



Voltha High Availability Design Document

Proxy & Core

Khenaidoo Nursimulu, Sergio Slobodrian

April 10, 2018

Contents

Contents

Fundamental Design Principles

Proxy Design

Core Addition and Removal

Algorithms / Flowcharts

Core Distribution Across HW

Voltha Core Design

Current Core Overview

Proposed Core Architecture

Success & Failure Examples

Fundamental Design Principles

Fundamental Design Principles

- **Ensure 5 or 6 9's availability**
 - 5 9's – Less than 5 min / system / year outage
 - 6 9's – Less than 30 sec / system / year outage
- **Support horizontal scaling**
- **Support 50ms recovery from failure in most cases**
- **Never lose a transaction**

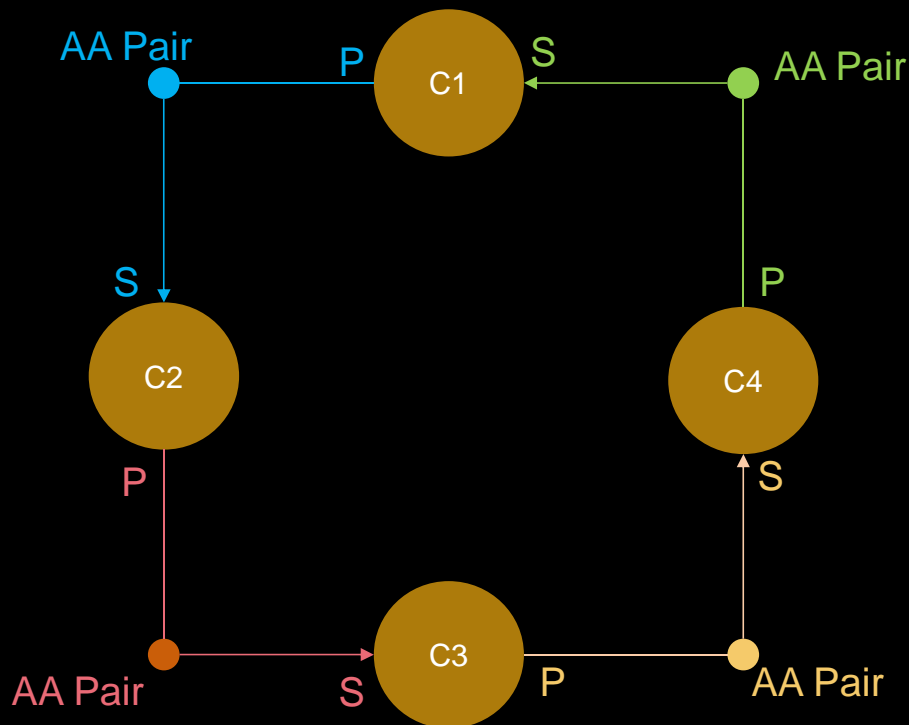
Fundamental principles → Design decisions

Principle	Design Decisions
Ensure 5 or 6 9's availability	Primary/Secondary Cores
Support horizontal scaling	Stateless cores
Support 50ms recovery from failure in most cases	Primary/Secondary cores, stateless cores
Never lose a transaction	Primary/Secondary cores

Proxy Design

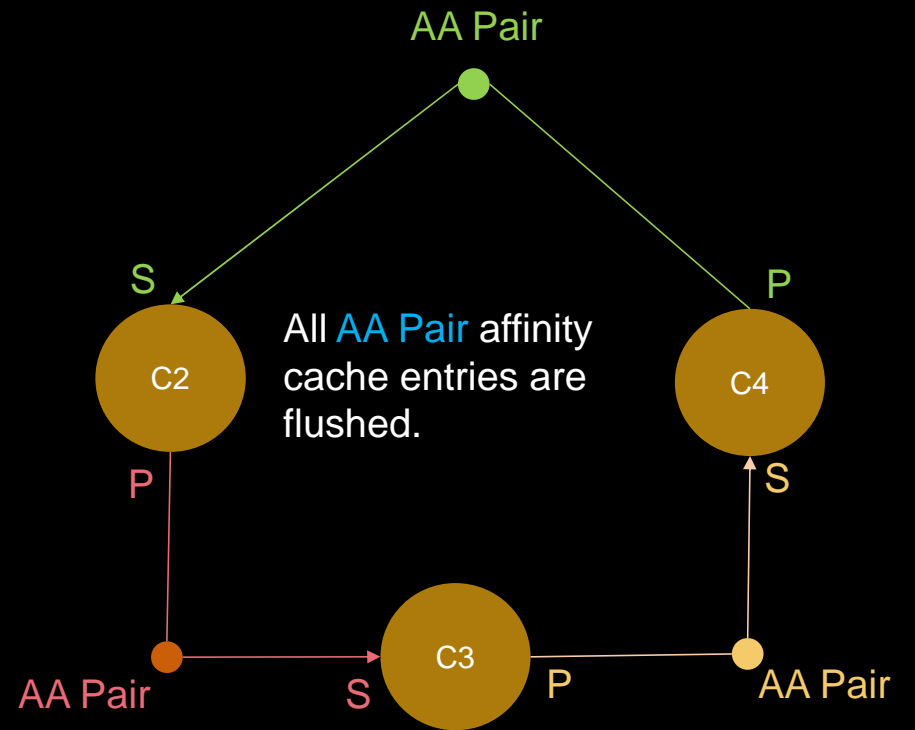
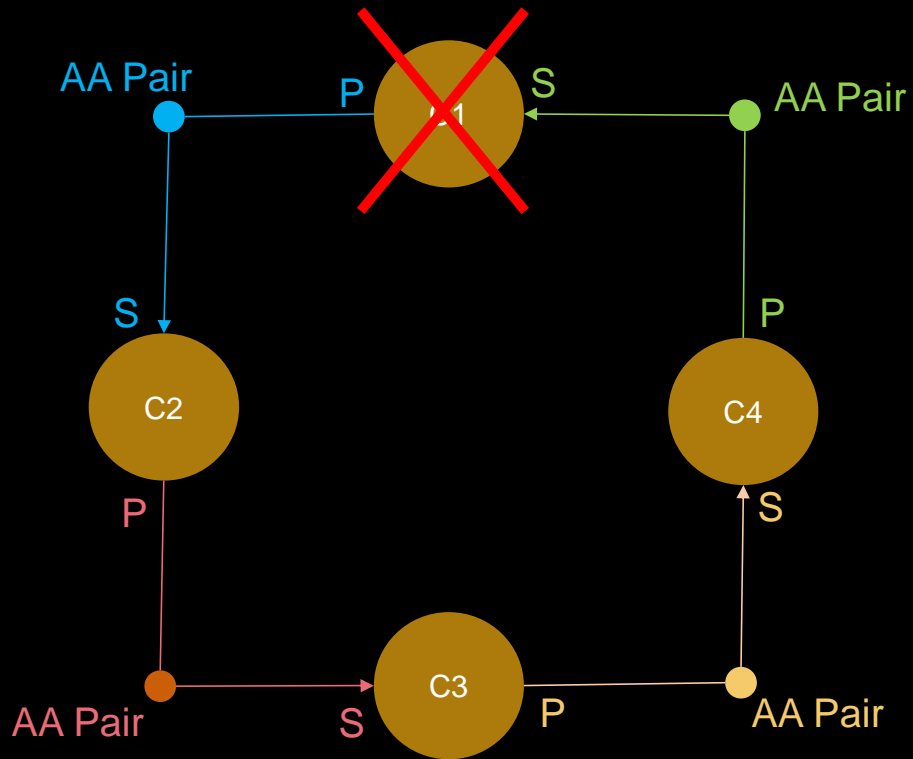
Core Addition and Removal

Active-Active Core Pairing

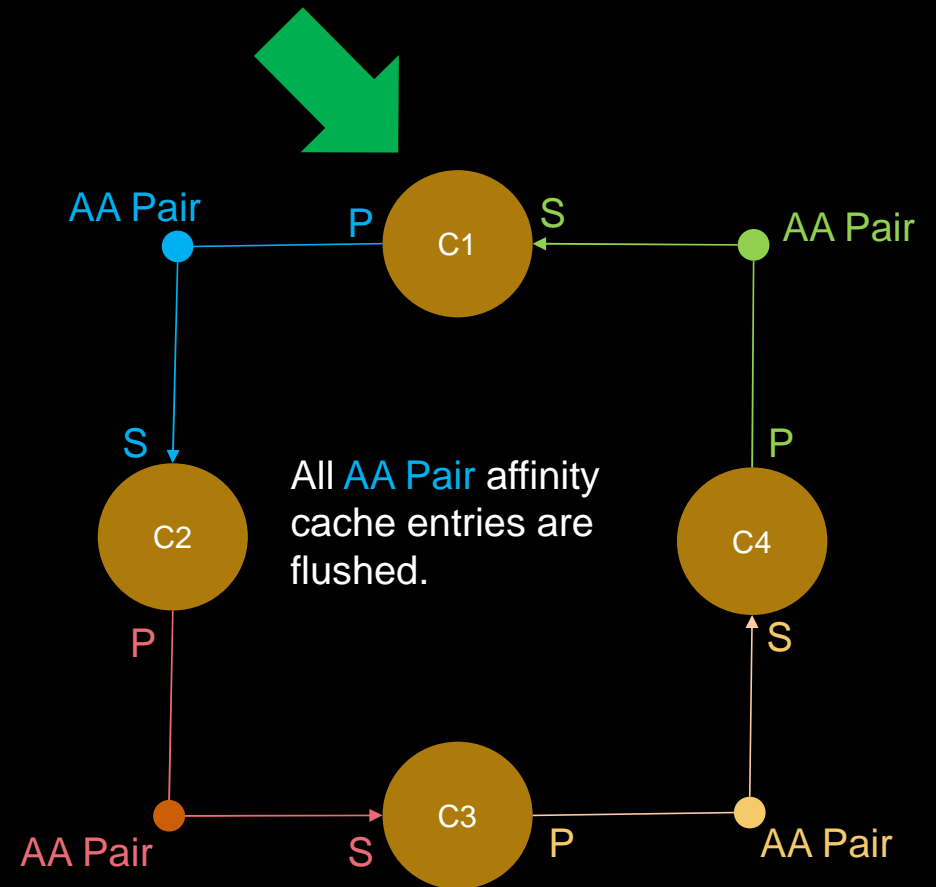
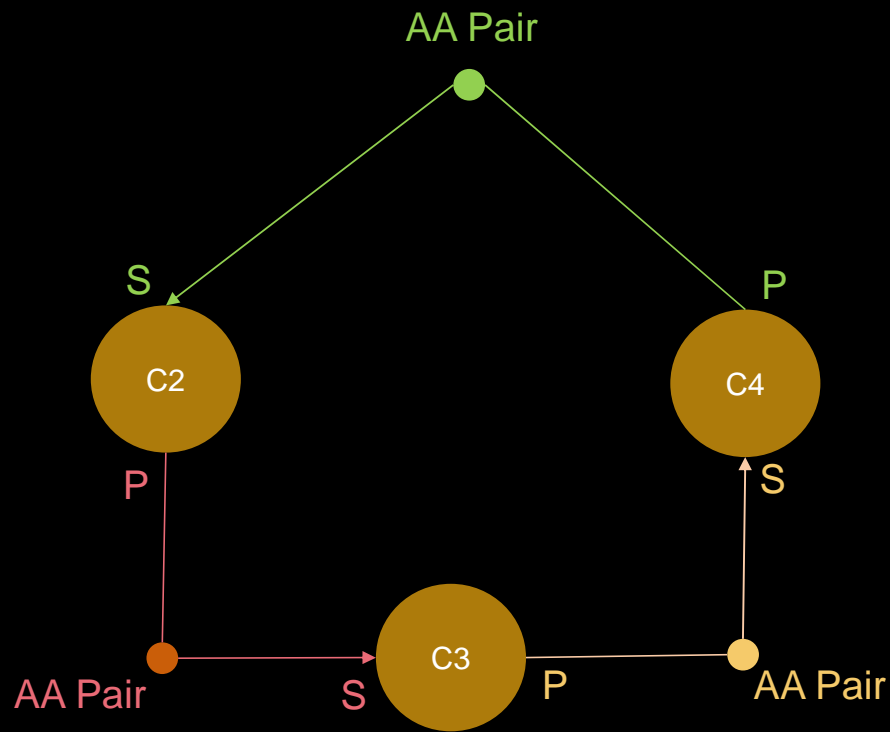


- For a non-id specific request the proxy will send it to the least loaded AA pair.
- For an id specific request where the id is in the affinity cache, the cached AAPair is used
- For an id specific request where the id isn't in the affinity cache the least loaded AA pair will be selected and that selection cached.
- For all requests going to an AA pair no matter how they got there
 - The AA pair will send it to the primary (P) core
 - The AA pair will immediately send it to the secondary (S) core.
- Core loading is based on the responses since only one of the 2 cores in a pair will process the request.

Active-Active Core Removal

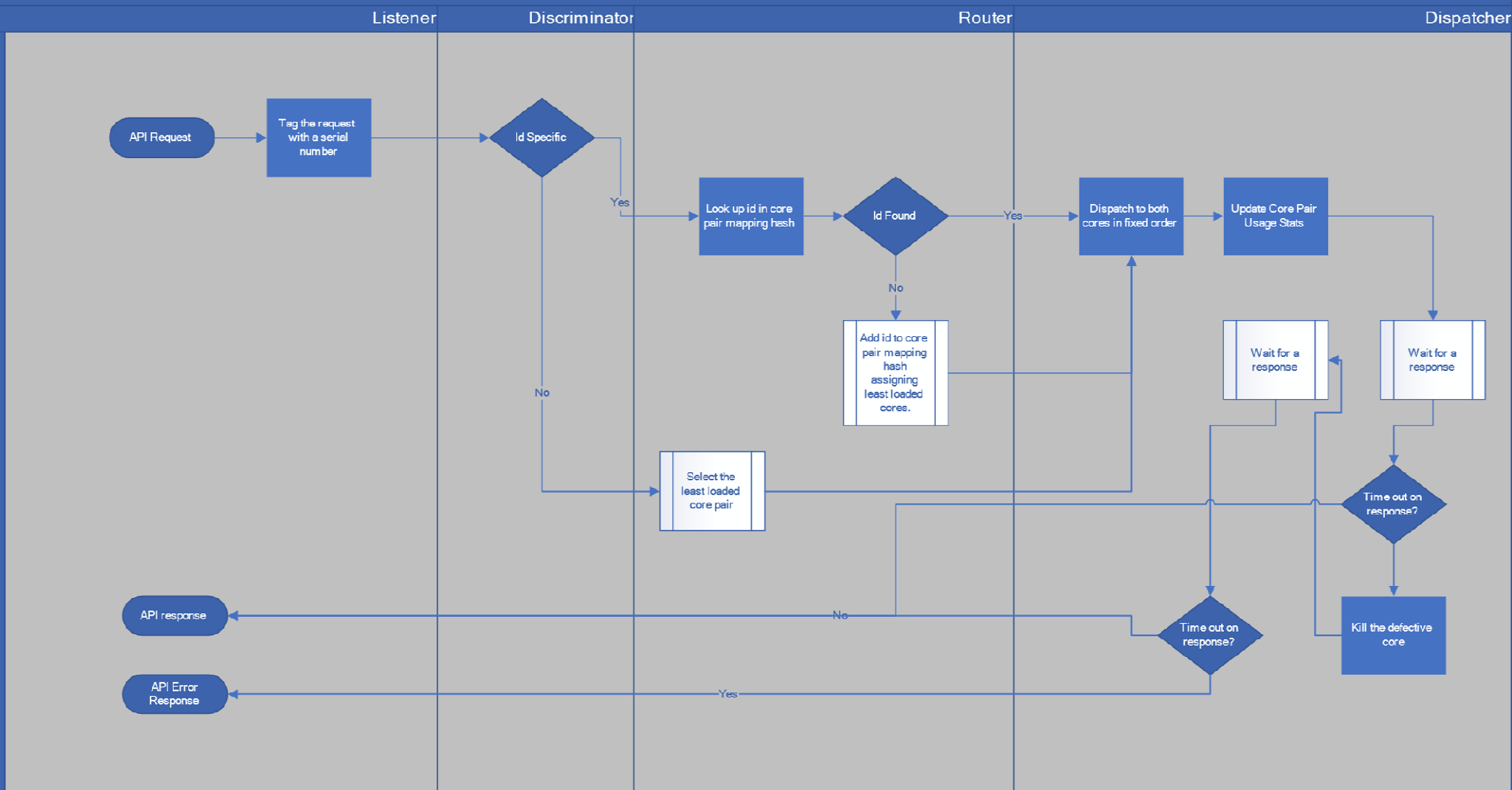


Active-Active Core Addition

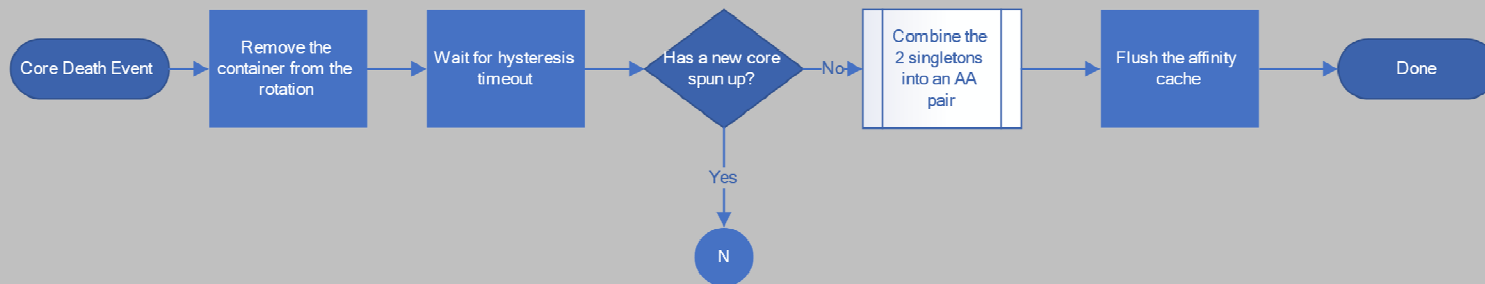


Algorithms / Flowcharts

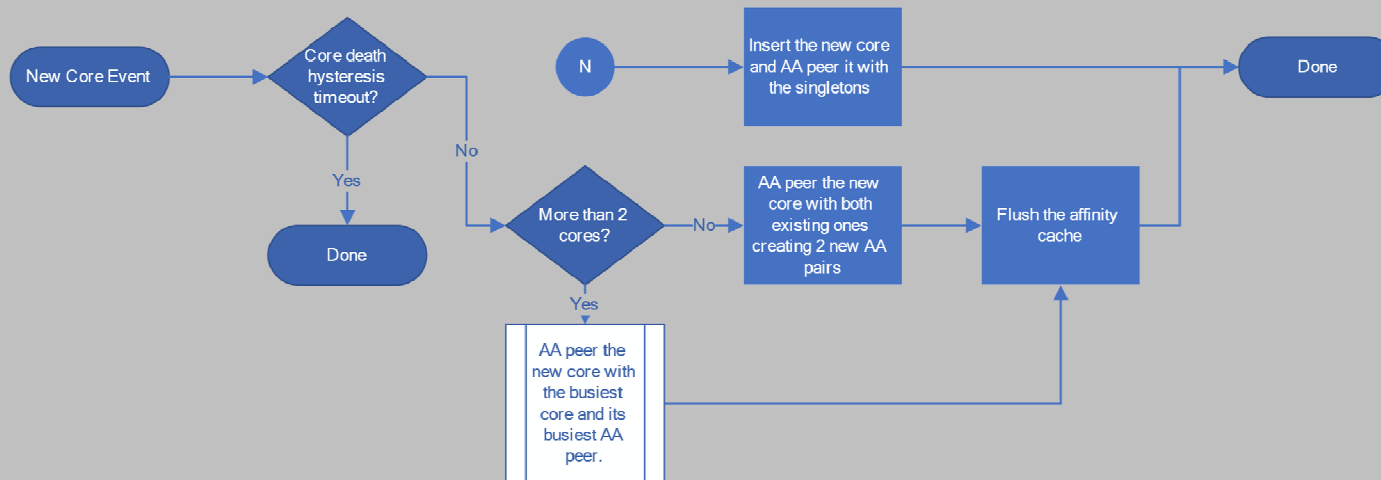
Proxy Request Handler Thread



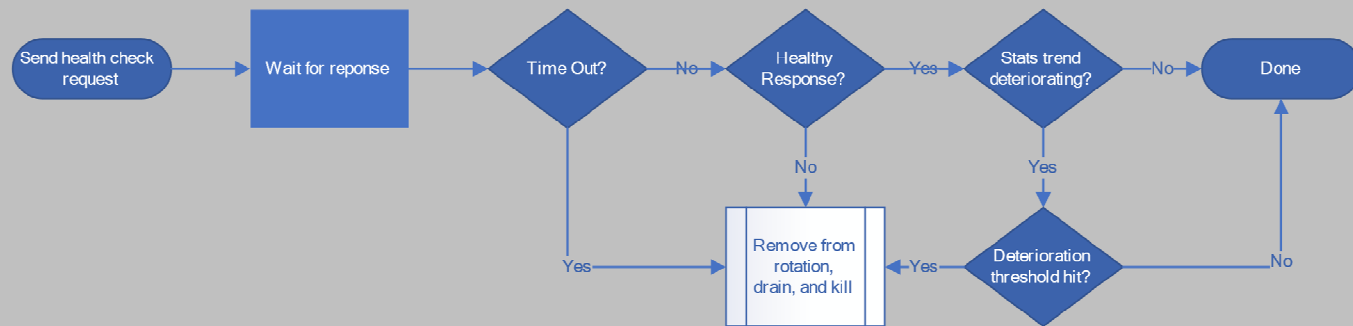
Core Death



New Core

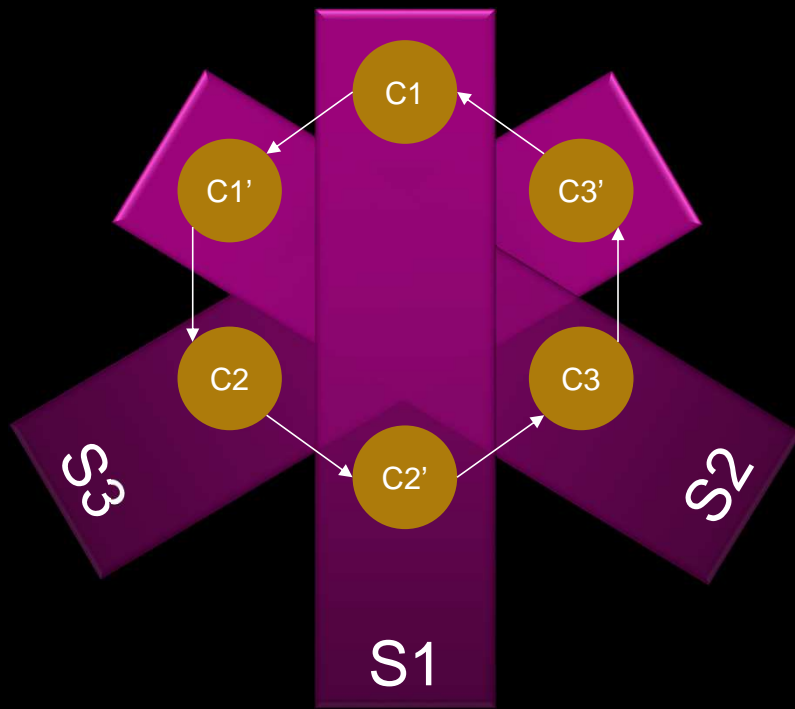


Health Check Thread



Core distribution across HW

Core distribution across HW



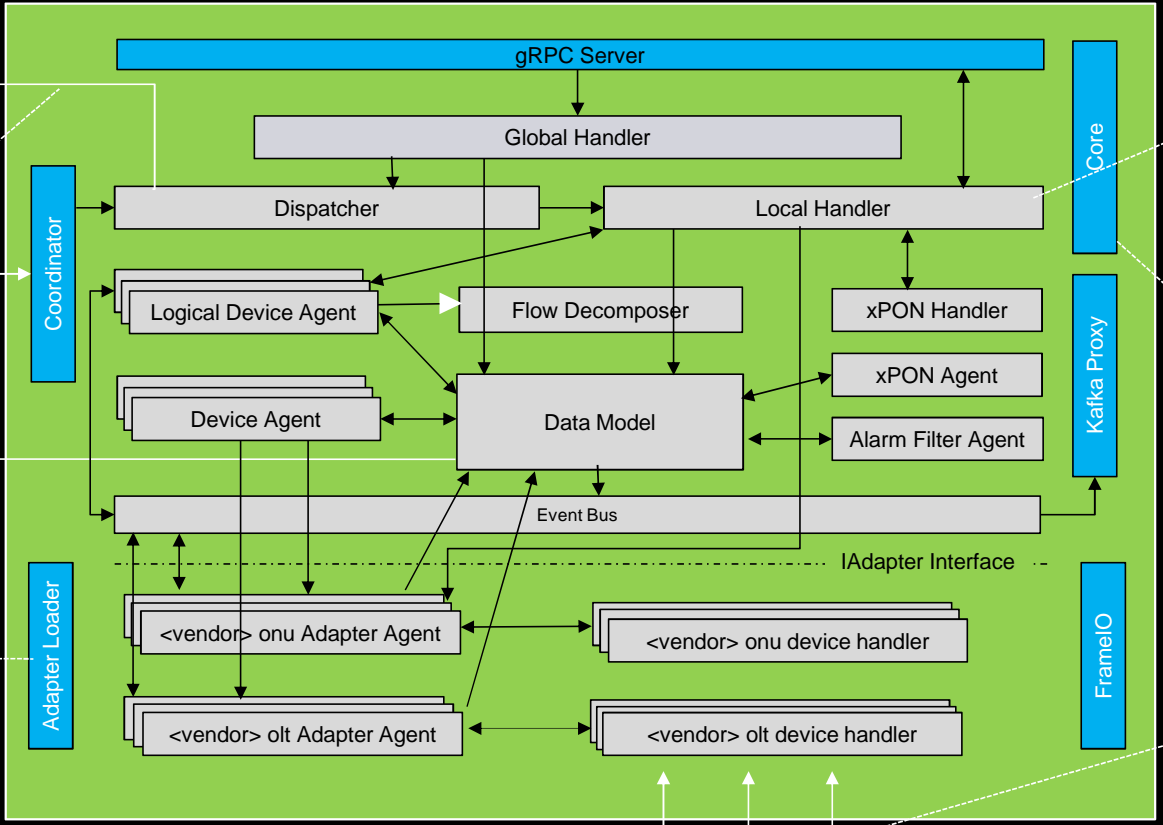
- One possible core distribution is shown on the left.
- This is the minimum HA configuration, 3 servers and 6 cores.
- A core and its peer should never be co-located on the same server
- For better load distributions, each server should have a non-prime and a prime not associated with the non-prime.
- Shown is $S1:\{C1,C2'\}$, $S2:\{C1',C3\}$, $S3:\{C2,C3'\}$
- Could have also distributed as: $S1:\{C1, C2\}$, $S2\{C2',C3'\}$, $S3\{C3,C1'\}$
 - this would have been less optimal since the non=prime cores are expected to be the busiest based the proxy behavior.

Voltha High Availability Design Document

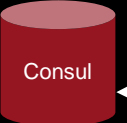
Voltha Core Design

Current Core Overview

Voltha Core - Current Design



Dispatch request to a specific Vcores if the request refers to devices of that core using grpc



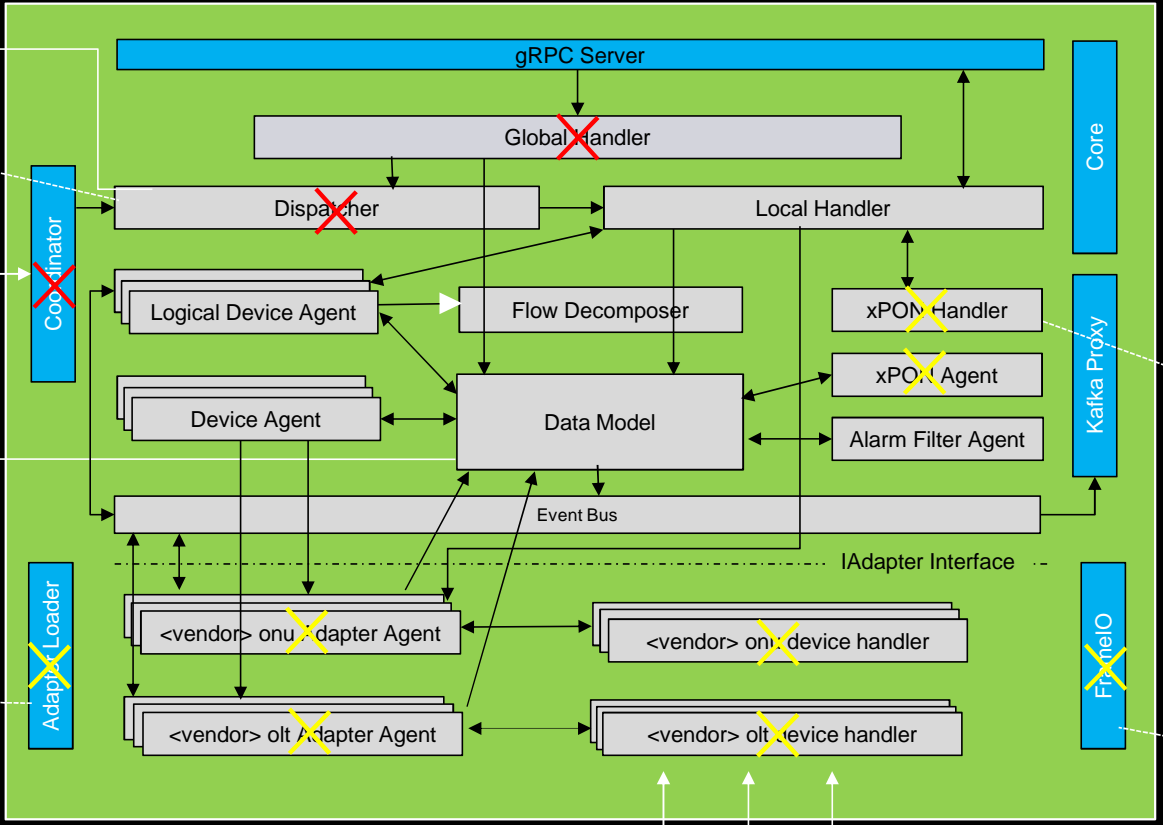
Dynamically loads adapters

- Handles requests for this core.
- Dispatches OF events (igmp, Eapol, etc) to OFAgent

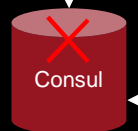
- Create global, local and xpon handler/agent
- Reconcile Data after a restart
- Proxy to the data model
- Handle add/remove device and logical device
- Create/remove xpon interface

Connection to vendor specific OLTs with varied protocols

Voltha Core 2.0 – Removed/Transferred out components



With affinity routing and stateless vcore there is no need of Dispatcher, Global Handler and Coordinator service.



Replaced by Etcd

Since adapters will be moved out of vcore then there are no need to have the adapters and the adapter loader.

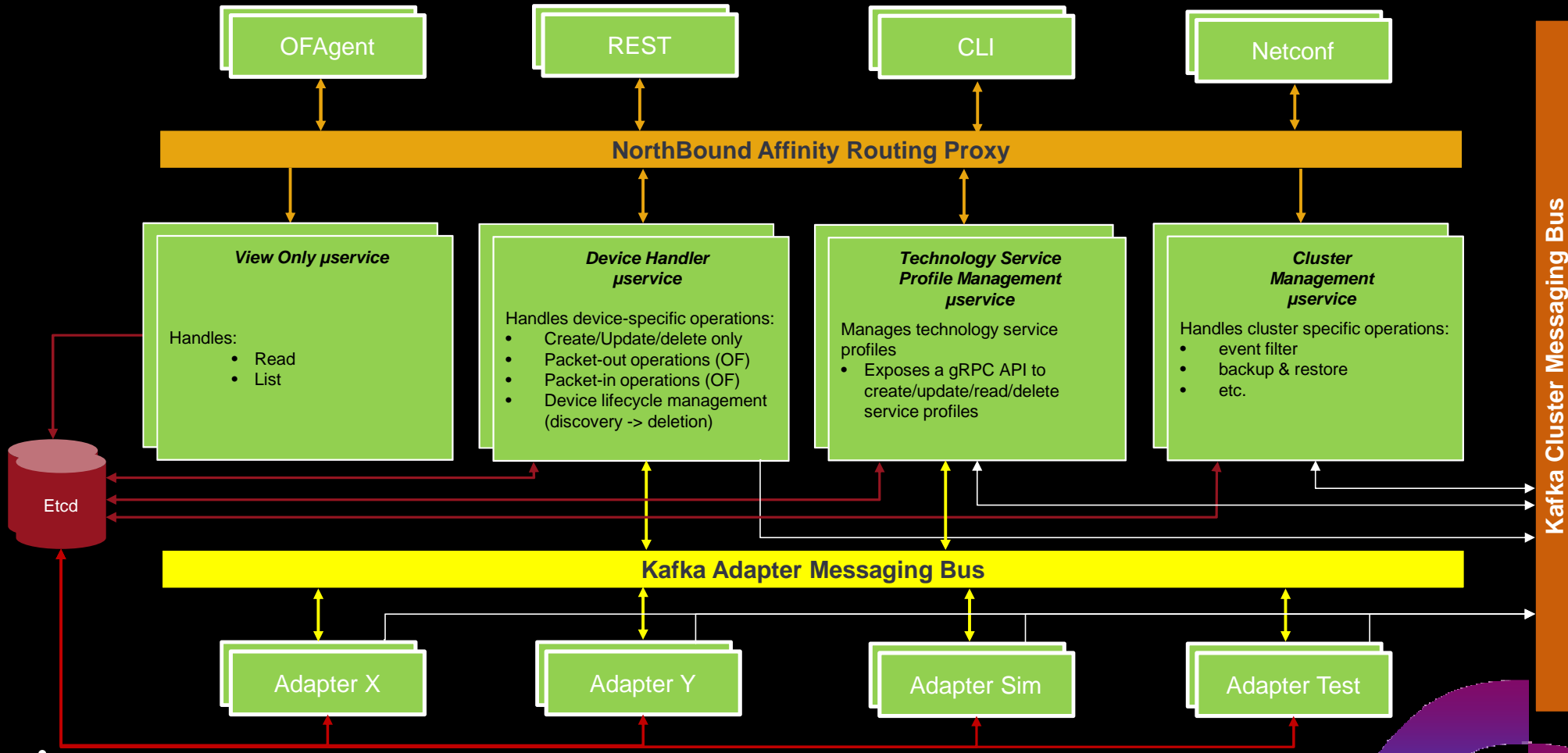
X Removed
Y Transferred out

• Technology specific components like xPON need to live in their own containers

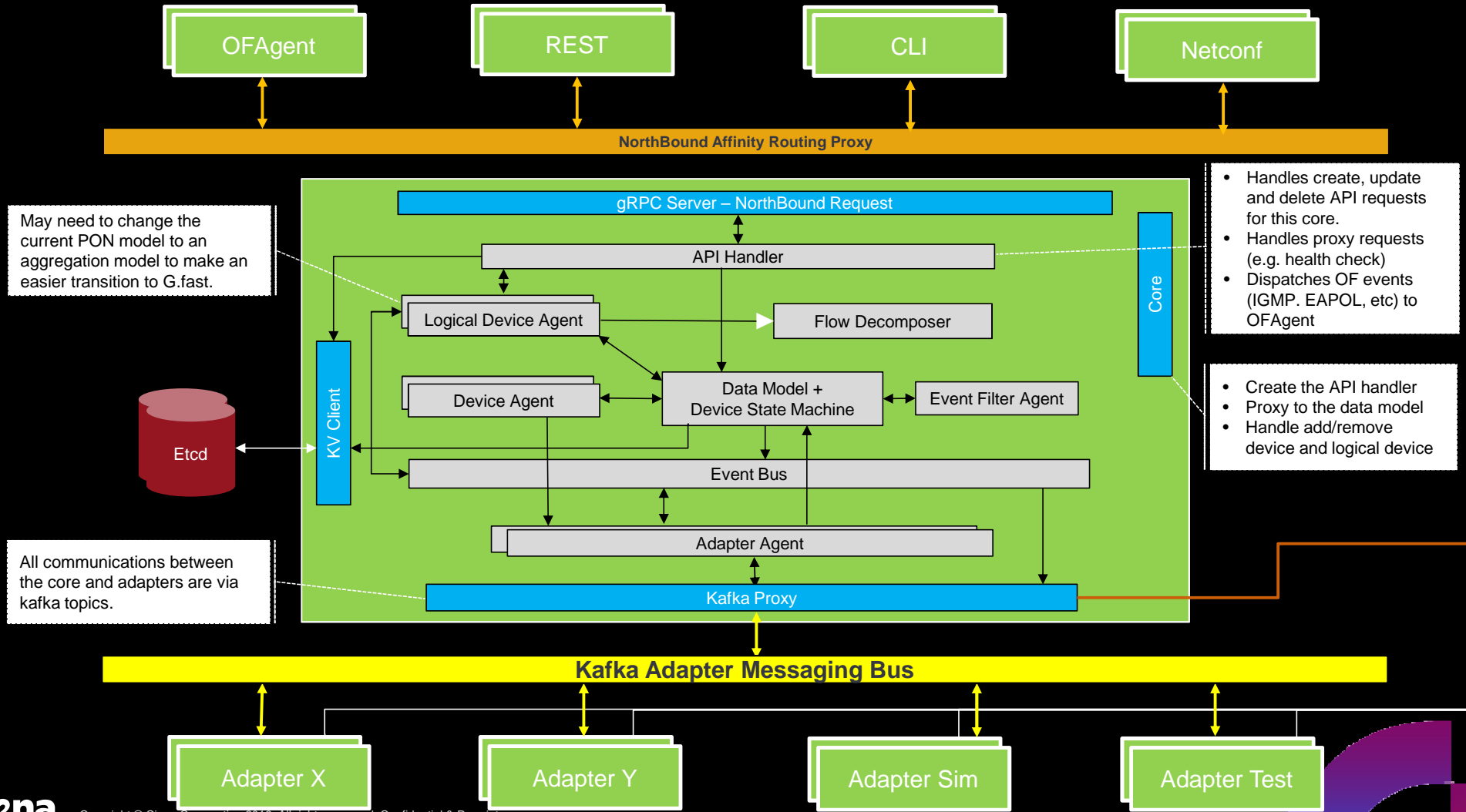
No need for FrameIO as the core does not need it. Only some adapters may still need it

Proposed Core Architecture

Voltha Core 2.0 – 4 μservices



Device Handler μservice



May need to change the current PON model to an aggregation model to make an easier transition to G.fast.

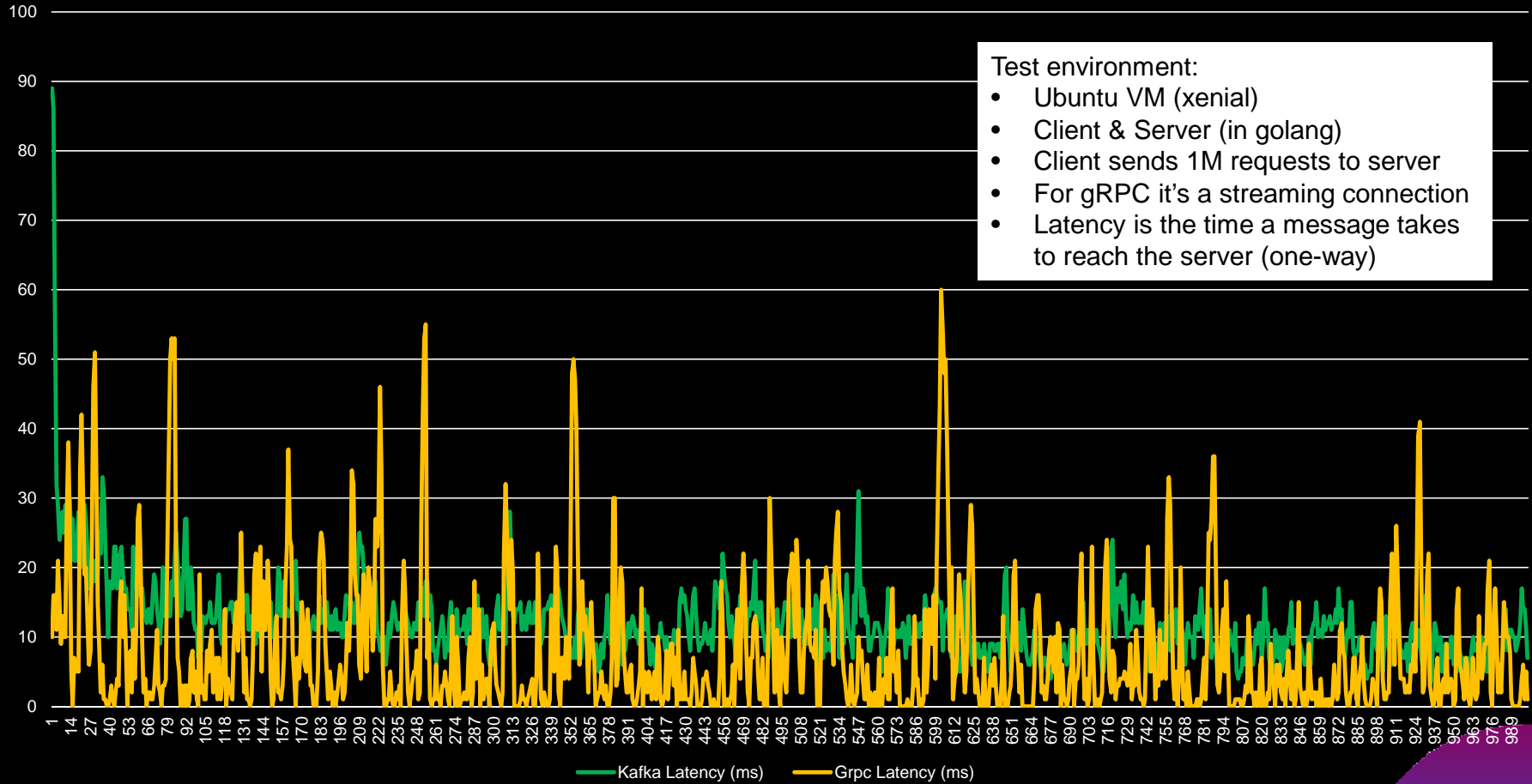
- Handles create, update and delete API requests for this core.
- Handles proxy requests (e.g. health check)
- Dispatches OF events (IGMP, EAPOL, etc) to OFAgent

- Create the API handler
- Proxy to the data model
- Handle add/remove device and logical device

All communications between the core and adapters are via kafka topics.

Message latency – Kafka vs gRPC

Kafka/gRPC latency (1 Million messages)

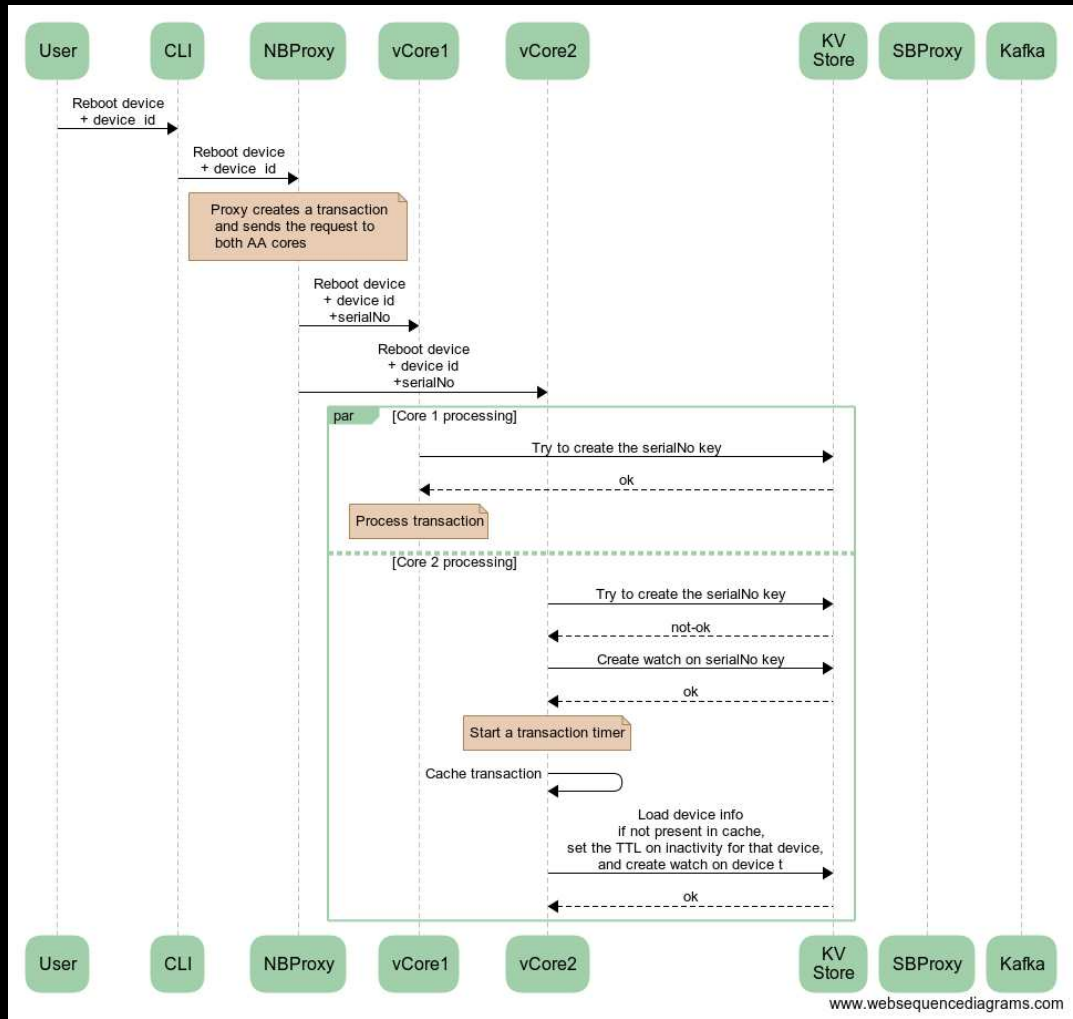


Test environment:

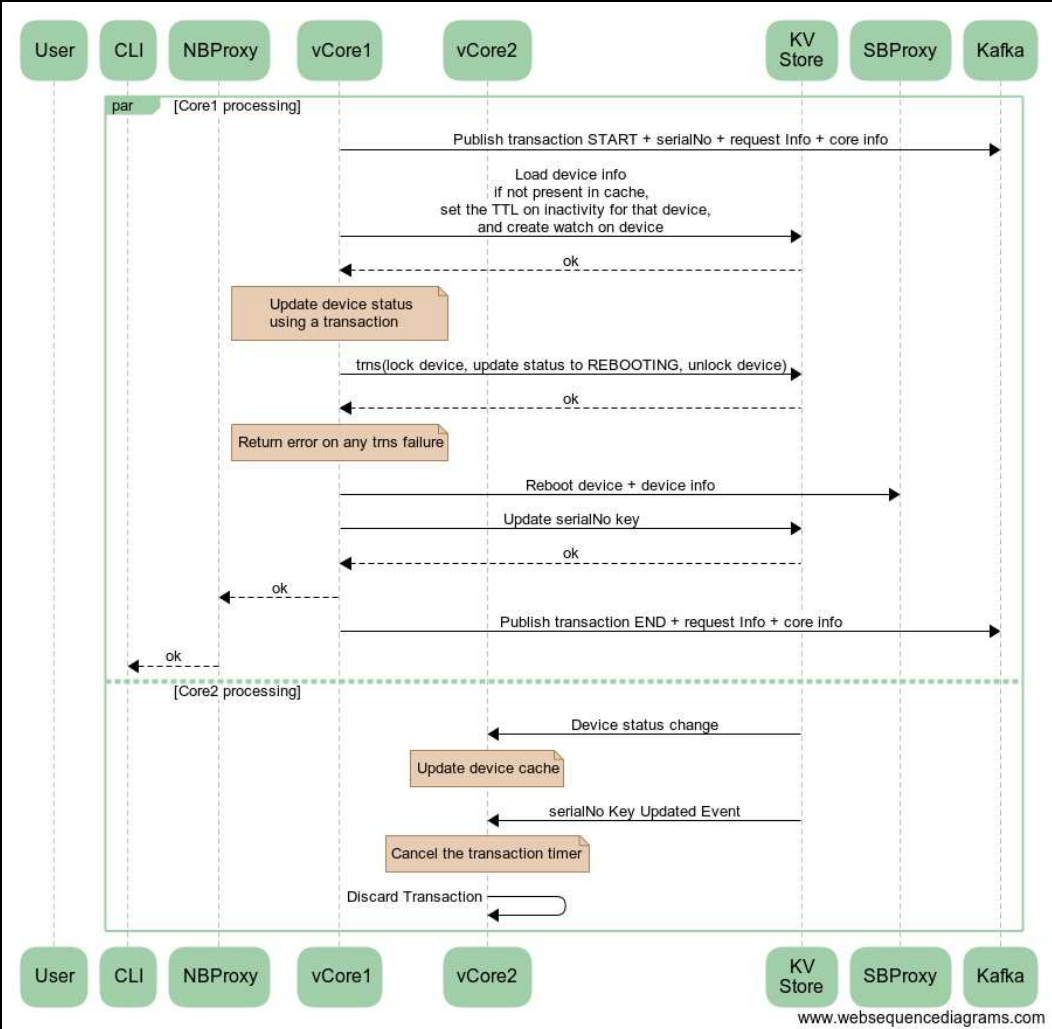
- Ubuntu VM (xenial)
- Client & Server (in golang)
- Client sends 1M requests to server
- For gRPC it's a streaming connection
- Latency is the time a message takes to reach the server (one-way)

Success & Failure Examples

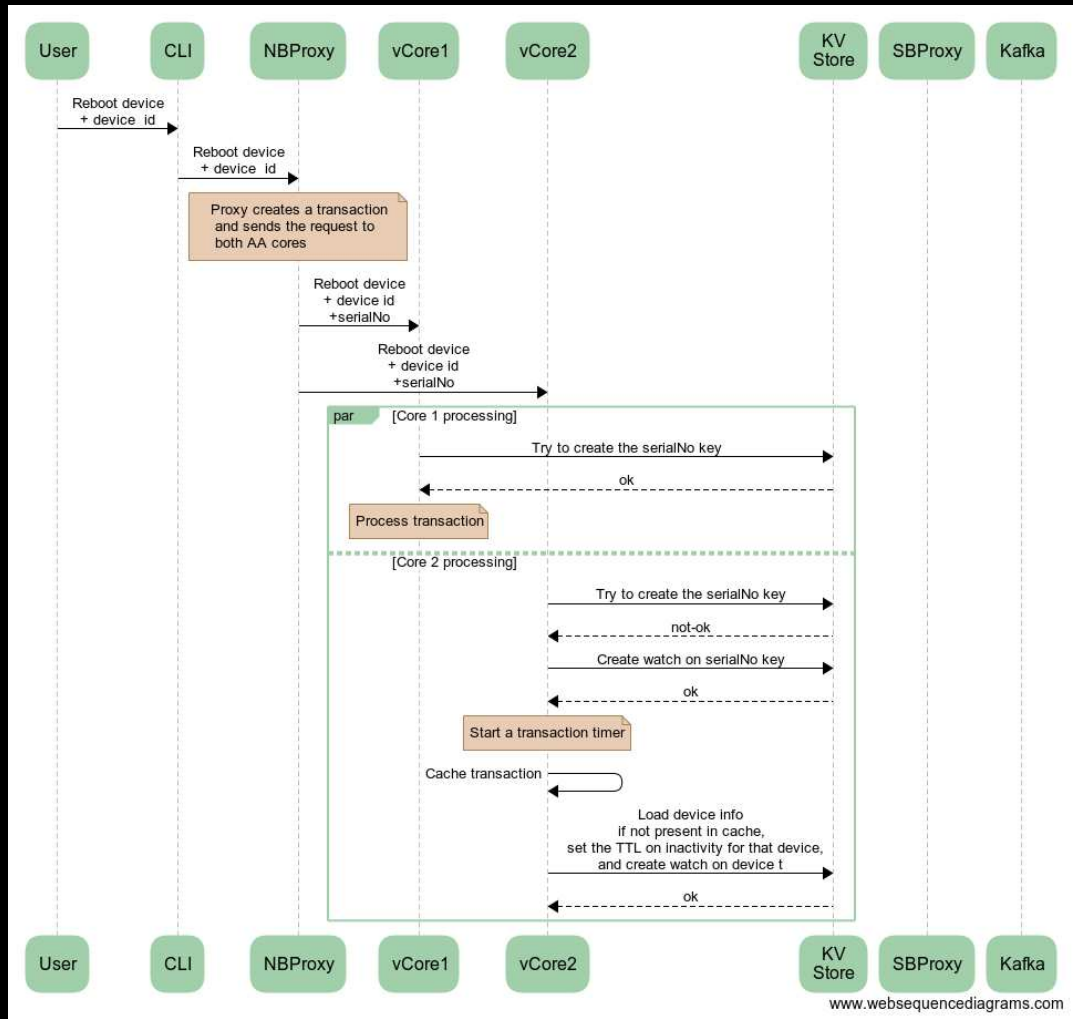
Example— Reboot request (no error) – part I



Example— Reboot request (no error) – part II



Example— Reboot request (error) – part I



Example— Reboot request (error) – part II

