

# Automated Disaster Recovery for Multi-Region GenAI applications using vector databases on Kubernetes



Sergey Pronin  
Group Product Manager



Paweł Bojanowski  
Principal Software Engineer

# Agenda

1. Intro - “AI is magic”
2. Problem space
  - a. pgvector on Kubernetes 101
  - b. Multi-region deployments
  - c. Disaster Recovery (DR) problems
3. Automating DR
4. Demo
5. Q&A

✨ AI is Magic

# Creativity and imagination



Pretty kitten working behind the laptop



Unicorn plays hockey



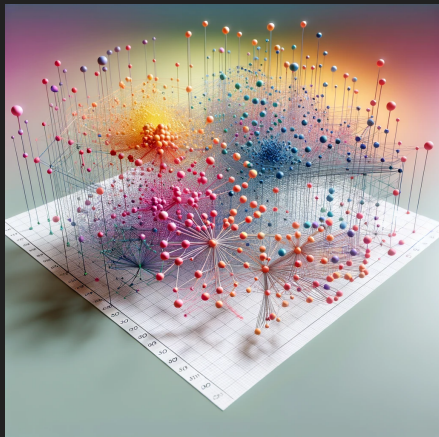
Dad became a monster, cause we broke  
his Lego

# AI is magic, but for Muggles

Various technologies power AI:

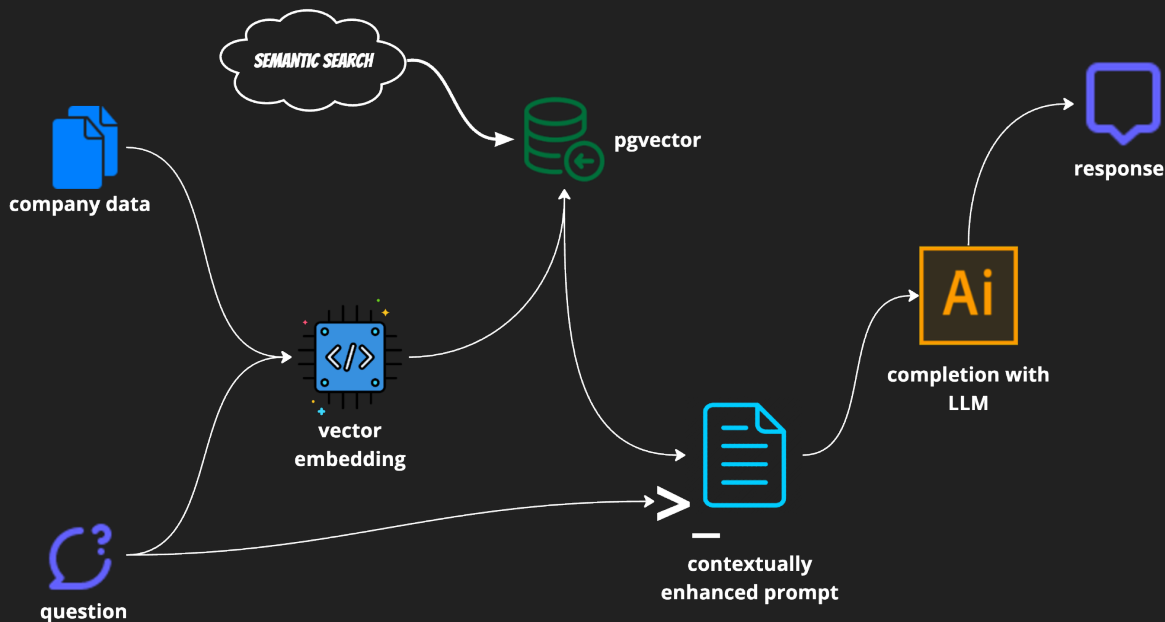
- Large Language Models - heart of modern Generative AI
  - Powered by Transformers (Attention is All You Need)
- GPUs for training
- Vector databases for semantic search, visual or multimodal search
  - Through vector embeddings

# Everything is a vector



This is how DALL-E sees vector embedding

Example of RAG - Retrieval Augmented Generation



# Vector databases

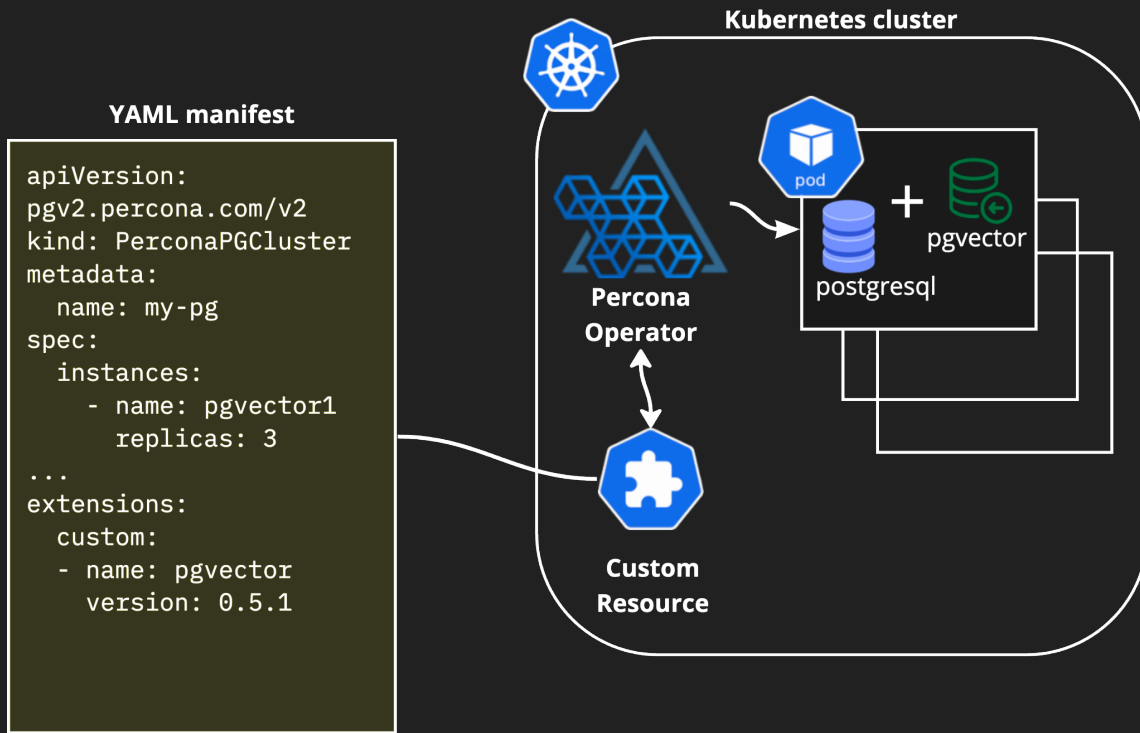
## Various options

- Pinecone
- Milvus
- Chroma
- FAISS
- pgvector (extension for PostgreSQL)



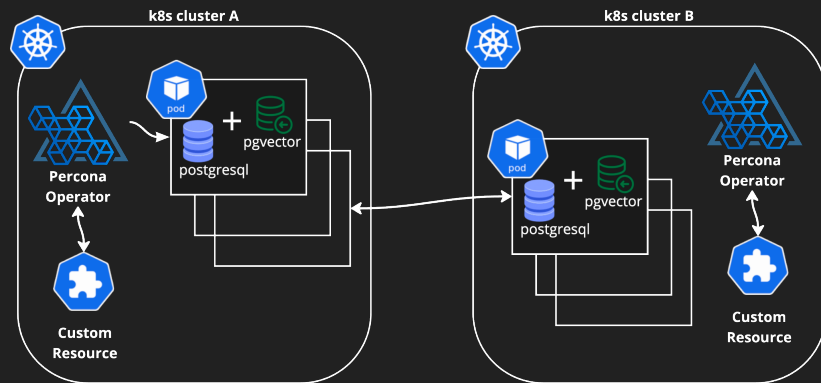
Problem space

# pgvector on Kubernetes



# Why multi-region

- Not only multi-region:
  - multi-cloud
  - hybrid-cloud
- Main reasons
  - Disaster recovery
  - Migration
  - No-vendor lock-in / cloud-ready



# Disaster recovery

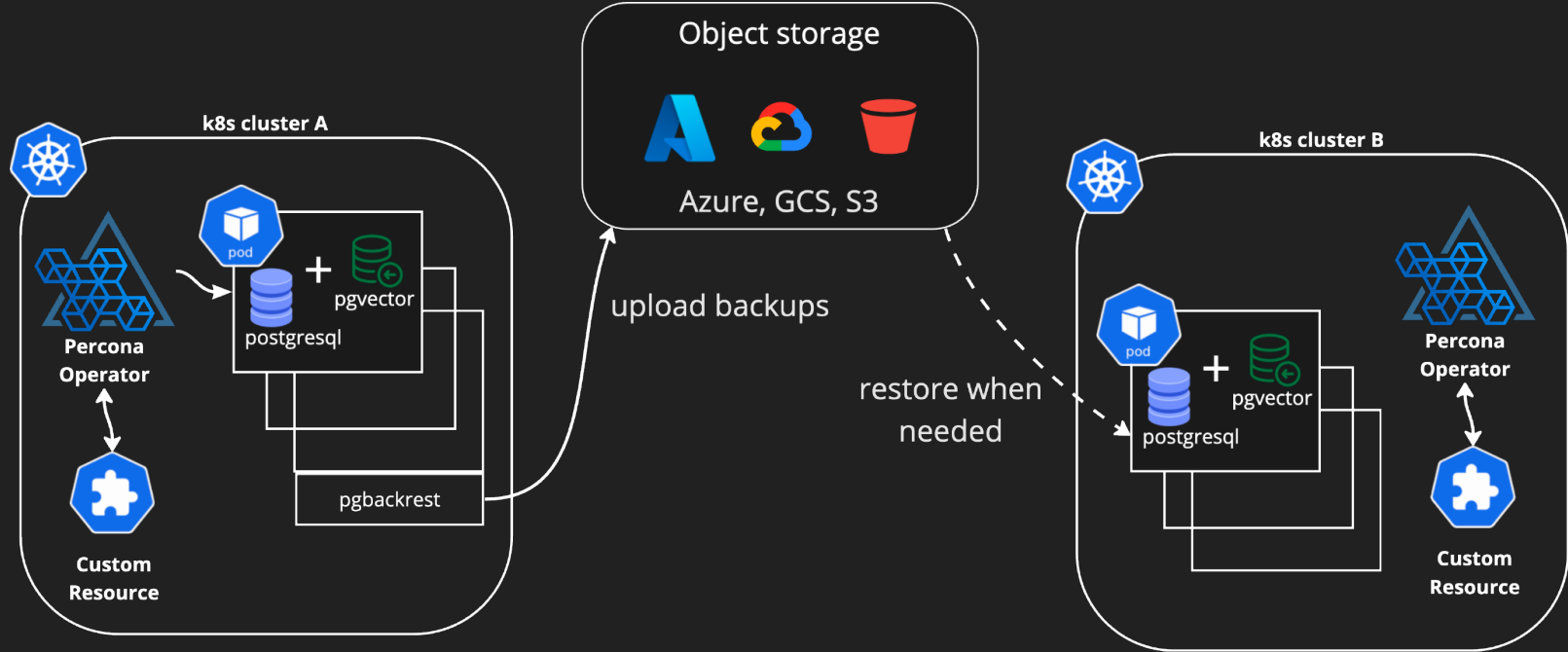
## Recovery Time Objective (RTO)

- How fast do you recover from failure?

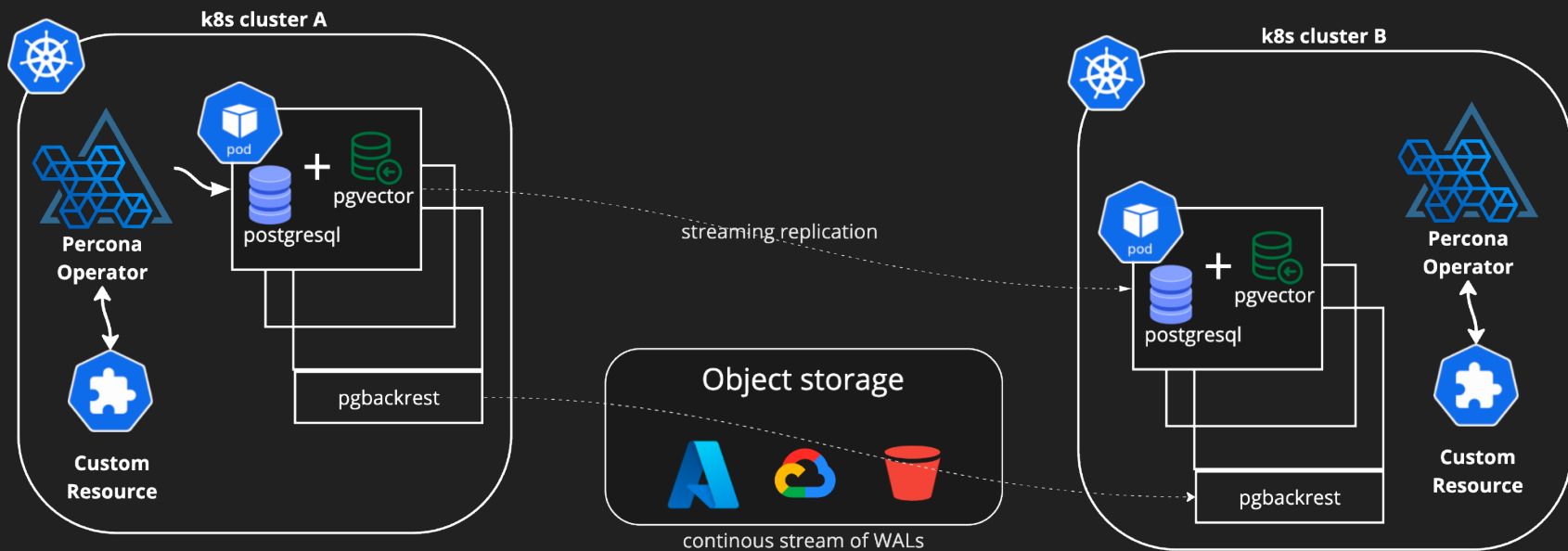
## Recovery Point Objective (RPO)

- What is your last transaction?

DR through Backups - high RTO



## DR through Replication - low RTO



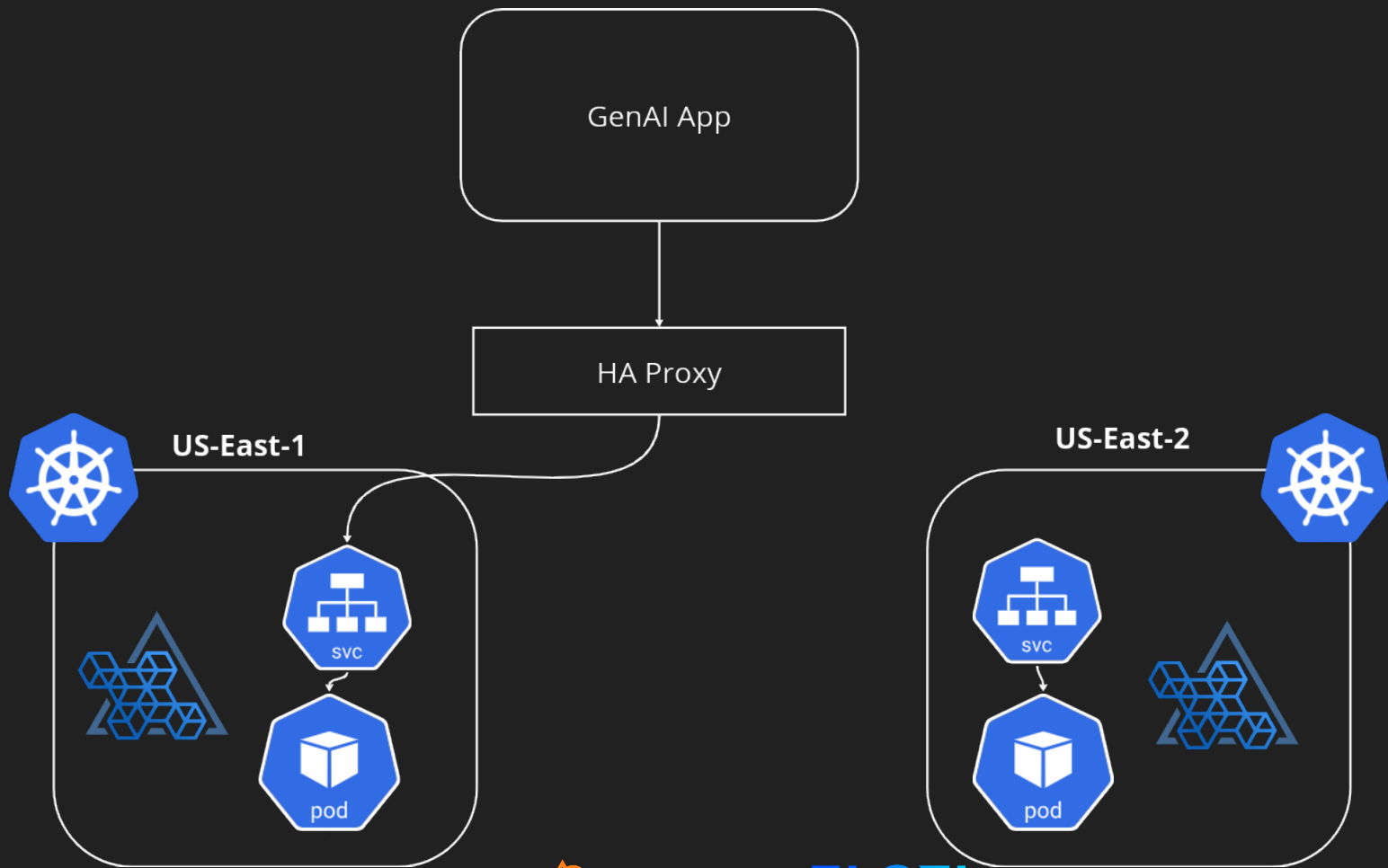
# Looks like problem solved?

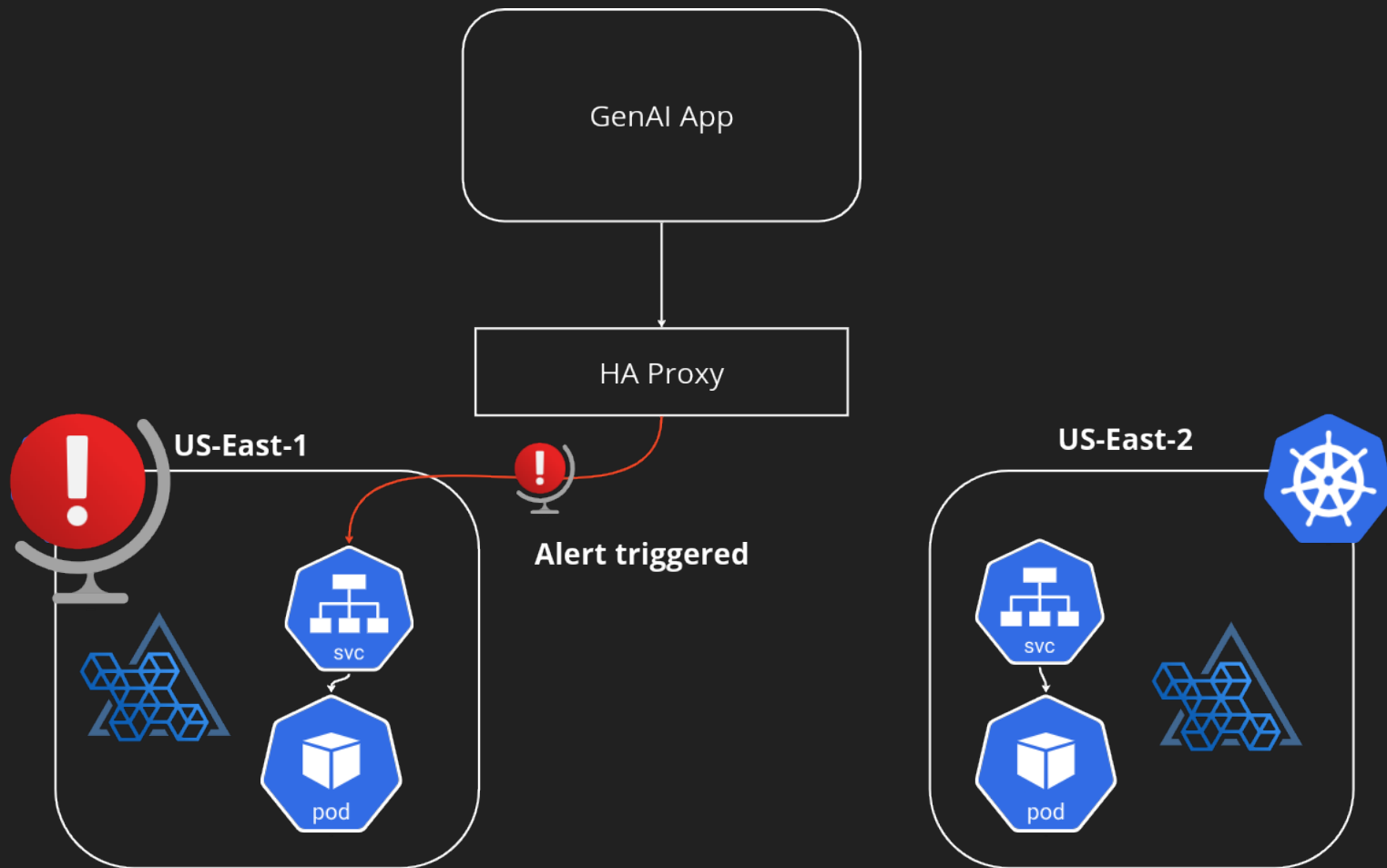
1. Recover fast from failure
2. Recovery Point Objective is solved on DB level
  - a. For streaming replication it can be down to almost 0

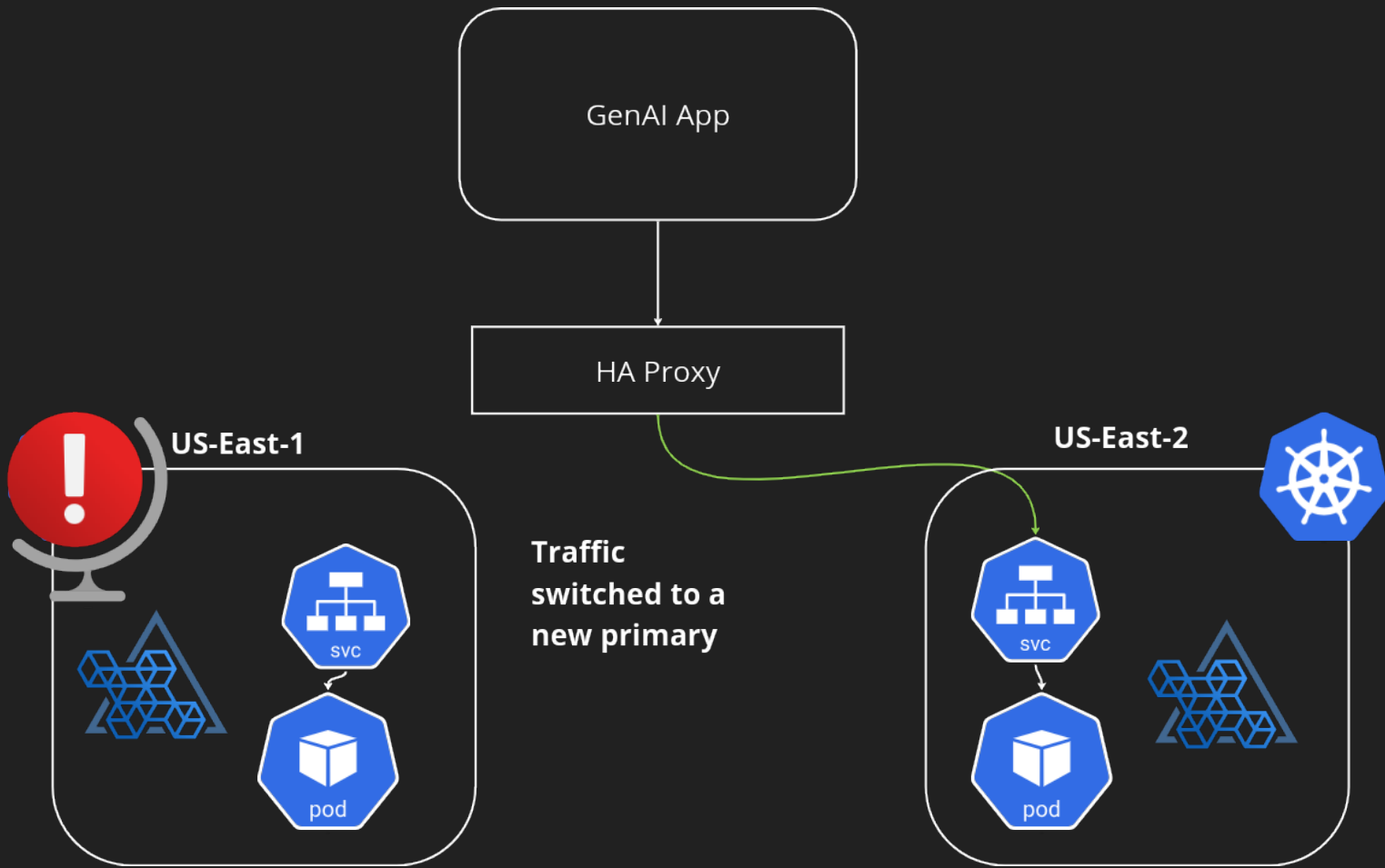
# What about the app and failover?

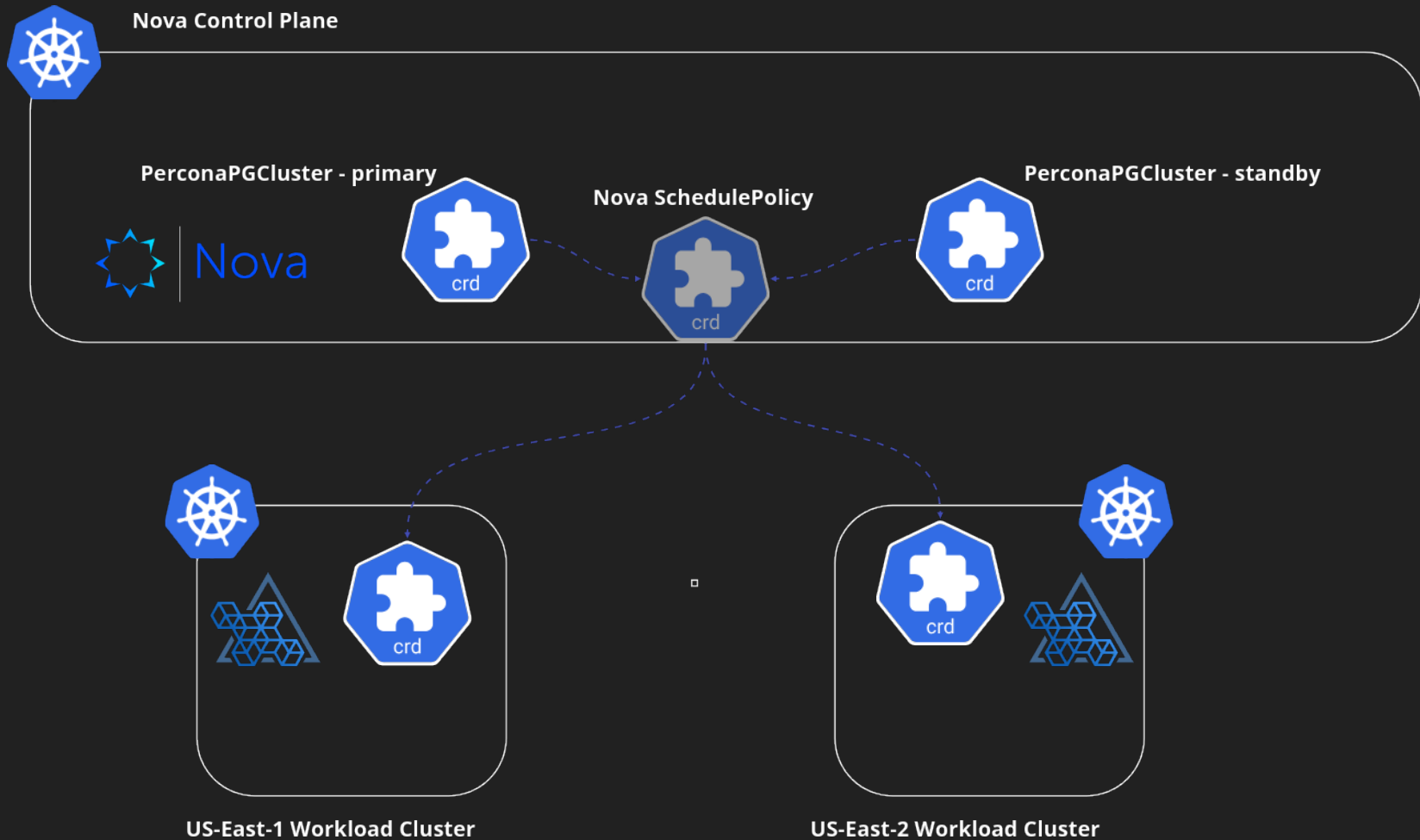
1. Failover is not automated
  - Operators can detect failure, but 3rd agent needed (quorum)
2. Point application to correct database
  - Which endpoint should the app use when failover happens?
  - Who is reconfiguring it and how?
3. Failback
  - What should happen when main site goes back up?

💡 Automating DR











Nova Control Plane

PerconaPGCluster - primary



crd

Nova SchedulePolicy



crd

PerconaPGCluster - standby



crd



crd



US-East-1 Workload Cluster



crd



US-East-2 Workload Cluster

Yaml manifest

```
apiVersion: pgv2.percona.com/v2
kind: PerconaPGCluster
metadata:
  name: cluster1
spec:
  standby:
    enabled: false
```

Yaml manifest

```
apiVersion: pgv2.percona.com/v2
kind: PerconaPGCluster
metadata:
  name: cluster2
spec:
  standby:
    enabled: true
```



PERCONA

ELOTL



```
→ nova git:(fixes-percona) X
```

I



PERCONA

ELOTL

```
apiVersion: recovery.elotl.co/v1alpha1
kind: RecoveryPlan
metadata:
  name: psql-primary-failover-plan
spec:
```

```
  alertLabels:
    app: example-app
```

```
  steps:
```

```
- type: patch
  patch:
    apiVersion: "pgv2.percona.com/v2"
    resource: "perconapgclusters"
    namespace: "psql-operator"
    name: "cluster1"
    override:
      fieldPath: "spec.standby.enabled"
      value:
        raw: false
    patchType: "application/merge-patch+json"
```

```
- type: patch
  patch:
    apiVersion: "pgv2.percona.com/v2"
    resource: "perconapgclusters"
    namespace: "psql-operator"
    name: "cluster2"
    override:
      fieldPath: "spec.standby.enabled"
      value:
        raw: true
    patchType: "application/merge-patch+json"
```

**Matches metadata of  
incoming alert (e.g. from  
prometheus)**

```
groups:
- name: example
  rules:
  - alert: HighRequestLatency
    expr: percona_pg_cluster_failure_total(cluster="cluster-1") > 0
    for: 10m
    labels:
      app: example-app
```

[https://prometheus.io/docs/alerting/latest/configuration/#webhook\\_config](https://prometheus.io/docs/alerting/latest/configuration/#webhook_config)



PERCONA

ELOTL

```
apiVersion: recovery.elotl.co/v1alpha1
kind: RecoveryPlan
metadata:
  name: psql-primary-failover-plan
spec:
  alertLabels:
    app: example-app
  steps:
```

```
- type: patch
  patch:
    apiVersion: "pgv2.percona.com/v2"
    resource: "perconapgclusters"
    namespace: "psql-operator"
    name: "cluster1"
    override:
      fieldPath: "spec.standby.enabled"
      value:
        raw: true
    patchType: "application/merge-patch+json"
```

Puts PerconaPGCluster  
primary into standby

```
- type: patch
  patch:
    apiVersion: "pgv2.percona.com/v2"
    resource: "perconapgclusters"
    namespace: "psql-operator"
    name: "cluster2"
    override:
      fieldPath: "spec.standby.enabled"
      value:
        raw: false
    patchType: "application/merge-patch+json"
```

```
- type: readField
```



PERCONA

ELOTL

```
apiVersion: recovery.elotl.co/v1alpha1
kind: RecoveryPlan
metadata:
  name: psql-primary-failover-plan
spec:
  alertLabels:
    app: example-app
  steps:
```

```
- type: patch
  patch:
    apiVersion: "pgv2.percona.com/v2"
    resource: "perconapgclusters"
    namespace: "psql-operator"
    name: "cluster1"
    override:
      fieldPath: "spec.standby.enabled"
      value:
        raw: true
    patchType: "application/merge-patch+json"
```

```
- type: patch
  patch:
    apiVersion: "pgv2.percona.com/v2"
    resource: "perconapgclusters"
    namespace: "psql-operator"
    name: "cluster2"
    override:
      fieldPath: "spec.standby.enabled"
      value:
        raw: false
    patchType: "application/merge-patch+json"
```

Turns PerconaPGCluster  
standby into primary

```
- type: readField
```



PERCONA

ELOTL

```
- type: patch
patch:
  apiVersion: "pgv2.percona.com/v2"
  resource: "perconapgclusters"
  namespace: "psql-operator"
  name: "cluster1"
  override:
    fieldPath: "spec.standby.enabled"
    value:
      raw: true
  patchType: "application/merge-patch+json"
```

```
- type: patch
patch:
  apiVersion: "pgv2.percona.com/v2"
  resource: "perconapgclusters"
  namespace: "psql-operator"
  name: "cluster2"
  override:
    fieldPath: "spec.standby.enabled"
    value:
      raw: false
  patchType: "application/merge-patch+json"
```

```
- type: readField
readField:
  apiVersion: "pgv2.percona.com/v2"
  resource: "perconapgclusters"
  namespace: "psql-operator"
  name: "cluster2"
  fieldPath: "status.host"
  outputKey: "Cluster2IP"
```

Reads address of a new  
primary

```
- patch:
  apiVersion: "v1"
  resource: "configmaps"
  namespace: "default"
  name: "haproxy-config"
```



PERCONA

ELOTL

```
- type: patch
patch:
  apiVersion: "pgv2.percona.com/v2"
  resource: "perconapgclusters"
  namespace: "psql-operator"
  name: "cluster2"
  override:
    fieldPath: "spec.standby.enabled"
    value:
      raw: false
  patchType: "application/merge-patch+json"
```

```
- type: readField
readField:
  apiVersion: "pgv2.percona.com/v2"
  resource: "perconapgclusters"
  namespace: "psql-operator"
  name: "cluster2"
  fieldPath: "status.host"
  outputKey: "Cluster2IP"
```

Updates HAProxy  
configuration with an  
address of a new  
primary

```
- patch:
  apiVersion: "v1"
  resource: "configmaps"
  namespace: "default"
  name: "haproxy-config"
  override:
    fieldPath: "data"
    value:
      raw: {"haproxy.cfg": "defaults\n  mode tcp\n  timeout connect 5000ms\n  timeout client 50000ms\n  timeout server\n50000ms\n\nfrontend fe_main\n  bind *:5432\n  default_backend be_db_2\n\nbackend be_db_2\n  server db2 {{ .Values.Cluster2IP\n}}:5432 check"}
  patchType: "application/merge-patch+json"
type: patch
```

# Demo

Untitled x +

localhost:8888/notebooks/Untitled.ipynb

jupyter Untitled Last Checkpoint: 2 days ago

File Edit View Run Kernel Settings Help Trusted

JupyterLab Python 3 (ipykernel)

```
HAPROXY_HOST = "a8f74e18d50d145f2baf78876061c21d-1459685038.us-east-2.elb.amazonaws.com:5432/zoo?sslrootcert=/tmp/tls/ca.crt"
```

```
[92]:
```

```
[99]: print(HAPROXY_HOST)
```

```
a8f74e18d50d145f2baf78876061c21d-1459685038.us-east-2.elb.amazonaws.com:5432/zoo?sslrootcert=/tmp/tls/ca.crt
```

```
*[101]: conn = psycopg2.connect(
        f"{CONN_CREDS}@{HAPROXY_HOST}"
    )
    cur = conn.cursor()

    while True:
        cur.execute("""
            select * from test.perconavec;
        """)
        print(cur.fetchmany(2))
        time.sleep(3)
    conn.close()
```

```
[(5, 'This framework generates embeddings for each input sentence', 'dummy1', array([-0.01371733, -0.04285153, -0.01562861, ..., 0.0250469
7,
      -0.05183711, -0.04365626], dtype=float32)), (6, 'Sentences are passed as a list of string.', 'dummy2', array([ 0.05645249, 0.0550024
3, 0.0313796, ..., 0.01458188,
      0.00743343, -0.03071318], dtype=float32))]
[(7, 'The quick brown fox jumps over the lazy dog.', 'dummy3', array([ 0.04393358, 0.05893442, 0.04817839, ..., -0.10047367,
      -0.03837524, -0.05880814], dtype=float32))]
[]
[]
[]
[]
[]
```

# Nova's secret sauce: Schedule Policies



**apiVersion: pgv2.percona.com/v2**

kind: PerconaPGCluster

metadata:

name: cluster1

namespace: psql-operator

labels:

psql-cluster: cluster-1

spec:

...



Nova

apiVersion: policy.elotl.co/v1alpha1

kind: SchedulePolicy

metadata:

name: psql-cluster-1

spec:

namespaceSelector:

matchLabels:

kubernetes.io/metadata.name: psql-operator

clusterSelector:

matchLabels:

kubernetes.io/metadata.name: us-east-1-wlc

resourceSelectors:

labelSelectors:

- matchLabels:

psql-cluster: cluster-1



PERCONA

ELOTL



```
kind: CustomResourceDefinition
metadata:
  labels:
    psql-cluster: all
names:
  kind: PerconaPGBackup
  ...
  ---
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
metadata:
  labels:
    psql-cluster: all
names:
  kind: PerconaPGCluster
  ...
  ---
apiVersion: apps/v1
kind: Deployment
metadata:
  labels:
    psql-cluster: all
  ...
```



Nova

```
apiVersion: policy.elotl.co/v1alpha1
kind: SchedulePolicy
metadata:
  name: psql-all-clusters
spec:
  namespaceSelector:
    matchLabels:
      kubernetes.io/metadata.name: psql-operator
  groupBy:
    labelKey: psql-cluster
    clusterSelector:
      matchExpressions:
        - key: kubernetes.io/metadata.name
          operator: In
          values:
            - us-east-1-wlc
            - us-east-2-wlc
      spreadConstraints:
        spreadMode: Duplicate
        topologyKey: kubernetes.io/metadata.name
    resourceSelectors:
      labelSelectors:
        - matchLabels:
            psql-cluster: all
```



PERCONA

ELOTL

# Thank you!

Learn more:

Percona Operators



Elotl Nova

