# Event-driven architecture

Unpacking Knative

Jeff Barnes
GC Micro-mission

# Event-driven (Serverless) Architecture

**Event-Driven**, sometimes also called serverless or functions as a service, is a computing execution model in which the infrastructure/provider dynamically manages the allocation of machine resources
- Code is deployed to a CSP or elsewhere
- Instance does not exist until invoked, spins up and scales as required, scales to zero when complete
- Changes IT requirements;
  - Security
  - Administration (patching, deployment)
  - Architecture
  - Cold start
- Currently various standards and little interoperability between CSP's
- Specific use cases include 'intents' for voice activated assistants (AWS Echo skills deployed in AWS Lambda)

IBM
OpenWhisk

AWS
Lambda

Google
Cloud Functions

Azure
Functions

• • • • •

# Event-driven (Serverless) Architecture

Originally found in public clouds, CSP instances are heading on-prem via appliance (Google Cloud Services Platform, Azure Stack, AWS Greengrass). A number of open-source versions are now available which can be deployed on-prem or in public cloud.

Examples include knative, OpenFaaS and Kubeless. This deck unpacks Knative, which can be used to auto-scale serverless-style functions, applications, and containers on Kubernetes



Knative          OpenFaaS          Kubeless

Initially serverless, Functions as a service (FaaS) and event-driven architecture referred to a microservice that is run 'on public cloud' compute engine only when invoked. Recently, serverless is also being used to describe managed container service (CaaS) and application deployments where the user is not responsible for the compute where/when the container runs (Azure AKS and ACI, AWS EKS and Fargate, Google GKE, Google App Engine, Ibm Kubernetes Service....)

# Knative

Knative can be used to deploy serverless-style functions, applications, and containers to an auto-scaling runtime on Kubernetes. Also has many other features: deploy multiple versions, perform custom tasks on your application's source code, build reusable templates…

- [Install Knative on an Istio Cluster](#) (Istio injection enabled)
- [Select and install demo app](#) ([helloworld-go](#))

Invoke the app

- Determine if running 'LoadBalancer' or 'NodePort'

```
export IP_ADDRESS=$(kubectl get svc istio-ingressgateway --namespace istio-system --output
'jsonpath={.status.loadBalancer.ingress[0].ip}')

OR

export IP_ADDRESS=$(kubectl get node  --output
'jsonpath={.items[0].status.addresses[0].address}'):$(kubectl get svc istio-ingressgateway --
namespace istio-system   --output 'jsonpath={.spec.ports[?(@.port==80)].nodePort}')

export HOST_URL=$(kubectl get ksvc helloworld-go  --output jsonpath='{.status.domain}')

curl -H "Host: ${HOST_URL}" http://${IP_ADDRESS}
Hello World: Go Sample v1!
```

# Knative – Cold Start

```
~/knative$ curl --header "Host: ${HOST_URL}" --write-out "
> lookup          %{time_namelookup}
> connect         %{time_connect}
> appconnect      %{time_appconnect}
> pretransfer     %{time_pretransfer}
> redirect        %{time_redirect}
> starttransfer  %{time_starttransfer}
> total           %{time_total}\n" \
> http://${IP_ADDRESS}
Hello Go Sample v1!

lookup          0.000067
connect         0.002420
appconnect      0.000000
pretransfer     0.002490
redirect        0.000000
starttransfer  15.767278
total          15.767717
```

Cold Start time 15 sec (need to fix it in my demo)

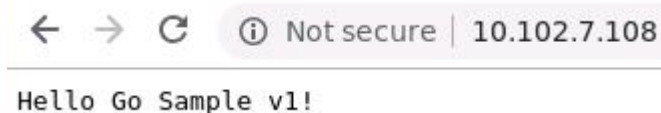# Knative – Subsequent Start

```
~/knative$ curl --header "Host: ${HOST_URL}" --write-out "
lookup         %{time_namelookup}
connect        %{time_connect}
appconnect     %{time_appconnect}
pretransfer    %{time_pretransfer}
redirect       %{time_redirect}
starttransfer  %{time_starttransfer}
total          %{time_total}\n" http://${IP_ADDRESS}
Hello Go Sample v1!

lookup         0.000073
connect        0.003975
appconnect     0.000000
pretransfer    0.004061
redirect       0.000000
starttransfer  0.062721
total          0.062793
```
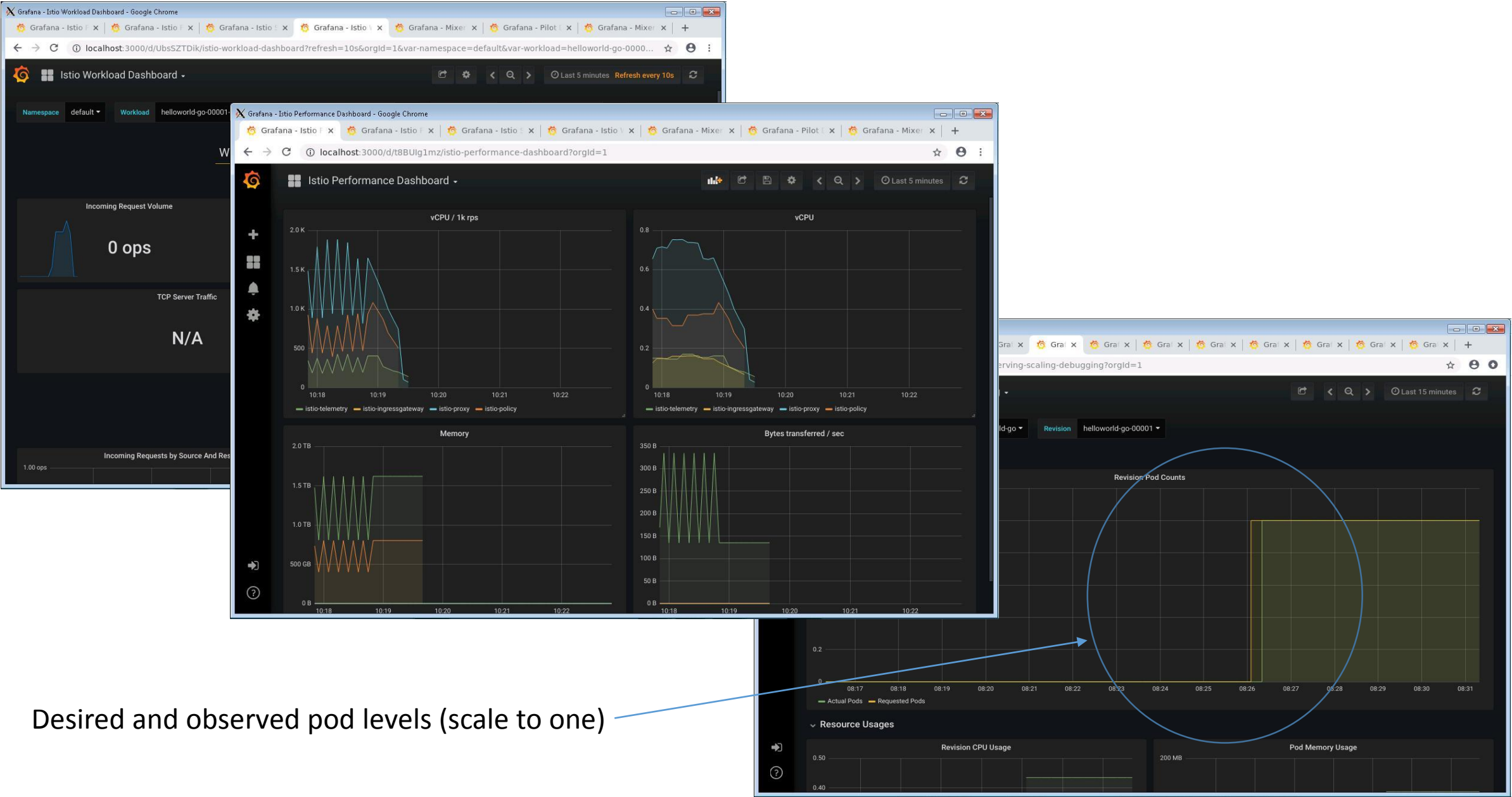
Subsequent response time < 1s

# Knative – Performance Monitoring Grafana



Desired and observed pod levels (scale to one)
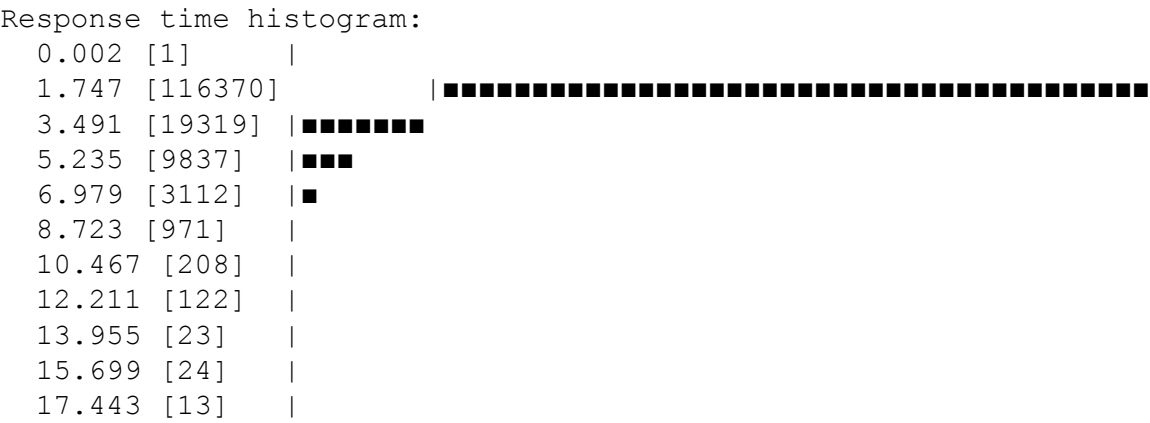
# Knative – Scale based on workload

[hey](#) runs number of requests in the provided concurrency level and prints stats.

-n  Number of requests to run. Default is 200.

-c  Number of requests to run concurrently. Total number of requests cannot
    be smaller than the concurrency level. Default is 50.

```
~/go/bin$ ./hey -host helloworld-go.default.example.com -c 500 -n
150000   "http://${IP_ADDRESS?}?sleep=1000"
```

Send 150,000 requests (with 500 requests in parallel), each taking 1 second

X 150,000

```
Summary:
  Total:         291.8676 secs
  Slowest:       17.4429 secs
  Fastest:       0.0025 secs
  Average:       0.9101 secs
  Requests/sec: 513.9316

  Total data:   3002220 bytes
  Size/request: 20 bytes

Response time histogram:
  0.002 [1]       |
  1.747 [116370]  |■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■■
  3.491 [19319]  |■■■■■■■
  5.235 [9837]   |■■■
  6.979 [3112]   |■
  8.723 [971]    |
  10.467 [208]   |
  12.211 [122]   |
  13.955 [23]    |
  15.699 [24]    |
  17.443 [13]    |
```

```
Latency distribution:
  10% in 0.0041 secs
  25% in 0.0050 secs
  50% in 0.0083 secs
  75% in 1.3462 secs
  90% in 3.4178 secs
  95% in 4.5164 secs
  99% in 6.8141 secs

Details (average, fastest, slowest):
  DNS+dialup:   0.0002 secs, 0.0025 secs, 17.4429 secs
  DNS-lookup:   0.0000 secs, 0.0000 secs, 0.0000 secs
  req write:    0.0002 secs, 0.0000 secs, 0.1051 secs
  resp wait:    0.9091 secs, 0.0024 secs, 17.4428 secs
  resp read:    0.0005 secs, 0.0000 secs, 0.1377 secs

Status code distribution:
  [200] 149940 responses
  [503] 60 responses
```
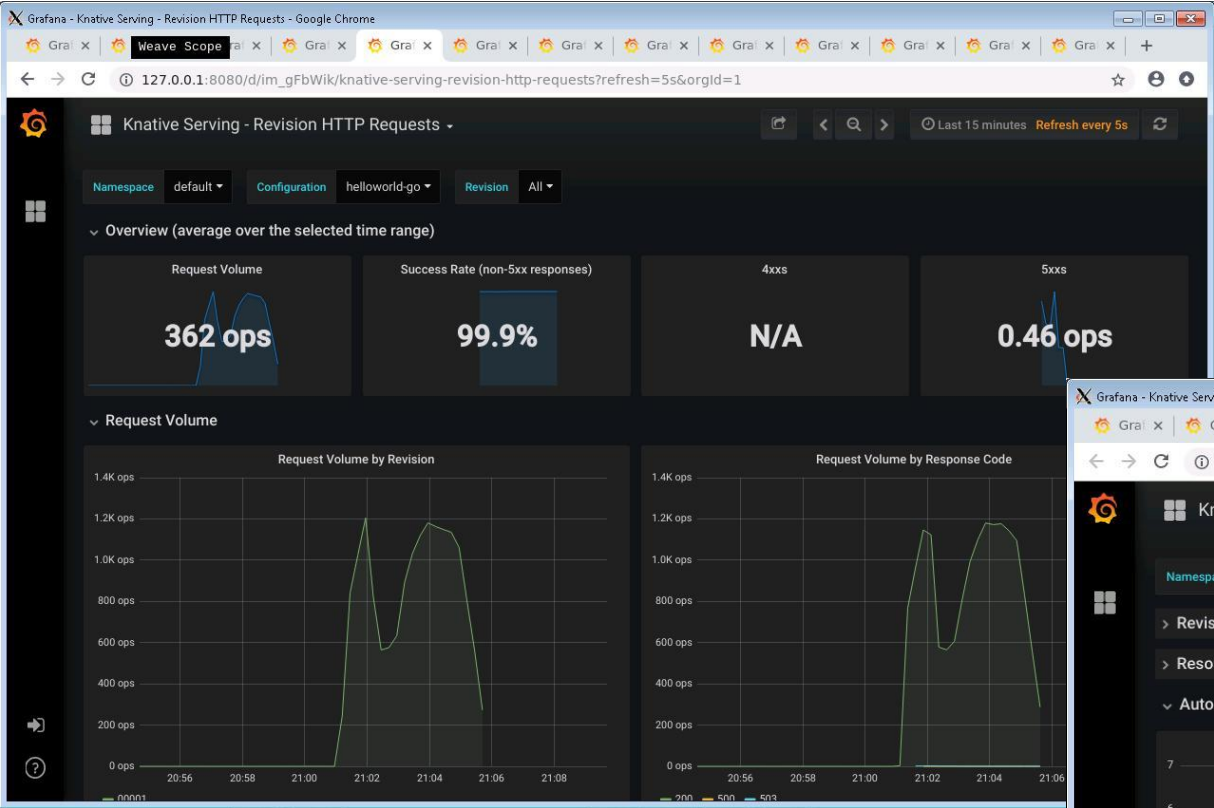
Not secure | 10.102.7.108

Hello Go Sample v1!

# Knative – Scale based on workload w/Performance Monitoring Grafana
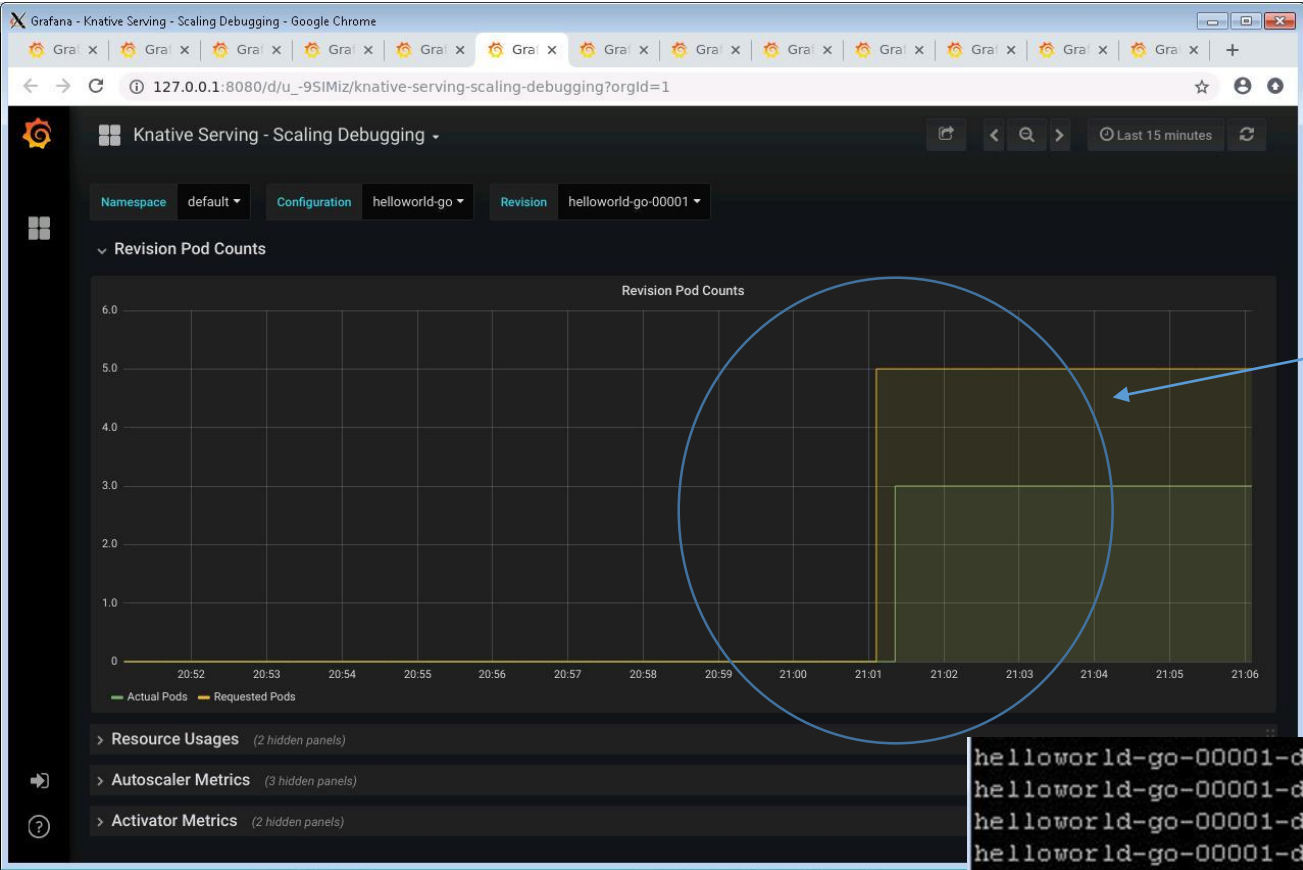
Desired and observed pod levels

# Knative – Scale based on workload

Knative Serving, by default, has a concurrent requests target of 100. Sending 500 concurrent requests causes autoscaling to note that it needs to run 5 Pods to satisfy this level



Desired and observed pod levels

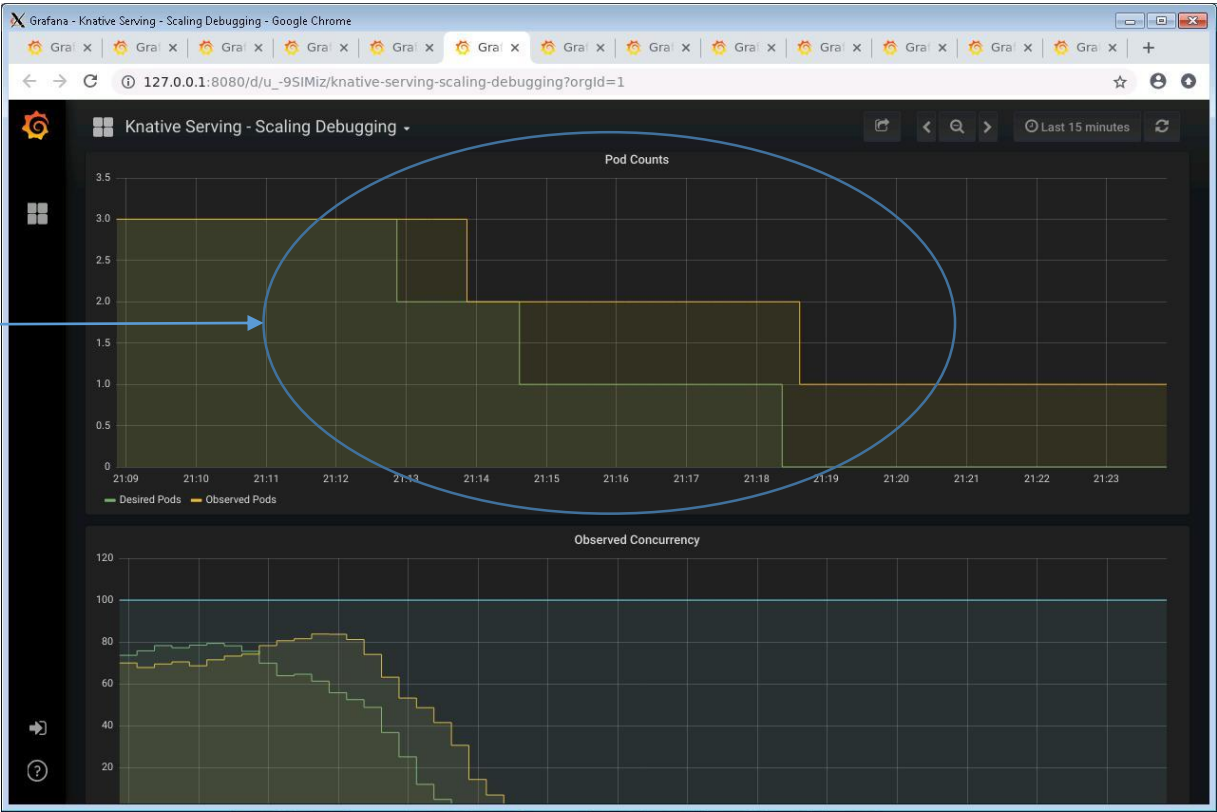Running and terminating Pods
(as required)

# Knative – Scale based on workload w/Performance Monitoring Grafana

```
helloworld-go-00001-deployment-94667496b-2npf2   0/3   Terminating   0   7m20s
helloworld-go-00001-deployment-94667496b-6tccn   3/3   Running       0   8m46s
helloworld-go-00001-deployment-94667496b-xrzjm   3/3   Running       0   9m6s
```

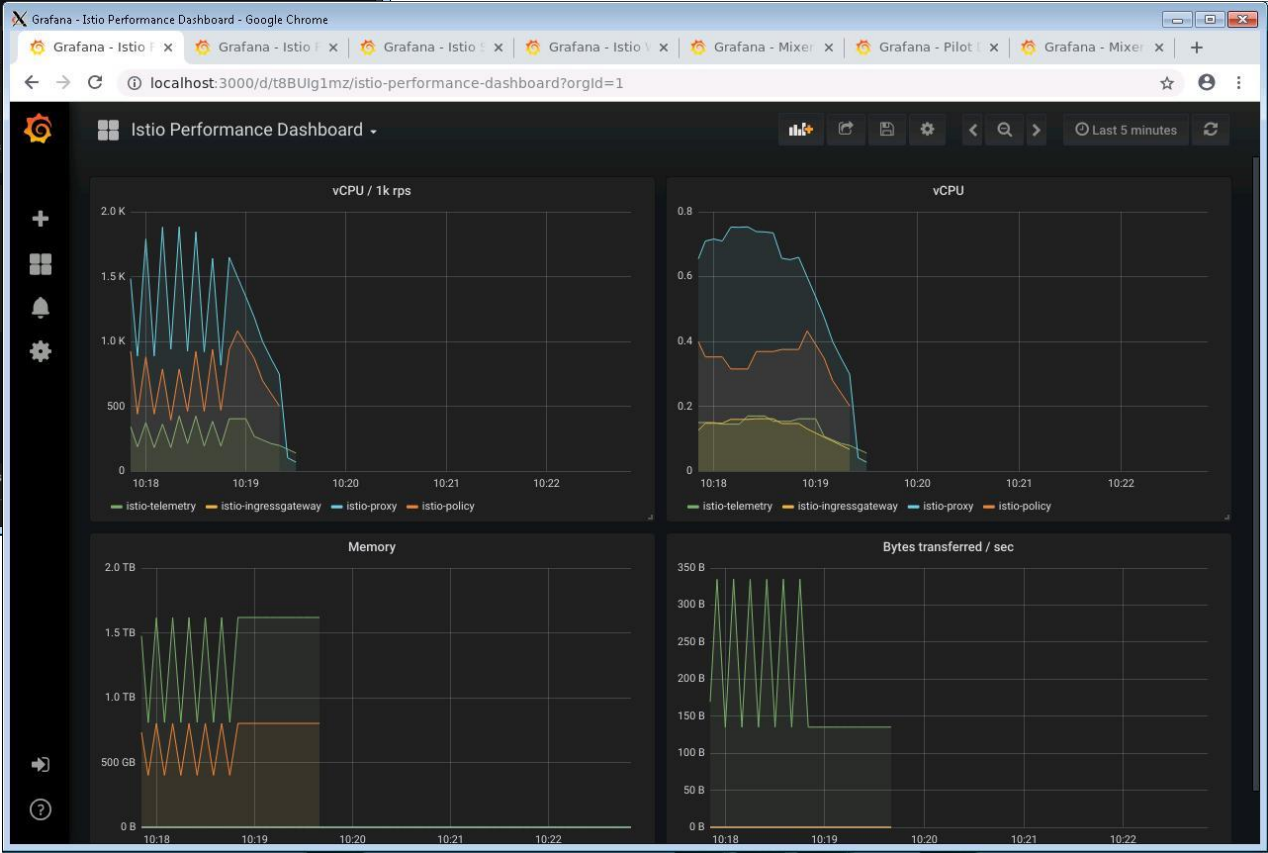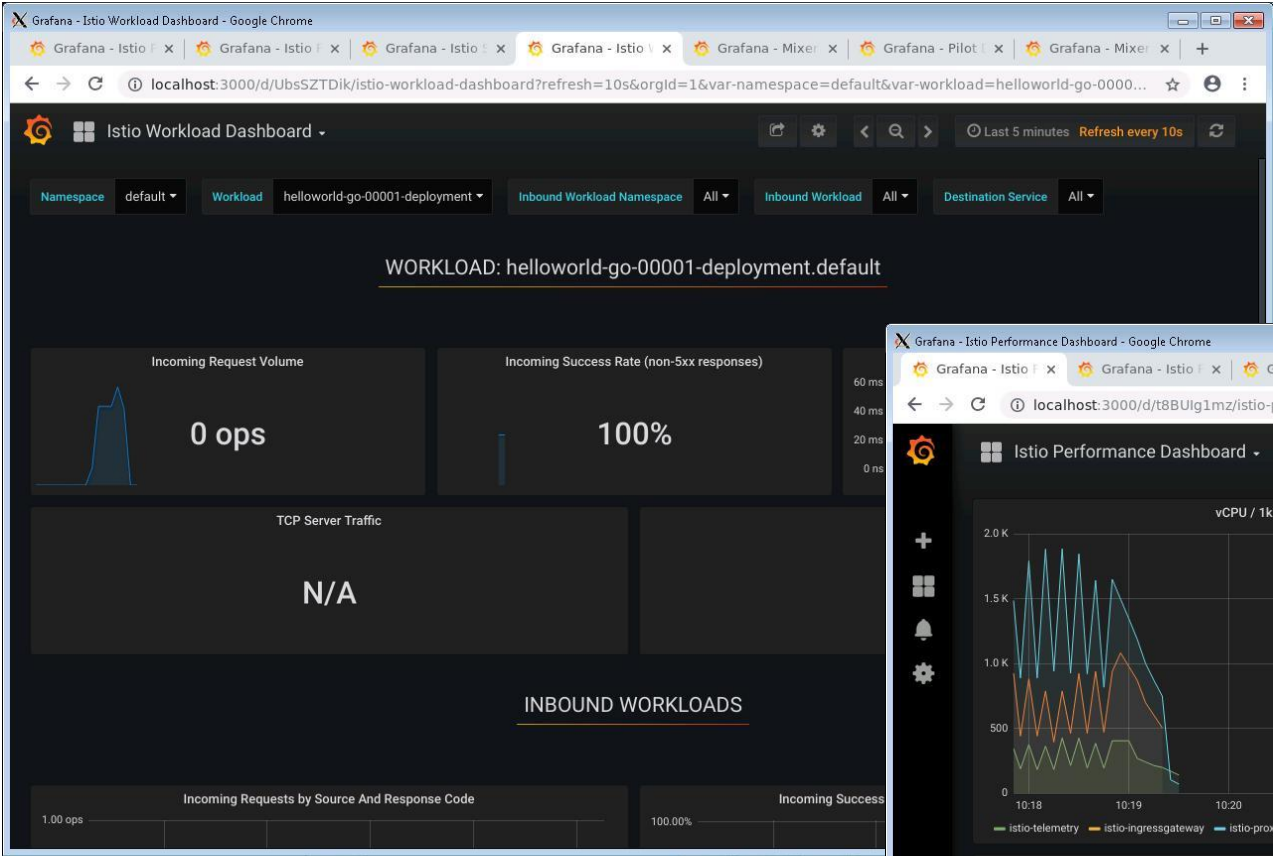Running and terminating pods as traffic decreases

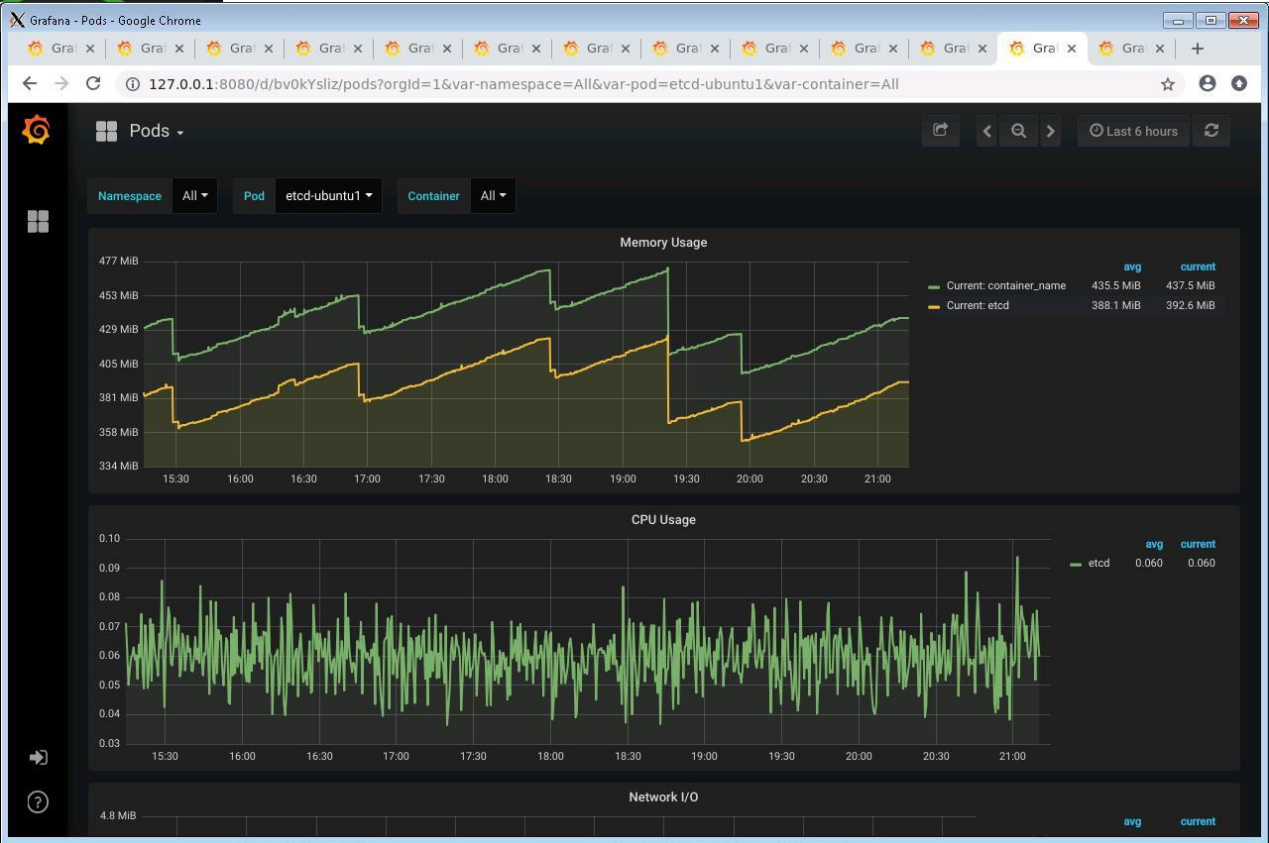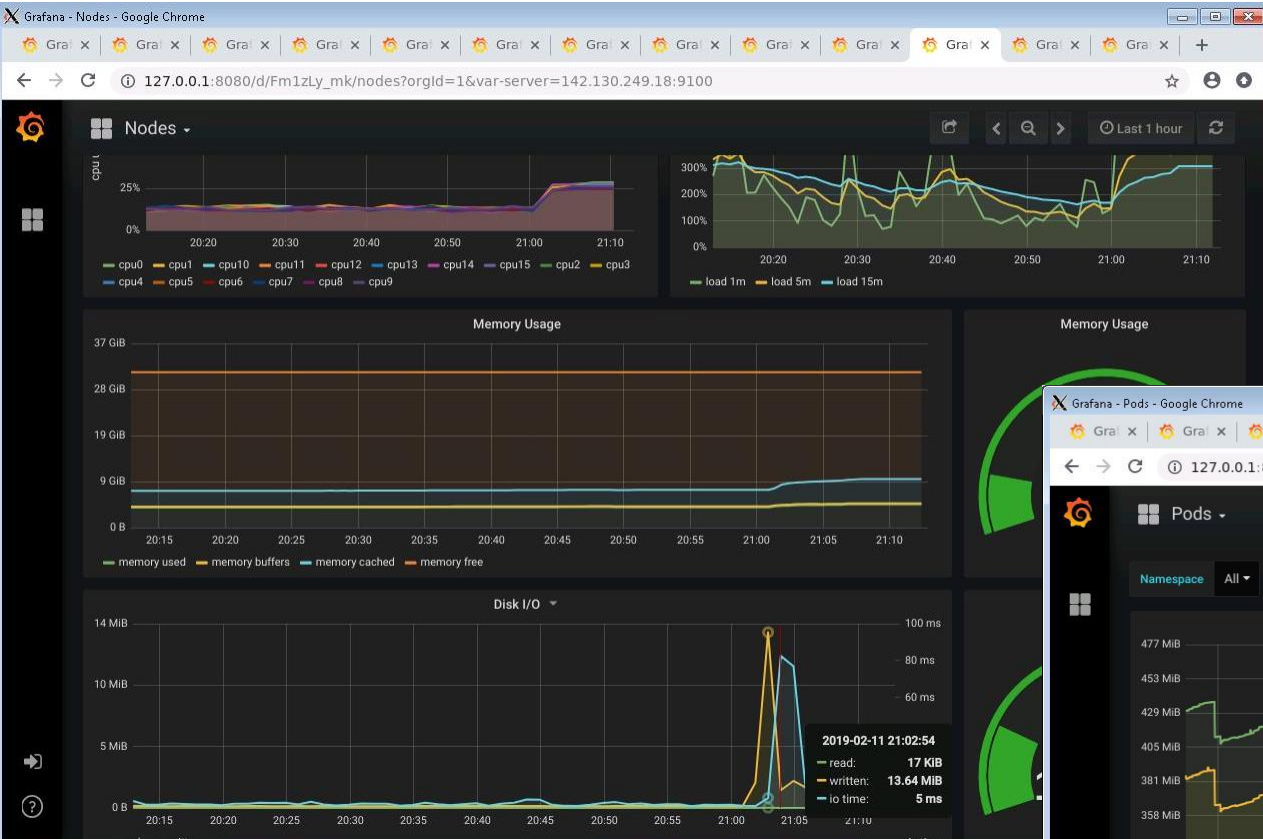Desired and observed pod levels as traffic decreases

Desired and observed concurrency



```
helloworld-go-00001-deployment-94667496b-xrzjm   2/3   Terminating   0   18m
```

# Knative – Performance Monitoring Grafana

# Knative – Performance Monitoring Grafana

# Knative – Demo

## $ kubectl get pods

Runs as required

```
helloworld-go-00001-deployment-94667496b-68x2q    0/3    Init:0/1            0    7s
helloworld-go-00001-deployment-94667496b-9p2cq    0/3    Init:0/1            0    5s
helloworld-go-00001-deployment-94667496b-9sff9    0/3    PodInitializing    0    7s
helloworld-go-00001-deployment-94667496b-jqbjp    0/3    PodInitializing    0    7s
helloworld-go-00001-deployment-94667496b-pb2zq    0/3    Init:0/1            0    7s


helloworld-go-00001-deployment-94667496b-68x2q    0/3    Init:0/1            0    11s
helloworld-go-00001-deployment-94667496b-9p2cq    0/3    PodInitializing    0    9s
helloworld-go-00001-deployment-94667496b-9sff9    2/3    Running            0    11s
helloworld-go-00001-deployment-94667496b-jqbjp    0/3    PodInitializing    0    11s
helloworld-go-00001-deployment-94667496b-pb2zq    0/3    Init:0/1            0    11s
```

Then self-terminates when no longer required

```
helloworld-go-00001-deployment-94667496b-68x2q    2/3    Terminating    0    95s
helloworld-go-00001-deployment-94667496b-9p2cq    2/3    Terminating    0    93s
helloworld-go-00001-deployment-94667496b-9sff9    2/3    Terminating    0    95s
helloworld-go-00001-deployment-94667496b-jqbjp    3/3    Running        0    95s
helloworld-go-00001-deployment-94667496b-pb2zq    2/3    Terminating    0    95s


helloworld-go-00001-deployment-94667496b-68x2q    2/3    Terminating    0    2m47s
helloworld-go-00001-deployment-94667496b-9p2cq    2/3    Terminating    0    2m45s
helloworld-go-00001-deployment-94667496b-9sff9    2/3    Terminating    0    2m47s
helloworld-go-00001-deployment-94667496b-jqbjp    2/3    Terminating    0    2m47s
helloworld-go-00001-deployment-94667496b-pb2zq    2/3    Terminating    0    2m47s
```

# Knative – Demo