# COST EFFECTIVE KUBERNETES
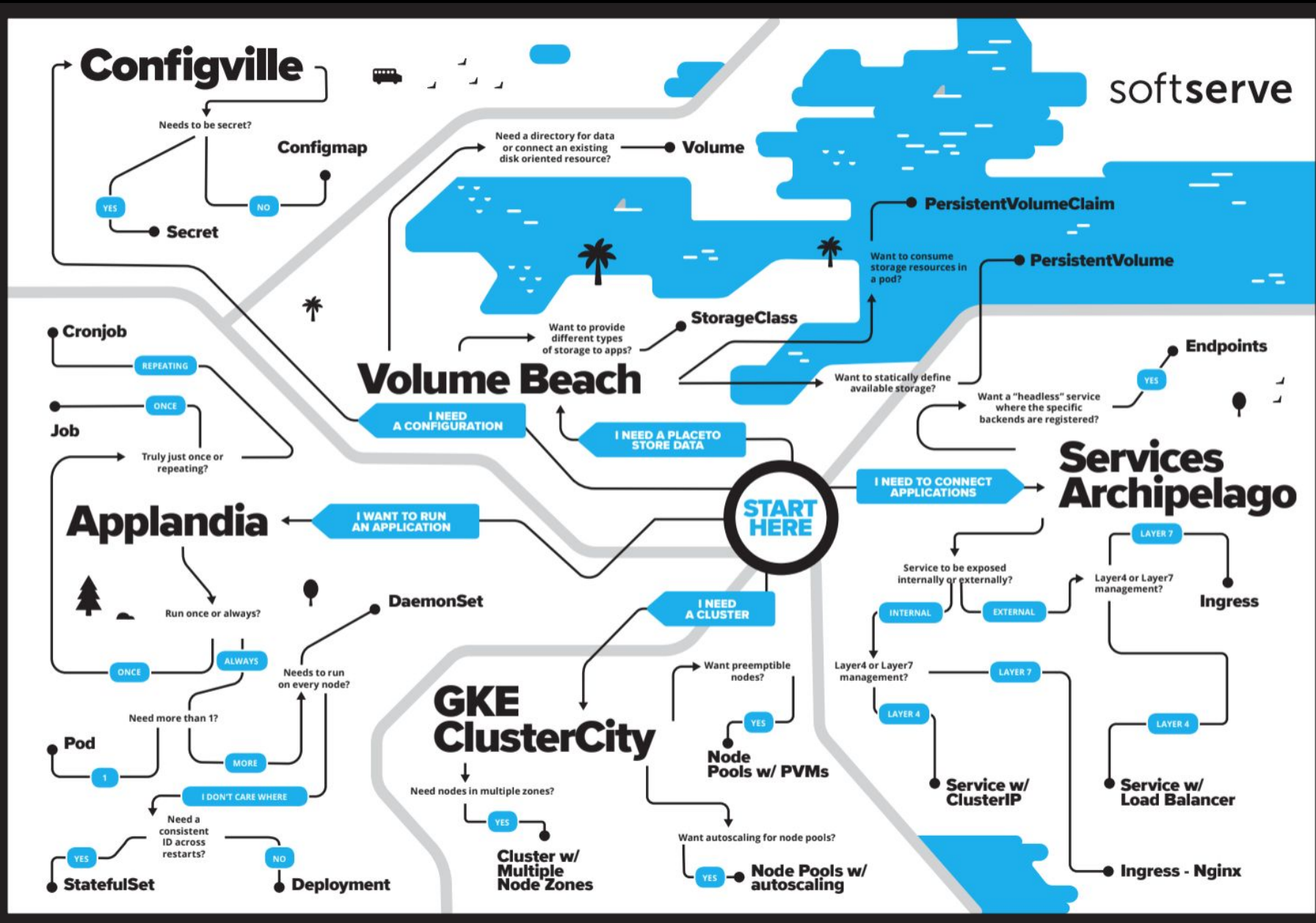
Myroslav Rys — Softserve
Ryan Richard — Google
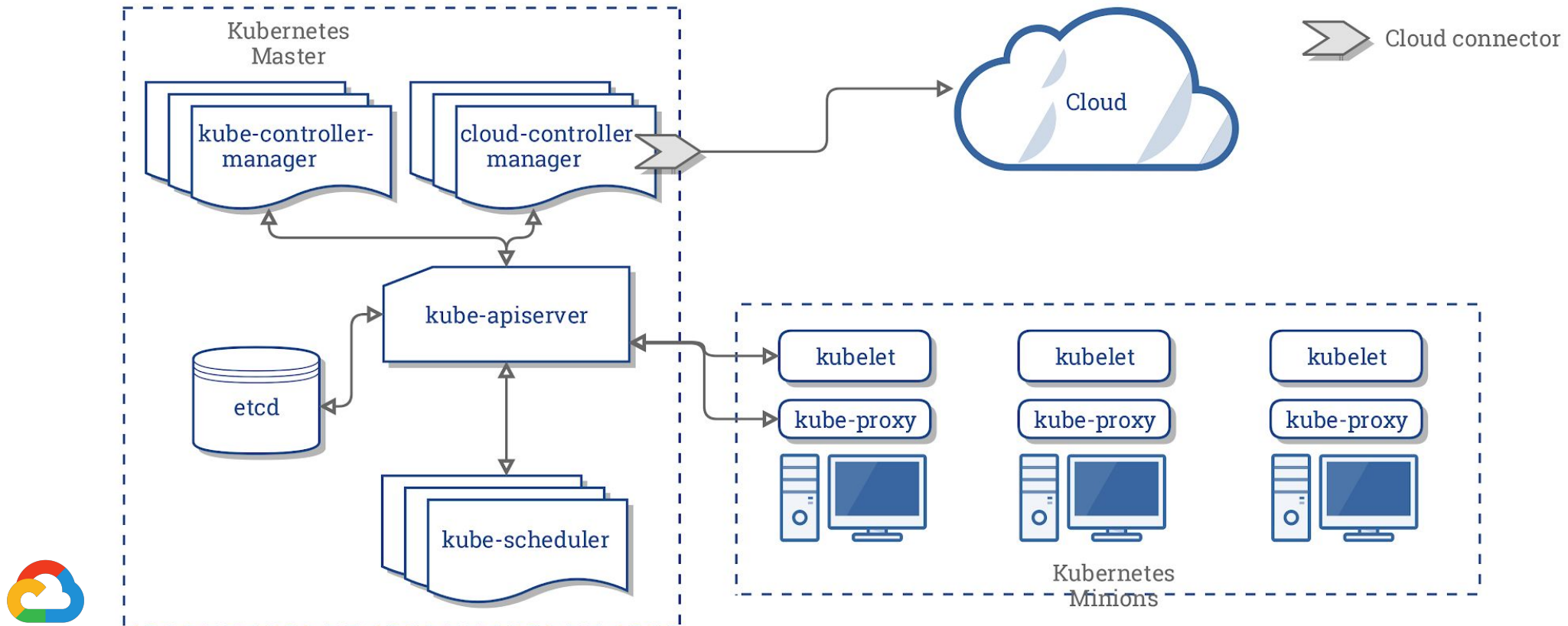
softserve
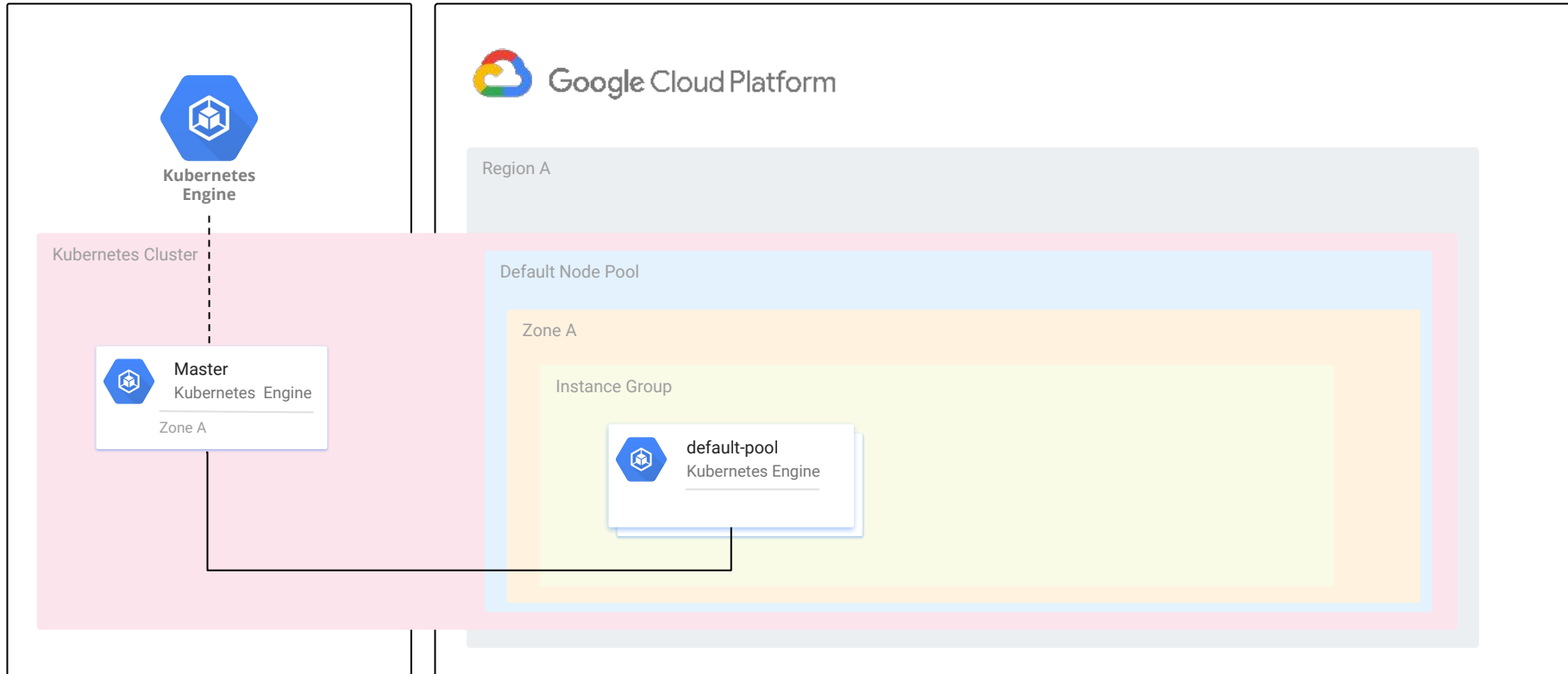
# KUBERNETES ARCHITECTURE FOUNDATION

# GKE — SINGLE ZONE / NODE POOL



Kubernetes Engine

Kubernetes Cluster

Google Cloud Platform

Region A

Default Node Pool

Zone A

Instance Group

Master
Kubernetes Engine
Zone A

default-pool
Kubernetes Engine

softserve

# GKE — SINGLE NODE POOL, MULTI-ZONE

Kubernetes Engine

Google Cloud Platform

Kubernetes Cluster

Region A

Default Node Pool

Master
Kubernetes Engine
Zone A

Zone A

Instance Group

default-pool
Kubernetes Engine

Zone B

Instance Group

default-pool
Kubernetes Engine

Zone C

Instance Group

default-pool
Kubernetes Engine

softserve

# GKE CLUSTER AUTOSCALER

Allows for Autoscaling of the compute nodes which make up the cluster.

```
gcloud beta container clusters create meetup \
--enable-autoscaling \
--min-nodes 2 \
--max-nodes 6
...
```

# CLUSTER AUTOSCALER VS CUSTOM CONTROLLER OR MANUAL MANAGEMENT

Scaling Approach
GKE

**Less 1000 Nodes 30 Pods each**
— YES → **Nodes in Pools same instance type**
— YES → **Pools size should even balanced across zone**
— YES → **Services are disruption tolerant**
— YES → Autoscaler GKE

**Nodes in Pools same instance type** — NO →

**Pools size should even balanced across zone** — NO →

**Services are disruption tolerant** — NO → **Less 20 Nodes. Static load**

**Less 20 Nodes. Static load** — YES → Manual Scheduler GKE

**Less 20 Nodes. Static load** — NO → Custom Scheduler GKE

**Less 1000 Nodes 30 Pods each** — NO → Custom Scheduler GKE

# GKE CLUSTER AUTOSCALER LIMITATIONS

- Cluster autoscaler works based on **Pod resource requests**

- Cluster autoscaler does **not track labels manually added** after initial cluster or node pool creation

- Cluster autoscaler considers the relative cost of each instance type in the node pool and attempts to expand the **least expensive** possible node pool

- Cluster with multiple node pools with the same instance type, cluster autoscaler will attempt to keep those node pools' sizes **balanced**

- Maximum period for **graceful termination** for a Pod up to 10 minutes

softserve

# CAPACITY VS COST (autoscaling)

Capacity

100%

50%  A  B  C

0  Zone

Cost

50%

100%

**BEST CASE:**

**50%** of spend,
**100%** of the time

**WORST CASE:**

**100%** of spend,
**100%** of the time

regular instances

softserve

# PREEMPTIBLE VMS

A preemptible VM (PVM) is an instance that you can create and run at a much lower price than normal instances. However, Compute Engine might terminate (preempt) these instances if it requires access to those resources for other tasks.

~**80%** Discount!

**24HR** Life (max)

# PREEMPTIBLE VMS STATS

**580,000**

**20,000**

**10% - 15%**

softserve

# PREEMPTIBLE VMS STATS

**580,000**
cores for 1 HPC  workload

**10% - 15%**
Average Preemption Rate*

**$20,000**
over a weekend for HPC workload

softserve

# EXTEND GKE CLUSTER WITH PVM NODE POOL DEMO

softserve

# RECOMMENDED WORKLOAD FOR PVM NODE POOL

```
gcloud container node-pools create pvm-pool \
    --cluster $CLUSTER_NAME \
    --zone $CLUSTER_ZONE \
    --scopes cloud-platform \
    --enable-autoupgrade \
    --preemptible \
    --num-nodes 1 --machine-type g1-small \
    --enable-autoscaling --min-nodes=1 --max-nodes=6
```

softserve

# GKE — MULTI-ZONE / MULTIPLE NODE POOLS

Kubernetes
Engine

Google Cloud Platform

Kubernetes Cluster

Default Pool

Zone A

Instance Group

default-pool
Kubernetes Engine

Zone B

Instance Group

default-pool
Kubernetes Engine

Zone C

Instance Group

default-pool
Kubernetes Engine

Master
Kubernetes  Engine

Zone A

PVM Node Pool

Zone A

Instance Group

PVM Pool
Kubernetes Engine

Zone B

Instance Group

PVM Pool
Kubernetes Engine

Zone C

Instance Group

PVM Pool
Kubernetes Engine

# GKE — PREEMPTIBLE POOL DECISION MAP

Sample Workloads
Recommendation

Default

GKE
Default Pool

Preemptible

GKE
Preemptible Pool

Stateful

Static Load

SQL/NoSQL Servers
Not a Cloud SQL

Stateless Applications

Continuous Integration

QA and Testing

Web Apps with load spikes

Backend Apps with load spikes

Batch Processing

Analytics
Weekly/Daily

ML

Additional Node Pools

Additional Node Pools

Additional Node Pools

GPU
Preemptible

# EXAMPLE 1

Simple App manager by Autoscaler

# SAMPLE APP

Sample App

Static Load

| Backend |
| Stateful |

| Frontend |
| Static Load |

Dynamic Load

| Frontend |
| Spike Load |

**Old-style Java monolithic** enterprise app splitted for two major parts

- **Backend Java** old-style stateful service
- New and shiny **Node.js Frontend** service

  - Predictable static load
  - High load spikes at EOD/EOW

Nodes in preemptible pool will have label

`cloud.google.com/gke-preeptible: true`

soft**serve**

```
┌─────────────────────────┐                                                    ┌─────────────────────────┐
│ Create Initial Default  │                                                    │   Scale up pvm-pool     │
│ Pool with 1             │                                                    └─────────────────────────┘
│ node in every zone and  │                                                                │
│ preemptible Pool with 0 │                                                                ▼
│ nodes                   │                                                    ┌─────────────────────────┐
└─────────────────────────┘                                                    │ Scale up Frontend       │
            │                                                                   │ Service Pods            │
            ▼                                                                   └─────────────────────────┘
┌─────────────────────────┐                                                                │
│ Deploy Backend Service  │                                                                │
│ Pods with restriction   │                                                                │
│ only for default-pool.  │                                                                │
│ Min scale 2.            │                                                                │
│ Autoscaler will         │                                                                │
│ autobalance Pods across │                                                                │
│ zones                   │                                                                │
└─────────────────────────┘                                                                │
            │                                                                               │
            ▼                                                                               │
┌─────────────────────────┐                                                                │
│ Deploy Frontend Service │                                                                │
│ Pods with restriction   │                                        YES                      │
│ only for default-pool.  │                                         │                       │
│ Min scale 8.            │                                         │                       │
│ Autoscaler will         │                                         │                       │
│ autobalance Pods across │                                         │                       │
│ zones                   │                                         │                       │
└─────────────────────────┘                                         │                       │
            │                                                        │                       │
            ▼                                                        │                       │
┌─────────────────────────┐       EOD Spike                         │                       │
│ Autoscaler will check   │◄──────  Load    ◄───  [people icon]     │                       │
│ for available Node      │                                         │                       │
│ resources               │                                         │                       │
└─────────────────────────┘                                         │                       │
            │                                                        │                       │
            ▼                                                        │                       │
        ◇─────────◇                                             ◇─────────◇                  │
       /           \                                           /           \                 ▼
      / Available   \              NO                         / Cheapest    \    NO  ┌──────────────────┐
     /  Node         \ ─────────────────────────────────────/  resources    \──────▶│ Scale up         │
  YES\  resources    /                                       \  (PVM)        /       │ default-pool     │
   │  \  is enough? /                                         \  available? /        └──────────────────┘
   │   \           /                                           \           /
   │    ◇─────────◇                                             ◇─────────◇
   │
   └──────────────────────────────────────────────────────────────┘
```
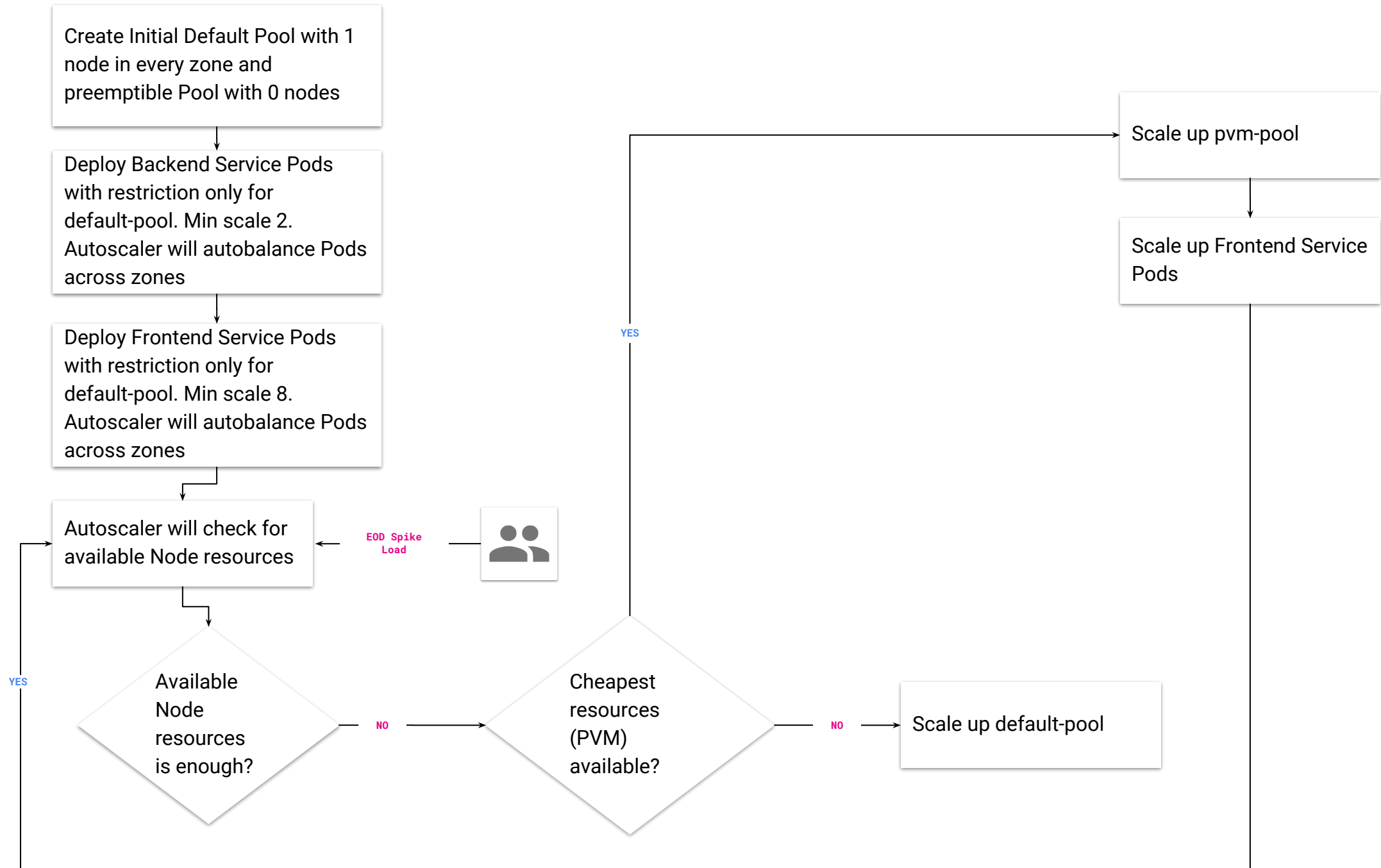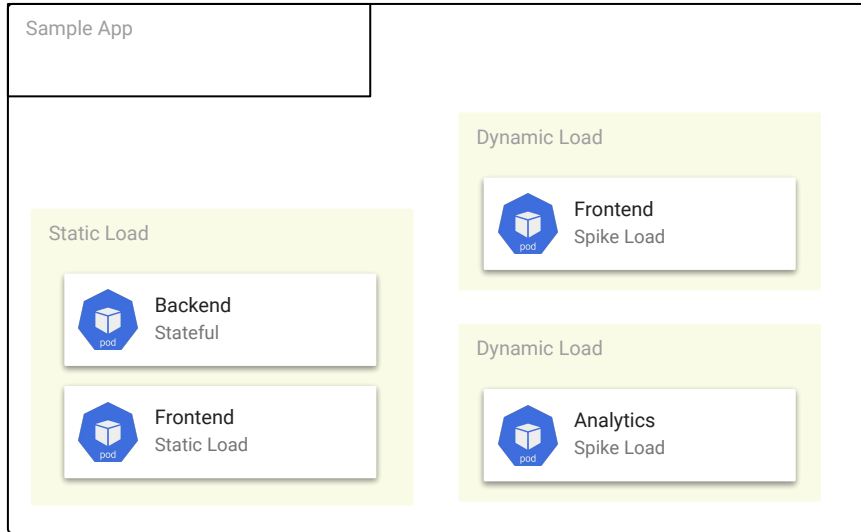
# EXAMPLE 2

Let's add Analytics

softserve

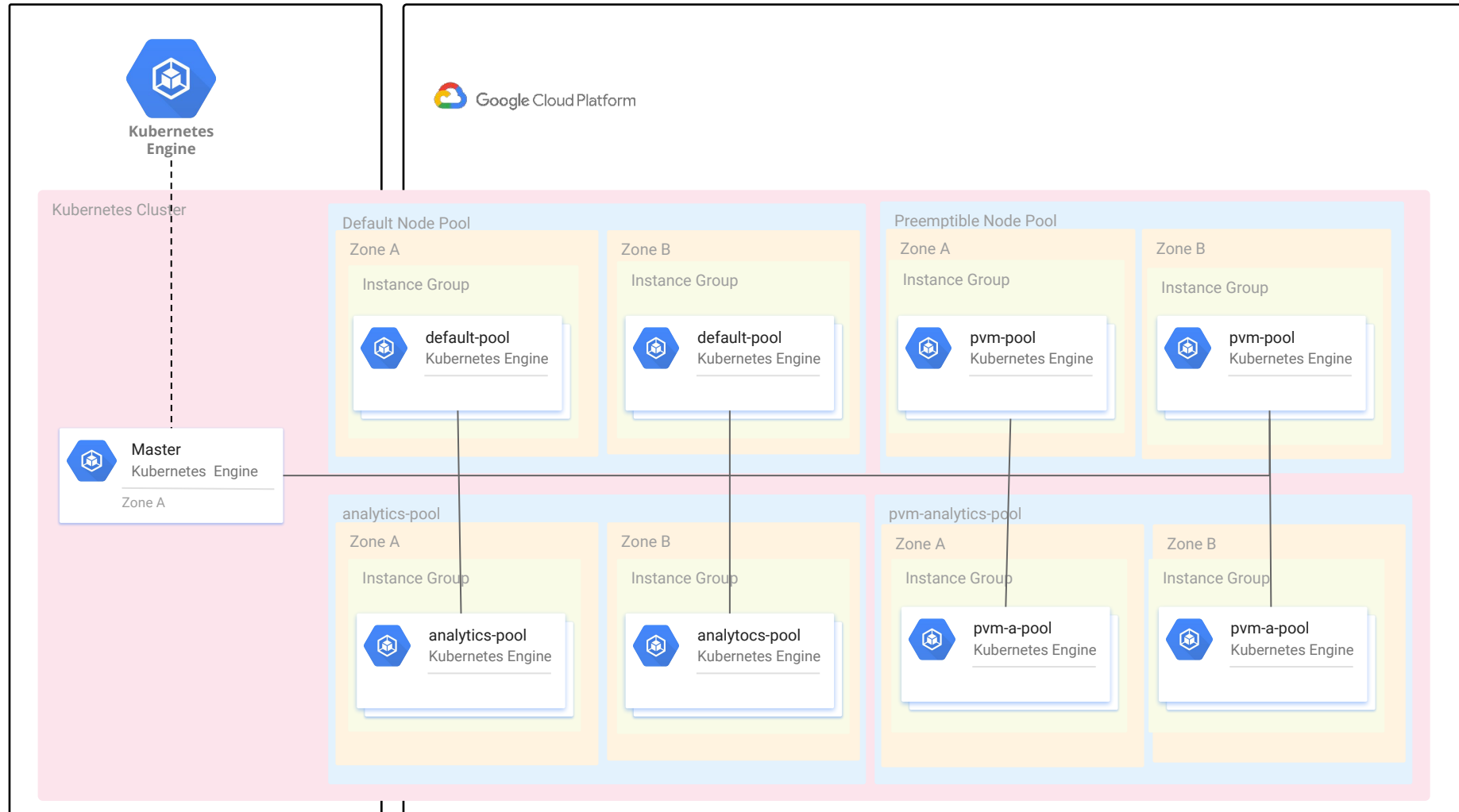# SAMPLE APP. ANALYTICS ADDED



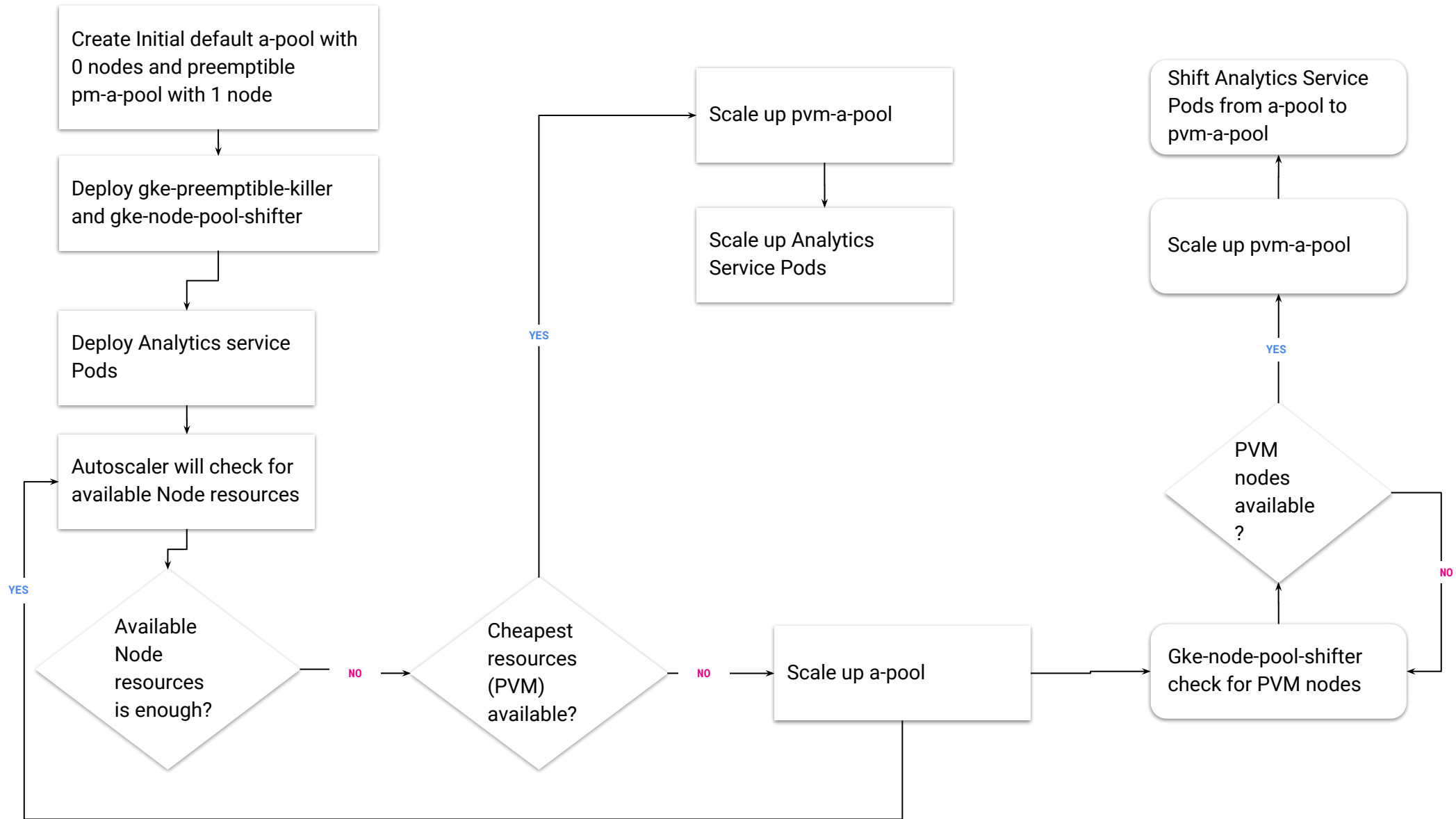Old-style Java monolithic enterprise app splitted for two major parts
- Backend Java old-style stateful service
- New and shiny Node.js Frontend service
  - Predictable static load
  - High load spikes at EOD/EOW

- **Daily/Weekly Analytics**
  - To decrease cost let's run heavy analytics only on PVM

softserve

# GKE — MULTI-ZONE DEFAULT NODE POOL / MULTI-ZONE PREEMPTIBLE NODE POOL

**Kubernetes Engine**

Google Cloud Platform

Kubernetes Cluster

**Default Node Pool**

Zone A

Instance Group

default-pool
Kubernetes Engine

Zone B

Instance Group

default-pool
Kubernetes Engine

**Preemptible Node Pool**

Zone A

Instance Group

pvm-pool
Kubernetes Engine

Zone B

Instance Group

pvm-pool
Kubernetes Engine

Master
Kubernetes  Engine

Zone A

**analytics-pool**

Zone A

Instance Group

analytics-pool
Kubernetes Engine

Zone B

Instance Group

analytocs-pool
Kubernetes Engine

**pvm-analytics-pool**

Zone A

Instance Group

pvm-a-pool
Kubernetes Engine

Zone B

Instance Group

pvm-a-pool
Kubernetes Engine

Create Initial default a-pool with 0 nodes and preemptible pm-a-pool with 1 node

Deploy gke-preemptible-killer and gke-node-pool-shifter

Deploy Analytics service Pods

Autoscaler will check for available Node resources

Available Node resources is enough?

**NO** →

Cheapest resources (PVM) available?

**YES** ↑ Scale up pvm-a-pool

Scale up Analytics Service Pods

**NO** →

Scale up a-pool →

Gke-node-pool-shifter check for PVM nodes

PVM nodes available ?

**NO** →

**YES** ↑

Scale up pvm-a-pool

Shift Analytics Service Pods from a-pool to pvm-a-pool

**YES** (from Available Node resources is enough? loops back to Autoscaler will check for available Node resources)

# SCHEDULE PODS ON PVM IF AVAILABLE

Modify your Pod or Deployment spec using

```
...
spec:
  affinity:
    nodeAffinity:
      preferredDuringSchedulingIgnoredDuringExecution:
      - preference:
          matchExpressions:
          - key: cloud.google.com/gke-preemptible
            operator: Exists
        weight: 100
...
```

# SCHEDULE PODS ON NON-PVM NODES

Modify your Pod or Deployment spec using

```
...
spec:
  affinity:
    nodeAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
        nodeSelectorTerms:
        - matchExpressions:
          - key: cloud.google.com/gke-preemptible
            operator: DoesNotExist
...
```

# CAPACITY VS COST(PVM)



**BEST CASE:**

**34%** of spend,
**100%** of the time

**WORST CASE:**

**120%** of spend,
**100%** of the time
**2x** capacity

Capacity

Cost

Zone

regular instances

PVM instances

softserve

# SHOW STACKDRIVER

# DEMO STRESS

# COST SAVING WITH PVM AS % OF CLUSTER

# COMMITTED USE DISCOUNTS

# SHOW STACKDRIVER FOR STRESS

softserve

Gke-preemptible-killer check if pvm-a-pool node run for 12h or more hours

pvm-a-pool node run for 12h or more hours?

NO

YES

Kill pvm-a-pool node that run for more than 12h

Kubernetes controller that ensures **deletion** of preemptible nodes in a GKE cluster **is spread out to avoid the risk** of all getting deleted at the same time after 24 hours

soft**serve**

# KUBERNETES CONTROLLERS FOR PREEMPTIBLE NODE POOL

- We recommend to randomly kill PVM in Preemptible Node Pool to avoid expire all nodes same time
  estafette-gke-preemptible-killer

- Another great tool will help to constantly monitor Node Pools and move nodes to preemptible PVM
  estafette-gke-node-pool-shifter

softserve

# APPENDIX

DEMO repository https://github.com/stonevil/gke-meetup-demo-project

softserve

# KUBERNETES AND EPHEMERAL COMPUTE

**Abstract**

In cloud computing, elasticity is defined as "the degree to which a system is able to adapt to workload changes by provisioning and de-provisioning resources in an autonomic manner, such that at each point in time the available resources match the current demand as closely as possible".

For those familiar with Kubernetes, this may seem like a solved problem. Not exactly. What about the underlying cluster resources? Node Autoscaling is a great feature of GKE but how can we take advantage of this in interesting, cost effective way?

In this interactive session, you will walk through few cases how-to cut Google Cloud Kubernetes Engine cluster cost with preemptible VM's, Stackdriver and Cluster Autoscaler.

softserve

# BIOS

Ryan Richard — Ryan is a Customer Engineer at Google working with enterprise customers and focused on GCP. He has background building and running services in Kubernetes and originally added the Rackspace deployment code to the repo in 2014.

Myroslav Rys — Myroslav is a Solution Architect at SoftServe Inc. More than 8 years experience in large scale enterprise solutions including SaaS / Clouds solutions. Experience building products and solutions with Kubernetes from 2015.

softserve

# UTILIZATION VS CAPACITY VS COST



Zone: A, B, C, D

Capacity: 6, 4, 2

**100%** of spend,
**100%** of the time

softserve

# SCALING (MYROSLAV)

## Scaling

- Manually scaling, autoscaling (Cluster Autoscaler)
- Better utilization but you're paying full price for these resources
- **what if we look at this graph, it seems that our capacity spikes only for a few hours a day (Stackdriver)**
- Or a known batch job that will increase usage for a known amount of time, want to pay the least for it.
- Is there a way to handle this capacity temporarily without paying full price?

**Cluster autoscaler considers the relative cost of each instance type in the node pool and attempts to expand the least expensive possible node pool. [1]**

**Cluster Autoscaler only takes requested resources into account for autoscaling. It does not take current utilization into account.**

**[1]  https://cloud.google.com/kubernetes-engine/docs/concepts/cluster-autoscaler**

softserve

# UTILIZATION VS CAPACITY VS COST



**BEST CASE:**

**50%** of spend,
**100%** of the time

**WORST CASE:**

**100%** of spend,
**100%** of the time

**Capacity**

6 —

4 —

2 —

Zone  **A**  **B**  **C**  **D**

100%

0%

Utilization

**BEST CASE:**

**36.5%** of spend,
**100%** of the time

**WORST CASE:**

**120%** of spend,
**100%** of the time
**2x** capacity

softserve