

# We wrote a Kubernetes operator

Here is why you should too

```
[arthurb@padok.fr ~]$ whoami
```

Site Reliability Engineer

Certified Kubernetes Administrator



# What is a Kubernetes operator?

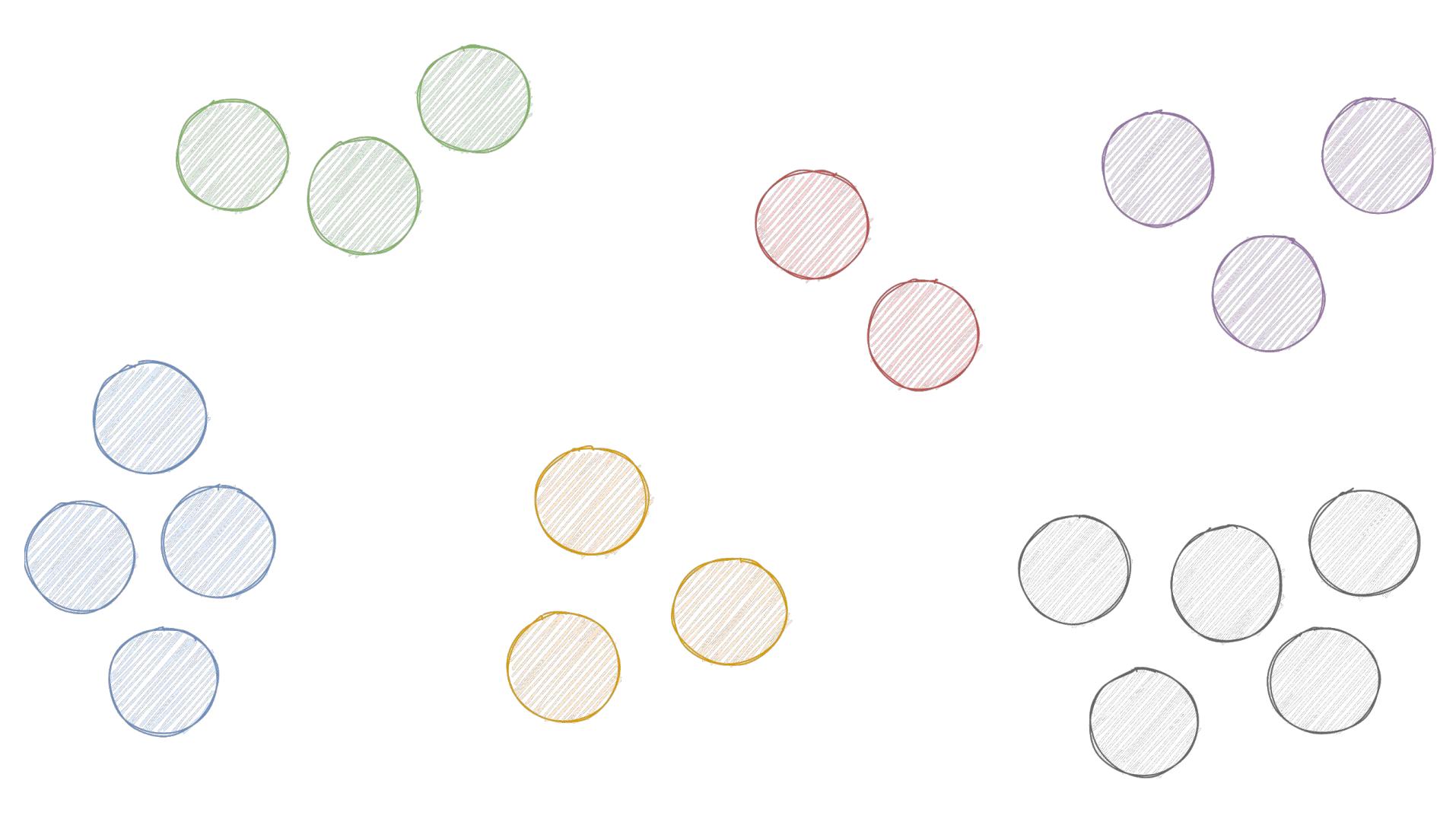
*An Operator is an application-specific controller that extends the Kubernetes API to create, configure, and manage instances of complex stateful applications on behalf of a Kubernetes user.*

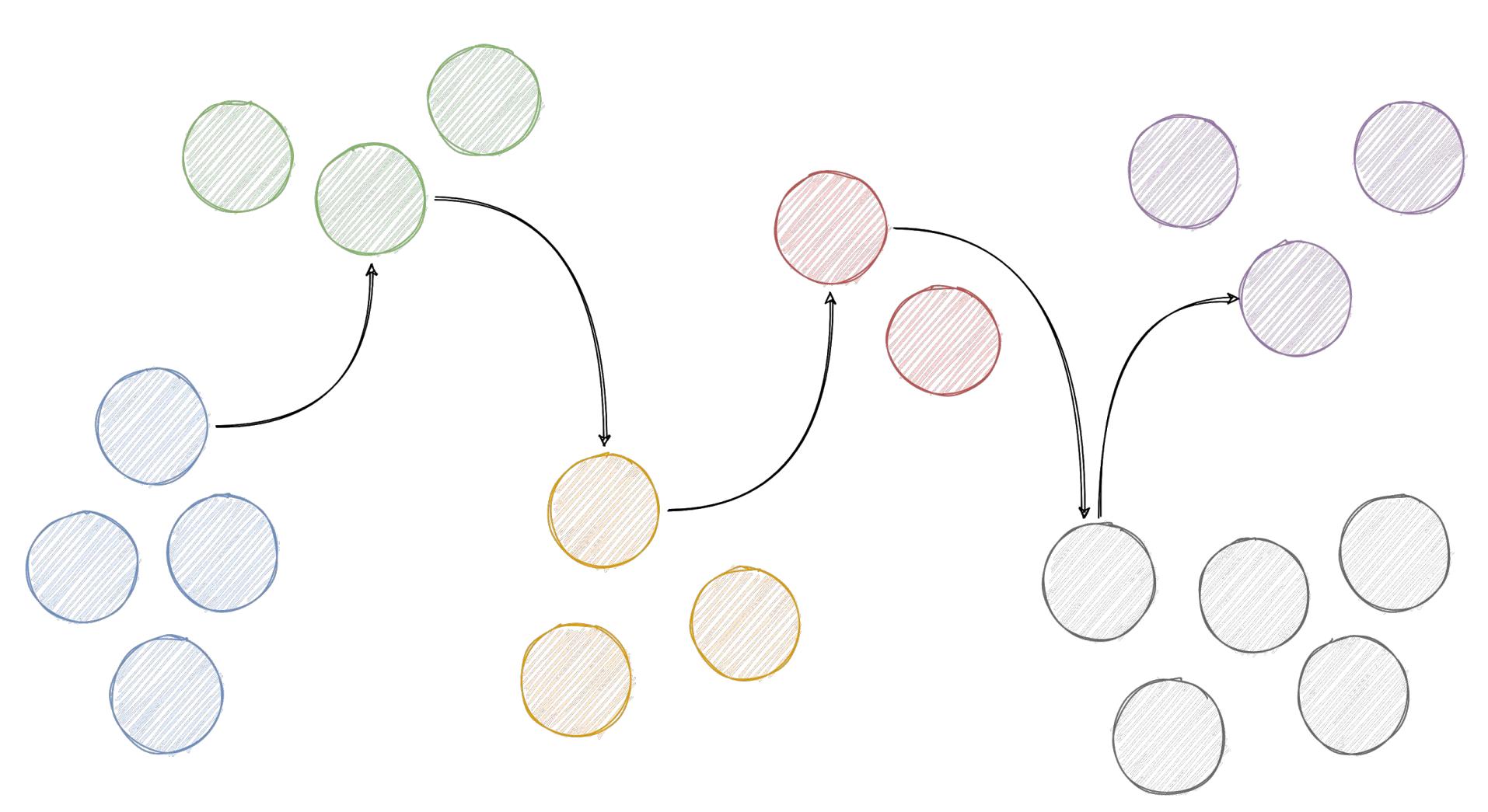
<https://coreos.com/blog/introducing-operators.html>

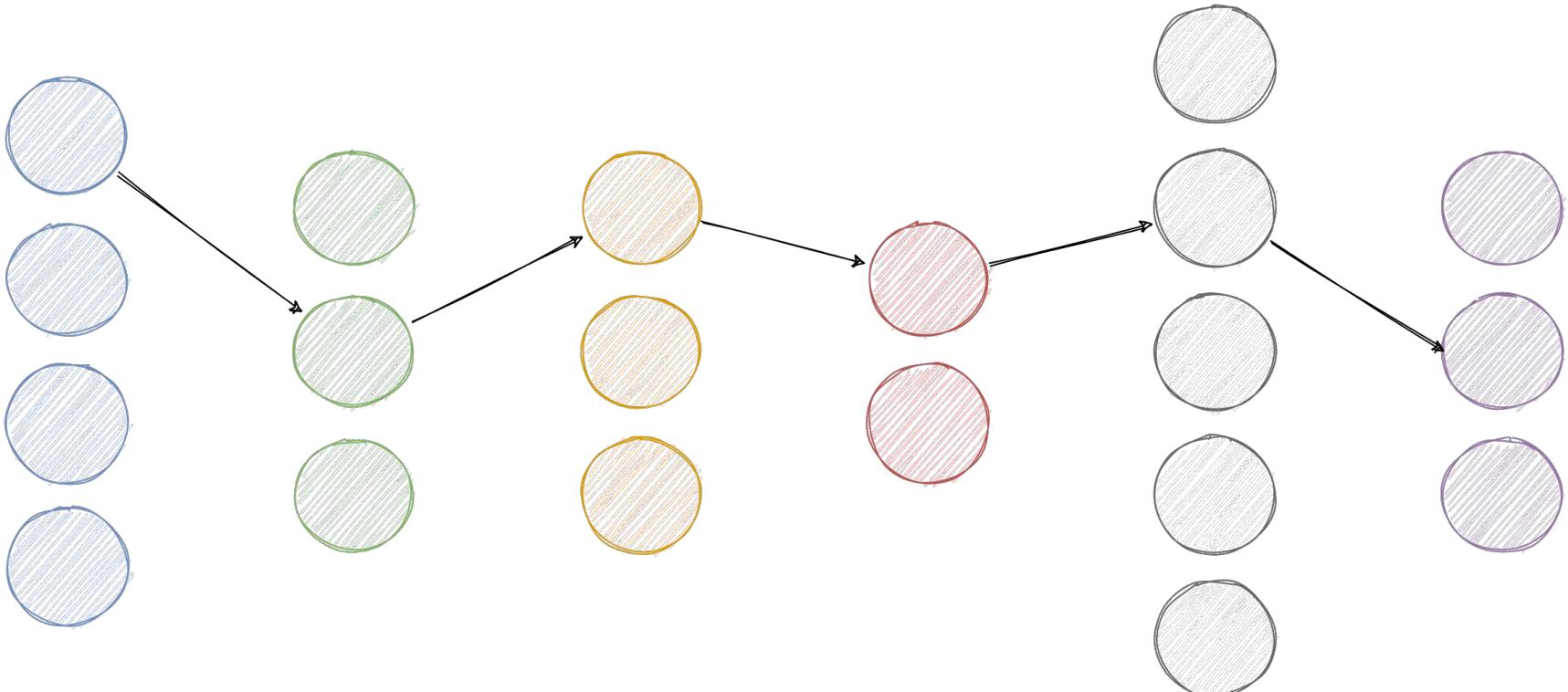
An operator consists of  
multiple controllers

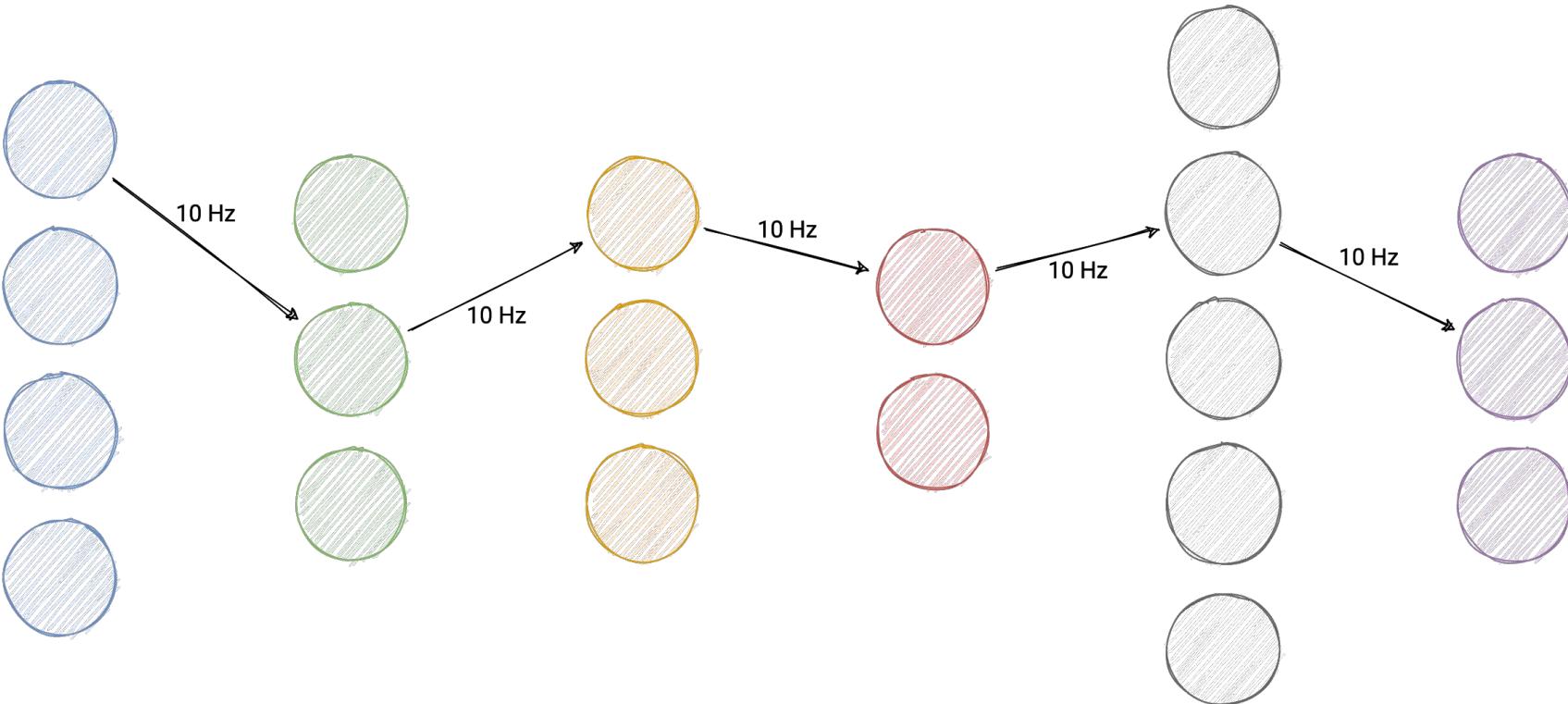
A Kubernetes controller is a finite state machine  
with a reconciliation loop that converges  
on desired state step by step

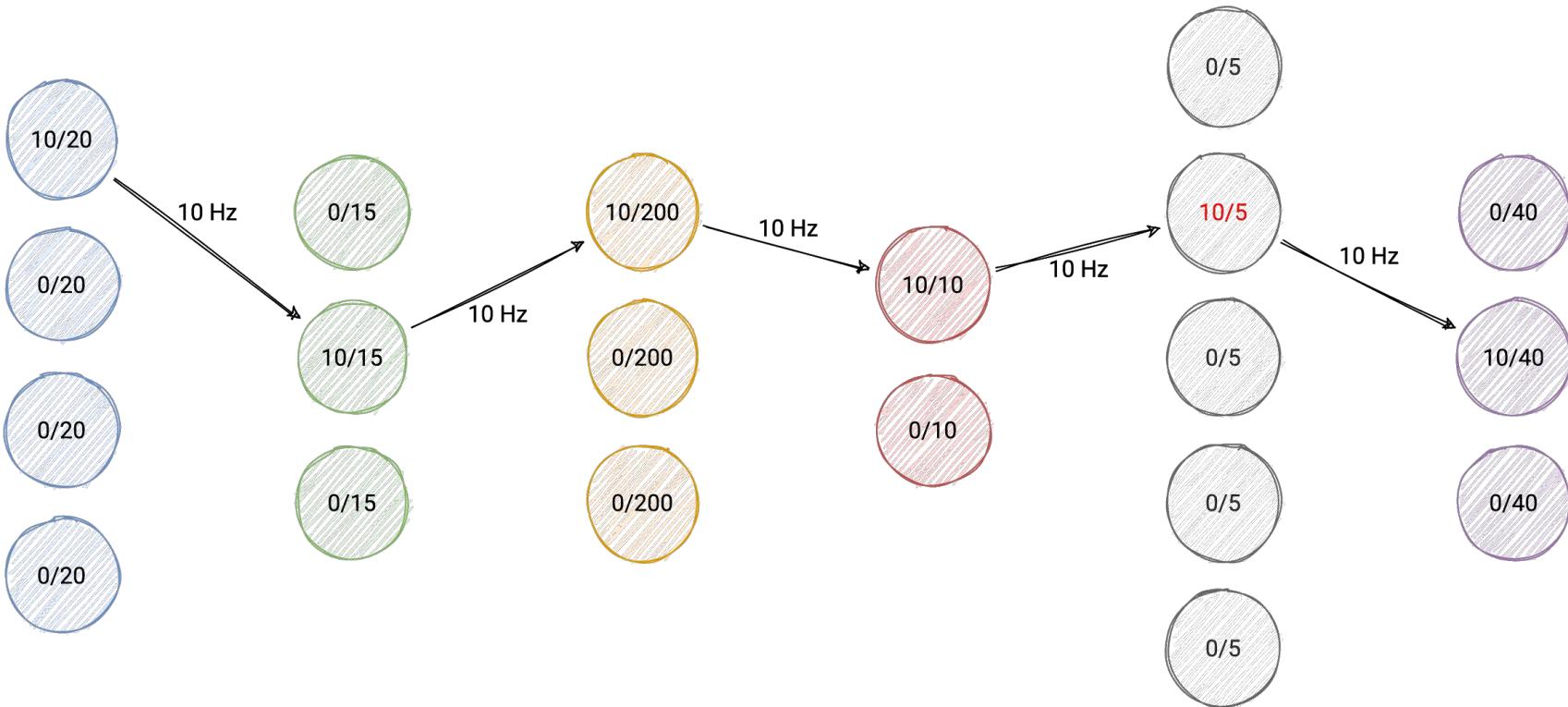
- **Why we built it**
- What we built it with
- How it works
- Lessons learned
- Overall experience

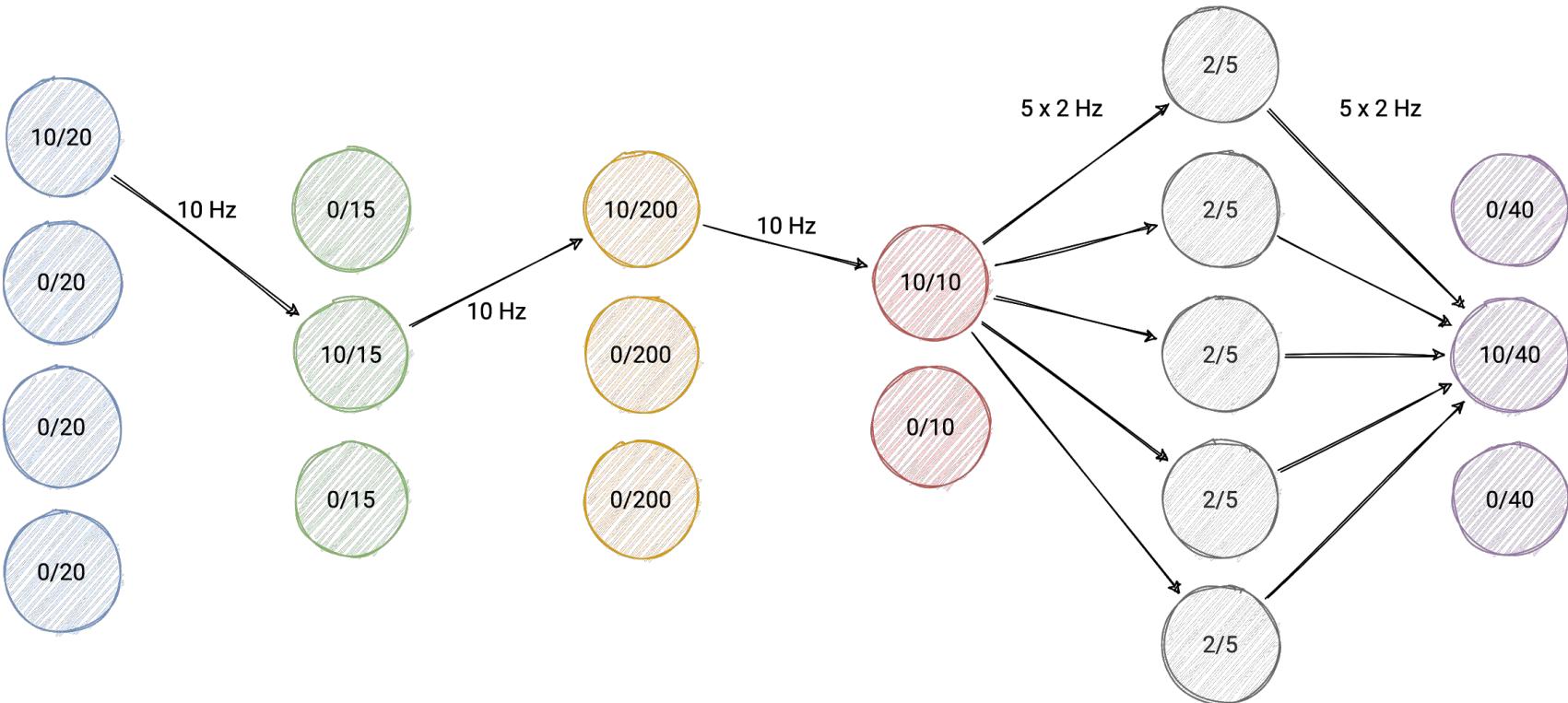


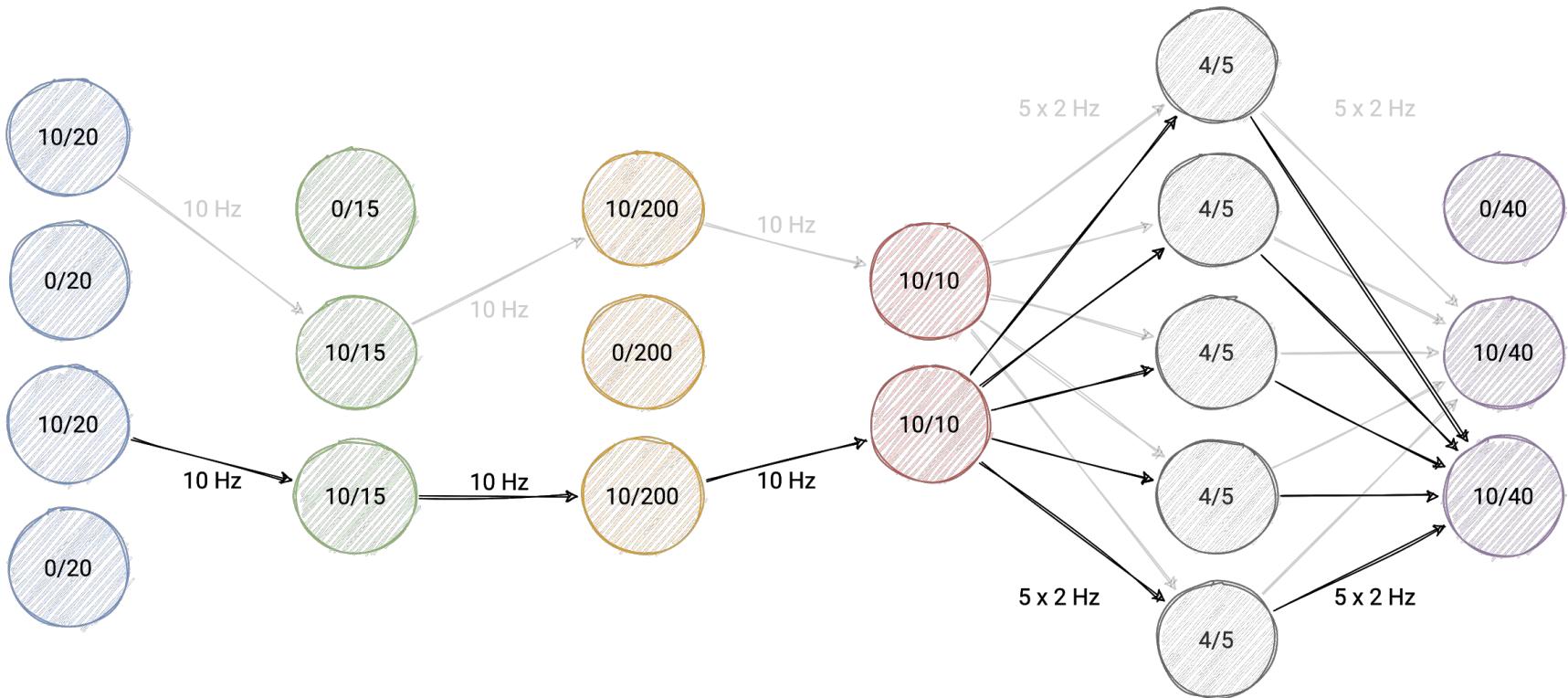


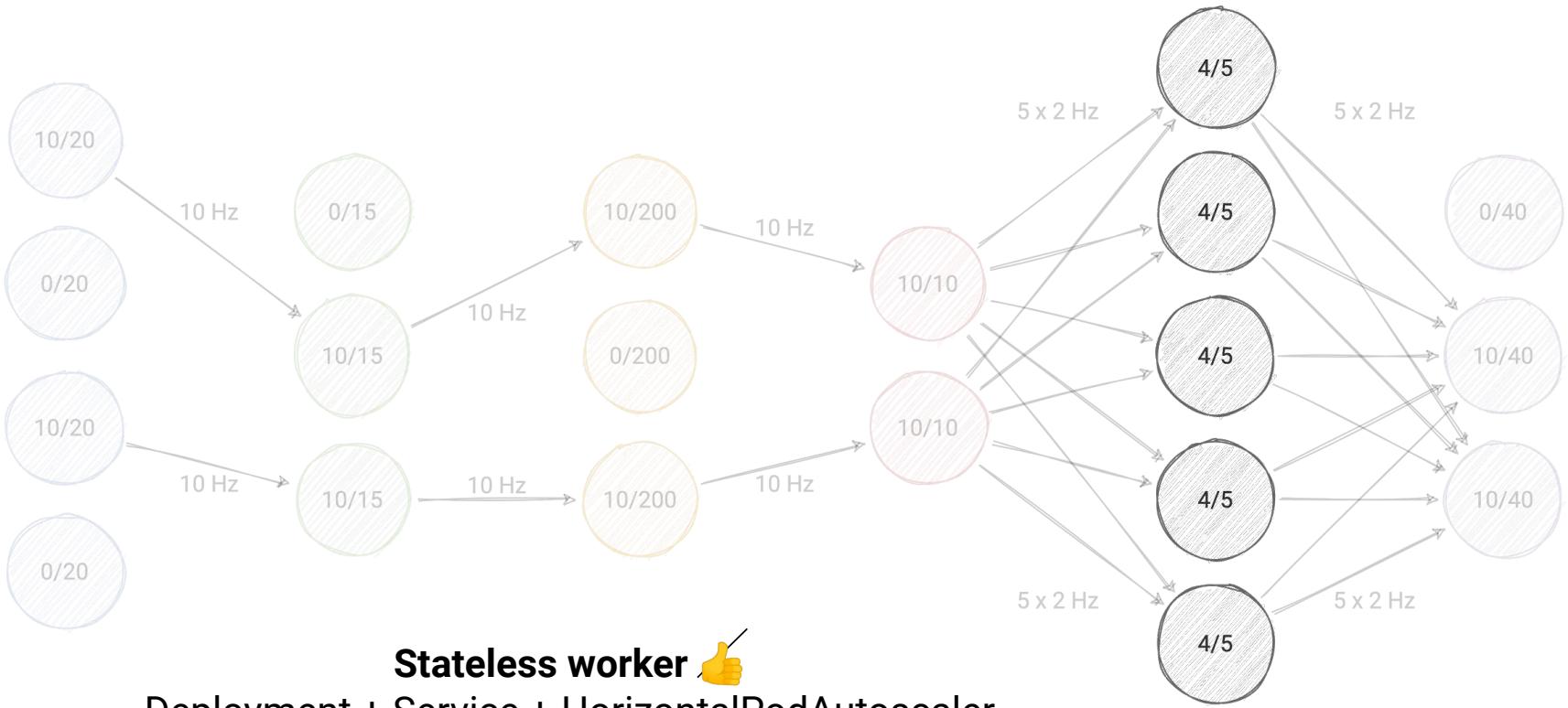


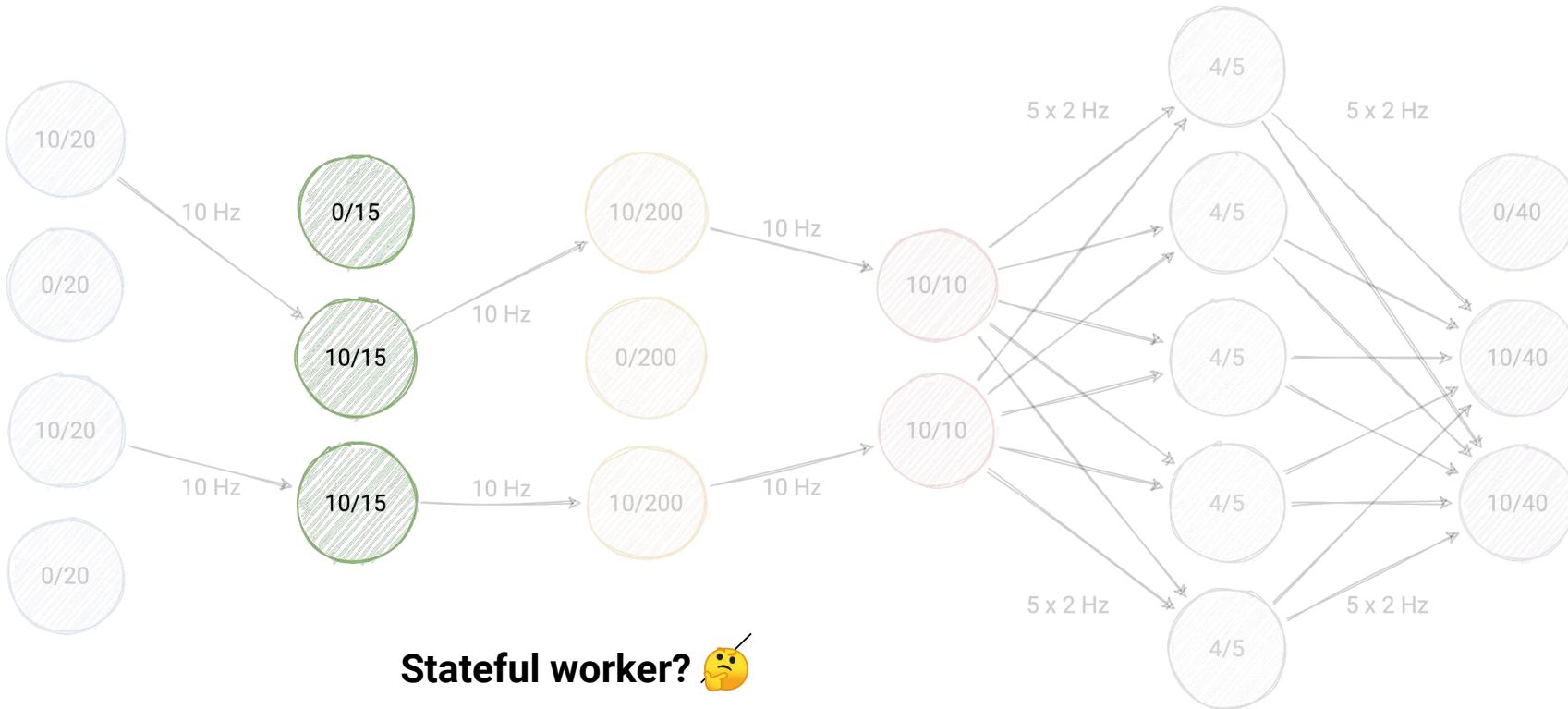




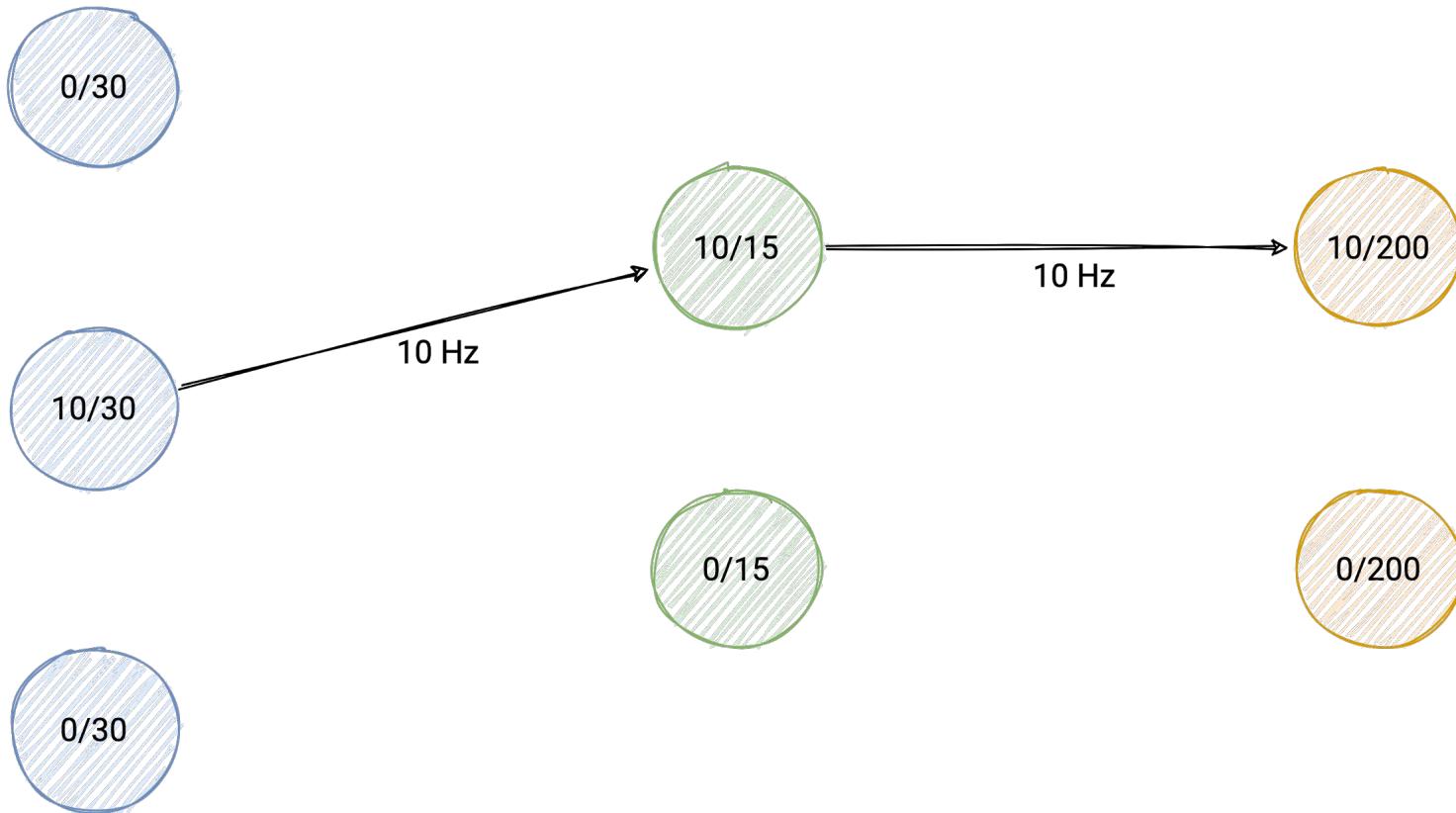


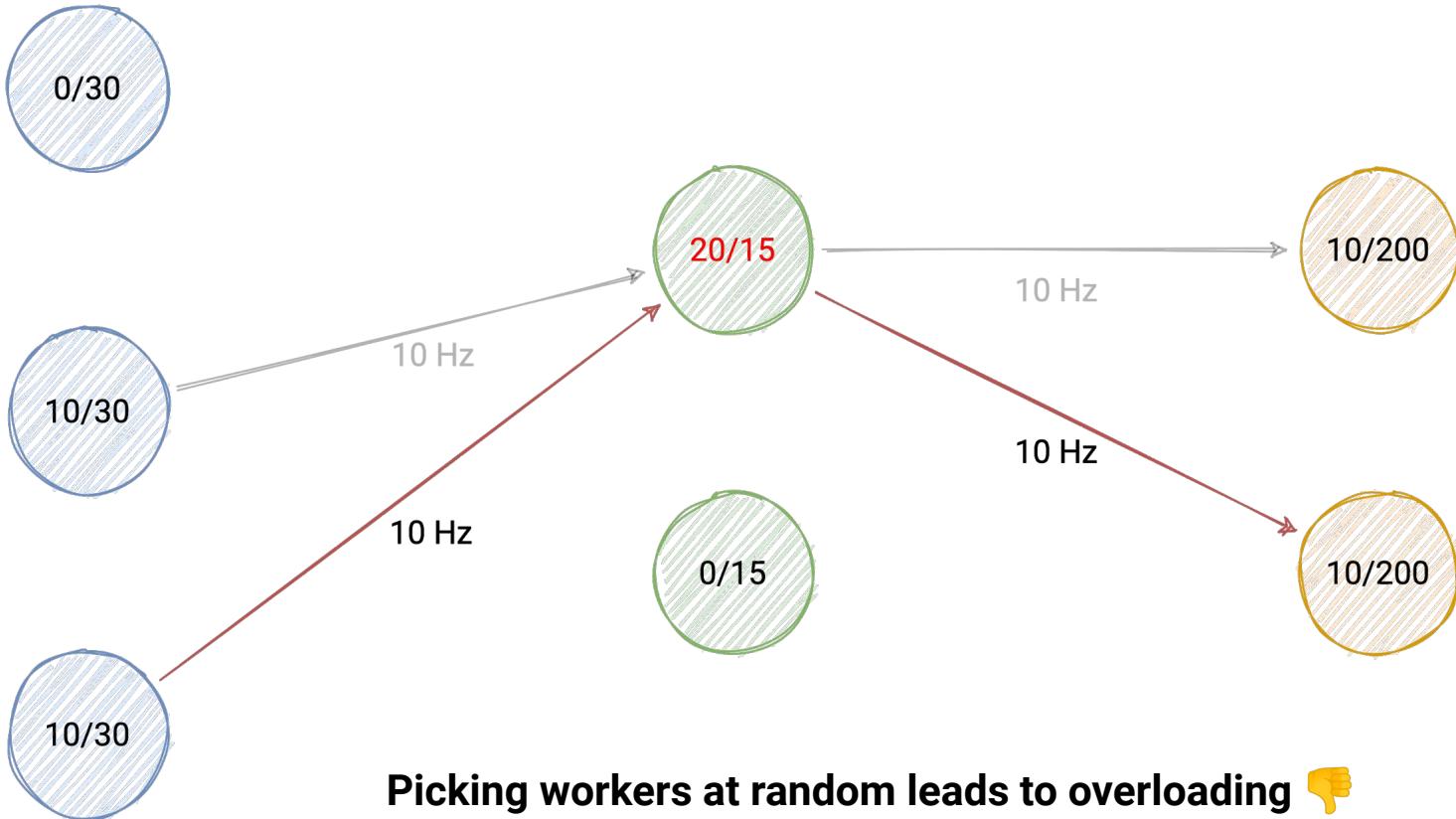




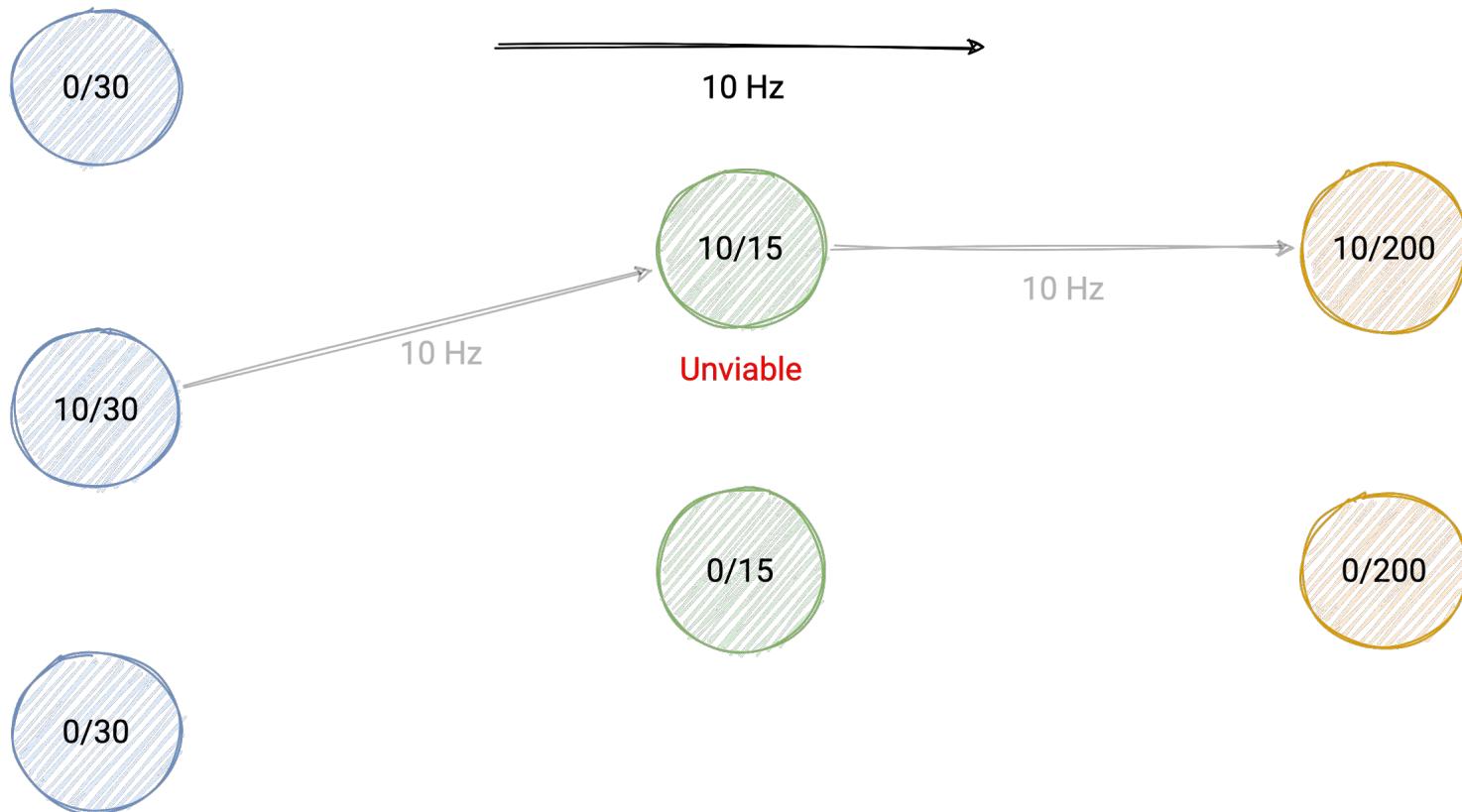


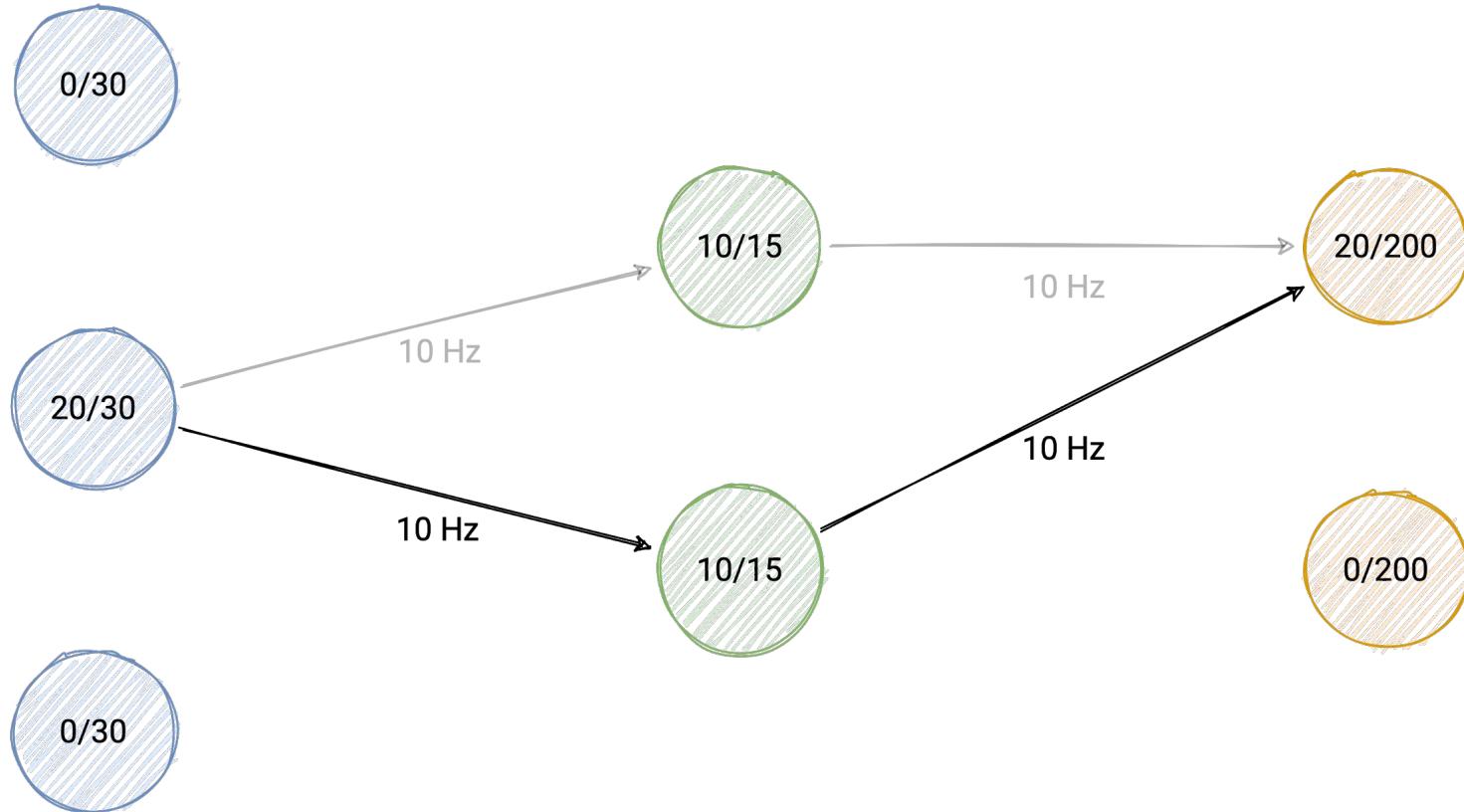


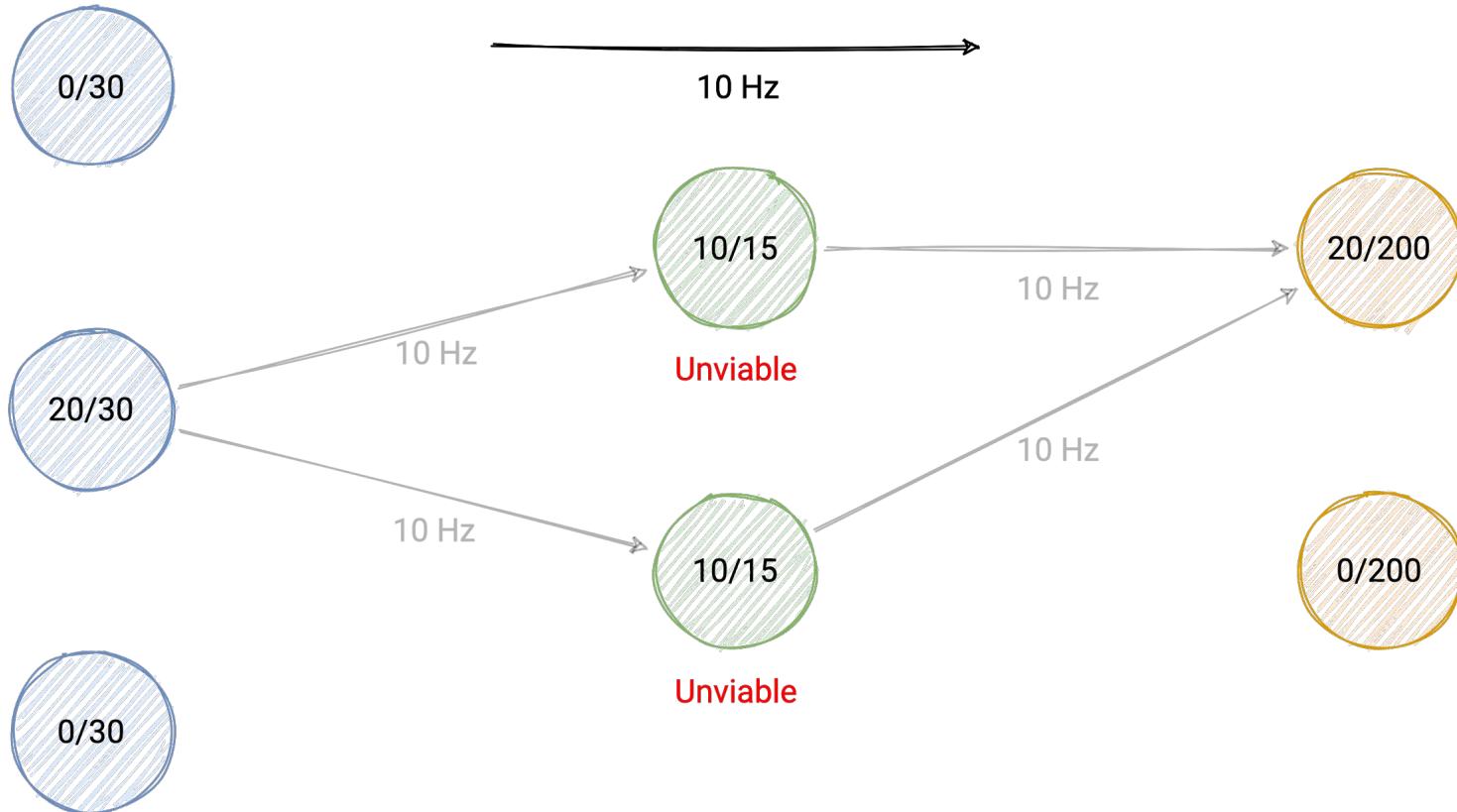


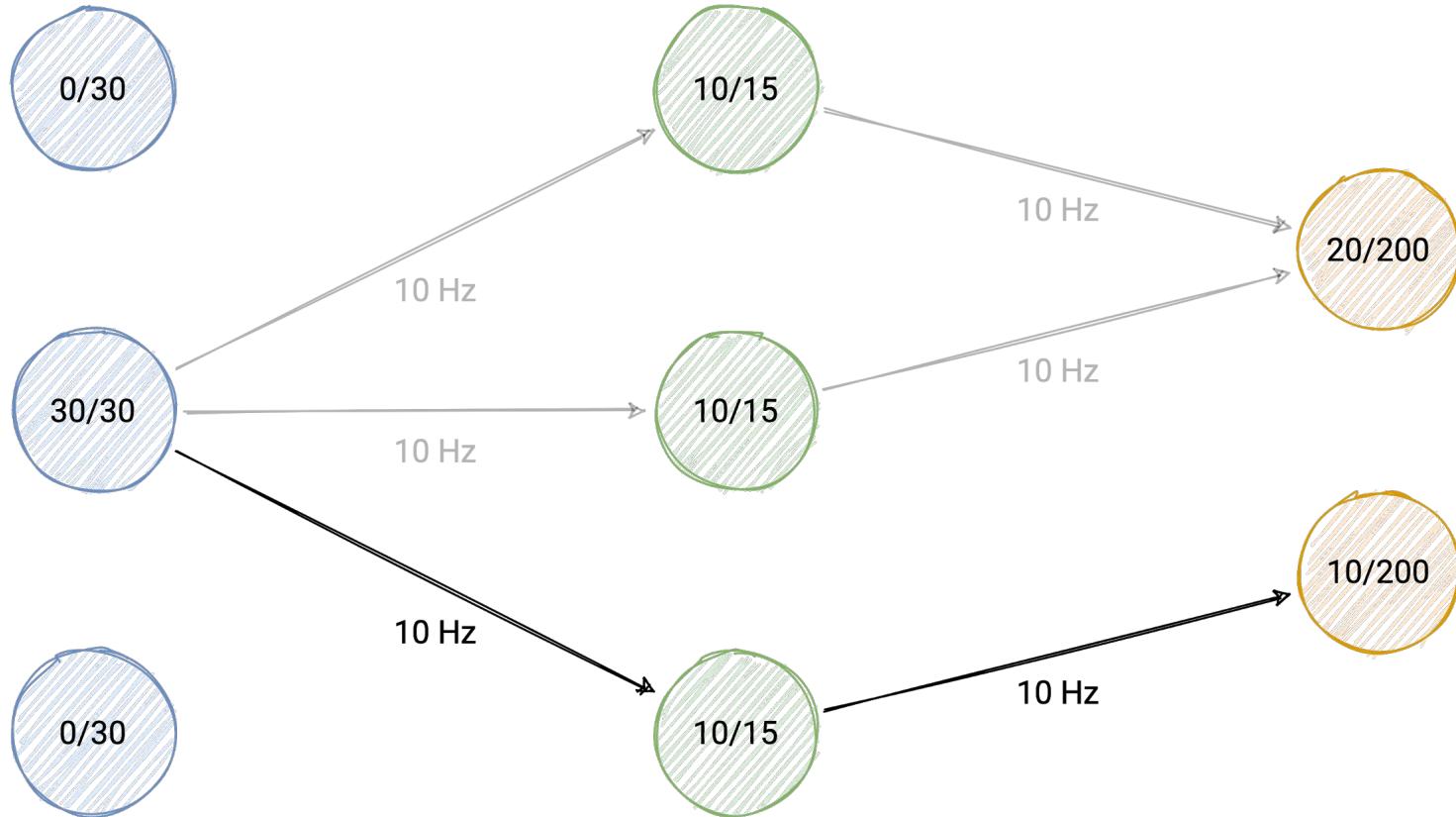


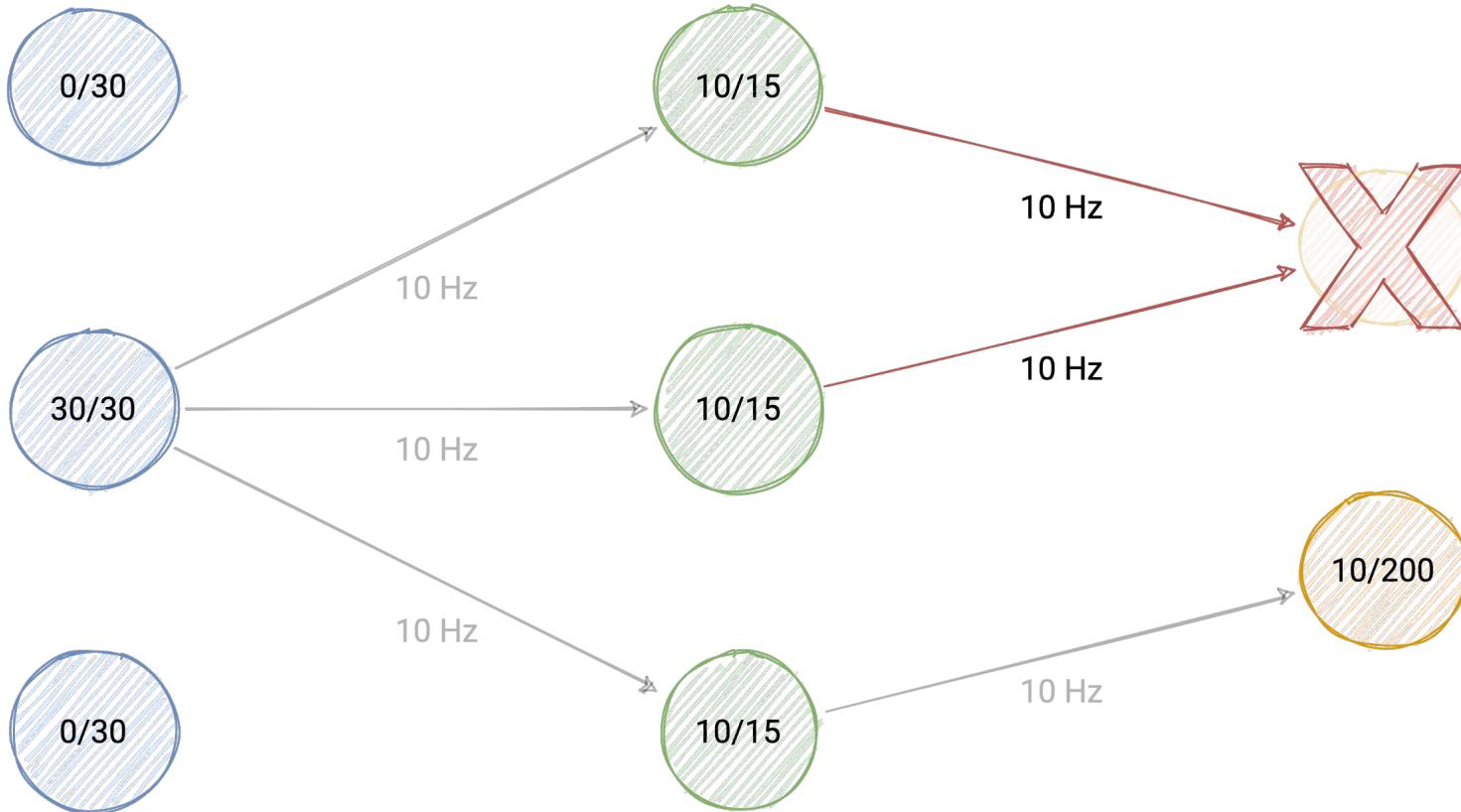
Picking workers at random leads to overloading

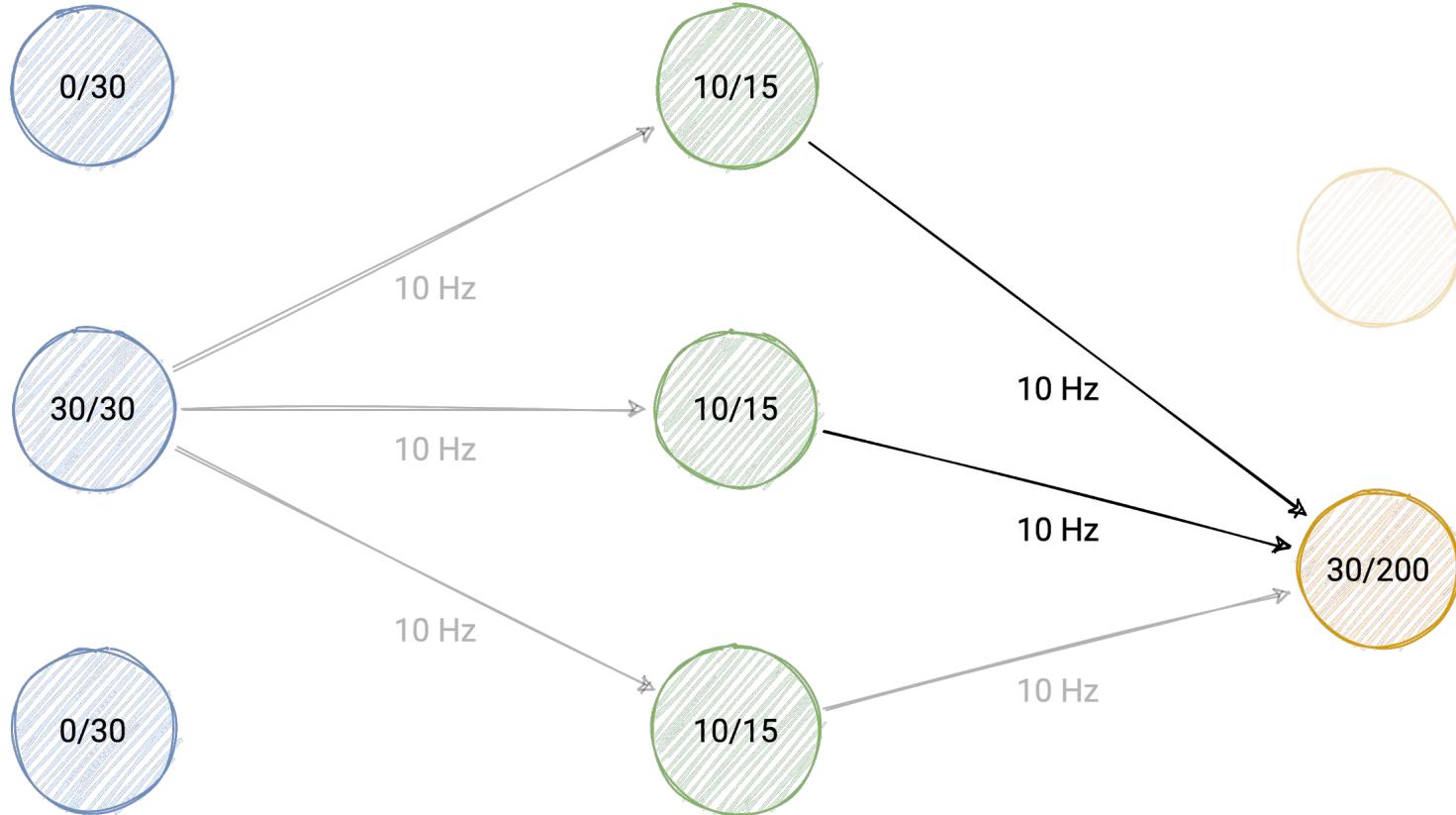












# Functional Requirements:

Select viable workers for each pipeline

Scale number of workers based on load

React to changes in worker status

- Why we built it
- **What we built it with**
- How it works
- Lessons learned
- Overall experience



```
kubebuilder init --domain=k8s.xxii.io
```

# Custom Resource Definitions

Pipeline

PipelineBinding

WorkerClass

```
kubebuilder create api --group=scheduling --version=v1alpha1 --kind=Pipeline  
kubebuilder create api --group=scheduling --version=v1alpha1 --kind=PipelineBinding  
kubebuilder create api --group=scheduling --version=v1alpha1 --kind=WorkerClass
```

```
// PipelineSpec defines the desired state of Route
type PipelineSpec struct {
    // Number of documents that will be sent along this Pipeline every second.
    // +kubebuilder:validation:Minimum=0
    DocumentsPerSecond int `json:"documentsPerSecond"`

    // The directed acyclic graph of workers to use in this Pipeline.
    Workers map[string]Worker `json:"workers"`
}

// Worker is a node in a Pipeline's graph of workers.
type Worker struct {
    // Name of the WorkerClass to get the worker's specification from.
    // +kubebuilder:validation:MinLength=1
    Class string `json:"class"`

    // Next workers in the the Pipelines's graph.
    Next []string `json:"next"`
}
```

make manifests  
make install

```
# Sample Pipeline
apiVersion: scheduling.k8s.xxii.io/v1alpha1
kind: Pipeline
metadata:
  name: my-pipeline
spec:
  documentsPerSecond: 10
  workers:
    first:
      class: blue
      next:
        - second
    second:
      class: green
      next:
        - third
    third:
      class: yellow
      next: []
status: {}
```

```
// Reconcile the cluster state based on a Pipeline's spec.
func (r *PipelineReconciler) Reconcile(req ctrl.Request) (ctrl.Result, error) {
    var pipeline scheduling.Pipeline
    r.get(req, &pipeline)

    r.updateStatus(&pipeline)

    ok := r.selectWorkers(&pipeline)
    if !ok {
        r.unsetBindings(&pipeline)
        return ctrl.Result{RequeueAfter: time.Second}, nil
    }

    r.createBindings(&pipeline)

    return ctrl.Result{}, nil
}
```

```
# HELP workerclass_request Load the WorkerClass needs to process
# TYPE workerclass_request gauge
workerclass_request{workerclass_name="blue" ,workerclass_namespace="default"} 30
workerclass_request{workerclass_name="green" ,workerclass_namespace="default"} 75
workerclass_request{workerclass_name="yellow",workerclass_namespace="default"} 110
workerclass_request{workerclass_name="red"   ,workerclass_namespace="default"} 0
workerclass_request{workerclass_name="grey"  ,workerclass_namespace="default"} 0
workerclass_request{workerclass_name="purple",workerclass_namespace="default"} 60
```

```
# Prometheus Adapter Rules
rules:
  custom:
    - seriesQuery: '{__name__=~"^workerclass_.*"}'
      name:
        matches: ^workerclass_.*$  

        as: ${1}
      resources:
        overrides:
          workerclass_namespace:
            resource: namespace
          workerclass_name:
            group: scheduling.k8s.xxii.io
            resource: workerclass
      metricsQuery: <<.Series>>{<<.LabelMatchers>>}
```

# Horizontal Pod Autoscaler

Core Kubernetes feature

```
apiVersion: autoscaling/v2beta2
kind: HorizontalPodAutoscaler
metadata: ...
spec:
  minReplicas: 0
  maxReplicas: 10
  scaleTargetRef:
    apiVersion: apps/v1
    kind: StatefulSet
    name: blue-workers
  metrics:
    - type: Object
      object:
        metric:
          name: workerclass_request
        describedObject:
          apiVersion: scheduling.k8s.xxii.io/v1alpha1
          kind: WorkerClass
          name: blue
        target:
          type: Value
          averageValue: '20'
```

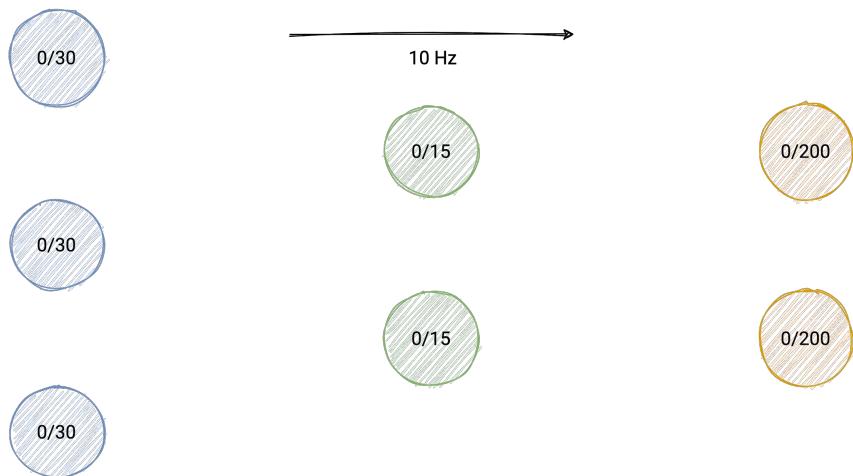
- Why we built it
- What we built it with
- How it works**
- Development
- Lessons learned
- Overall experience

## Select viable workers for each pipeline

Scale number of workers based on load

React to changes in worker status

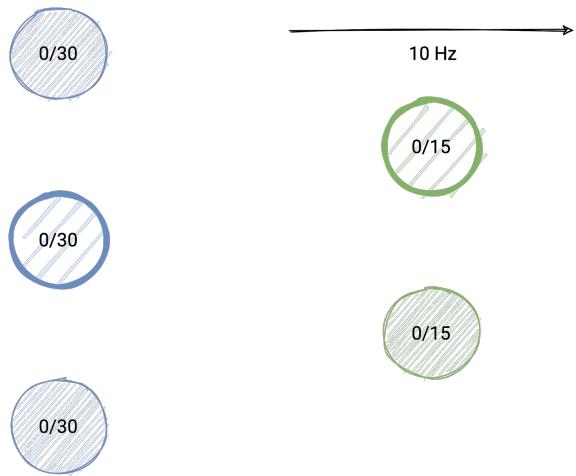




WorkerClass	NAME	CAPACITY
blue	30	
green	15	
yellow	200	

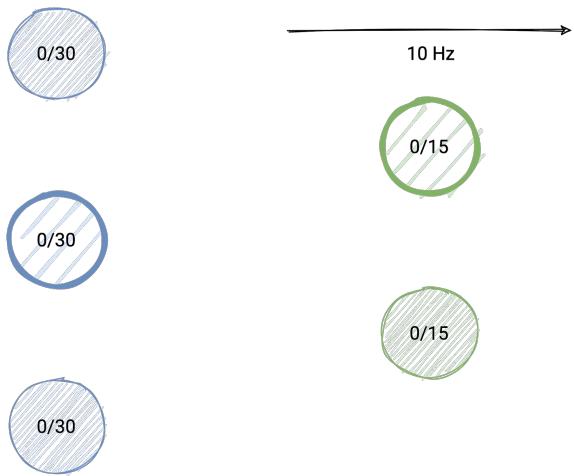
Pipeline	NAME	LOAD	STATUS
first	10		Pending



WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

Pipeline	NAME	LOAD	STATUS
	first	10	<b>Plotted</b>

PipelineBinding	NAME	STATUS
	first-blue	<b>Pending</b>
	first-green	<b>Pending</b>
	first-yellow	<b>Pending</b>

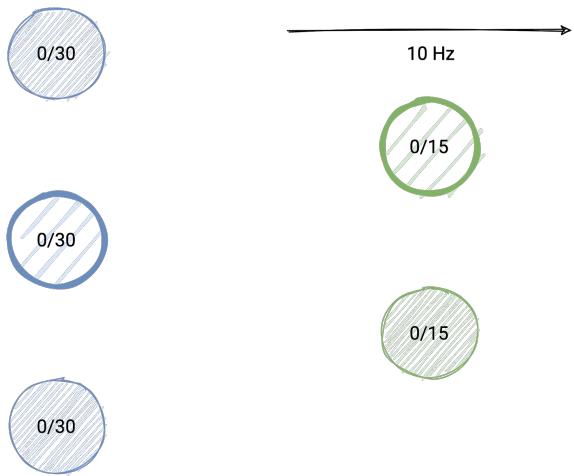


WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

Pipeline	NAME	LOAD	STATUS
	first	10	Plotted

PipelineBinding	NAME	STATUS
	first-blue	Pending
	first-green	Pending
	first-yellow	Pending

Service	NAME	ENDPOINTS
	first-blue	blue-1
	first-green	green-0
	first-yellow	yellow-0

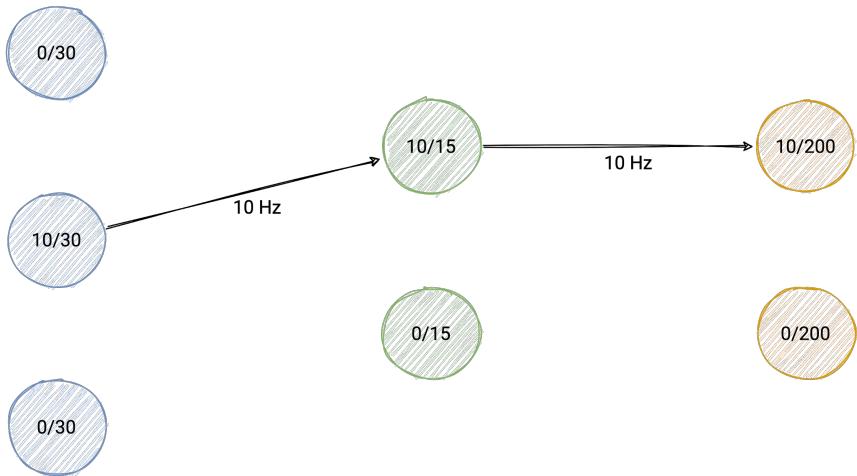


WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

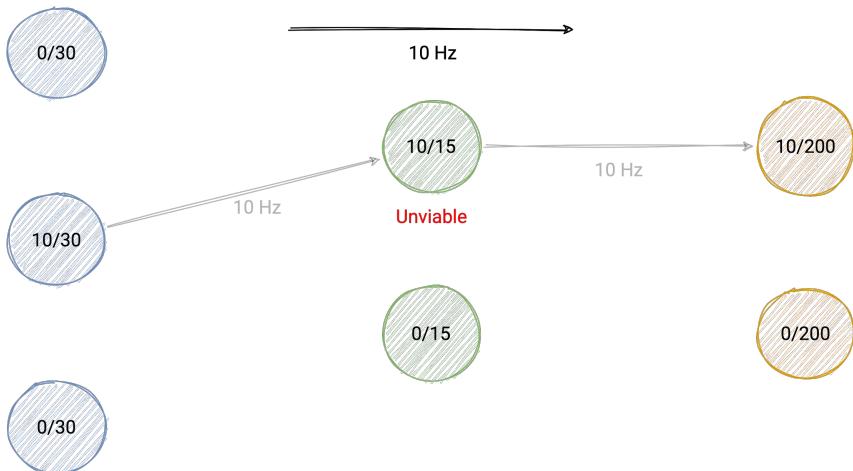
Pipeline	NAME	LOAD	STATUS
	first	10	Plotted

PipelineBinding	NAME	STATUS
	first-blue	<b>Ready</b>
	first-green	<b>Ready</b>
	first-yellow	<b>Ready</b>

Service	NAME	ENDPOINTS
	first-blue	blue-1
	first-green	green-0
	first-yellow	yellow-0



WorkerClass	NAME	CAPACITY	
	blue	30	
	green	15	
	yellow	200	
Pipeline	NAME	LOAD	STATUS
	first	10	<b>Ready</b>
PipelineBinding	NAME	STATUS	
	first-blue	Ready	
	first-green	Ready	
	first-yellow	Ready	
Service	NAME	ENDPOINTS	
	first-blue	blue-1	
	first-green	green-0	
	first-yellow	yellow-0	

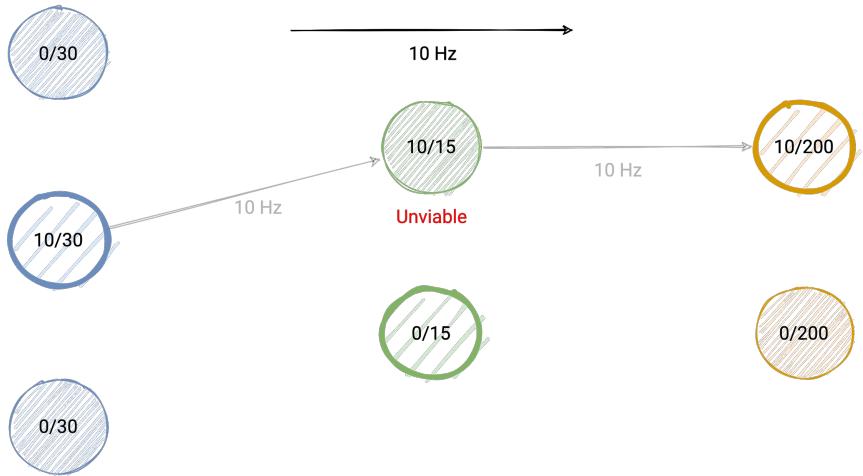


WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

Pipeline	NAME	LOAD	STATUS
	first	10	Ready
	<b>second</b>	<b>10</b>	<b>Pending</b>

PipelineBinding	NAME	STATUS
	first-blue	Ready
	first-green	Ready
	first-yellow	Ready

Service	NAME	ENDPOINTS
	first-blue	blue-1
	first-green	green-0
	first-yellow	yellow-0

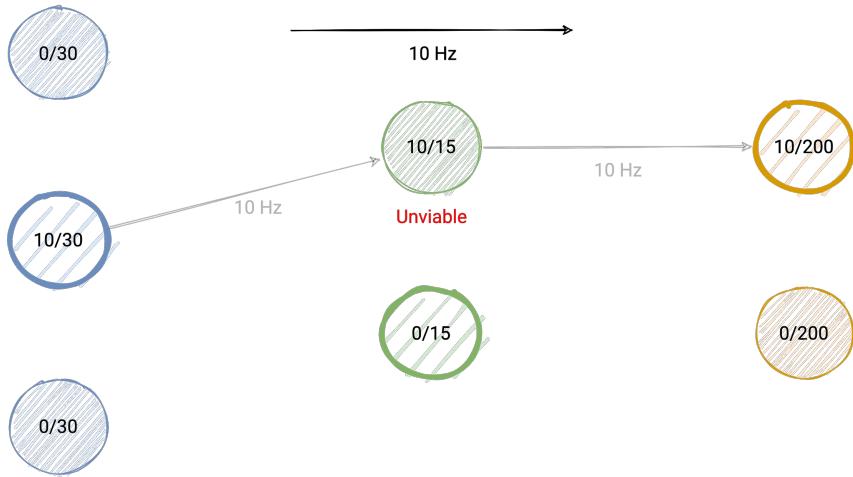


WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

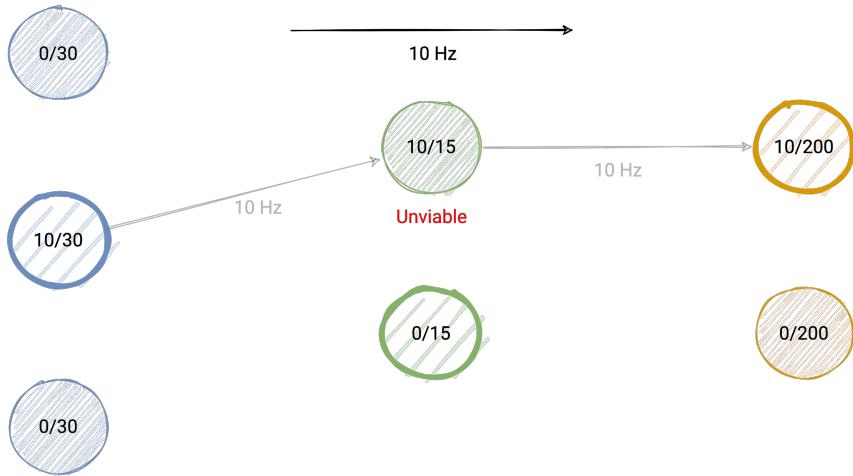
Pipeline	NAME	LOAD	STATUS
	first	10	Ready
	second	10	<b>Plotted</b>

PipelineBinding	NAME	STATUS
	first-blue	Ready
	first-green	Ready
	first-yellow	Ready
	<b>second-blue</b>	<b>Pending</b>
	<b>second-green</b>	<b>Pending</b>
	<b>second-yellow</b>	<b>Pending</b>

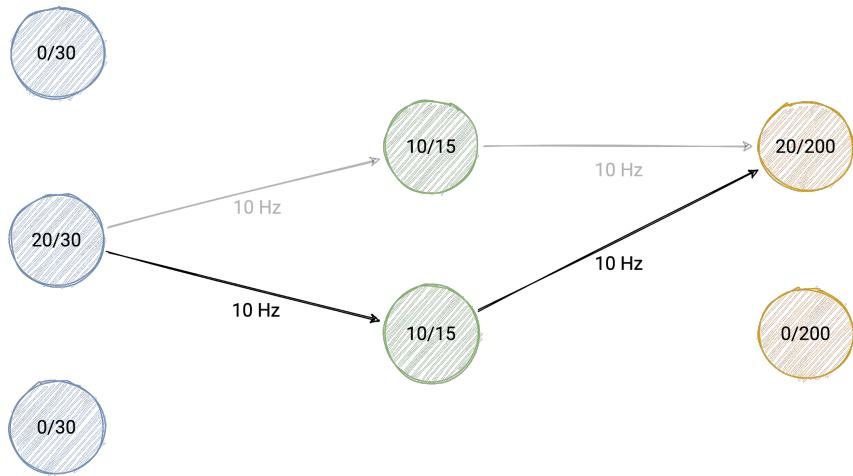
Service	NAME	ENDPOINTS
	first-blue	blue-1
	first-green	green-0
	first-yellow	yellow-0



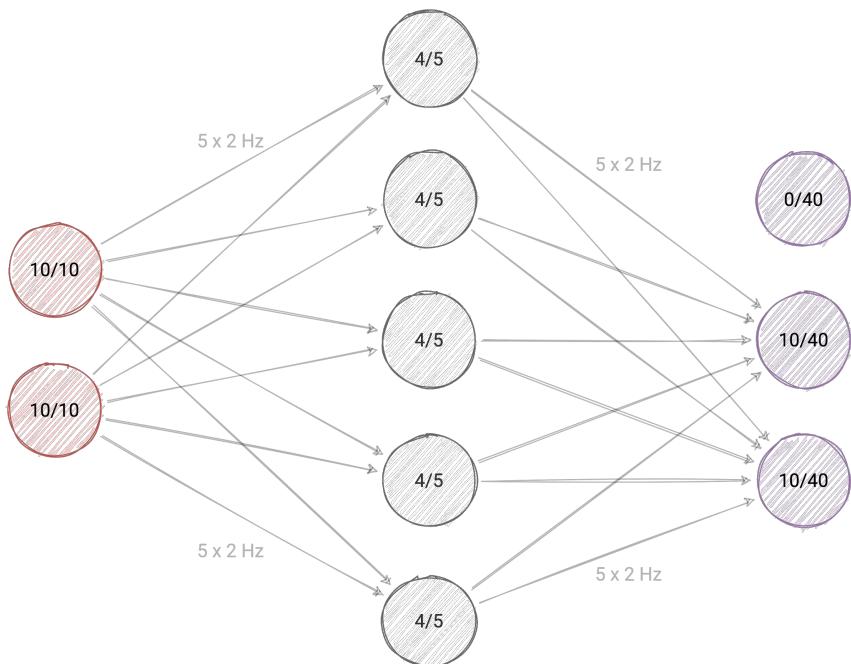
WorkerClass	NAME	CAPACITY	
	blue	30	
	green	15	
	yellow	200	
Pipeline	NAME	LOAD	STATUS
	first	10	Ready
	second	10	Plotted
PipelineBinding	NAME		STATUS
	first-blue		Ready
	first-green		Ready
	first-yellow		Ready
	second-blue		Pending
	second-green		Pending
	second-yellow		Pending
Service	NAME	ENDPOINTS	
	first-blue	blue-1	
	first-green	green-0	
	first-yellow	yellow-0	
	<b>second-blue</b>	<b>blue-1</b>	
	<b>second-green</b>	<b>green-1</b>	
	<b>second-yellow</b>	<b>yellow-0</b>	



WorkerClass	NAME	CAPACITY	
blue		30	
green		15	
yellow		200	
Pipeline	NAME	LOAD	STATUS
first		10	Ready
second		10	Plotted
PipelineBinding	NAME		STATUS
first-blue			Ready
first-green			Ready
first-yellow			Ready
second-blue			<b>Ready</b>
second-green			<b>Ready</b>
second-yellow			<b>Ready</b>
Service	NAME	ENDPOINTS	
first-blue		blue-1	
first-green		green-0	
first-yellow		yellow-0	
second-blue		blue-1	
second-green		green-1	
second-yellow		yellow-0	



WorkerClass	NAME	CAPACITY	
blue		30	
green		15	
yellow		200	
Pipeline	NAME	LOAD	STATUS
first		10	Ready
second		10	<b>Ready</b>
PipelineBinding	NAME		STATUS
first-blue		Ready	
first-green		Ready	
first-yellow		Ready	
second-blue		Ready	
second-green		Ready	
second-yellow		Ready	
Service	NAME	ENDPOINTS	
first-blue		blue-1	
first-green		green-0	
first-yellow		yellow-0	
second-blue		blue-1	
second-green		green-1	
second-yellow		yellow-0	

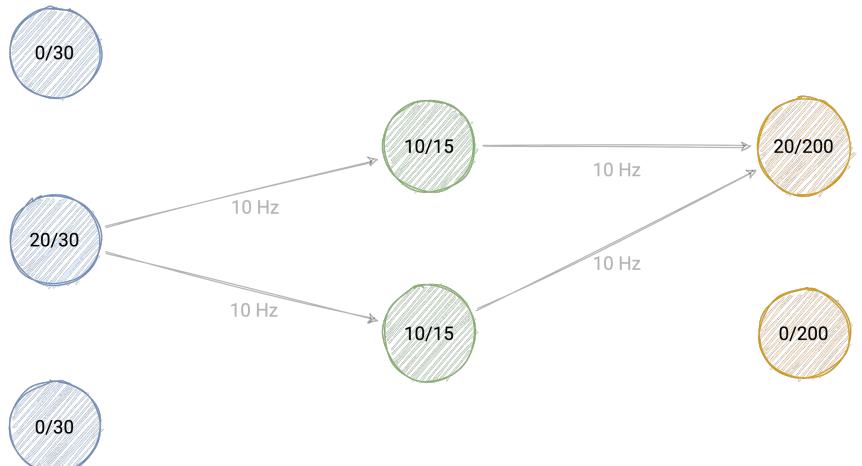


WorkerClass	NAME	CAPACITY	
	red	10	
	grey	5	
	purple	40	
Pipeline	NAME	LOAD	STATUS
	first	10	Ready
	second	10	Ready
PipelineBinding	NAME	STATUS	
	first-red	Ready	
	first-grey	Ready	
	first-purple	Ready	
	second-red	Ready	
	second-grey	Ready	
	second-purple	Ready	
Service	NAME	ENDPOINTS	
	first-red	red-0	
	<b>first-grey</b>	<b>grey-0,1,2,3,4</b>	
	first-purple	purple-1	
	second-red	red-1	
	<b>second-grey</b>	<b>grey-0,1,2,3,4</b>	
	second-purple	purple-2	

Select viable workers for each pipeline

**Scale number of workers based on load**

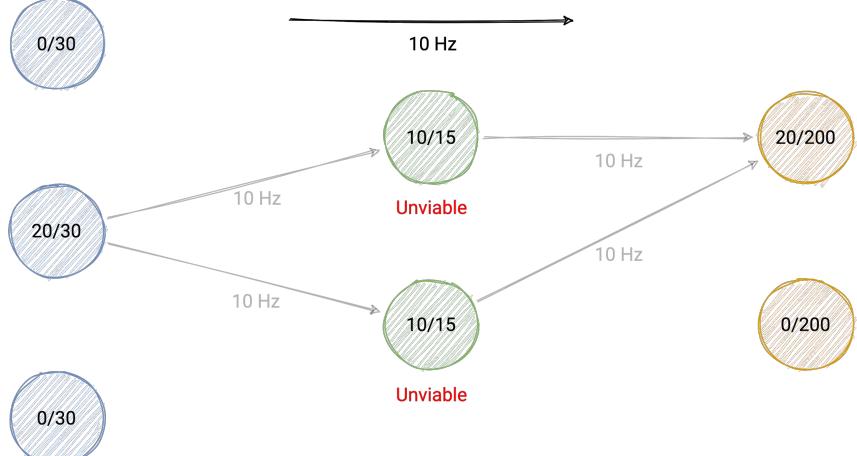
React to changes in worker status



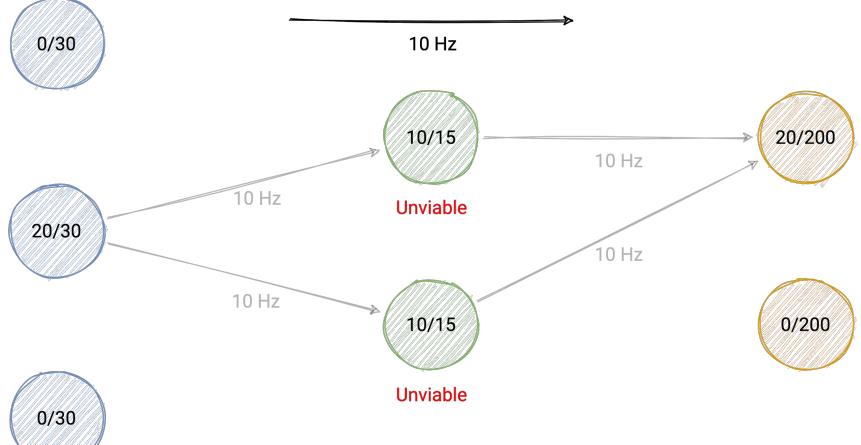
WorkerClass		NAME	CAPACITY	
		blue	30	
		green	15	
		yellow	200	
HPA	NAME	TARGET	ACTUAL	REPLICAS
	blue	70%	22%	3
	green	70%	66%	2
	yellow	70%	5%	2

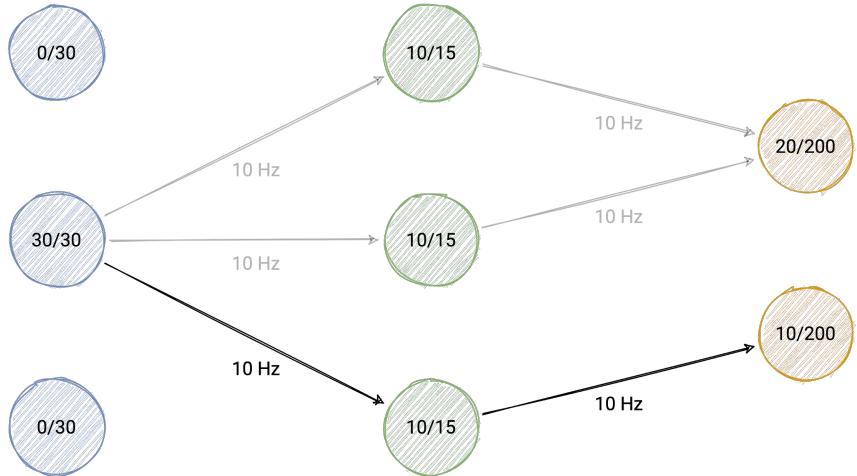
Pipeline		NAME	LOAD	STATUS
		first	10	Ready
		second	10	Ready



WorkerClass		NAME	CAPACITY	
	blue	30		
	green	15		
	yellow	200		
HPA	NAME	TARGET	ACTUAL	REPLICAS
	blue	70%	22%	3
	green	70%	66%	2
	yellow	70%	5%	2
Pipeline		NAME	LOAD	STATUS
	first	10	Ready	
	second	10	Ready	
	<b>third</b>	<b>10</b>	<b>Pending</b>	



WorkerClass		NAME	CAPACITY	
	blue	30		
	green	15		
	yellow	200		
HPA	NAME	TARGET	ACTUAL	REPLICAS
	blue	70%	<b>33%</b>	3
	green	70%	<b>100%</b>	2
	yellow	70%	<b>7%</b>	2
Pipeline		NAME	LOAD	STATUS
	first	10	Ready	
	second	10	Ready	
	third	10	Pending	

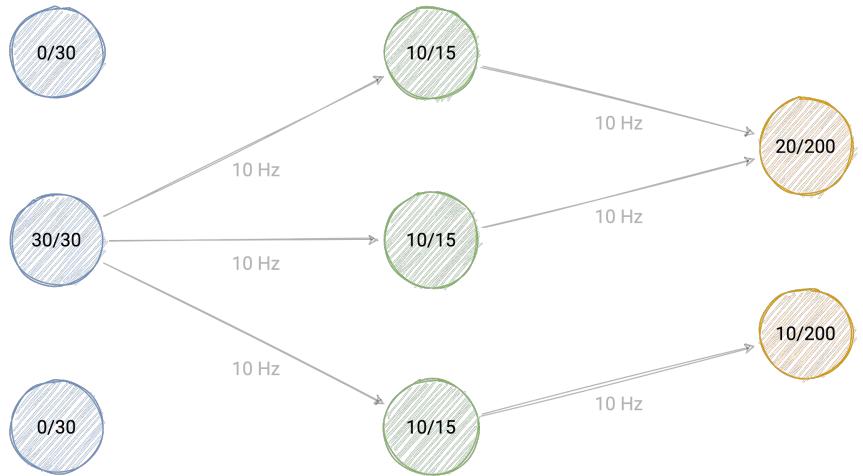


WorkerClass		NAME	CAPACITY	
	blue	30	30	
	green	15	15	
	yellow	200	200	
HPA	NAME	TARGET	ACTUAL	REPLICAS
	blue	70%	33%	3
	green	70%	<b>66%</b>	<b>3</b>
	yellow	70%	7%	2
Pipeline		NAME	LOAD	STATUS
	first	10	Ready	
	second	10	Ready	
	third	10	<b>Ready</b>	

Select viable workers for each pipeline

Scale number of workers based on load

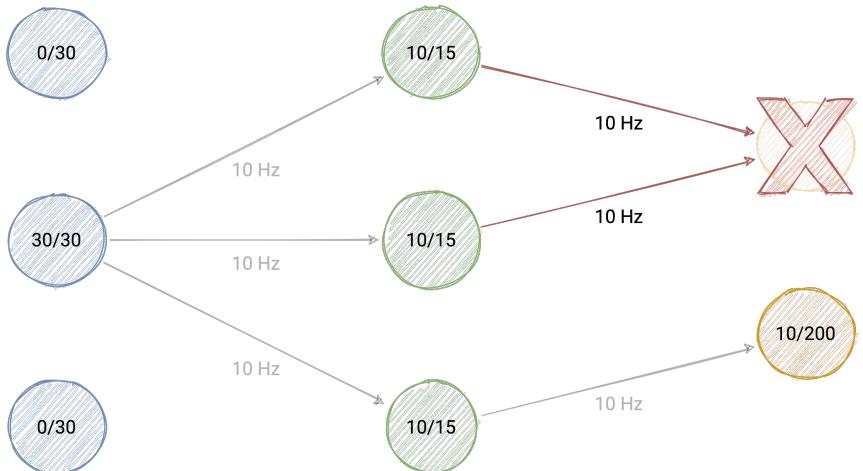
**React to changes in worker status**



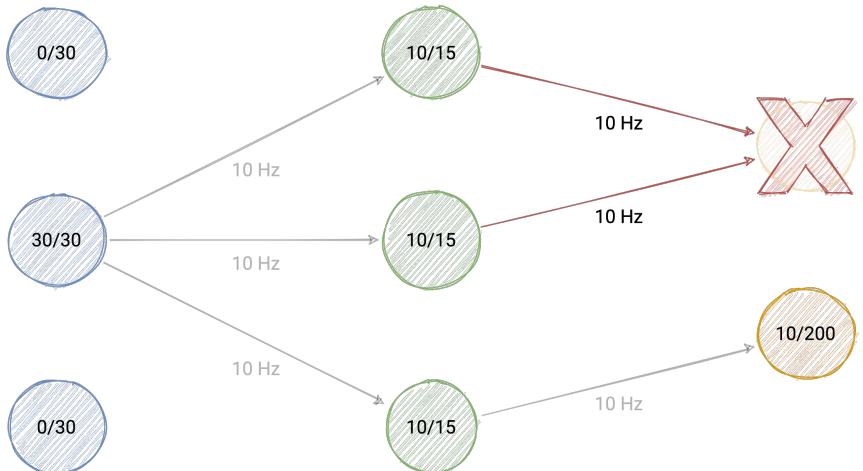
WorkerClass	NAME	CAPACITY
blue	30	
green	15	
yellow	200	

Pipeline	NAME	LOAD	STATUS
first	10	Ready	
second	10	Ready	
third	10	Ready	

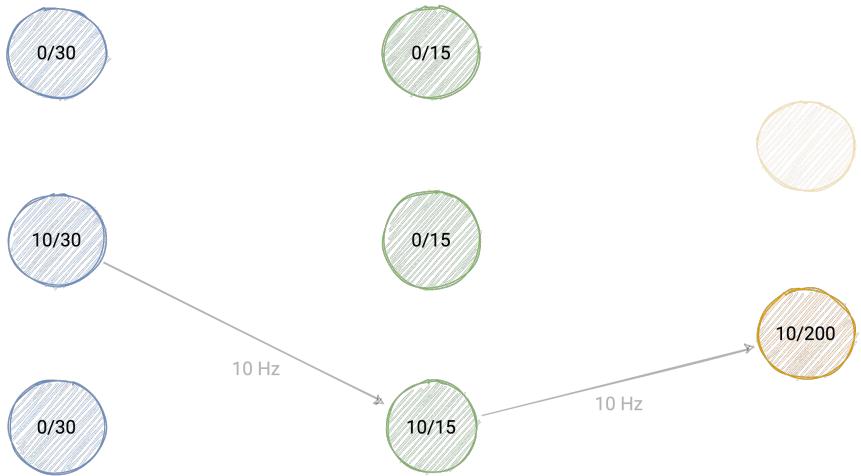
PipelineBinding	NAME	STATUS
first-blue	Ready	
first-green	Ready	
first-yellow	Ready	
second-blue	Ready	
second-green	Ready	
second-yellow	Ready	
third-blue	Ready	
third-green	Ready	
third-yellow	Ready	



WorkerClass	NAME	CAPACITY	
	blue	30	
	green	15	
	yellow	200	
Pipeline	NAME	LOAD	STATUS
	first	10	Ready
	second	10	Ready
	third	10	Ready
PipelineBinding	NAME		STATUS
	first-blue		Ready
	first-green		Ready
	first-yellow		<b>DeadEnd</b>
	second-blue		Ready
	second-green		Ready
	second-yellow		<b>DeadEnd</b>
	third-blue		Ready
	third-green		Ready
	third-yellow		Ready



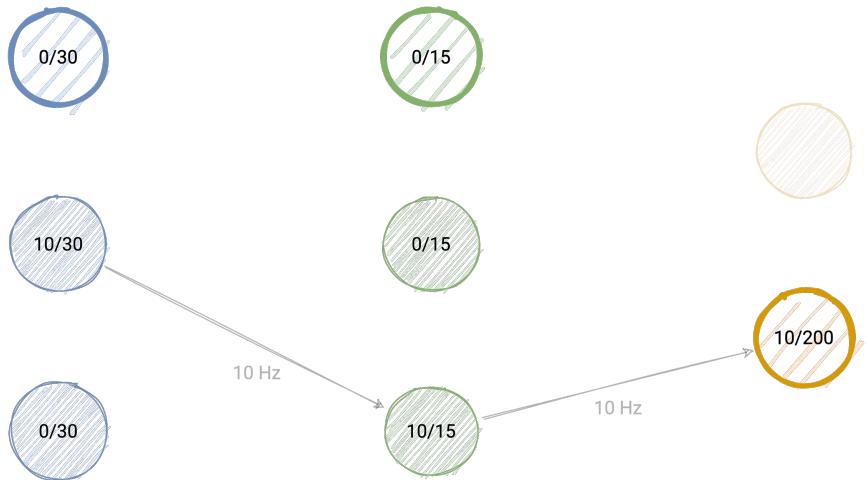
WorkerClass	NAME	CAPACITY	
	blue	30	
	green	15	
	yellow	200	
Pipeline	NAME	LOAD	STATUS
	first	10	<b>Broken</b>
	second	10	<b>Broken</b>
	third	10	Ready
PipelineBinding	NAME	STATUS	
	first-blue	Ready	
	first-green	Ready	
	first-yellow	DeadEnd	
	second-blue	Ready	
	second-green	Ready	
	second-yellow	DeadEnd	
	third-blue	Ready	
	third-green	Ready	
	third-yellow	Ready	



WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

Pipeline	NAME	LOAD	STATUS
	first	10	Broken
	second	10	Broken
	third	10	Ready

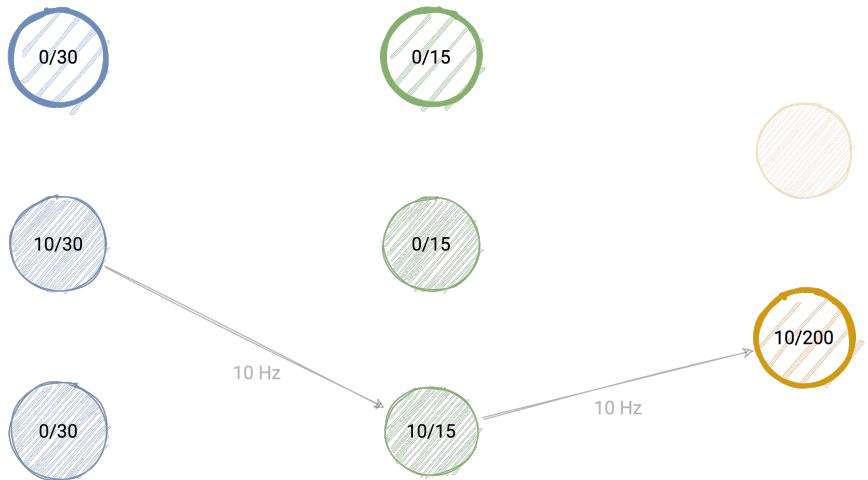
PipelineBinding	NAME	STATUS
	first-blue	Unset
	first-green	Unset
	first-yellow	Unset
	second-blue	Unset
	second-green	Unset
	second-yellow	Unset
	third-blue	Ready
	third-green	Ready
	third-yellow	Ready



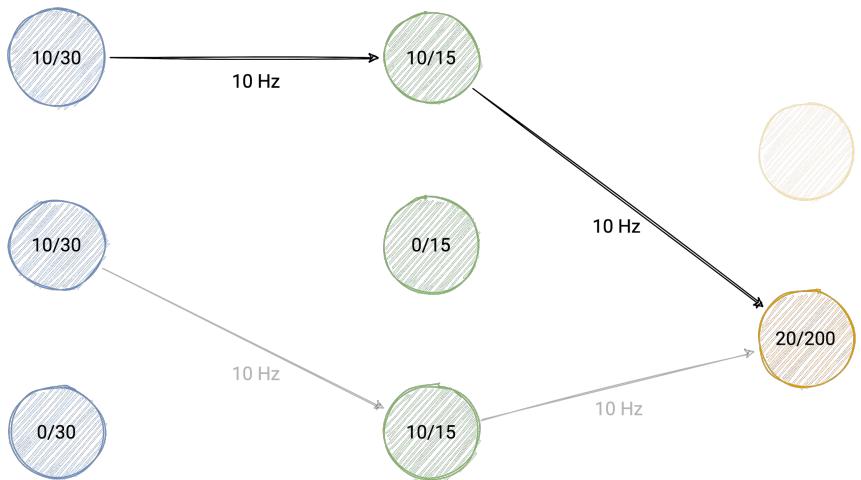
WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

Pipeline	NAME	LOAD	STATUS
	first	10	<b>Plotted</b>
	second	10	Broken
	third	10	Ready

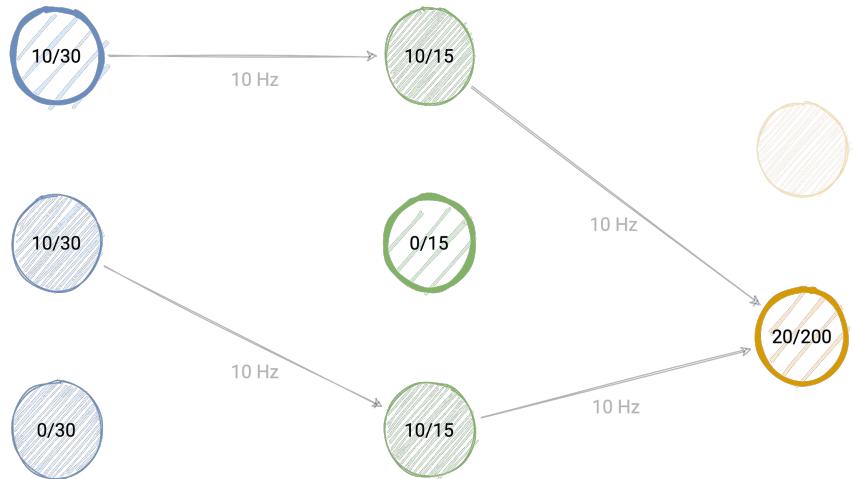
PipelineBinding	NAME	STATUS
	first-blue	<b>Pending</b>
	first-green	<b>Pending</b>
	first-yellow	<b>Pending</b>
	second-blue	Unset
	second-green	Unset
	second-yellow	Unset
	third-blue	Ready
	third-green	Ready
	third-yellow	Ready



WorkerClass	NAME	CAPACITY	
	blue	30	
	green	15	
	yellow	200	
Pipeline	NAME	LOAD	STATUS
	first	10	Plotted
	second	10	Broken
	third	10	Ready
PipelineBinding	NAME		STATUS
	first-blue		<b>Ready</b>
	first-green		<b>Ready</b>
	first-yellow		<b>Ready</b>
	second-blue		Unset
	second-green		Unset
	second-yellow		Unset
	third-blue		Ready
	third-green		Ready
	third-yellow		Ready



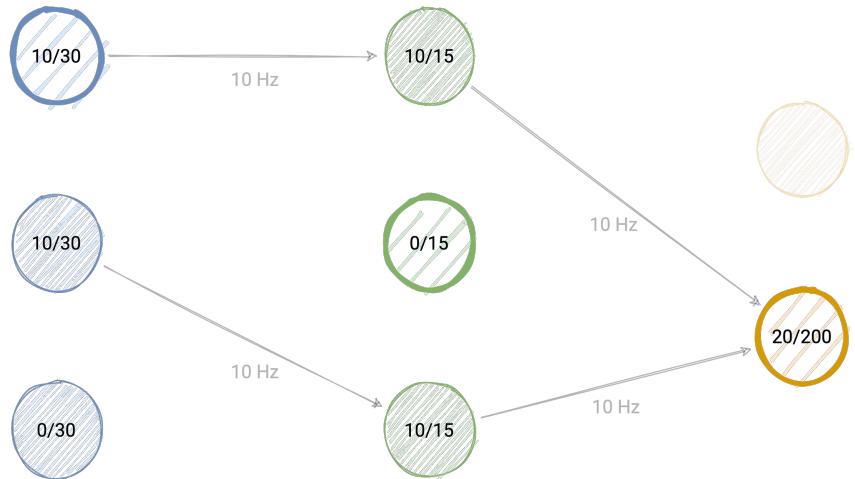
WorkerClass	NAME	CAPACITY	
	blue	30	
	green	15	
	yellow	200	
Pipeline	NAME	LOAD	STATUS
	first	10	<b>Ready</b>
	second	10	Broken
	third	10	Ready
PipelineBinding	NAME	STATUS	
	first-blue	Ready	
	first-green	Ready	
	first-yellow	Ready	
	second-blue	Unset	
	second-green	Unset	
	second-yellow	Unset	
	third-blue	Ready	
	third-green	Ready	
	third-yellow	Ready	



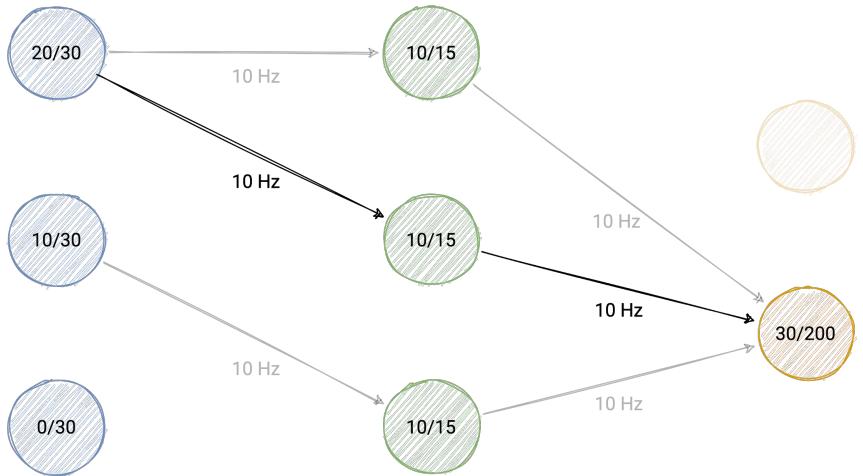
WorkerClass	NAME	CAPACITY
blue	30	
green	15	
yellow	200	

Pipeline	NAME	LOAD	STATUS
first	10	Ready	
second	10	<b>Plotted</b>	
third	10	Ready	

PipelineBinding	NAME	STATUS
first-blue	Ready	
first-green	Ready	
first-yellow	Ready	
second-blue	<b>Pending</b>	
second-green	<b>Pending</b>	
second-yellow	<b>Pending</b>	
third-blue	Ready	
third-green	Ready	
third-yellow	Ready	



WorkerClass	NAME	CAPACITY	
blue	30		
green	15		
yellow	200		
Pipeline	NAME	LOAD	STATUS
first	10	Ready	
second	10	Plotted	
third	10	Ready	
PipelineBinding	NAME	STATUS	
first-blue	Ready		
first-green	Ready		
first-yellow	Ready		
second-blue	<b>Ready</b>		
second-green	<b>Ready</b>		
second-yellow	<b>Ready</b>		
third-blue	Ready		
third-green	Ready		
third-yellow	Ready		



WorkerClass	NAME	CAPACITY
	blue	30
	green	15
	yellow	200

Pipeline	NAME	LOAD	STATUS
	first	10	Ready
	second	10	<b>Ready</b>
	third	10	Ready

PipelineBinding	NAME	STATUS
	first-blue	Ready
	first-green	Ready
	first-yellow	Ready
	second-blue	Ready
	second-green	Ready
	second-yellow	Ready
	third-blue	Ready
	third-green	Ready
	third-yellow	Ready

- Why we built it
- What we built it with
- How it works
- Lessons learned**
- Overall experience

# Operator does not scale horizontally

Only one instance can create/update resources at once

We have active-passive high availability

Services can't select Pods by name

Have to use StatefulSet to have **pod-name** label

Could we implement **WorkerSet** custom resource?

# Beware of race conditions

Resources of different kind will be handled concurrently

Resources of same kind will be handled sequentially

# Beware of deadlocks

Pipeline can remain Pending forever if HPA does not scale

Solution: constraints on Pipeline spec  
depending on which workers are used

- Why we built it
- What we built it with
- How it works
- Lessons learned
- **Overall experience**

Was it worth it? Absolutely.

The operator is fast and reliable

Adding new behavior is a breeze

# Kubebuilder is great

Handles all the plumbing

95% of our time was spent implementing features

# Significant design phase

Choice of custom resources

Finite state machine design, possible transitions, etc.

# Should you build your own?

Yes, if you have specific scheduling needs

Perfect for stateful applications running in Kubernetes

Thank you

We are hiring



PADOK

XXII



Arthur Busser



arthurb@padok.fr



ArthurBusser



busser