**KubeCon** | **CloudNativeCon**

Europe 2019

# Kubernetes Scalability Definition Evolution

Wojciech Tyczyński, Staff Software Engineer, Google

# Scalability - what does it mean?

"**Scalability** is the property of a system to handle a growing amount of work by adding resources to the system.[1]"

"In computing, scalability is a characteristic of computers, networks, algorithms, networking protocols, programs and applications. An example is a search engine, which must support increasing numbers of users, and the number of topics it indexes.[3]"

Wikipedia contributors, "Scalability," *Wikipedia, The Free Encyclopedia,* https://en.wikipedia.org/w/index.php?title=Scalability&oldid=892100604 (accessed May 13, 2019).
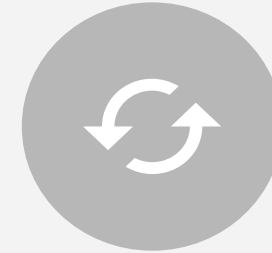
# Scalability - what does it mean?

**Scalability definition**

**Driving improvements**

**Testing infrastructure**

**Tests & guarding against regressions**

# Scalability - what does it mean?

**Scalability definition**

Driving improvements

Testing infrastructure

Tests & guarding against regressions

# Scalability - how to define it?

SLI - Service Level Indicator

SLO - Service Level Objective

# Cluster scales

# =

# **all** SLOs are satisfied

# Scalability SLOs

2015 SLOs:

**API Responsiveness**: 99% of all API calls return in less than 1s

**Pod startup latency**: 99% of pods and their containers (with pre-pulled images) start within 5s

Poor coverage

Lack of precision

# Scalability SLOs - coverage

April 2017: First attempt to improve coverage:

[Target SLIs and SLOs in Kubernetes](#)

Failed due to high scope

What about other issues?

# SLI/SLO principles

- precise and well-defined

- consistent

- user-oriented

- testable

# How to provide SLOs?

- cluster configuration

- Kubernetes extensibility

- load in the cluster

# Defining Kubernetes limits

- ## scalability dimension

- ## scalability envelope

# Secrets

# Namespaces

# Nodes

# Services

# Ingresses

# Backends/service

# Net LBs

# Pods/node

Pod Churn

Source of hypercube image: http://www.gregegan.net/APPLETS/29/29.html

# "You promise

- correctly configure cluster
- keeping load within the limits

# we promise"

- satisfied SLOs

# Refining SLIs/SLOs

2015:
- **SLO**: 99% of all API calls return in less than 1s

2017:
- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes
- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s

# Refining SLIs/SLOs

2015:
- **SLO**: 99% of all API calls return in less than 1s

2017:
- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes
- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s

**Explicit SLI/SLO split**

# Refining SLIs/SLOs

2015:

- **SLO**: 99% of all API calls return in less than 1s

2017:

**What is measured?**

- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes
- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s

# Refining SLIs/SLOs

2015:

- **SLO**: 99% of all API calls return in less than 1s

**How it is grouped?**

2017:

- **SLI**: Latency of mutating API calls for single objects **for every** **(resource, verb) pair,** measured as 99th percentile over last 5 minutes
- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s

# Refining SLIs/SLOs

2015:
- **SLO**: 99% of all API calls return in less than 1s

2017:
- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes
- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s

**How it is aggregated?**

# Refining SLIs/SLOs

2015:
- **SLO**: 99% of all API calls return in less than 1s

2017:
- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes

- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s

**What has guarantees?**

# Refining SLIs/SLOs

2015:
- **SLO**: 99% of all API calls return in less than 1s

2017:

- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes

- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s

**What is excluded?**

# Refining SLIs/SLOs

2015:

- **SLO**: 99% of all API calls return in less than 1s

2017:

- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes
- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions,

**What is guaranteed?**

99th percentile per cluster-day <= 1s

# Refining SLIs/SLOs

2015:

- **SLO**: 99% of all API calls return in less than 1s

**Still missing bits?**

2017:

- **SLI**: Latency of mutating API calls for single objects for every (resource, verb) pair, measured as 99th percentile over last 5 minutes
- **SLO**: In default Kubernetes installation, for every (resource, verb) pair, excluding virtual and aggregated resources and Custom Resource Definitions, 99th percentile per cluster-day <= 1s
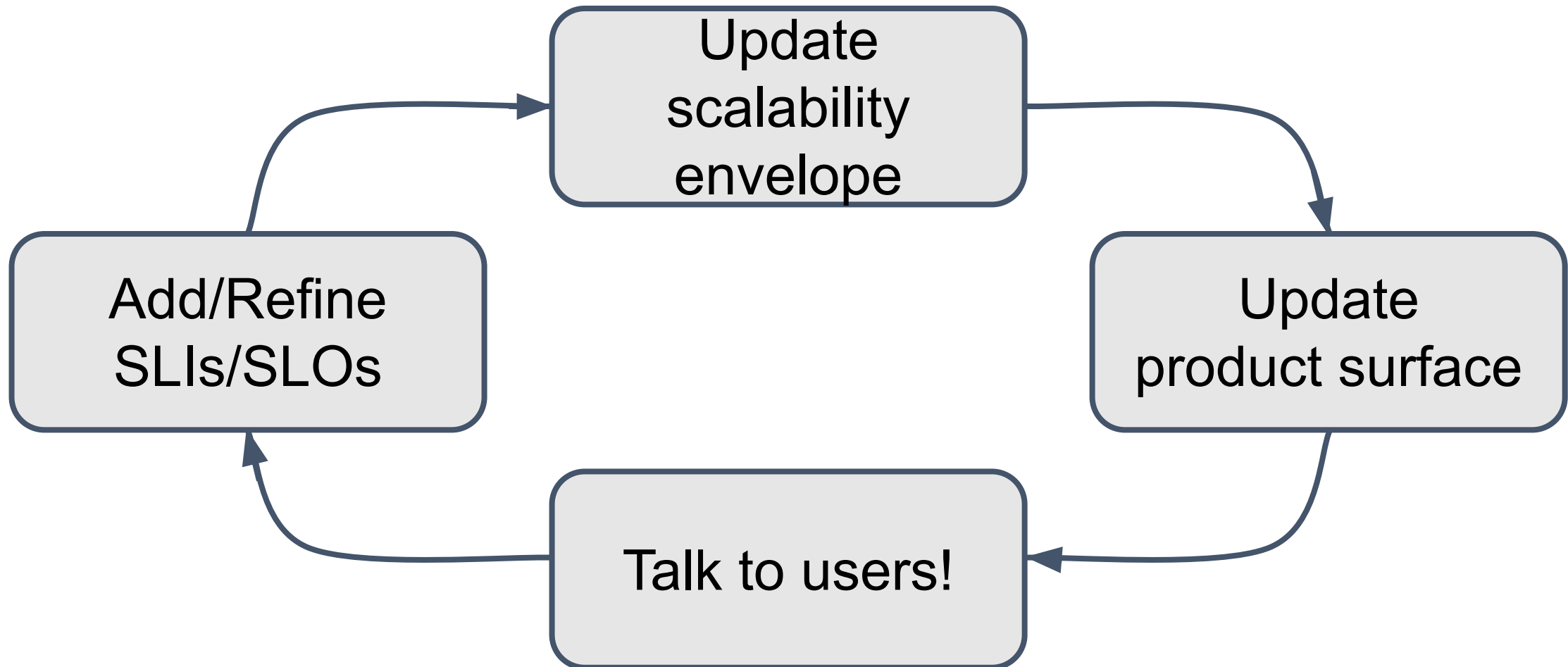
# Defining Scalability

# Current state

- 3 official SLIs/SLOs

- 5 more WIP SLIs/SLOs

- a lot of work to do :)
  - e.g. around apps concepts

# Join SIG Scalability

# BACKUP SLIDES