# K8s Nemesis

Requirements and analysis model

# Product description

**Product Description:** A cloud-native service designed to enhance application scalability and performance in Kubernetes (K8s) environments, specifically targeting machine learning (ML) applications with unique scaling and resource requirements.

Team: Sergey Lokhmatikov, Andrey Tamplon, Roman Kuzmenko, Vasiliy Lopatkin

Repo: https://github.com/Lokhmat/k8s_nemesis

Report: https://github.com/Lokhmat/k8s_nemesis/blob/main/final_slides.pdf
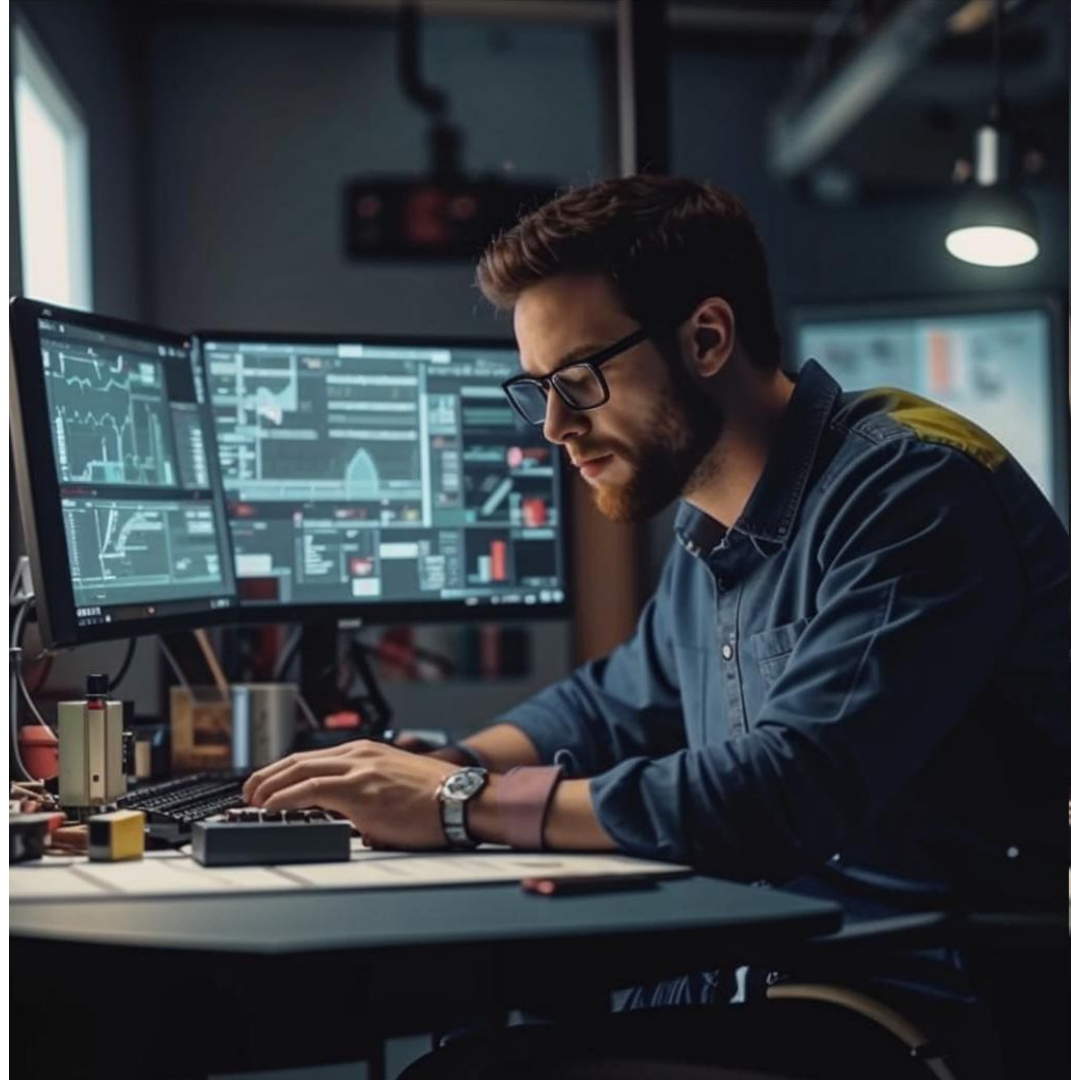
# MLOps engineer

**Background:** Focuses on deploying, monitoring, and maintaining ML models in production. He's skilled in cloud platforms, CI/CD pipelines, and container tools like Docker and Kubernetes.

**Needs:** Tools that help with continuous integration and deployment of ML models make it easier to quickly update and improve models without interrupting current services.

**Challenges:** Making sure models are reliable, keeping consistency across environments, and handling complex ML pipelines and dependencies.

**Goal:** Automate deployment and monitoring for optimal model performance.

## DevOps engineer

**Background:** Engineer with a deep understanding of a distributed system. High level expertise of working with Kubernetes, continuous integration and deployment. The one who is responsible for cloud architecture in company

**Needs:** Elegant and customizable tool for managing k8_s cluster. Tool that can reduce cost of maintenance of a large scale architecture while ensuring high SLA and good performance in a complicated systems.

**Challenges:** Balancing infrastructure maintenance cost vs service uptime. Managing distributed cluster, and event driven scaling based on customizable set of rules and metrics. Ensuring security and stability of a system.

**Goal:** Deployment and maintenance of a reliable distributed system. Increasing usage efficiency of a current organisation's cloud architecture.

# Story Map

## Storymap

| Persona goal | Equal load distribution across cluster | Save money and cluster resources | Integrate with k8s cluster |
|---|---|---|---|

**Persona task**

- Ensure minimal latency
- Guarantee high level availability of cluster
- Count of pods depends on current load
- Fair requests proxying

- Stop using resources if application is not processing requests
- Allow to define the most suitable settings for k8s cluster
- Activate application if it needs to handle requests but not active

- Integration with k8s cluster

miro

# Story Map

**Product feature**

| Automaticaly scale cluster up | Request proxying with different policies | Automatically scale cluster down | Disable application if there are no requests | Activate application if request comes in and system is disabled | Define scalling rules based on cluster load metrics | Define proxying rules based on cluster load metrics | Fetch k8s cluster metrics | Monitor current pods liveliness state |

**Stories**

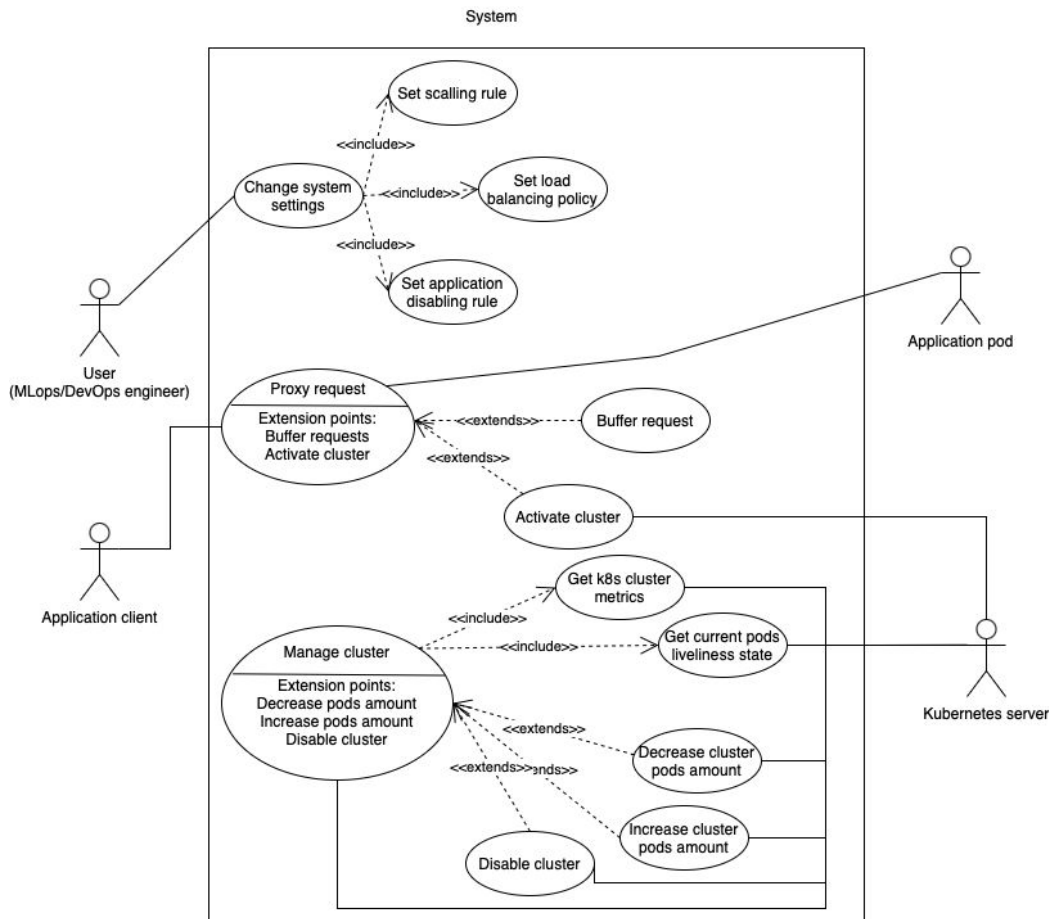| Increase cluster pods amount | Proxy requests | Decrease cluster pods amount | Disable cluster | Activate cluster | Set activation/de activation rule | Set load balancing policy | Get CPU/RAM/DIS K metrics from k8s | Get current pods liveliness state |

| | Set load balancing policy | | | Buffer request | Set service scalling rule | | | |

miro

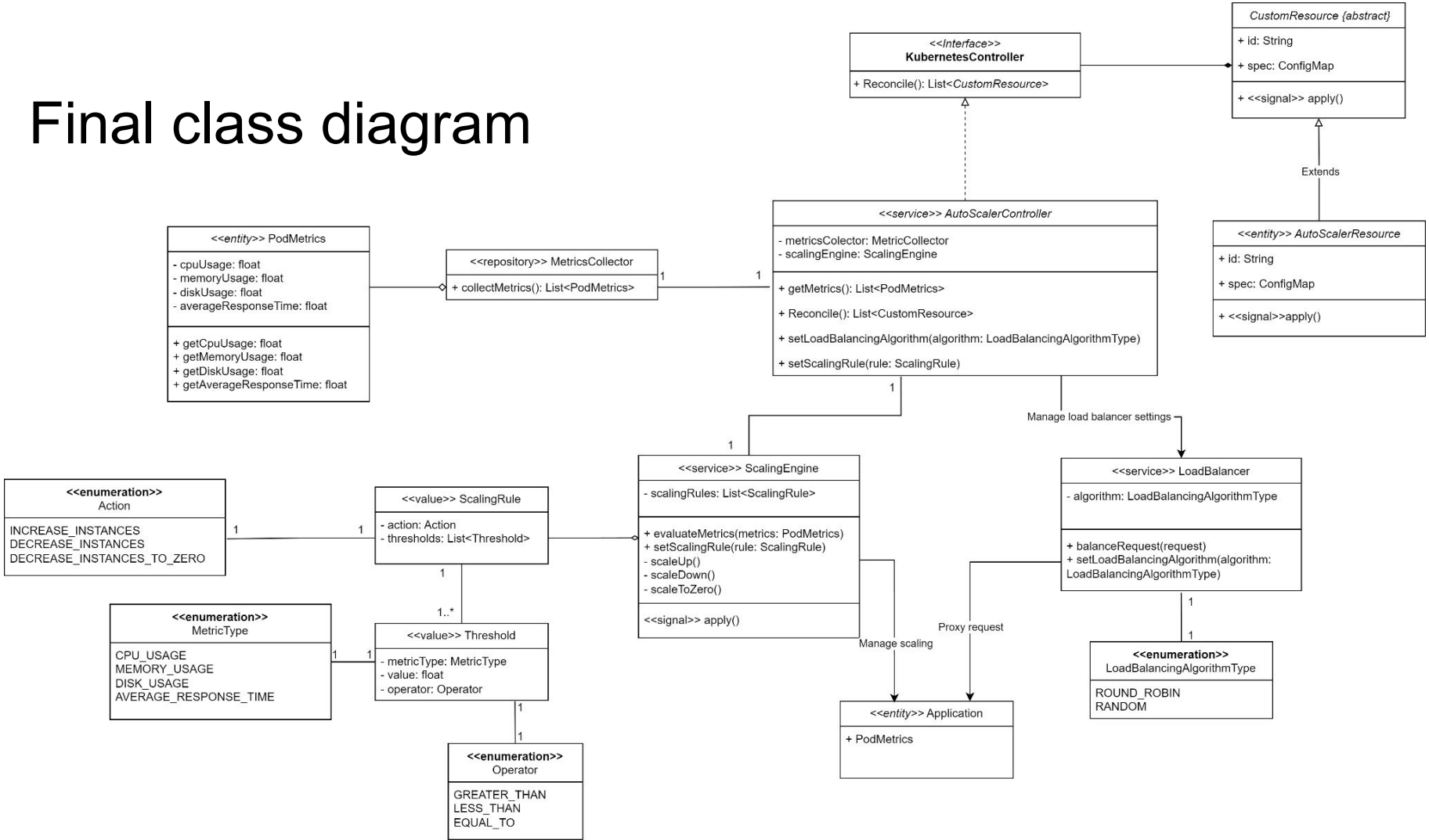# Use case diagram

Textual use case scenarios:

https://github.com/Lokhmat/
k8s_nemesis/blob/main/fina
l_task_materials/textual_us
e_cases.md

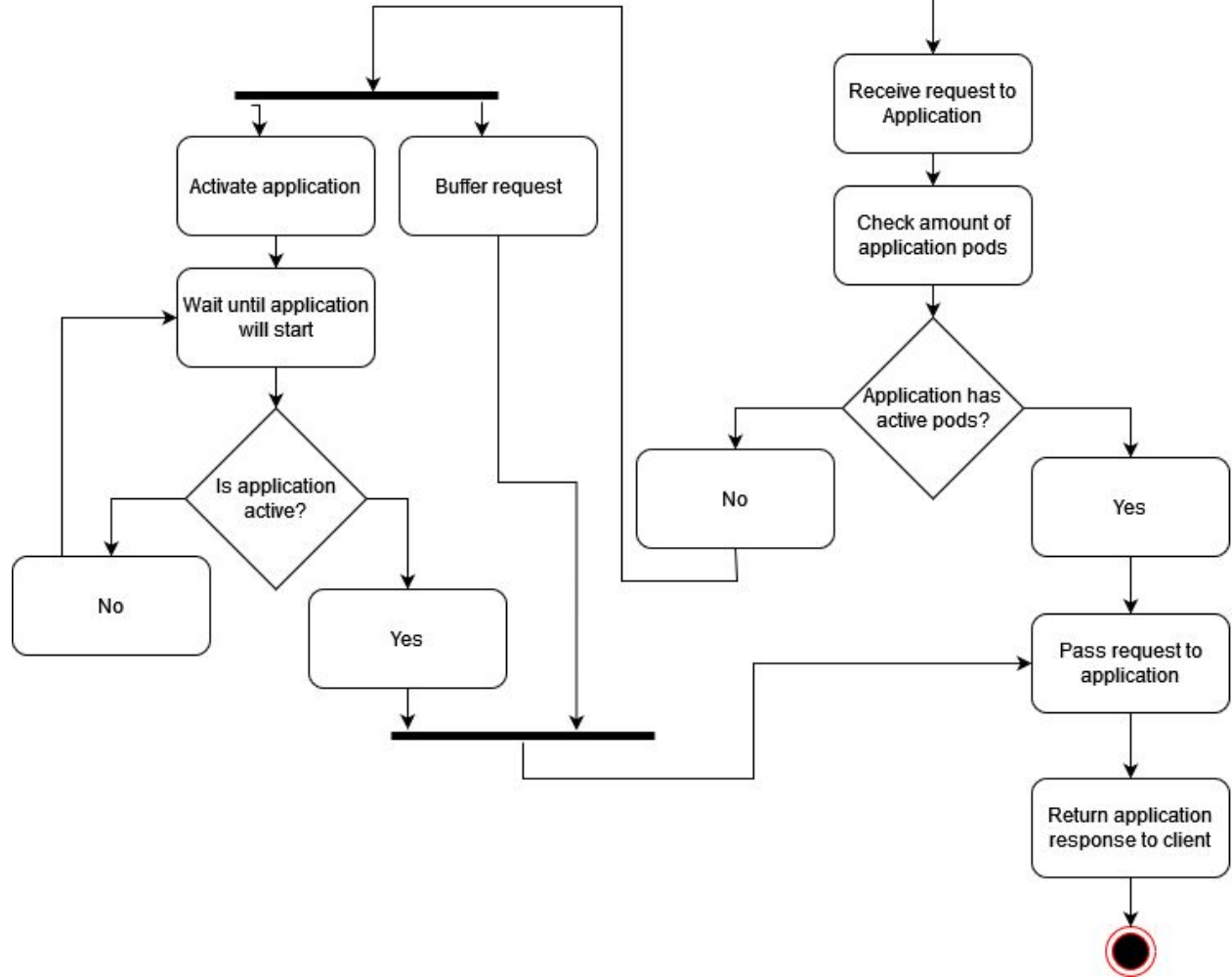# Interaction analysis

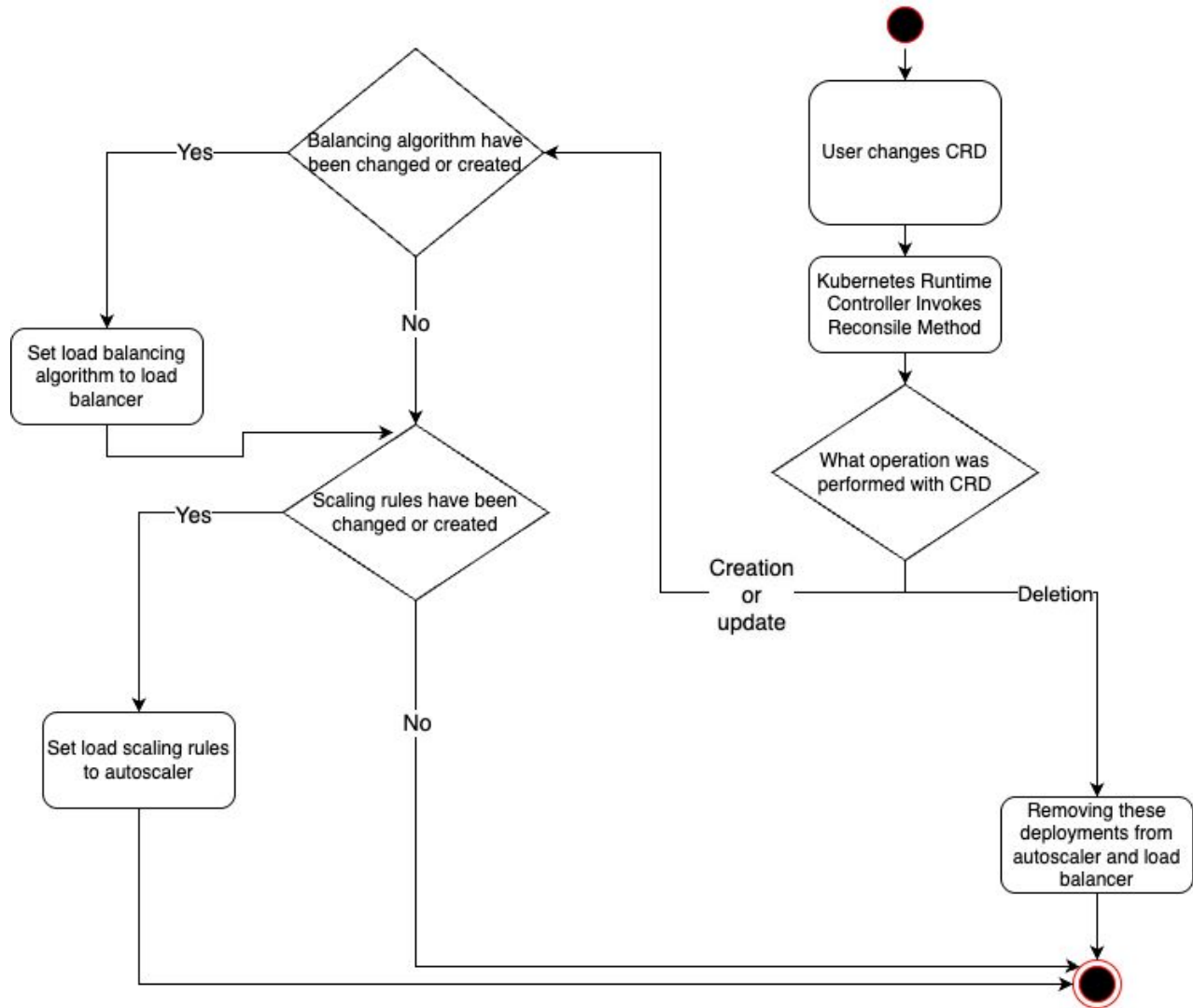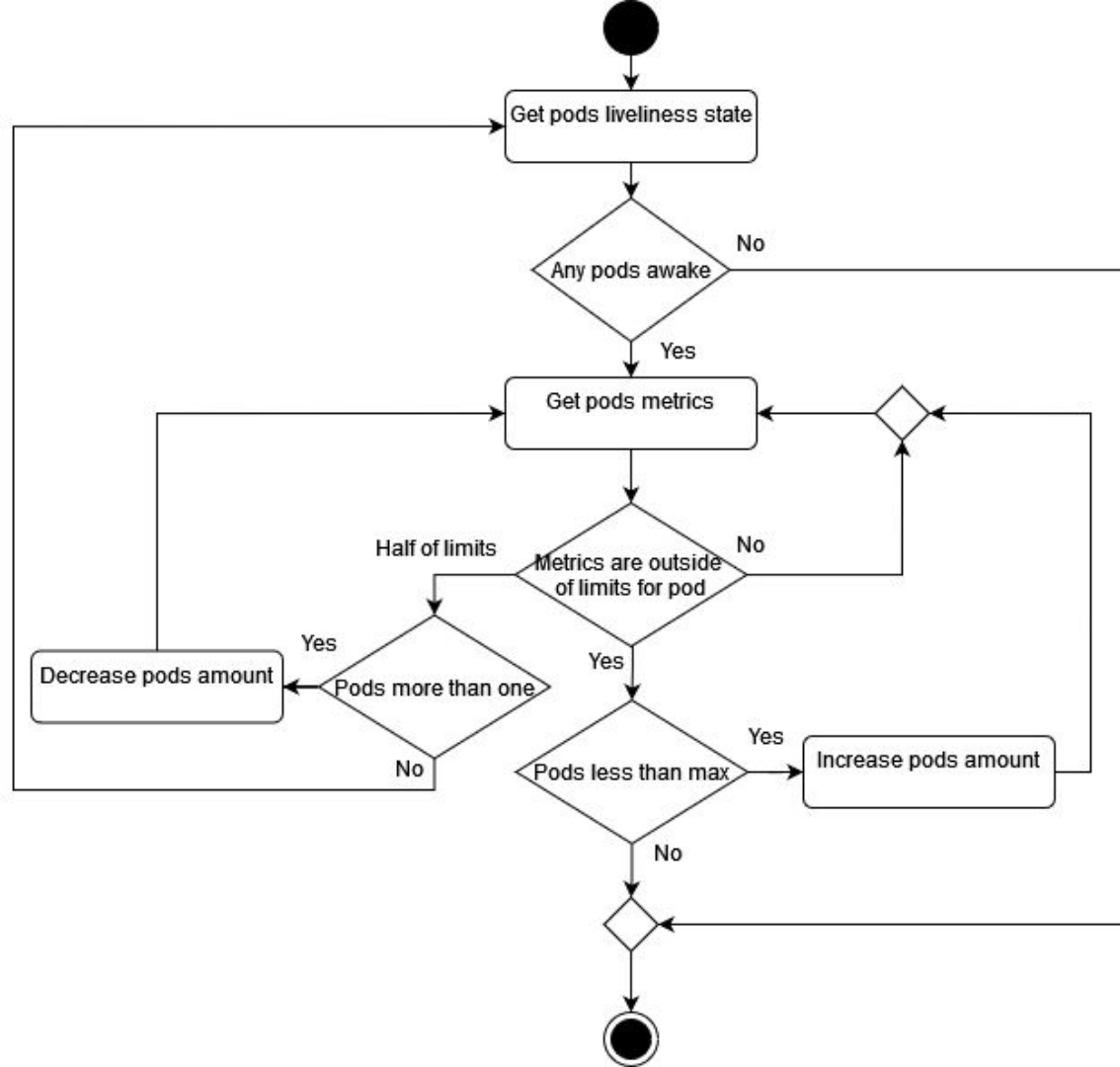| Use Case | Cooperation Type | Used Roles | Candidate Classes |
|---|---|---|---|
| Set scaling rule | Edit | User, ScalingEngine | ScalingEngine, ScalingRule, AutoScalerController, Treshhold |
| Set load balancing policy | Edit | User, LoadBalancer | LoadBalancer, AutoScalerController |
| Set application disabling rule | Edit | User, ScalingEngine | ScalingEngine, ScalingRule, AutoScalerController, Treshhold |
| Proxy request | Proxying | Application client, LoadBalancer, Application pod | LoadBalancer |
| Buffer request | Proxying | Application client, LoadBalancer | LoadBalancer |
| Activate cluster | Manage cluster | Application client, LoadBalancer, ScalingEngine | LoadBalancer |
| Get k8s cluster metrics | Validate | Kubernetes server, AutoScalerController | PodMetrics, MetricsCollector, AutoScalerController |
| Get current pods liveliness state | Validate | Kubernetes server, AutoScalerController, Application pod | PodMetrics, MetricsCollector, AutoScalerController |
| Decrease cluster pods amount | Manage cluster | Kubernetes server, AutoScalerController, Application pod | ScalingEngine, AutoScalerResource, AutoScalerController |
| Increase cluster pods amount | Manage cluster | Kubernetes server, AutoScalerController, Application pod | ScalingEngine, AutoScalerResource, AutoScalerController |
| Disable cluster | Manage cluster | Kubernetes server, AutoScalerController, Application pod | ScalingEngine, AutoScalerResource, AutoScalerController |

# Final class diagram



**CustomResource {abstract}**
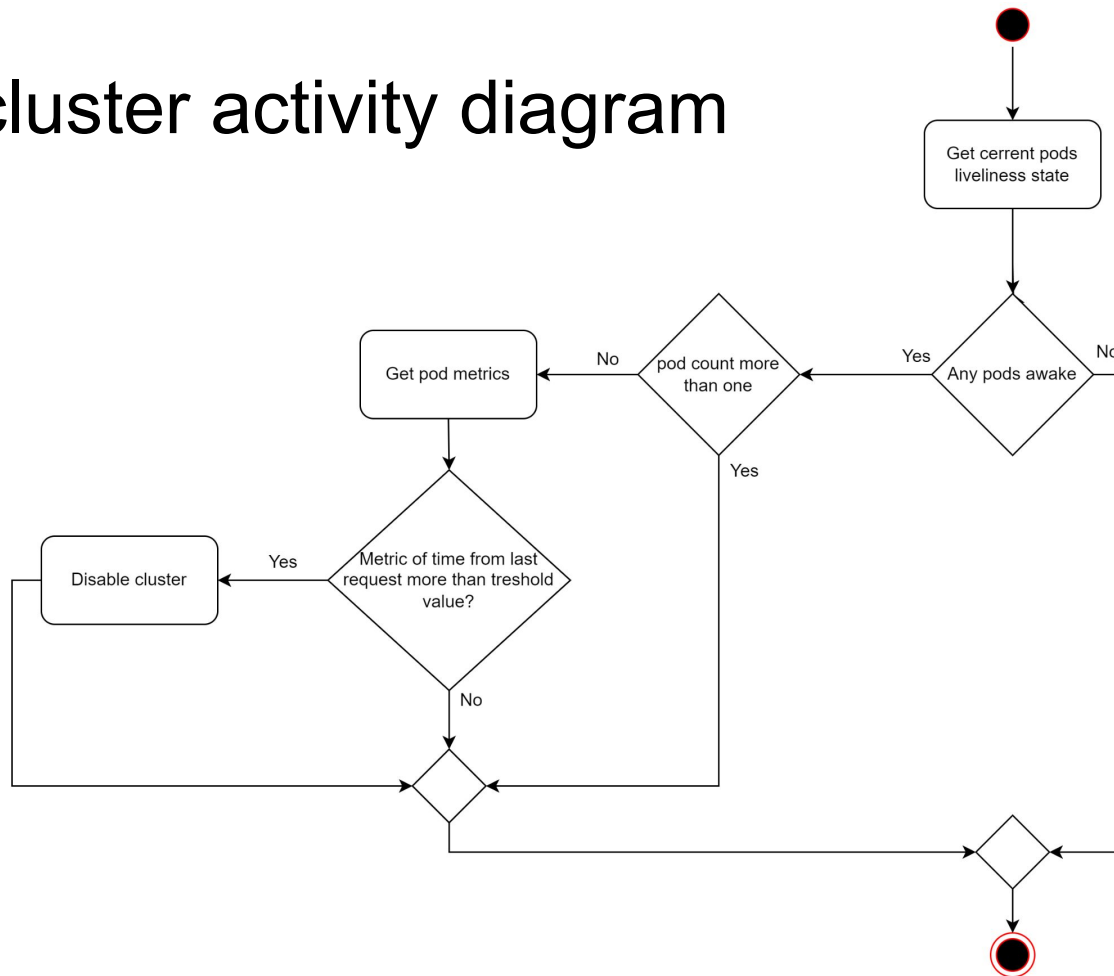+ id: String
+ spec: ConfigMap
+ <<signal>> apply()

**<<Interface>>**
**KubernetesController**
+ Reconcile(): List<CustomResource>

Extends

**<<service>> AutoScalerController**
- metricsColector: MetricCollector
- scalingEngine: ScalingEngine
+ getMetrics(): List<PodMetrics>
+ Reconcile(): List<CustomResource>
+ setLoadBalancingAlgorithm(algorithm: LoadBalancingAlgorithmType)
+ setScalingRule(rule: ScalingRule)

**<<entity>> AutoScalerResource**
+ id: String
+ spec: ConfigMap
+ <<signal>>apply()

**<<entity>> PodMetrics**
- cpuUsage: float
- memoryUsage: float
- diskUsage: float
- averageResponseTime: float
+ getCpuUsage: float
+ getMemoryUsage: float
+ getDiskUsage: float
+ getAverageResponseTime: float

**<<repository>> MetricsCollector**
+ collectMetrics(): List<PodMetrics>

Manage load balancer settings

**<<enumeration>>**
Action
INCREASE_INSTANCES
DECREASE_INSTANCES
DECREASE_INSTANCES_TO_ZERO

**<<value>> ScalingRule**
- action: Action
- thresholds: List<Threshold>

**<<service>> ScalingEngine**
- scalingRules: List<ScalingRule>
+ evaluateMetrics(metrics: PodMetrics)
+ setScalingRule(rule: ScalingRule)
- scaleUp()
- scaleDown()
- scaleToZero()
<<signal>> apply()

**<<service>> LoadBalancer**
- algorithm: LoadBalancingAlgorithmType
+ balanceRequest(request)
+ setLoadBalancingAlgorithm(algorithm: LoadBalancingAlgorithmType)

**<<enumeration>>**
MetricType
CPU_USAGE
MEMORY_USAGE
DISK_USAGE
AVERAGE_RESPONSE_TIME

1..*

**<<value>> Threshold**
- metricType: MetricType
- value: float
- operator: Operator

Manage scaling

Proxy request

**<<entity>> Application**
+ PodMetrics

**<<enumeration>>**
LoadBalancingAlgorithmType
ROUND_ROBIN
RANDOM

**<<enumeration>>**
Operator
GREATER_THAN
LESS_THAN
EQUAL_TO

# Activity diagram

# Configuration

Decrease and Increase cluster pods amount Activity diagram

# Disable cluster activity diagram

# Repository structure

# Team and roles



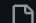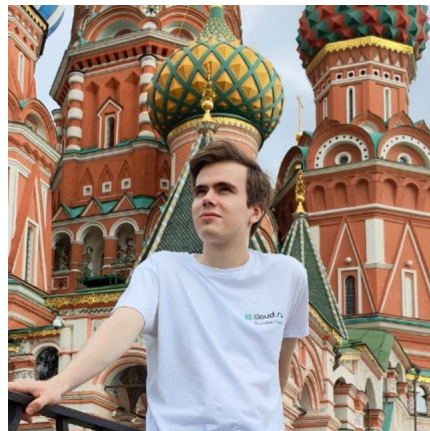Vasiliy Lopatkin
@DFCZok
Final class diagram
Interaction analysis

Sergey
Lokhmatikov
@Lohmat_Sergey
Story map
Use case diagram

Roman Kuzmenko
@definitely_not_rk
Personas
Behavioral diagram

Andrey Tamplon
@andreytamplon
Final class diagram
Interaction analysis