Fine-grained control of service-to-service communication including load balancing, retries, failovers, and rate limits.

Availability Zone 1
Availability Zone 2
Availability Zone 3

**API Gateway**

| Web | Job | Redis |

**Management Pod**

- VPC Manager (VNI/Route/ACL/ Security Group)
- Virtual Mac Allocator
- Virtual IP Allocator
- Private IP Allocator
- DNS Manager

**Infrastructure Pod**

- Node Manager
- Resource Pre-provisioning Manager
- Goal State Programming Manager

**Messaging Pod**

- Apache Kafka
- Zookeeper
- ...

**Monitoring Plane**

- NetData
- Prometheus
- Grafana

**Database Pod**

- Apache Ignite Node
- Apache Ignite Node
- Apache Ignite Node

**Compute Host**

| Host 1 | Host 2 | • • • | Host N |

# Throughput-Optimal Design

Focus on throughout optimization on every system layer

| API | Controller | Controller-Agent-Comm | Host Agent |
|-----|-----------|----------------------|-----------|
| • Allow group port creation for one network (e.g. create 1,000 ports) with one POST call<br>• Unified APIs to support network resource management for both VMs and containers | • Implicit batching for database write and network programming<br>• Insert data to database in a bulk mode using JDBC driver<br>• Bundle network configuration updates for the same host and drive them down to host agents in one shot | • Bundle network configuration update in the same host<br>• One configuration message could include various combinations of resource updates.<br>  ○ Multiple instances of resources (e.g. creating 10 ports for one VPC)<br>  ○ Multiple types of resources (e.g. updating 1 port+ creating 2 security groups)<br>  ○ Across VPC/subnet boundaries | • Parallel network setup on the host (creation of veth pairs and namespaces) and port programming to data plane<br>• Achieve 1000+ port RPM on the host with Mizar data plane |

Not Secure    10.213.43.166:19999/#menu_cgroup_endpoint_host_submenu_...

fw0009099

# endpoint host

Container resource utilization metrics. Netdata reads this information from **cgroups** (abbreviated from **control groups**), a Linux kernel feature that limits and accounts resource usage (CPU, memory, disk I/O, network, etc.) of a collection of processes. cgroups together with **namespaces** (that offer isolation between processes) provide what we usually call: **containers**.

| CPU | Memory | Read Disk I/O | Write Disk I/O | Received eth0 | Sent eth0 |
|-----|--------|---------------|----------------|---------------|-----------|
| 0.0 | 0.0126 | 0 | 0 | 432.1 | 592.1 |

## cpu

CPU Usage within the limits for cgroup endpoint_host (cgroup_endpoint_host.cpu_limit)

Thu, Nov 14, 2019
14:55:50

percentage
— used          0.0000

First endpoint host to be programmed.

CPU Usage (4000% = 40...                    ...point_host.cpu)

Thu, Nov 14, 2019
14:55:50

percentage

fw0009099

## net eth0

Bandwidth (cgroup_endpoint_host.net_eth0)

Thu, Nov 14, 2019
14:55:44

bits/s
— received    431
— sent        -591

Packets (cgroup_endpoint_host.net_packets_eth0)

Thu, Nov 14, 2019
14:55:44

pps
— received    1.00
— sent        -1.00

---

ericli — root@fw0009099: ~ — ssh root@10.213.43.166 — 89×27

```
root@172.17.0.7 / $ tail -f /var/log/syslog
Nov 14 14:43:34 2b29233632af kernel: [2073374.065071] docker0: port 47(vethd350d1c) enter
ed forwarding state
Nov 14 14:43:37 2b29233632af kernel: [2073377.375785] device veth401d3ff entered promiscu
ous mode
Nov 14 14:44:05 2b29233632af kernel: [2073405.948551] docker0: port 56(veth00463fa) enter
ed blocking state
Nov 14 14:44:09 2b29233632af kernel: [2073409.771339] docker0: port 57(vethd3ff720) enter
ed forwarding state
Nov 14 14:44:13 2b29233632af kernel: [2073413.145394] device veth18f431b entered promiscu
ous mode
Nov 14 14:45:04 2b29233632af kernel: [2073464.599646] docker0: port 72(vethe6577f0) enter
ed blocking state
Nov 14 14:45:52 2b29233632af kernel: [2073512.743385] docker0: port 85(veth1dec06f) enter
ed forwarding state
Nov 14 14:47:08 2b29233632af kernel: [2073588.157111] IPv6: ADDRCONF(NETDEV_UP): veth2019
d16: link is not ready
Nov 14 14:47:53 2b29233632af AliothControlAgent[140]: Network Control Agent started...
Nov 14 14:47:53 2b29233632af AliothControlAgent[140]: Server listening on 0.0.0.0:50001
```

ericli — root@fw0009099: ~ — ssh root@10.213.43.166 — 80×24

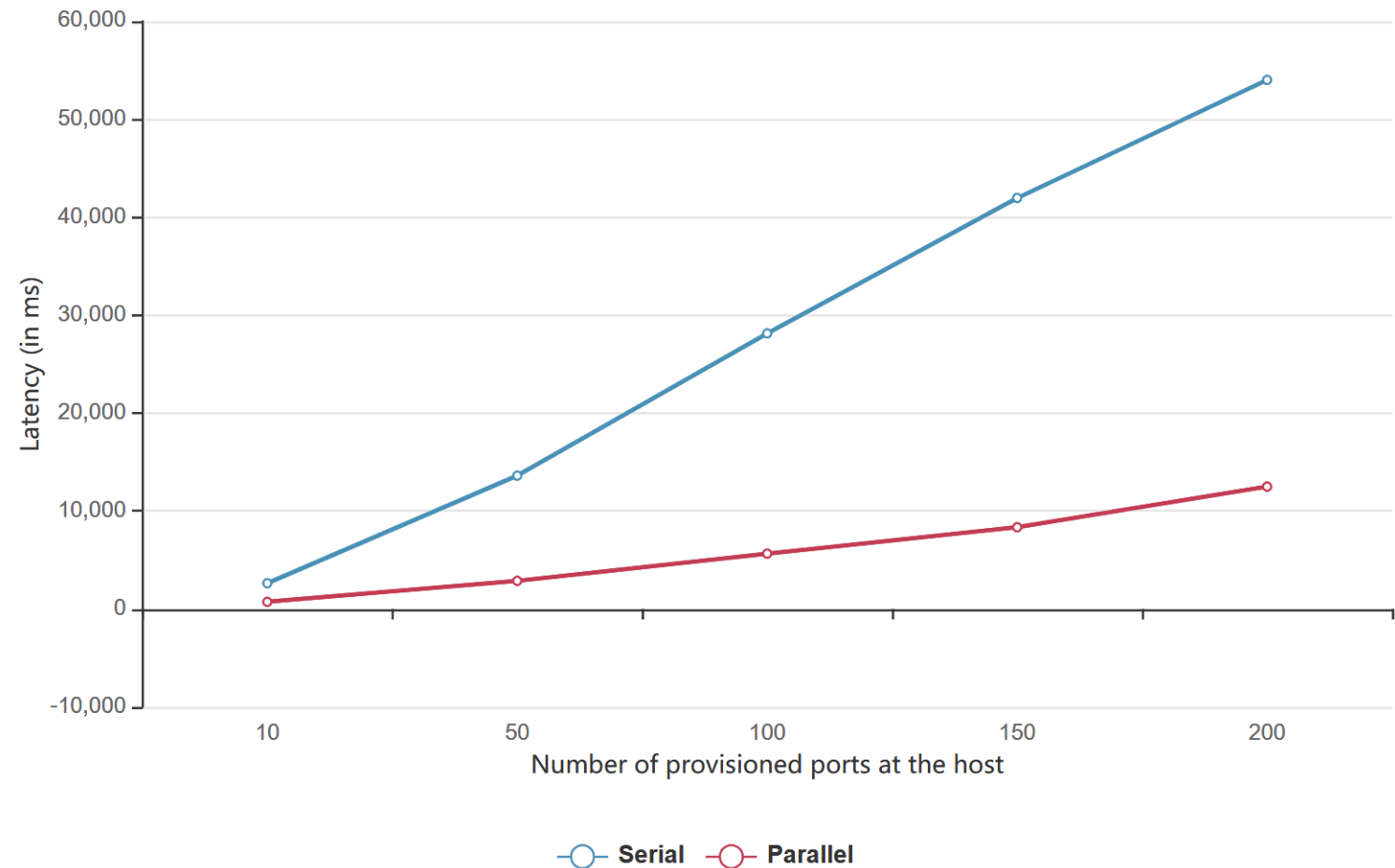```
root@172.17.0.7 / $
```

# Agent Programming Latency

**Serial Host Programming**

- Create ports one by one including network configuration and data plane programming
- Latency increases significantly for a large number of ports

**Parallel Host Programming**

- Create multiple threads for network configuration (veth pair and namespace creation etc.)
- Support programming of data plane in a single-threaded mode or multi-threaded mode

## Provisioning Latency Improvement at the Host



Latency (in ms) vs Number of provisioned ports at the host
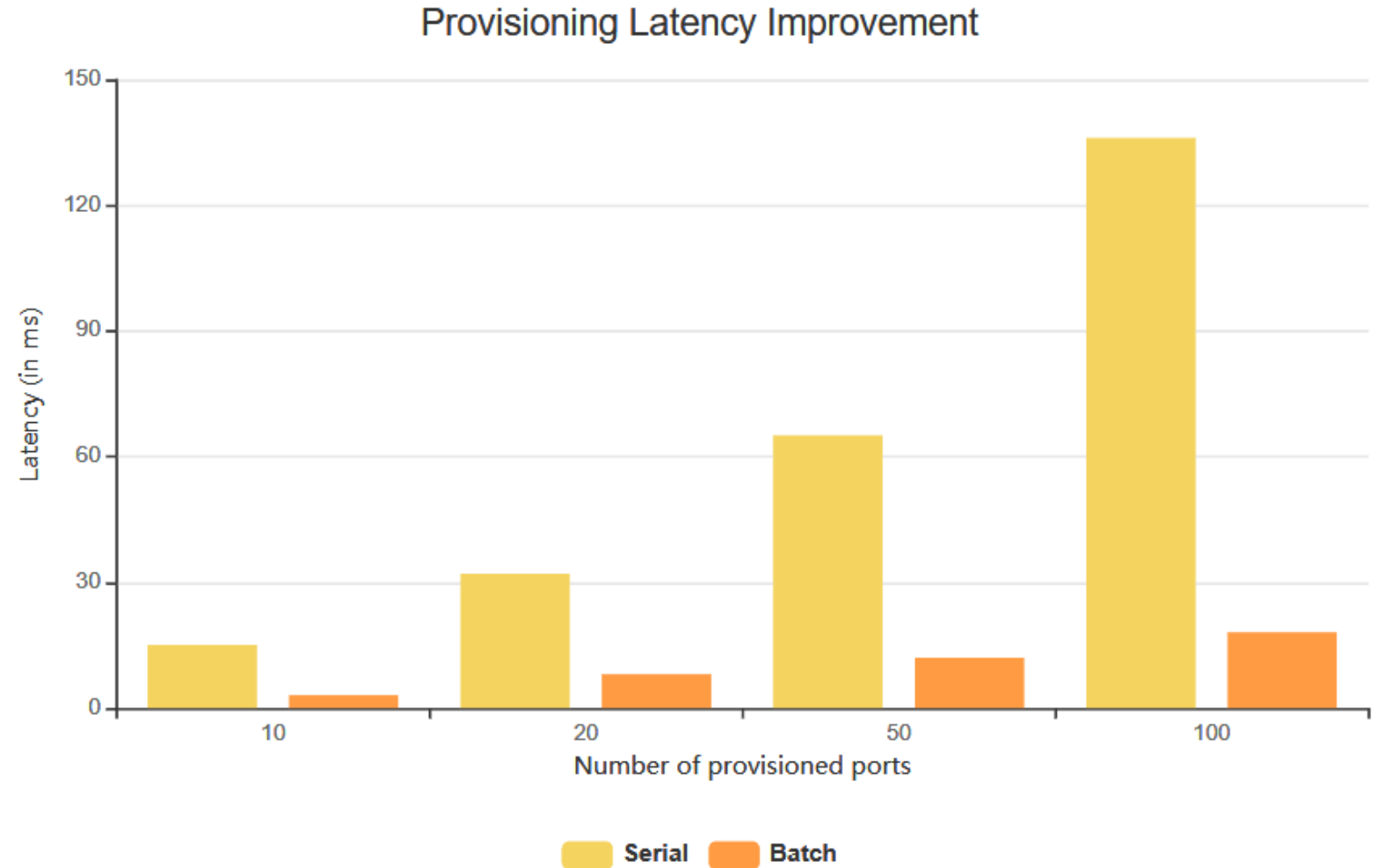
- Serial
- Parallel

# E2E Provisioning Latency

### Serial Provisioning

- Ports are created and distributed one by one
- Unable to scale to a large number of ports

### Batch Provisioning

- Created with a single API call (post portgroup)
- Distributed to hosts in a batched and parallel manner



Provisioning Latency Improvement

Latency (in ms) vs Number of provisioned ports (10, 20, 50, 100)
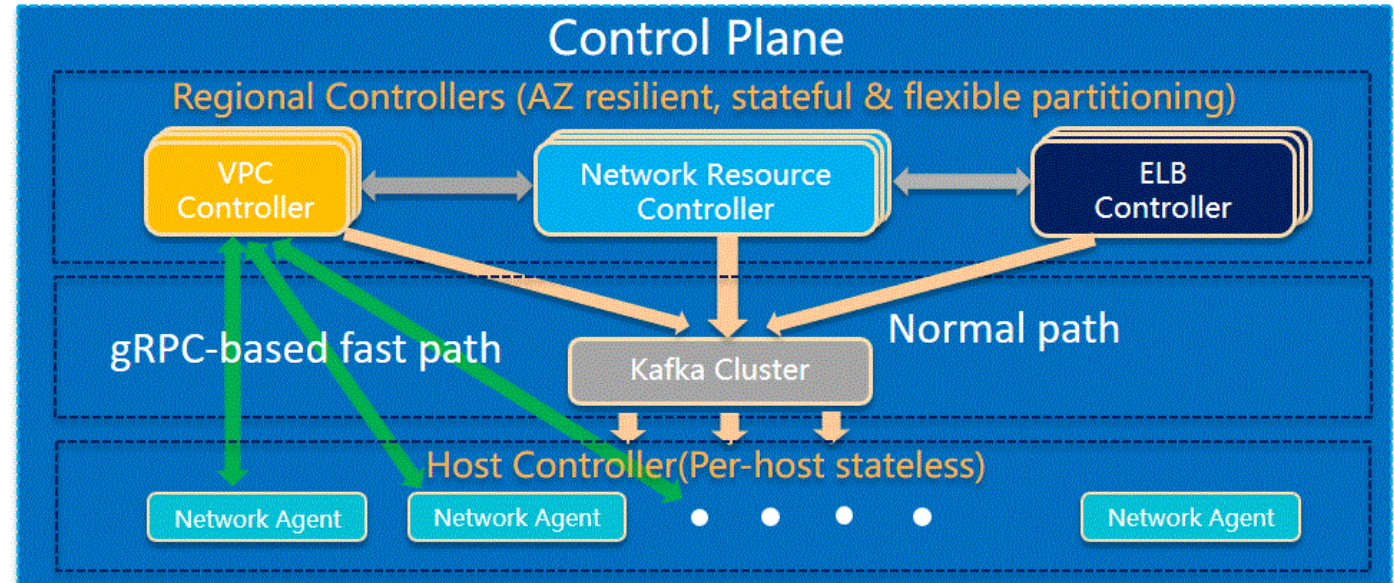
Legend: Serial, Batch

# Fast Port Provisioning

**Fast path design**

- Direct communication channel from controller to agent
- Alternative provisioning path for control plane reliability

**Background**

- Customer scenarios requires ultra-low latency for E2E network configuration provisioning (in a few 10 *ms* or 100 *ms*)
- Message queue subsystem, usually adopted as a high-throughput and scalable solution for network configuration updates, may not fit into time-critical customer scenarios

# Video

To be uploaded…