# AI Nativity

The path forward

# Pre-AI World

## Microservices Today

- **Key components**

- **Key Workflows**

- **Key Solutions**

  - **Routing**

  - **Policy Enforcement**

  - **Resiliency**

  - **Observability**

  - **More Complex Solutions**

    - WCP

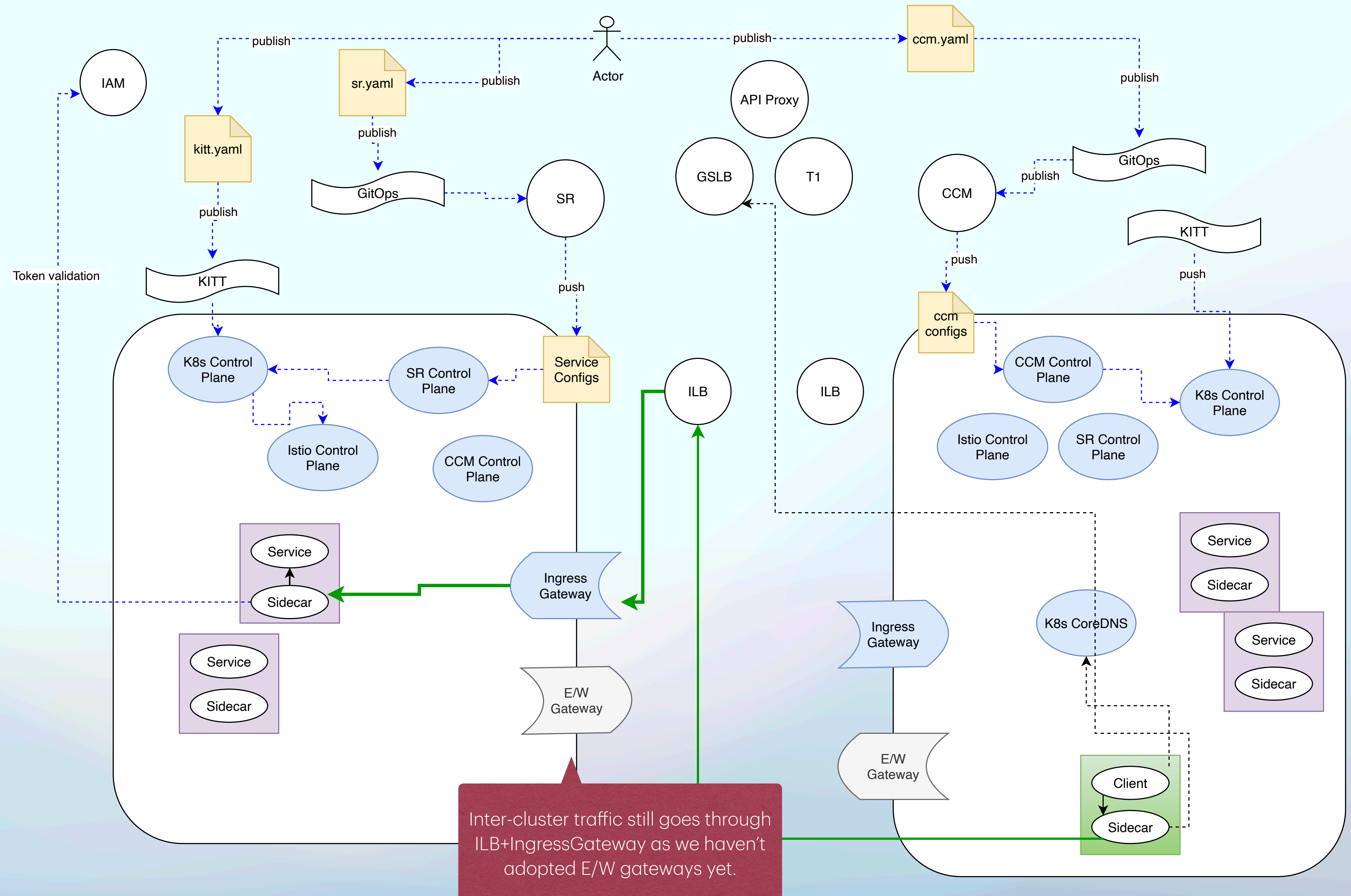    - API Gateway, Bubble

# AI World

## Where we want to be tomorrow

- **The Journey**

- **New players**

- **New/Adjusted Workflows**

- **Routing**

- **Policies**

- **Resiliency**

# Microservices World Today

- **Key components**
  - Services
  - WCNP
  - GSLB, ILB
  - Ingress Gateway (Envoy)
  - Sidecar (Envoy)
  - Istio Control Plane (IstioD)
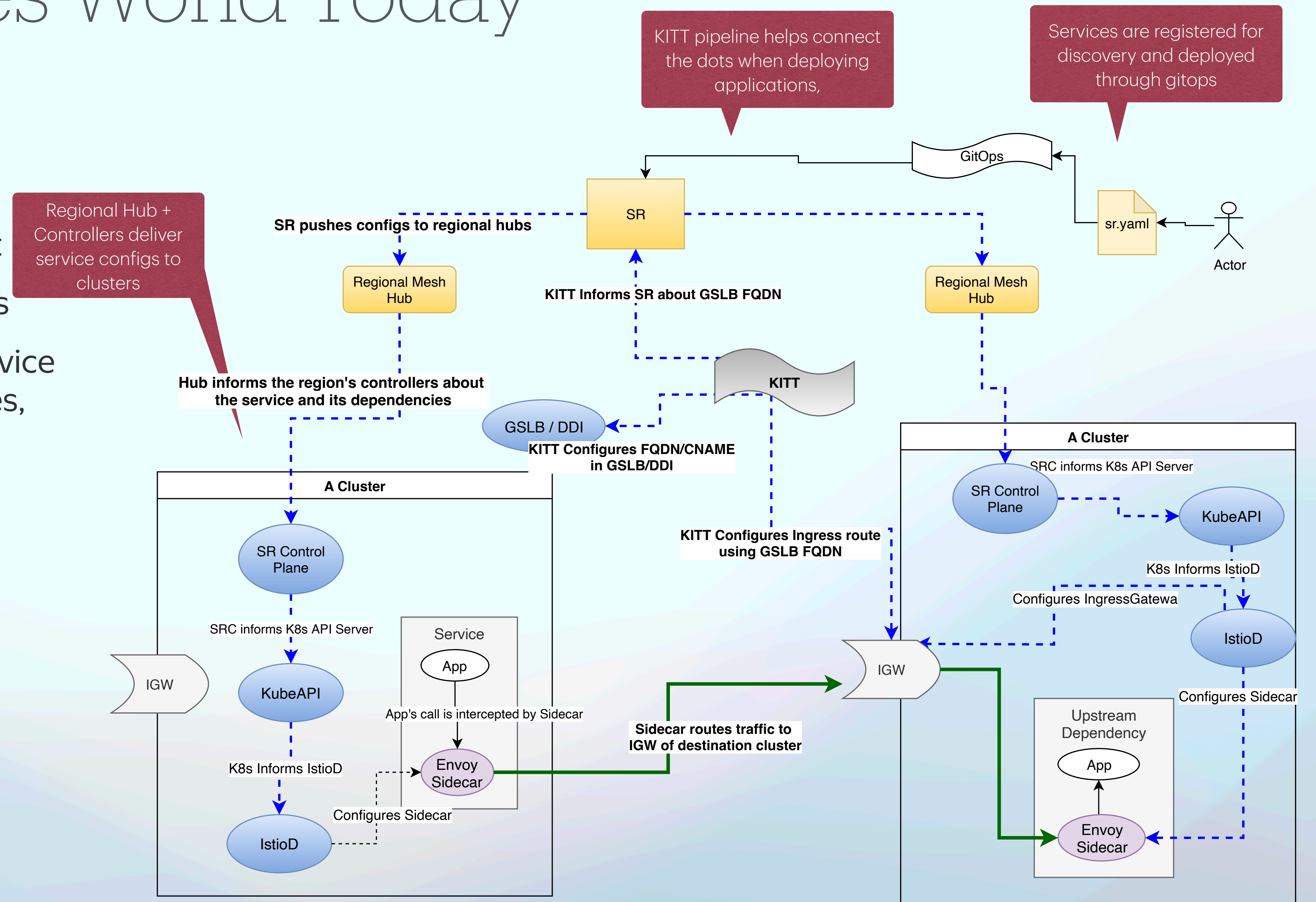  - K8s Control Plane (API Server, CoreDNS)
  - SR + Controller



Inter-cluster traffic still goes through ILB+IngressGateway as we haven't adopted E/W gateways yet.

# Microservices World Today
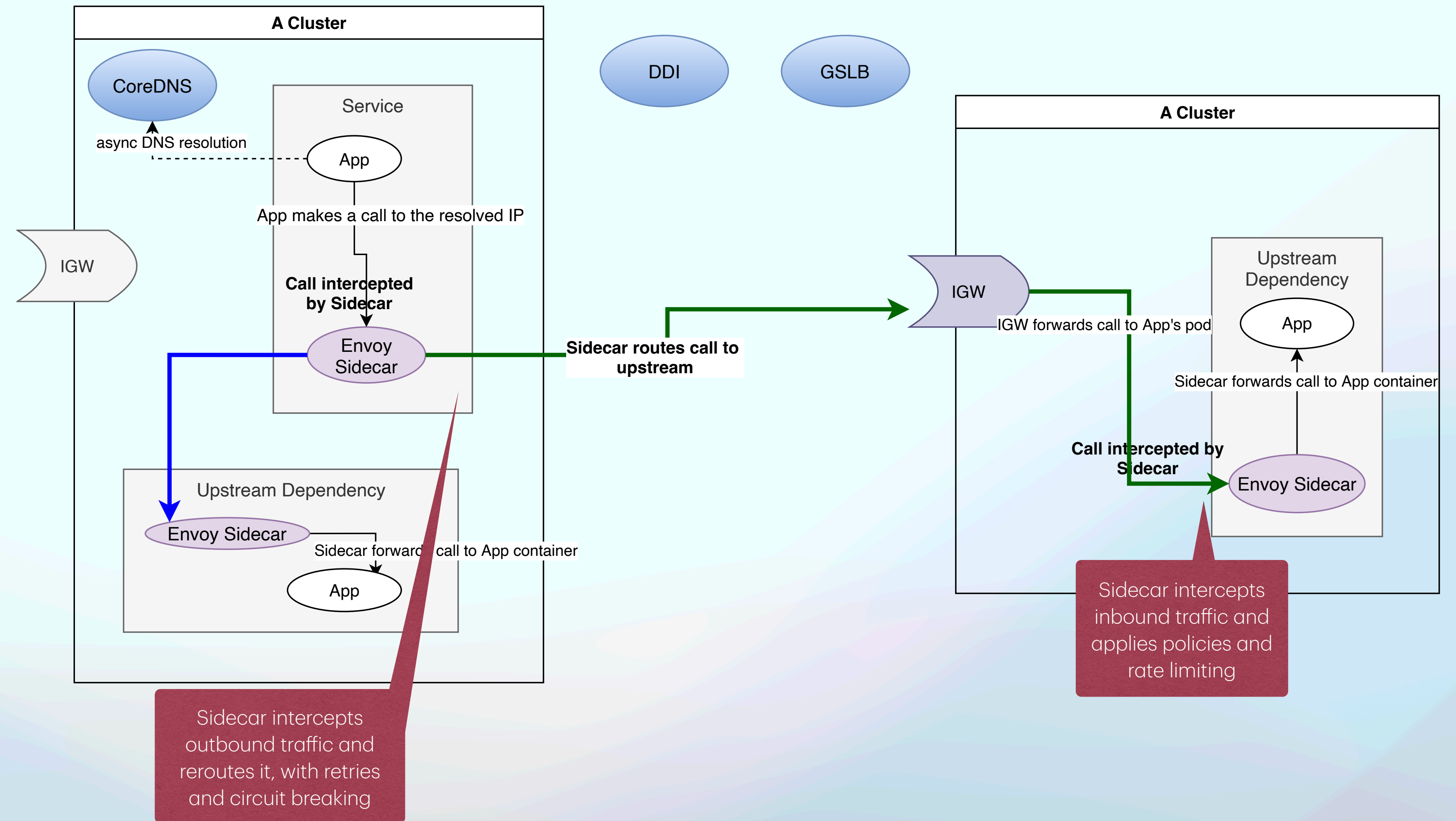
- **Key components**

- **Key Workflows**

  - CI/CI: build application artifact

  - KITT: deploy application to K8s

  - SR Gitops: Platform confis: service dependencies, routing, policies, resiliency
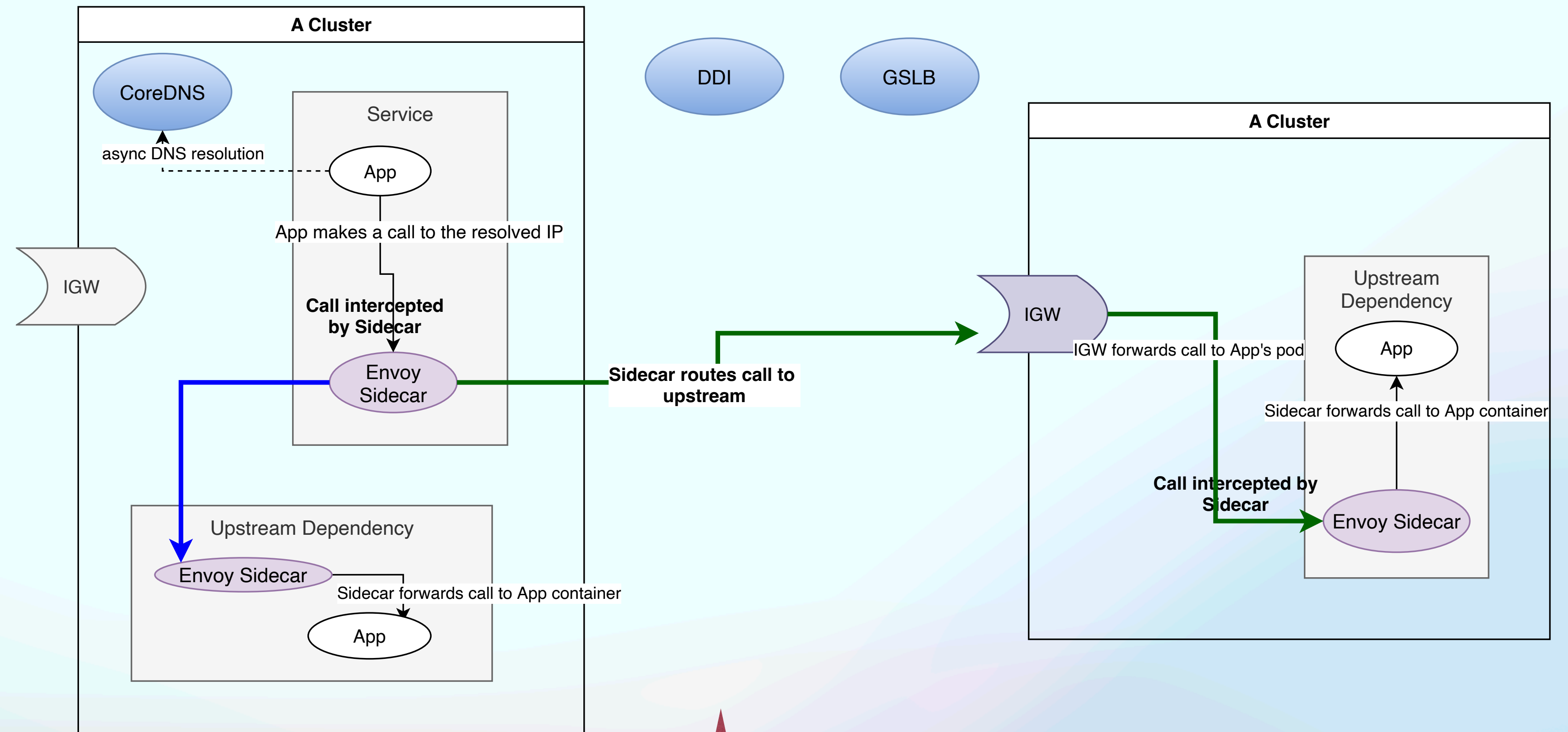
  - CCM Gitops: App's configs

# Microservices World Today

- **Key components**

- **Key Workflows**

- **Routing, Policies, Resiliency, Observability**
  - Let's zoom-in Into Service Mesh routing.



**A Cluster**

CoreDNS

async DNS resolution

**Service**

App

App makes a call to the resolved IP

**Call intercepted by Sidecar**

Envoy Sidecar

Sidecar routes call to upstream

Upstream Dependency

Envoy Sidecar

Sidecar forwards call to App container

App

Sidecar intercepts outbound traffic and reroutes it, with retries and circuit breaking

DDI

GSLB

IGW

**A Cluster**

IGW

IGW forwards call to App's pod

Upstream Dependency

App

Sidecar forwards call to App container

**Call intercepted by Sidecar**

Envoy Sidecar

Sidecar intercepts inbound traffic and applies policies and rate limiting

# Microservices World Today

- **Key components**

- **Key Workflows**

- **Routing, Policies, Resiliency, Observability**

  - Let's zoom-in Into Service Mesh routing.

  - The hops penalty question

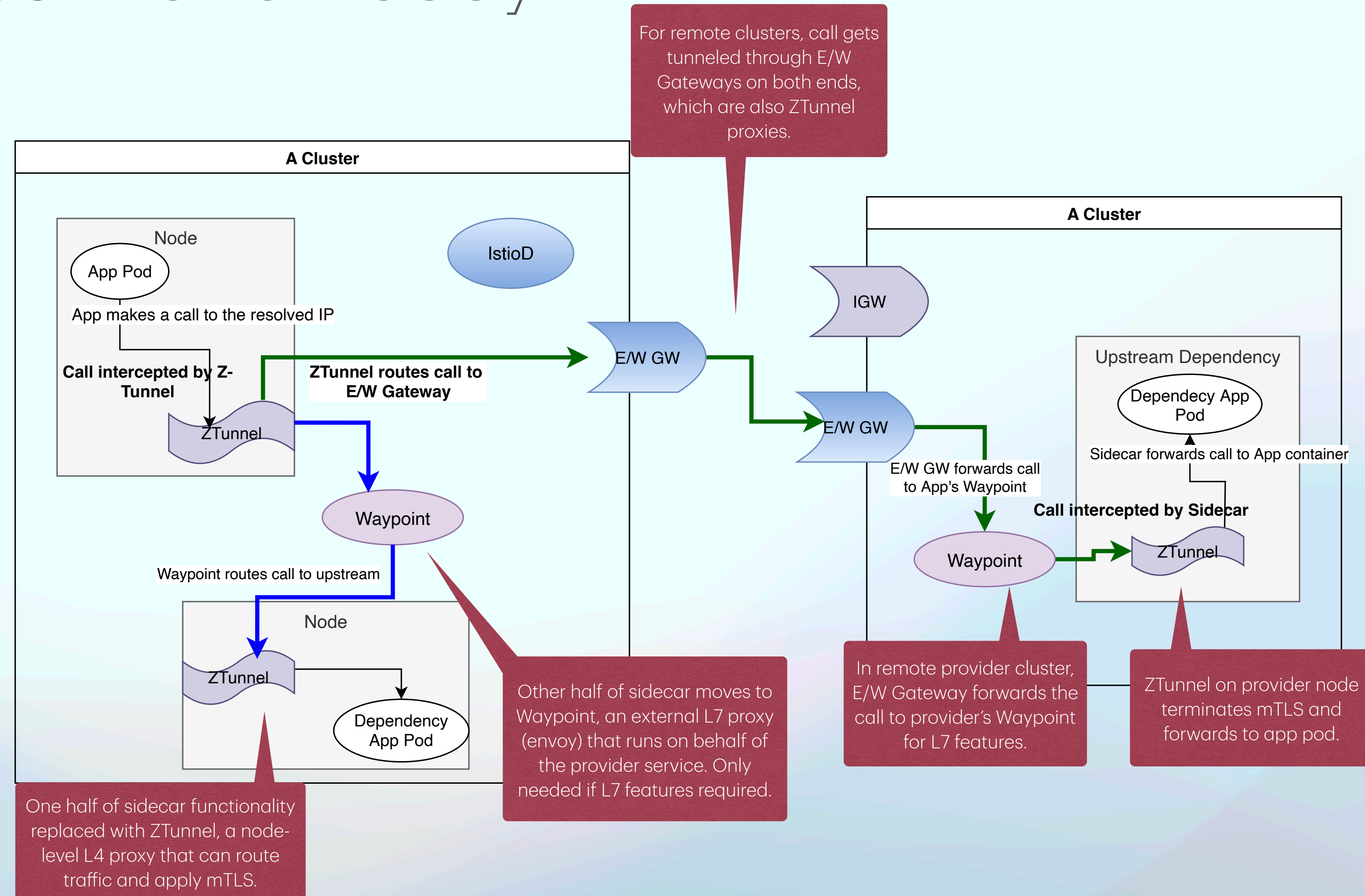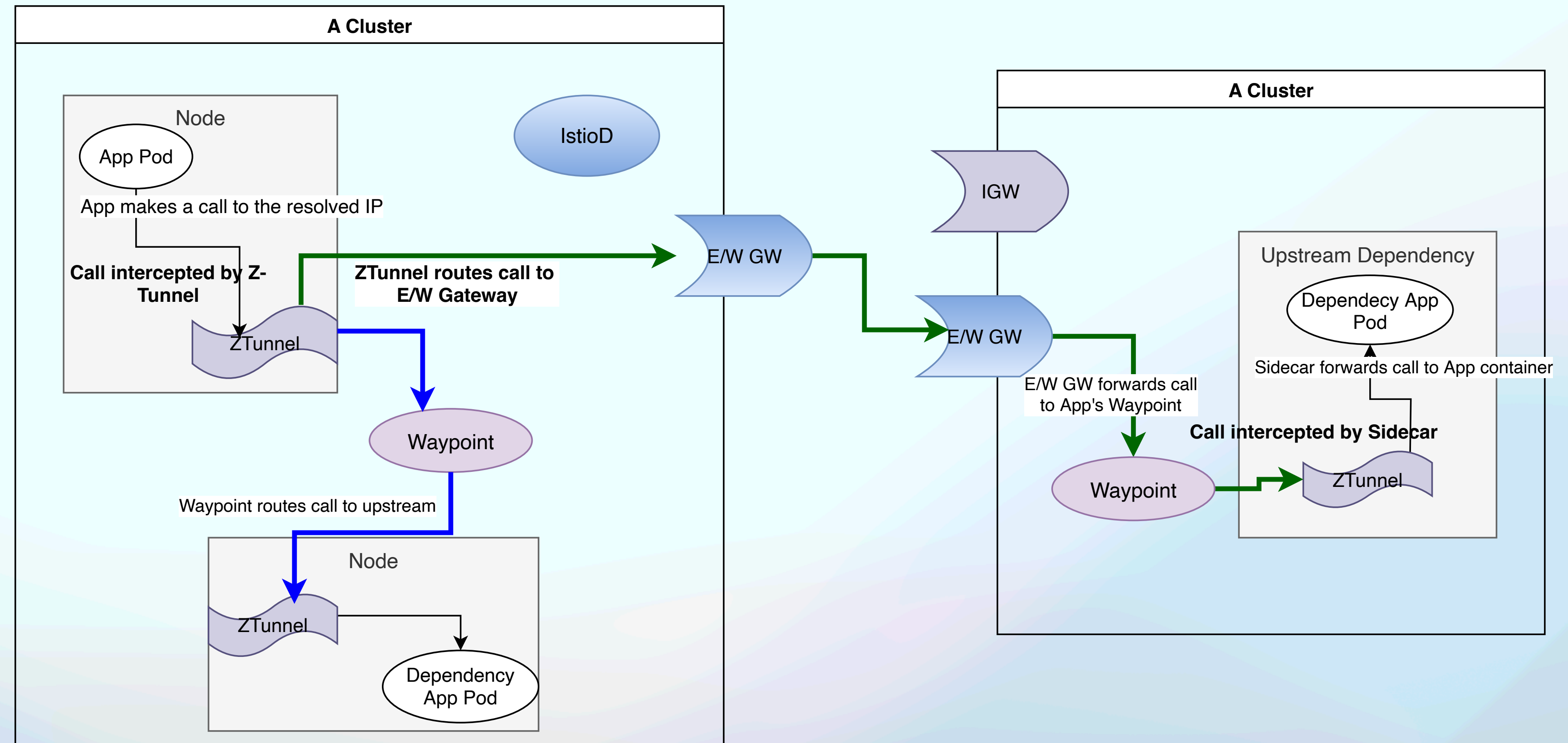  - Before that, let's look at the future of Service Mesh
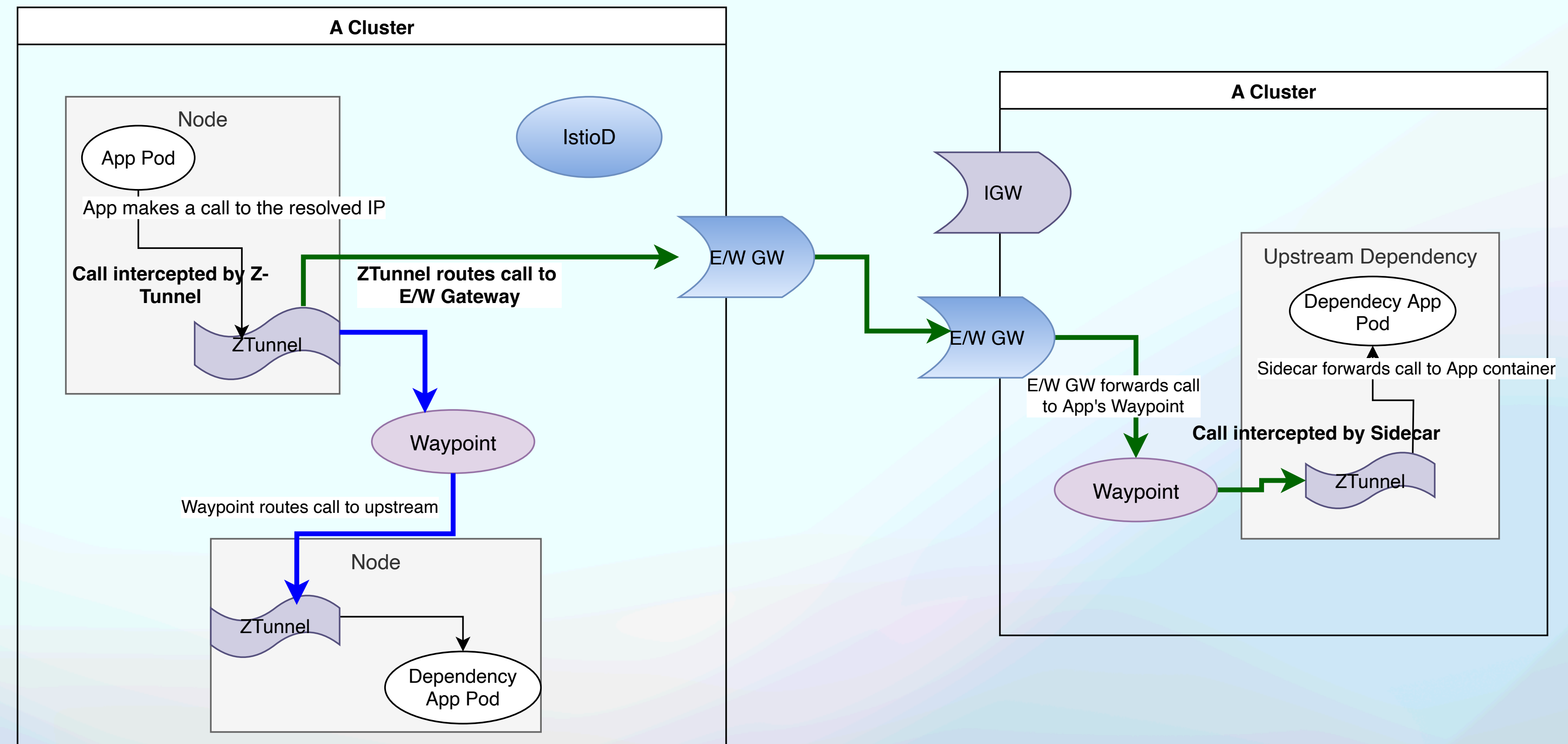
# Microservices World Today

## Ambient Mesh

- The logical evolution of Service Mesh.

- Why do we care? It helps lower the operational cost by cutting down the number of sidecar cores/memory needed (likely 1/10th).

**A Cluster**

For remote clusters, call gets tunneled through E/W Gateways on both ends, which are also ZTunnel proxies.

**A Cluster**

Node

App Pod

App makes a call to the resolved IP

**Call intercepted by Z-Tunnel**

**ZTunnel routes call to E/W Gateway**

ZTunnel

IstioD

E/W GW

IGW

E/W GW

Upstream Dependency

Dependecy App Pod

Sidecar forwards call to App container

E/W GW forwards call to App's Waypoint

**Call intercepted by Sidecar**

Waypoint

Waypoint routes call to upstream

ZTunnel

Node

ZTunnel

Dependency App Pod

In remote provider cluster, E/W Gateway forwards the call to provider's Waypoint for L7 features.

ZTunnel on provider node terminates mTLS and forwards to app pod.

Other half of sidecar moves to Waypoint, an external L7 proxy (envoy) that runs on behalf of the provider service. Only needed if L7 features required.

One half of sidecar functionality replaced with ZTunnel, a node-level L4 proxy that can route traffic and apply mTLS.

# Microservices World Today

## Ambient Mesh

- The logical evolution of Service Mesh.

- Why do we care? It helps lower the operational cost by cutting down the number of sidecar cores/memory needed (likely 1/10th).

- But the hop count increased!

- The hops equation is tricky. Let's visit it.

- And then we'll see how Ambient plays a role in the future AI world.

# Microservices World Today

## The Hops Equation

- Network call is much cheaper than L7 proxy processing overhead.

- Assume net cost=**1**, L7 cost=**2**

- Local Cluster Traffic:
  - Sidecar Local: (1x net) + (2x L7) = (1 + 4) = **5**
  - Ambient Local: (2x net) + (1x L7) = (2 + 2) = **4**

- Remote Cluster Traffic:
  - Sidecar Remote: (2x net) + (3x L7) = (2+6) = **8**
  - Ambient Remote: (4x net + 1x L7) = (4+2) = **6**

- The higher the L7 processing cost goes, the lesser the relative impact of network hops.


- We're almost ready to look at the AI space now…

**A Cluster**

Node

App Pod

App makes a call to the resolved IP

**Call intercepted by Z-Tunnel**

ZTunnel

**ZTunnel routes call to E/W Gateway**

IstioD

E/W GW

Waypoint

Waypoint routes call to upstream

Node

ZTunnel

Dependency App Pod

IGW

**A Cluster**

E/W GW

E/W GW forwards call to App's Waypoint

Waypoint

Upstream Dependency

Dependecy App Pod

Sidecar forwards call to App container

**Call intercepted by Sidecar**

ZTunnel

# Microservices World Today

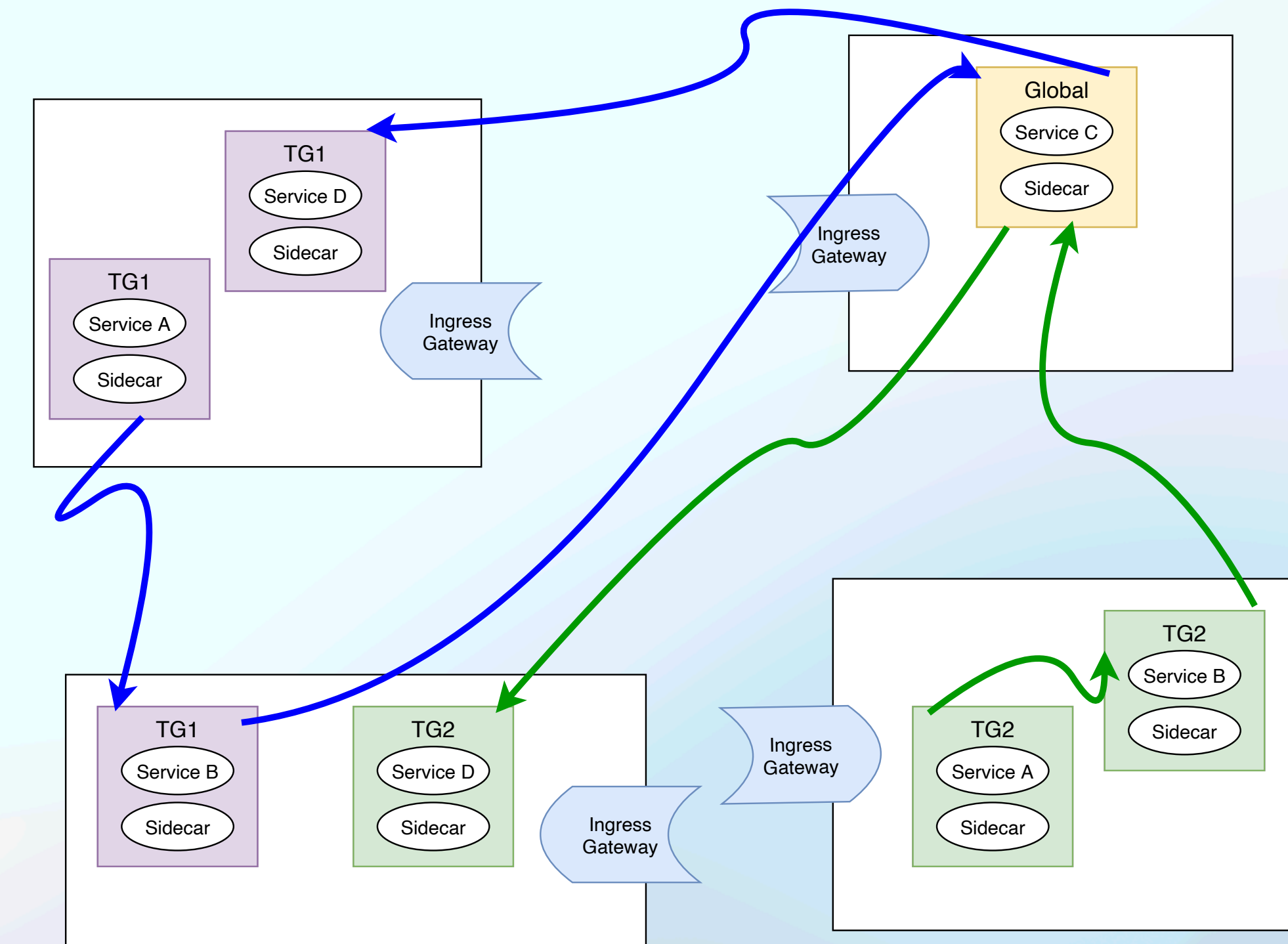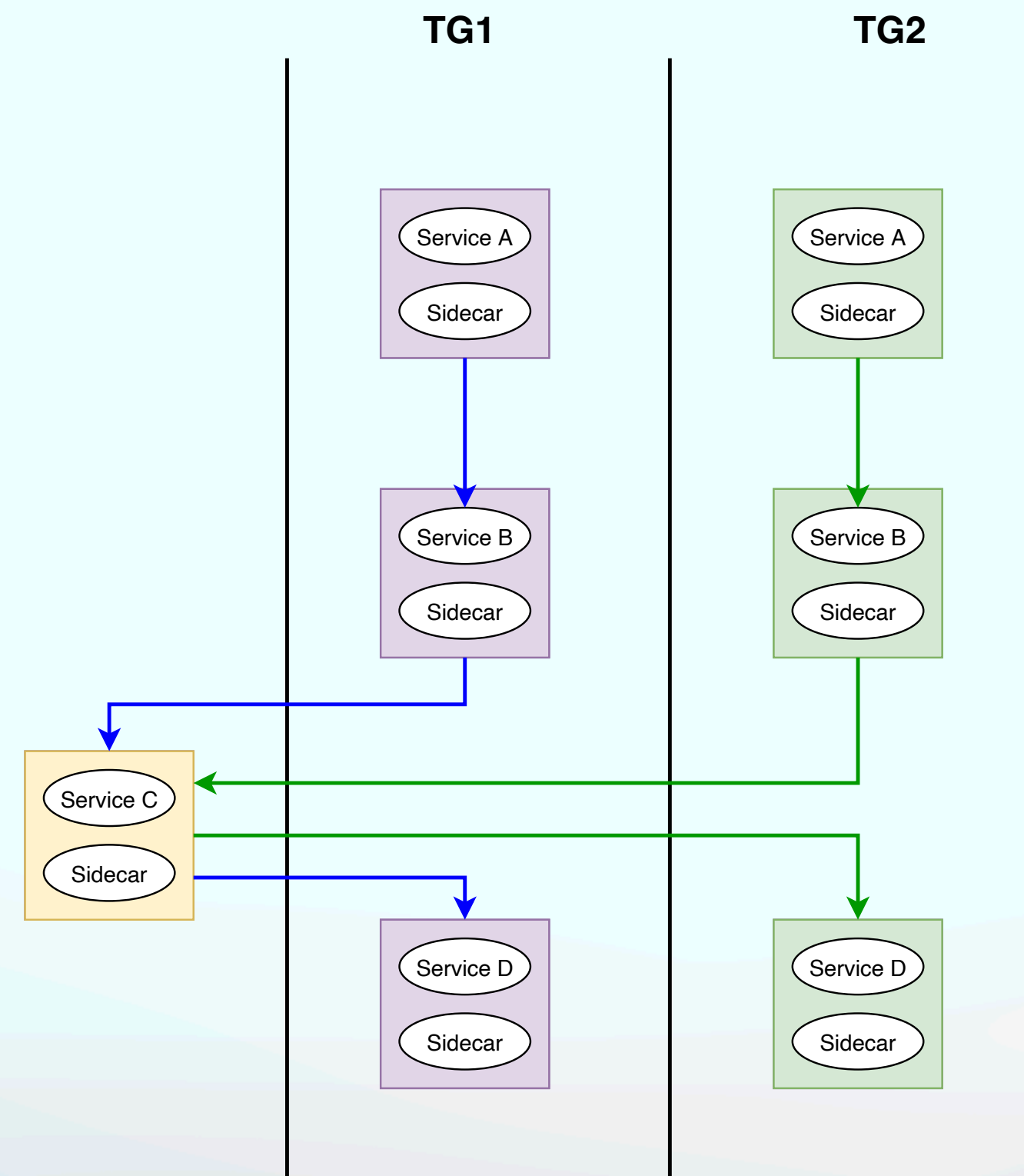**Advanced Use Cases**

- WCP

- API Gateway, Bubble

**Opportunities**

- **Proximity based routing**

  - Solved at GSLB via regional proximity. Service Mesh has a "WCNP Proximity" feature too.

  - Still, we could do more: capacity-aware, AZ-aware, load-aware, weight-based, and other advanced routing

- **API-based routing**

  - Mesh supports URI and header based routing, but we don't support a formal notion of APIs. If we did, we could make mesh aware of API-based capacity, weights, topology. proximity/colocation, availability, failovers.

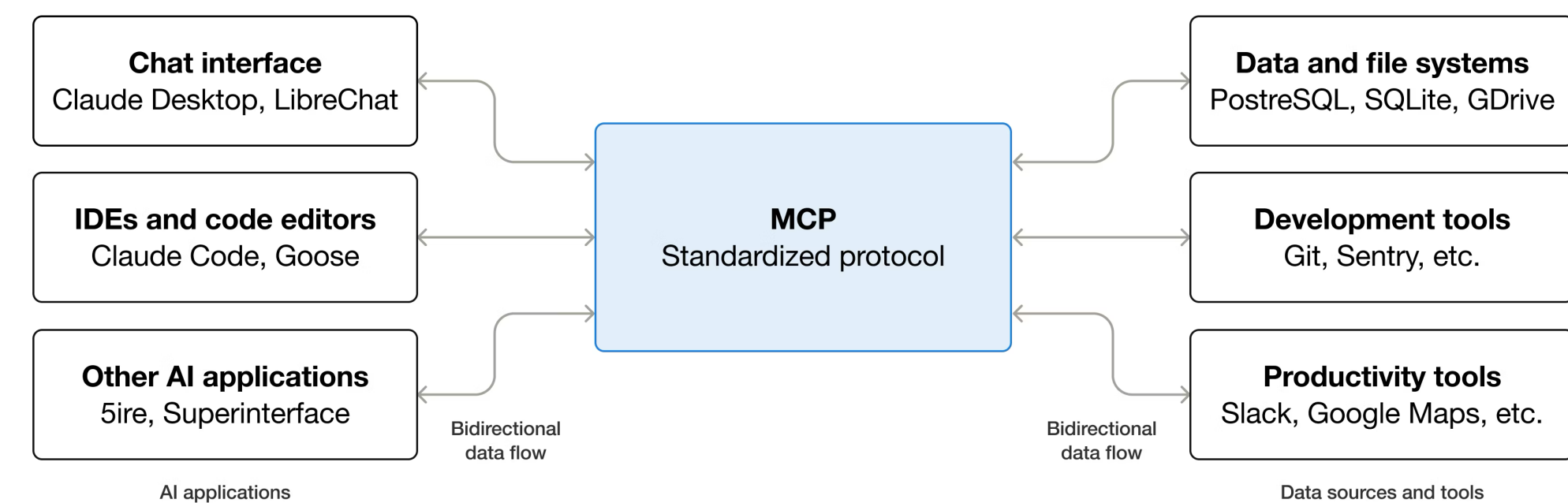  - We'll be circling back to this one in the AI world.
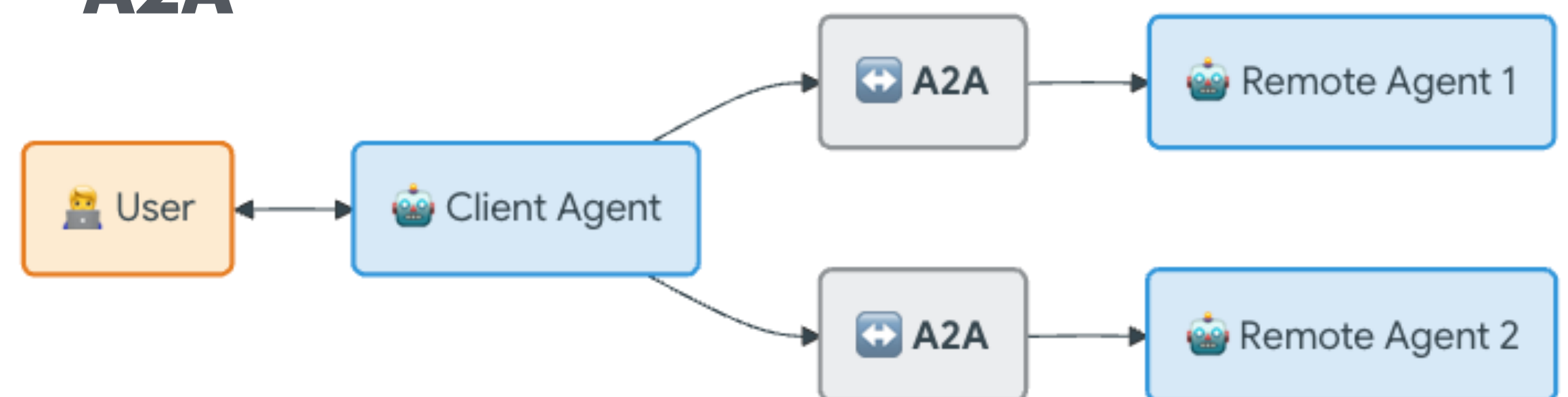
# AI World

## The Journey

- LLMs

- LLM Client Apps (e.g. ChatGPT, CoPilot)

- LLM Apps in enterprises: Need for integration with enterprise data

- **Enter MCP**

- The need for smarter clients: Agents

- The need for agents to converse

- **Enter A2A**

- LLM <- Agent <-> [A2A] <-> Agent -> [MCP] -> Services



**MCP**

| | | |
|---|---|---|
| **Chat interface** Claude Desktop, LibreChat | | **Data and file systems** PostreSQL, SQLite, GDrive |
| **IDEs and code editors** Claude Code, Goose | **MCP** Standardized protocol | **Development tools** Git, Sentry, etc. |
| **Other AI applications** 5ire, Superinterface | | **Productivity tools** Slack, Google Maps, etc. |

Bidirectional data flow

AI applications

Bidirectional data flow
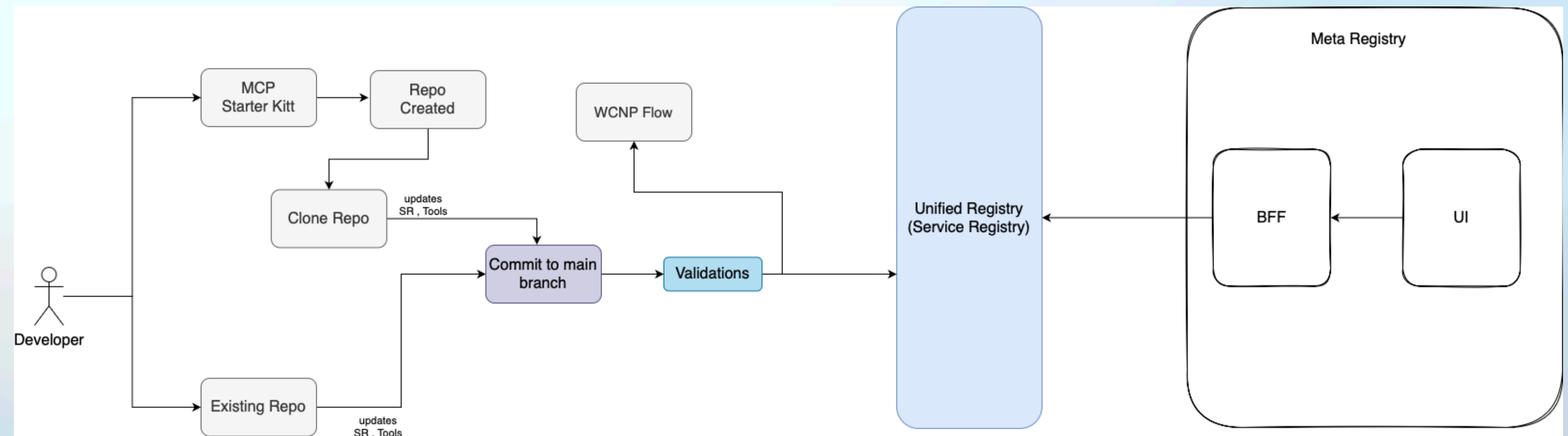
Data sources and tools



**A2A**

# AI World

## MCP Workflow

- MCP servers will mostly be the existing microservices in new clothing.

  - Some may be deployed as external servers that access the micro service via API to deliver tool functionality.

  - Others may be built embedded within the microservice as a new API/URI/Port.

- As such, it only follows that the MCP servers are registered into SR as microservices using the existing workflow.

- Gitops will register the server details in SR under a newly introduced "MCP" capability.

- Service Mesh will pick up the SR configs and enable routing/policies/resiliency/metrics for the MCP servers automatically.
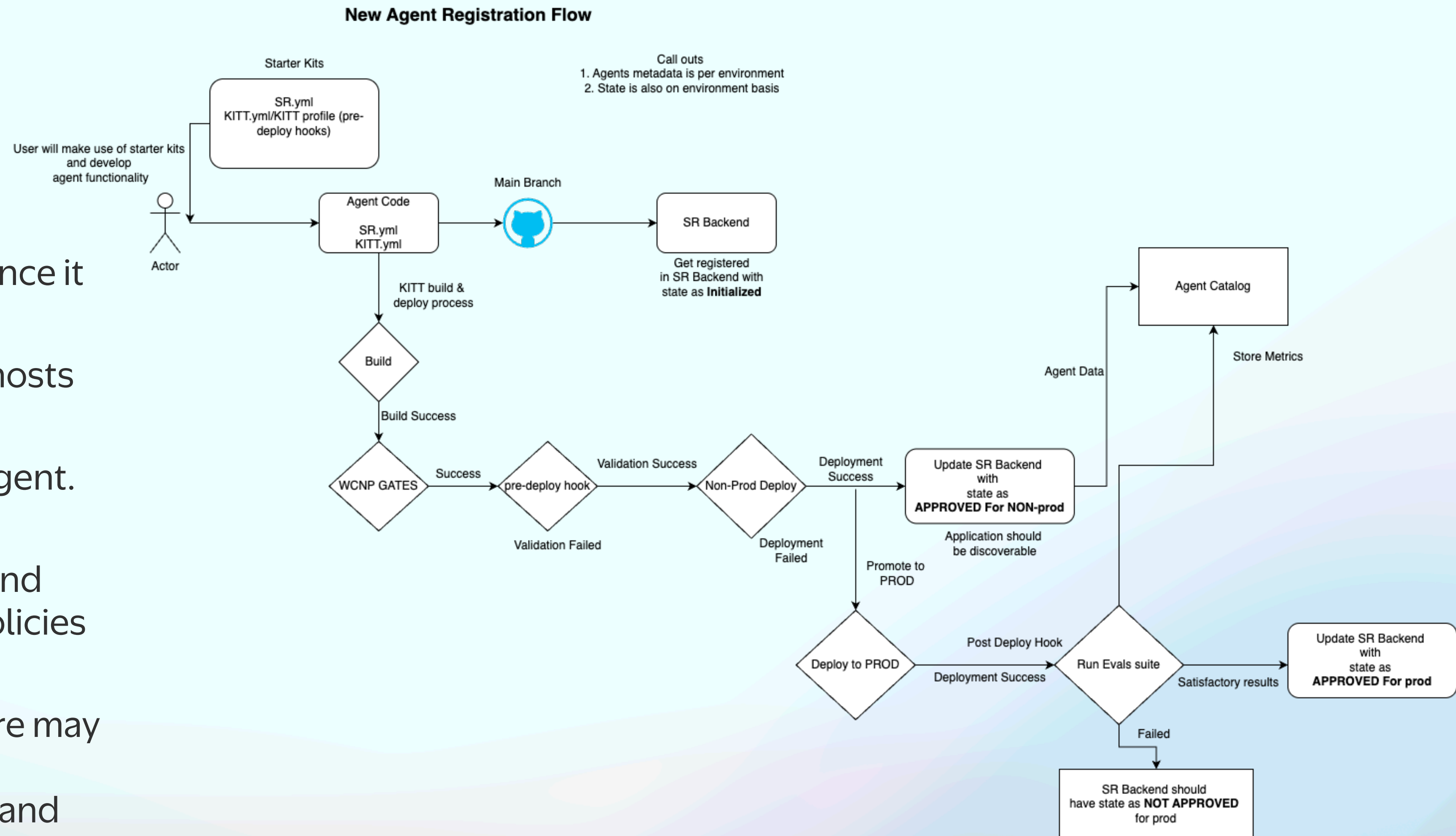
```
serviceType: MCP
wcnpProximity: true
mcp:
  displayName: Testhub MCP Server
  version: 0.0.1
  url: https://testhub-mcp-server.dev.walmart.com
  endpoint: /mcp
  healthcheck: /health
  category: CI/CD - Test Engineering
  publisher: Developer Experience
  language: Typescript
  auth:
    type: PingFed Token
    header: Authorization
  protocol:
    message: JSON-RPC
    transport: STREAMABLE_HTTP
    version: 2025-06-18
  disabled: false
```
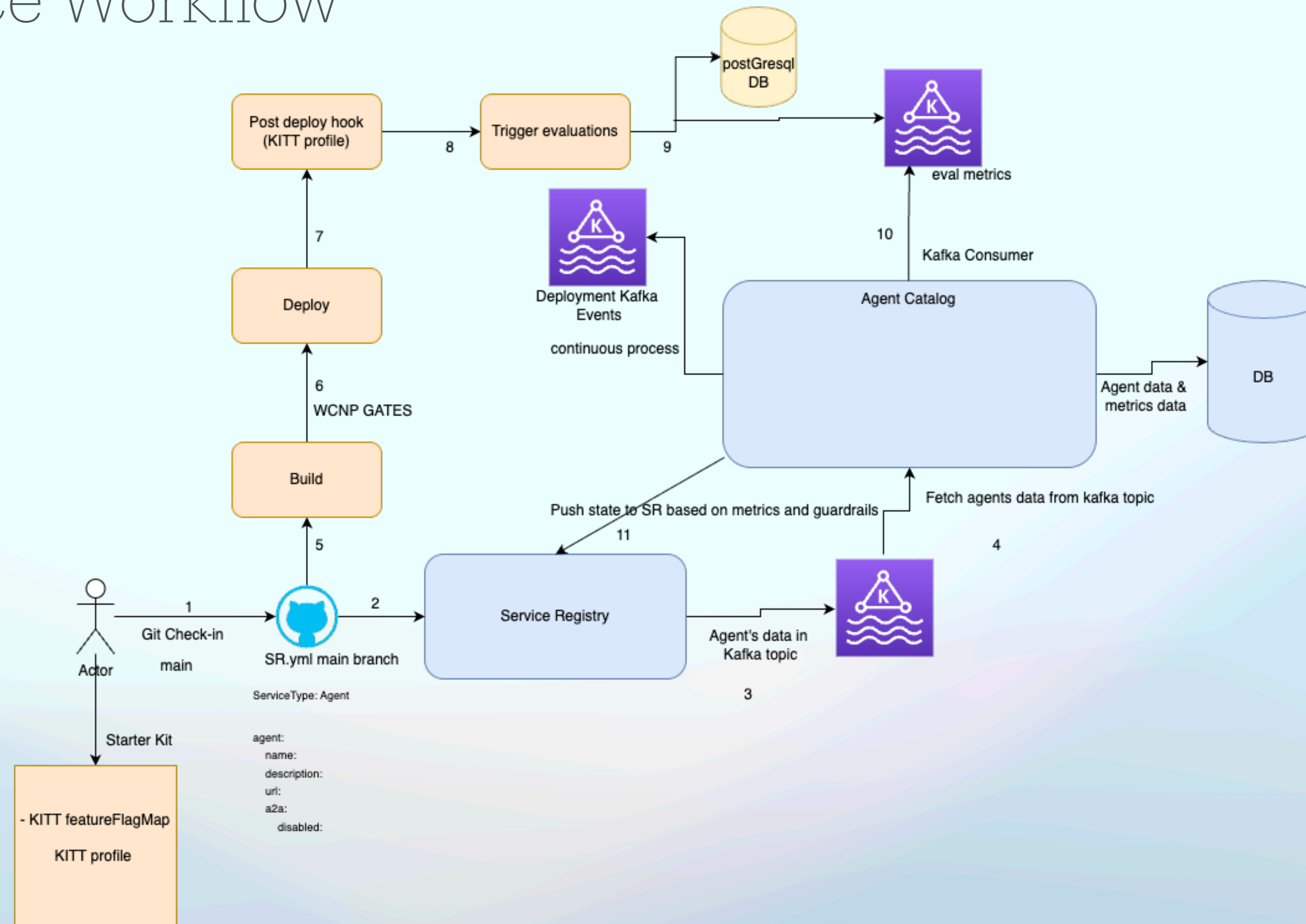
# AI World

## Agents Workflow

- Agent definition is more nuanced since it may or may not be a "service".

- A2A has a notion of a "Server" that hosts the agents.

- A server may not be 1-to-1 with an agent.

- It's the container server that will be registered in SR as a microservice, and service mesh will enable routing/policies for it.

- Since the space is still evolving, there may be additional needs in the future to perform more fine-grained routing and policy enforcement (at tools/tasks level).
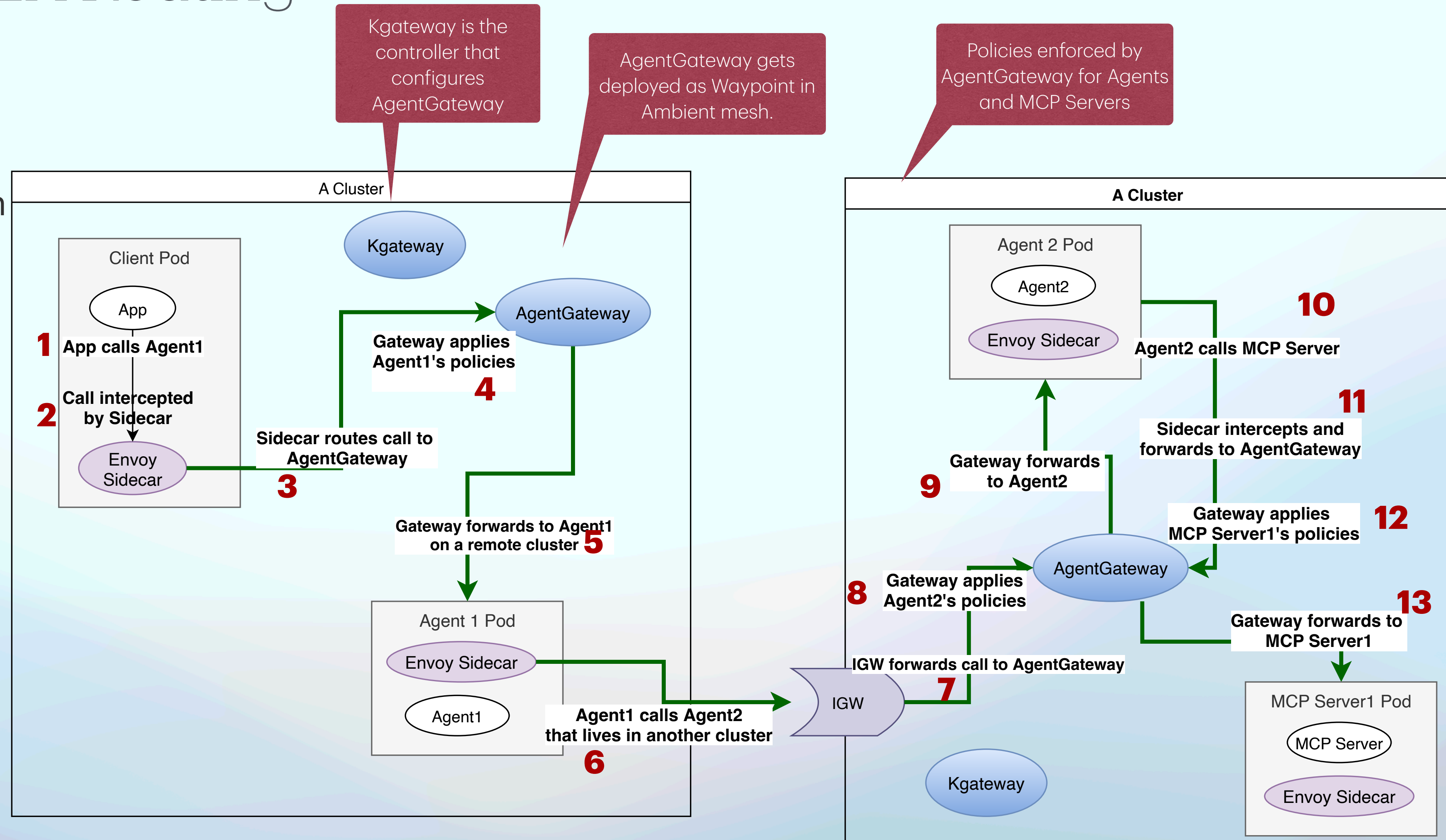


**New Agent Registration Flow**

# AI World

## A2A Governance Workflow

# AI World

## MCP and A2A Routing

- Demo
- QnA