



KubeCon



CloudNativeCon

— North America 2018 —

Performance Testing Ingress for Internet-Scale Workloads

Alexander Brand, Heptio



Kubernetes Ingress



- Organizations typically have a portfolio of applications that are offered to end users over the Internet
- Kubernetes supports multiple ways of exposing applications to the outside world
- Ingress is the answer if you are looking for layer 7 load balancing



heptio
Contour

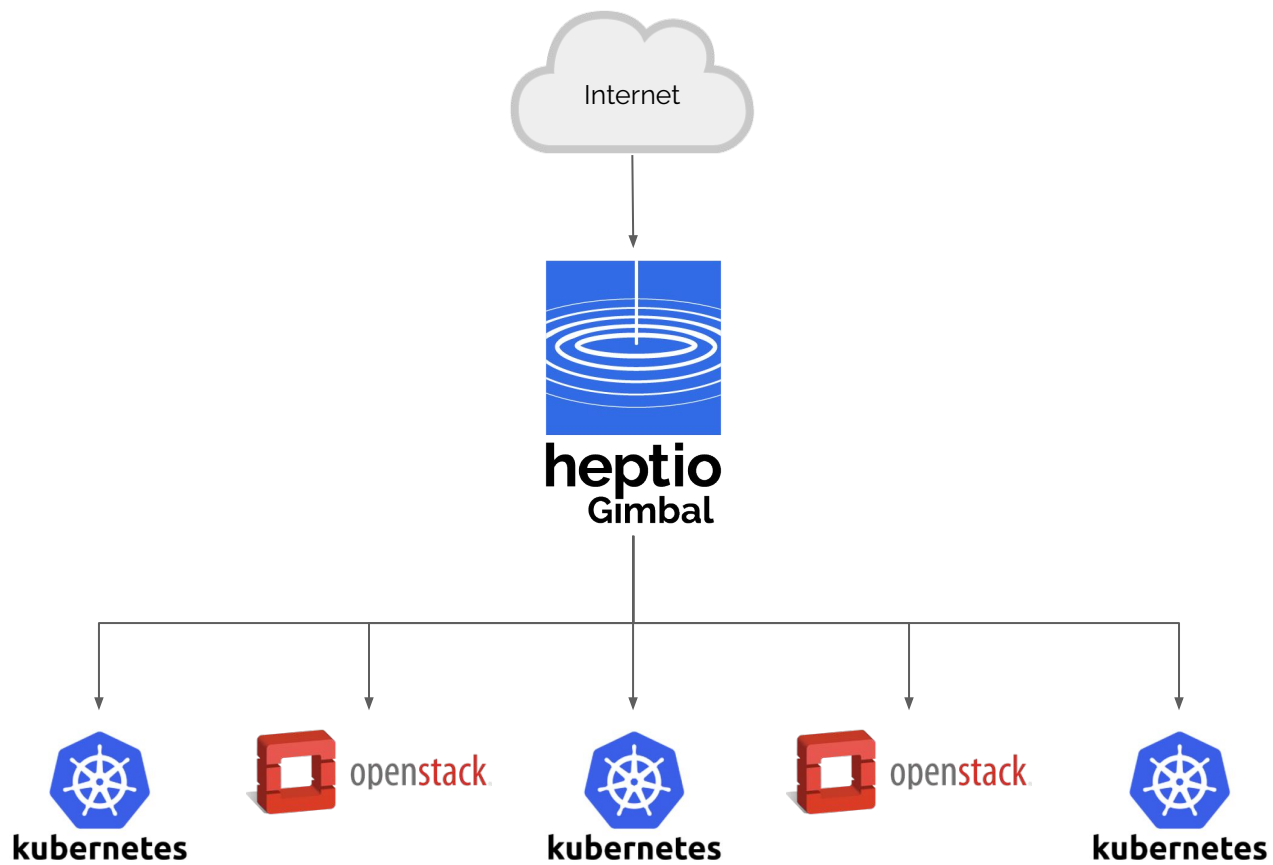


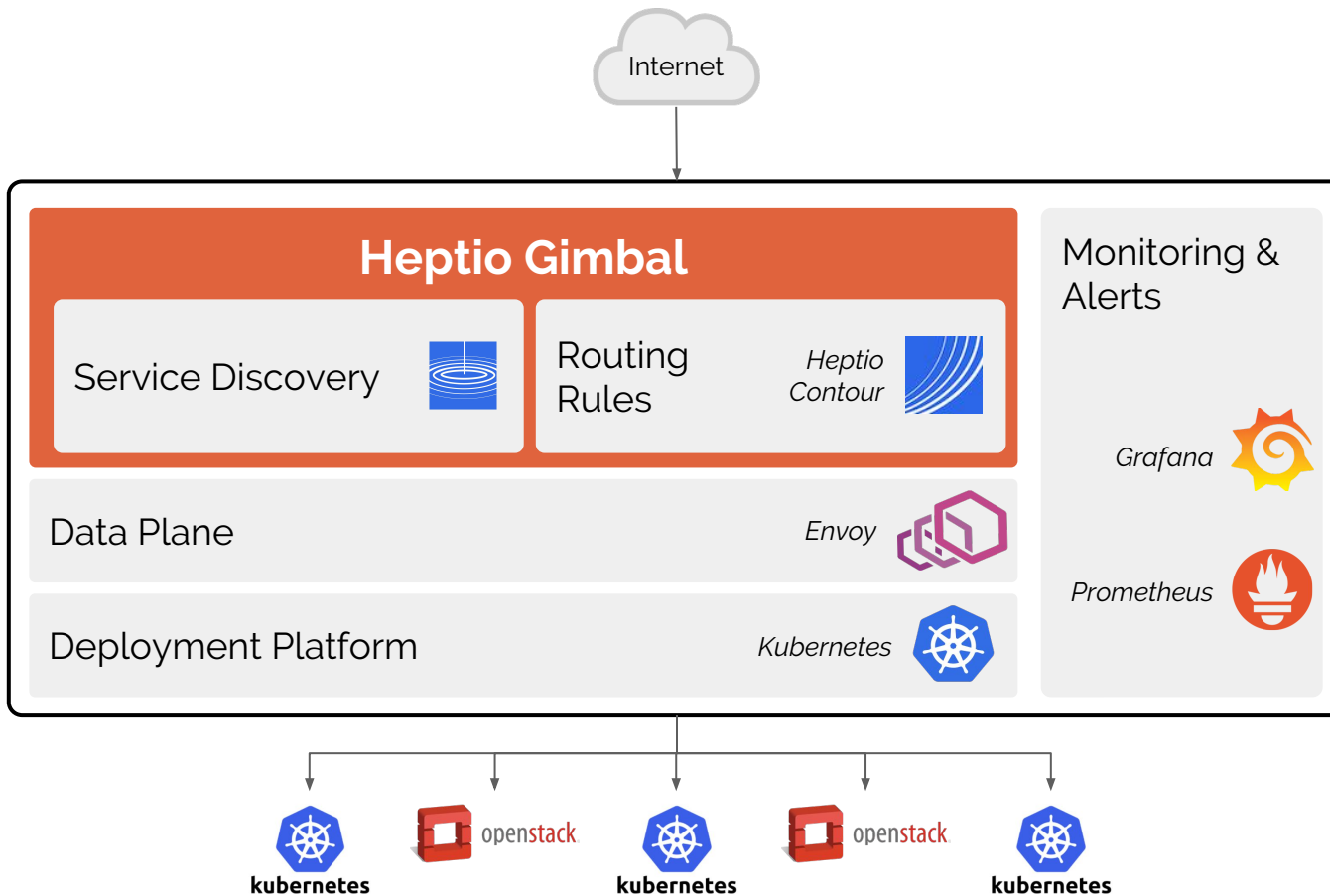
heptio
Gimbal

Why did we create Gimbal?



- Open sourced from co-development project with Yahoo Japan Corporation subsidiary, Actapio
- Operate hundreds of services that are exposed to the Internet
- Run multiple OpenStack and Kubernetes clusters





Let's talk performance

Requirements



- Millions of concurrent connections
- Thousands of services per datacenter
- Tens of thousands of endpoints per datacenter
- < 30ms P99 round-trip time latency

Multiple variables



- Requests per second
- Concurrent connections
- Response payload size
- Number of services
- Number of endpoints
- Number of Ingress / IngressRoutes
- Number of proxy (Envoy) pods

Multiple Subsystems



- Data plane
- Control plane
- Discovery system
- Monitoring system

Two pronged approach



Micro-benchmarks

Understand the impact of a single variable on a specific subsystem.

For example: “Understand the effect of number of concurrent connections on response latency”

Macro-benchmarks

Understand the impact of “realistic” load on the system (or subsystem)

For example: “Understand the effect of 100k CC and 30k RPS on the performance of the system”

Micro-benchmarks

Example: Concurrent connections vs Latency

Test Method: Adjust the number of concurrent connections in wrk2, and observe effect on latency

Variable Under Test: # of Concurrent Connections

Test Cases: 10k, 25k, 50k, 100k, 250k, 500k, 1m

Expected: < 30 ms P99 Latency

Micro-benchmarks



Pros

- Helped us identify bottlenecks/bugs in Contour and the discovery subsystem

panic in contour when there is a high rate of churn of IngressRoute resources #523

Closed alexbrand opened this issue on Jul 6 · 3 comments

alexbrand commented on Jul 6

While running performance tests for Gimbal, I noticed that Contour was crashing with an index out of range runtime error.

I was able to reproduce the issue using the internal/e2e test framework. The test can be found here: <https://github.com/alexbrand/contour/tree/rate-bug>

I have seen a couple of different stack traces, but they are all index out of range inside the route.pb.go file from envoyproxy/go-control-plane. This is one of them:

```
panic: runtime error: index out of range

goroutine 396 [running]:
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2/route.
  /Users/abrand/go/src/github.com/heptio/contour/vendor/github.com/envoyproxy/go-cont
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2. (#Rout
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2. (#Rout
github.com/heptio/contour/vendor/github.com/envoyproxy/go-control-plane/envoy/api/v2. (#Rout
github.com/heptio/contour/vendor/github.com/none/protobuf/protom.Marshal@x1cfa860...@xc42849
```

Assignees: davecheney

Labels: bug, p0 - Hair On Fire

Projects: None yet

Milestone: 0.6.0-beta.1

Notifications: Unsubscribe

3 participants

Lock conversation



Openstack Discoverer logs incorrect update #95

stevesloka opened this issue on May 8 · 3 comments

stevesloka commented on May 8

The openstack discoverer logs updates for services when no updates have happened. The following example show the lb being added, then shows an update on each cycle when it wasn't updated in openstack.

```
time="2018-05-08T13:10:17Z" level=info msg="Successfully handled: add service 'demo/lb1-2ad
time="2018-05-08T13:10:17Z" level=info msg="Successfully handled: add endpoints 'demo/lb1-2
time="2018-05-08T13:10:47Z" level=info msg="Successfully handled: update service 'demo/lb1-
time="2018-05-08T13:11:17Z" level=info msg="Successfully handled: update service 'demo/lb1-
time="2018-05-08T13:11:47Z" level=info msg="Successfully handled: update service 'demo/lb1-
time="2018-05-08T13:12:19Z" level=info msg="Successfully handled: update service 'demo/lb1-
time="2018-05-08T13:12:47Z" level=info msg="Successfully handled: update service 'demo/lb1-
time="2018-05-08T13:13:21Z" level=info msg="Successfully handled: update service 'demo/lb1-
time="2018-05-08T13:13:47Z" level=info msg="Successfully handled: update service 'demo/lb1-
time="2018-05-08T13:14:17Z" level=info msg="Successfully handled: update service 'demo/lb1-
```

stevesloka added the bug label on May 8

stevesloka self-assigned this on May 8

alexbrand commented on Jun 6 • edited

While running performance tests, this issue was very prominent. The discoverer is enqueueing updates at a faster rate than the workers are dequeuing items (Related to #140). This results in memory consumption and an ever growing queue.

Memory:

Assignees: alexbrand, stevesloka

Labels: bug, discoverer, p1 - Important

Projects: None yet

Milestone: v0.3-beta.1

Notifications: Unsubscribe

2 participants

Lock conversation

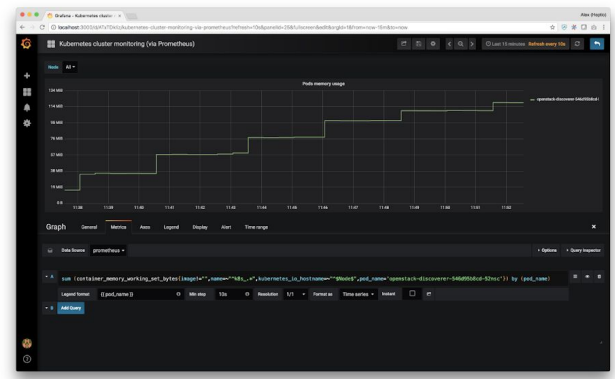




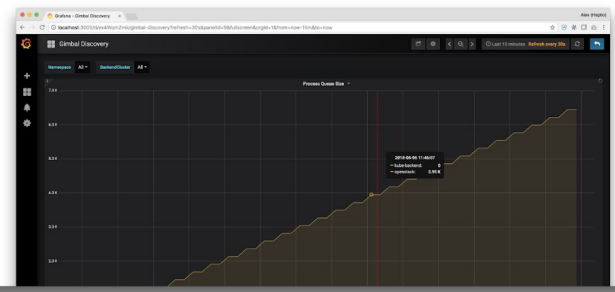
alexbrand commented on Jun 6 • edited ▾ Member + 👤 ⋮

While running performance tests, this issue was very prominent. The discoverer is enqueueing updates at a faster rate than the workers are dequeuing items (Related to #140). This results in memory consumption and an ever growing queue.

Memory:



Queue:



You're receiving notifications because you were assigned.

2 participants

- 🔒 Lock conversation
- ➡ Transfer this issue Beta

Micro-benchmarks



Pros

- Helped us identify bottlenecks/bugs in Contour and the discovery subsystem
- Gave us confidence that the control plane (Contour) could handle a large number of Services, Endpoints, Ingress and IngressRoute

Micro-benchmarks



Pros

- Helped us identify bottlenecks/bugs in Contour and the discovery subsystem
- Gave us confidence that the control plane (Contour) could handle a large number of Services, Endpoints, Ingress and IngressRoute
- Time is your friend - Less time needed to setup and run

Micro-benchmarks



Cons

- Evaluates the system through a narrow lens

Micro-benchmarks



Cons

- Evaluates the system through a narrow lens
- Doesn't necessarily reflect real world usage

Macro-benchmarks

Macro-benchmarks



- Test the system under “realistic” load
- Measure and evaluate multiple metrics
- Gives you an idea of where the bottlenecks are, and how the system should be scaled to handle more load
- Depending on hardware availability, budget, etc, you might have to scale the test down

Our approach

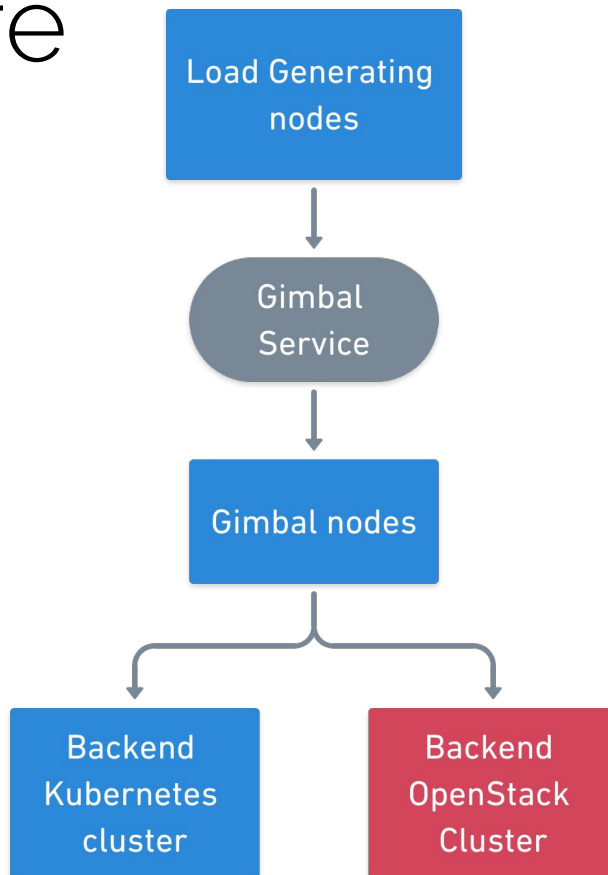


Run tests at three different scales and ensure resource utilization scales linearly

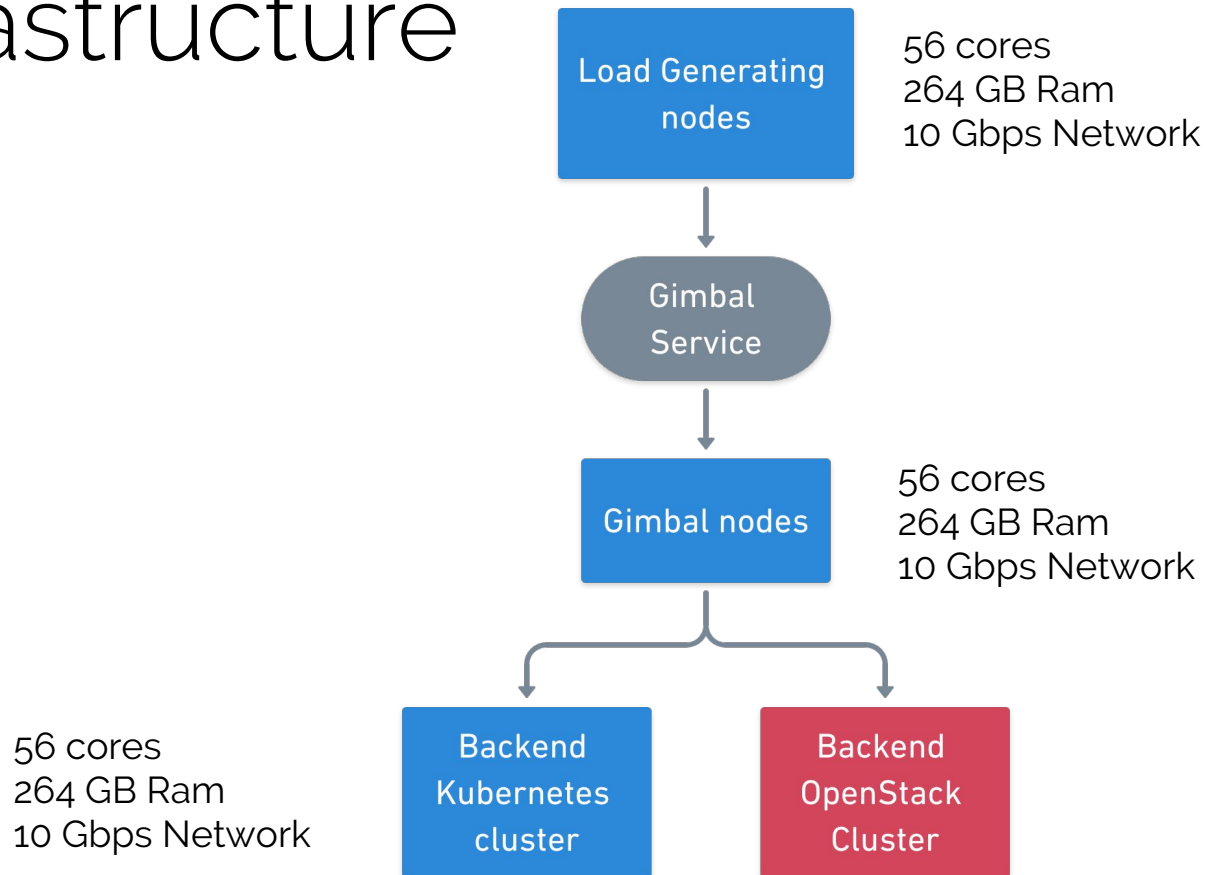
100k CC & 10k RPS	200k CC & 10k RPS	300k CC & 10k RPS
100k CC & 20k RPS	200k CC & 20k RPS	300k CC & 20k RPS
100k CC & 30k RPS	200k CC & 30k RPS	300k CC & 30k RPS

CC = concurrent connections
RPS = requests per second

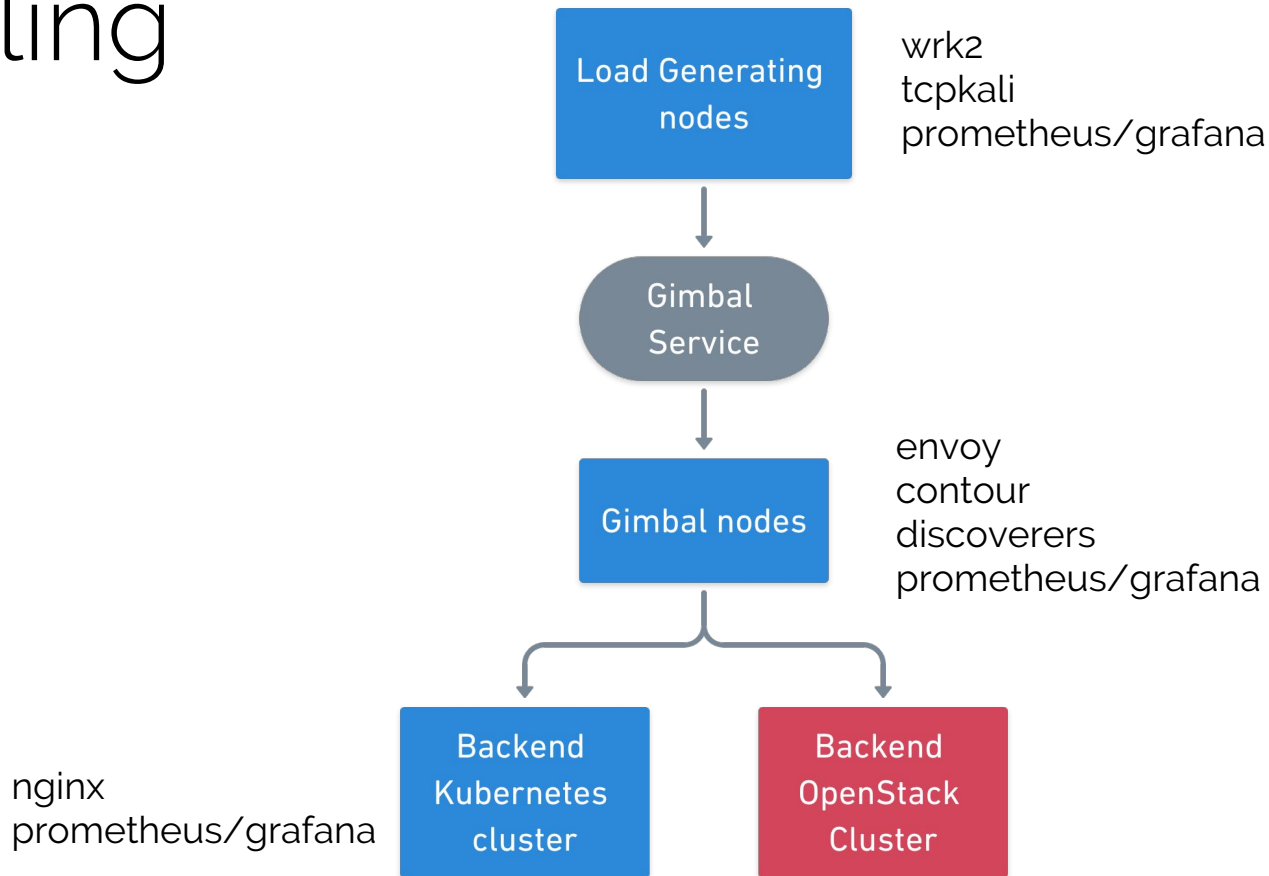
Infrastructure



Infrastructure



Tooling



wrk2



- HTTP benchmarking tool with accurate latency measurements
- Can generate significant load from a single, multi-core machine
- Deployed as a Kubernetes Job
- <https://github.com/giltene/wrk2>

abc apiVersion

```
1  apiVersion: batch/v1
2  kind: Job
3  metadata:
4    labels:
5      workload: wrk2
6    generateName: wrk2-
7  spec:
8    backoffLimit: 6
9    completions: 1
10   parallelism: 1
11   template:
12     metadata:
13       labels:
14         workload: wrk2
15     spec:
16       affinity:
17         podAntiAffinity:
18           preferredDuringSchedulingIgnoredDuringExecution:
19             - weight: 1
20               podAffinityTerm:
21                 labelSelector:
22                   matchExpressions:
23                     - key: workload
24                       operator: In
25                   values:
```

```
{}
```

```
7 spec:
8   backoffLimit: 6
9   completions: 1
10  parallelism: 1
11  template:
12    metadata:
13      labels:
14        workload: wrk2
15    spec:
16      affinity:
17        podAntiAffinity:
18          preferredDuringSchedulingIgnoredDuringExecution:
19            - weight: 1
20              podAffinityTerm:
21                labelSelector:
22                  matchExpressions:
23                    - key: workload
24                      operator: In
25                      values:
26                        - wrk2
27                topologyKey: kubernetes.io/hostname
28      initContainers:
29        - command:
30          - sh
31          - -c
```

```
{ } spec
20     podAffinityTerm:
21       labelSelector:
22         matchExpressions:
23           - key: workload
24             operator: In
25             values:
26               - wrk2
27       topologyKey: kubernetes.io/hostname
28     initContainers:
29       - command:
30         - sh
31         - -c
32         - sysctl -w net.ipv4.ip_local_port_range="1024 65535"
33         image: alpine:3.6
34         imagePullPolicy: IfNotPresent
35         name: sysctl-set
36         securityContext:
37           privileged: true
38     containers:
39       - command:
40         - wrk
41         - --threads
42         - "4"
43         - --latency
44         - --duration
```

```
{}
```

 spec

```
38     containers:
39     - command:
40       - wrk
41       - --threads
42       - "4"
43       - --latency
44       - --duration
45       - "300"
46       - --connections
47       - "600"
48       - --rate
49       - "30000"
50       - --header
51       - "Host: example.com"
52       - http://envoy.gimbal-contour.svc.cluster.local
53     image: bootjp/wrk2
54     imagePullPolicy: Always
55     name: wrk2
56   nodeSelector:
57     workload: wrk
58   restartPolicy: OnFailure
```


tcpkali



- TCP load generator used to open thousands of connections
- Opens connections and keeps them open
- Much better at opening connections than wrk2
- Deployed as a Kubernetes Job
- <https://github.com/satori-com/tcpkali>

```
{ } spec ▸ { } template ▸ { } spec ▸ abc restartPolicy
```

```
1  apiVersion: batch/v1
2  kind: Job
3  metadata:
4    labels:
5      workload: tcpkali
6    generateName: tcpkali-
7  spec:
8    completions: 6
9    parallelism: 6
10   template:
11     metadata:
12       labels:
13         workload: tcpkali
14     spec:
15       affinity:
16         podAntiAffinity:
17           preferredDuringSchedulingIgnoredDuringExecution:
18             - weight: 1
19               podAffinityTerm:
20                 labelSelector:
21                   matchExpressions:
22                     - key: workload
23                       operator: In
24                 values:
25                   - tcpkali
```

```
{ } spec ▸ { } template ▸ { } spec ▸ abc restartPolicy
```

```
29     - sh
30     - -c
31     - sysctl -w net.core.somaxconn=65535; sysctl -w net.ipv4.ip_local_port_range="1024 65535"
32     image: alpine:3.6
33     imagePullPolicy: IfNotPresent
34     name: sysctl-set
35     securityContext:
36       privileged: true
37     containers:
38     - args:
39       - tcpkali
40       - --connections
41       - "50000"
42       - --duration
43       - "420"
44       - --connect-rate
45       - "2000"
46       - envoy.gimbal-contour.svc.cluster.local:80
47     image: jojiisacth/tcpkali
48     imagePullPolicy: Always
49     name: tcpkali
50 restartPolicy: Never
```

Nginx



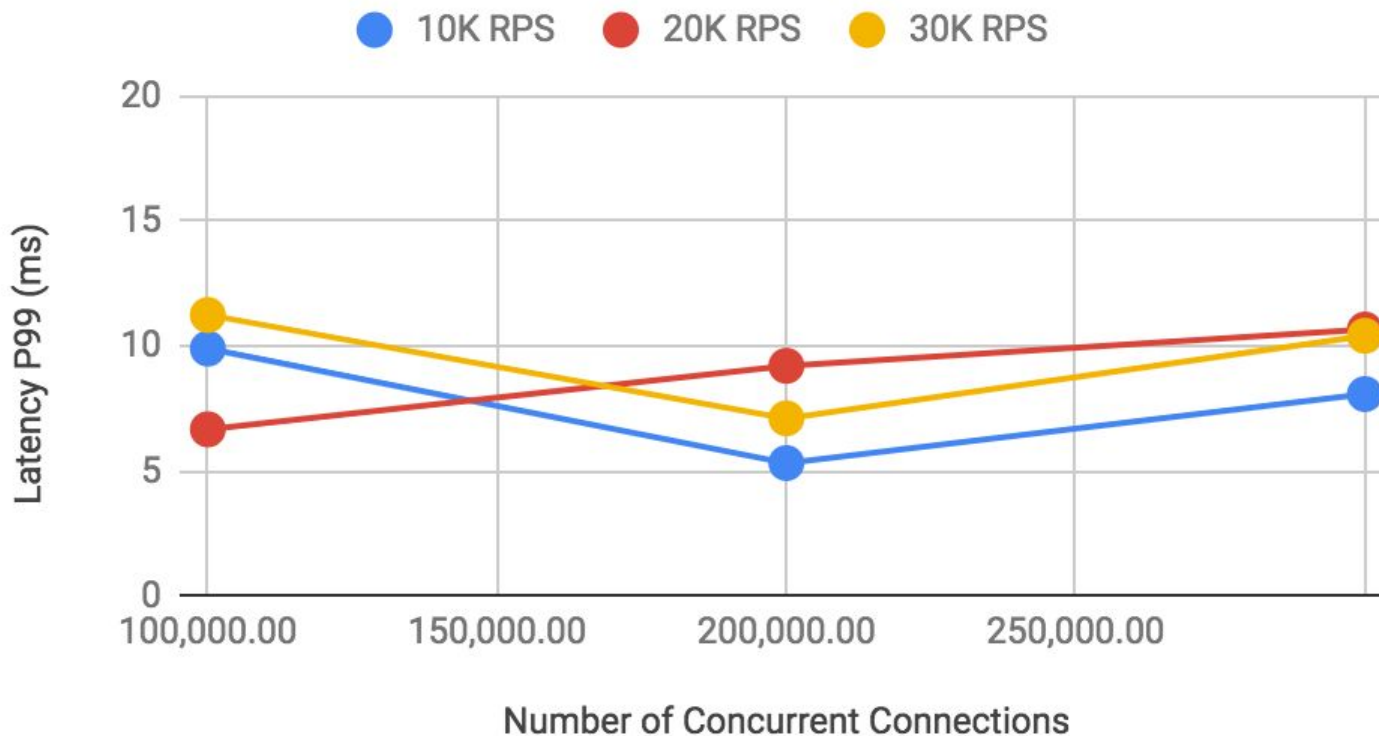
- Used nginx as the target service
- Ran as a Kubernetes Deployment
- Ran tests against two variants: "vanilla" (600 bytes) and custom (22 kilobytes)
- Default configuration was inadequate for our load test

```
{ } spec ▶ abc type
```

```
1  apiVersion: v1
2  data:
3    nginx.conf: |2
4      user root;
5      worker_processes auto;
6      worker_rlimit_nofile 262144;
7
8      error_log /var/log/nginx/error.log warn;
9      pid /var/run/nginx.pid;
10
11
12     events {
13       use epoll;
14       worker_connections 65536;
15       multi_accept on;
16     }
17
18
19     http {
20       include /etc/nginx/mime.types;
21       default_type application/octet-stream;
22
23       log_format main '$remote_addr - $remote_user [$time_local] "$request" '
24                       '$status $body_bytes_sent "$http_referer" '
25                       '"$http_user_agent" "$http_x_forwarded_for";
```

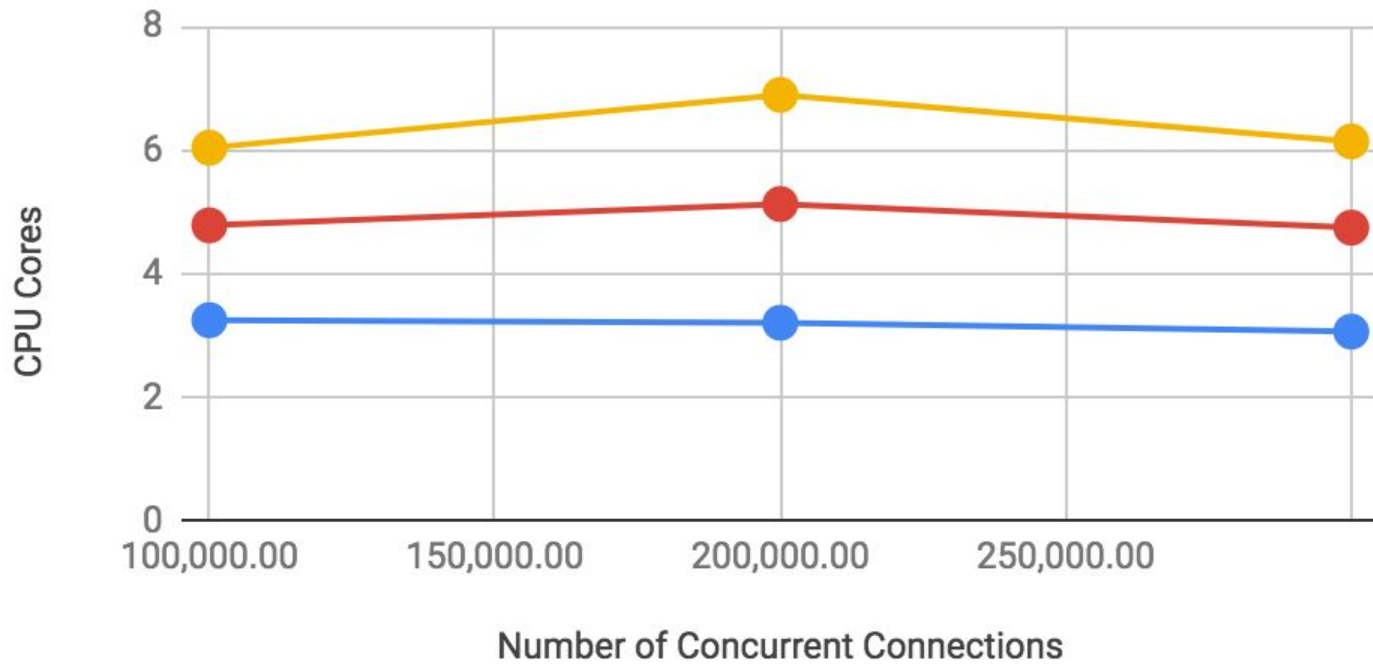
Results

Gimbal: Impact of Concurrency & RPS on Latency

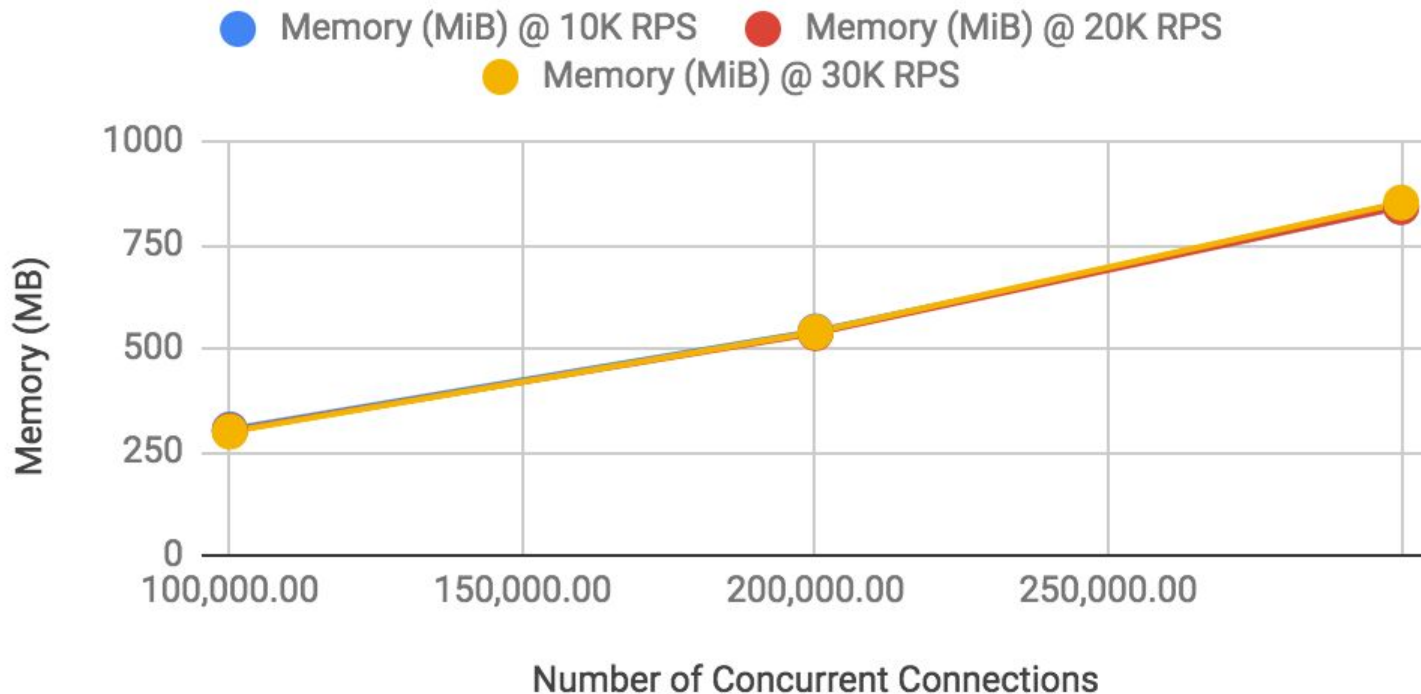


Gimbal: Impact of Concurrency & RPS on CPU

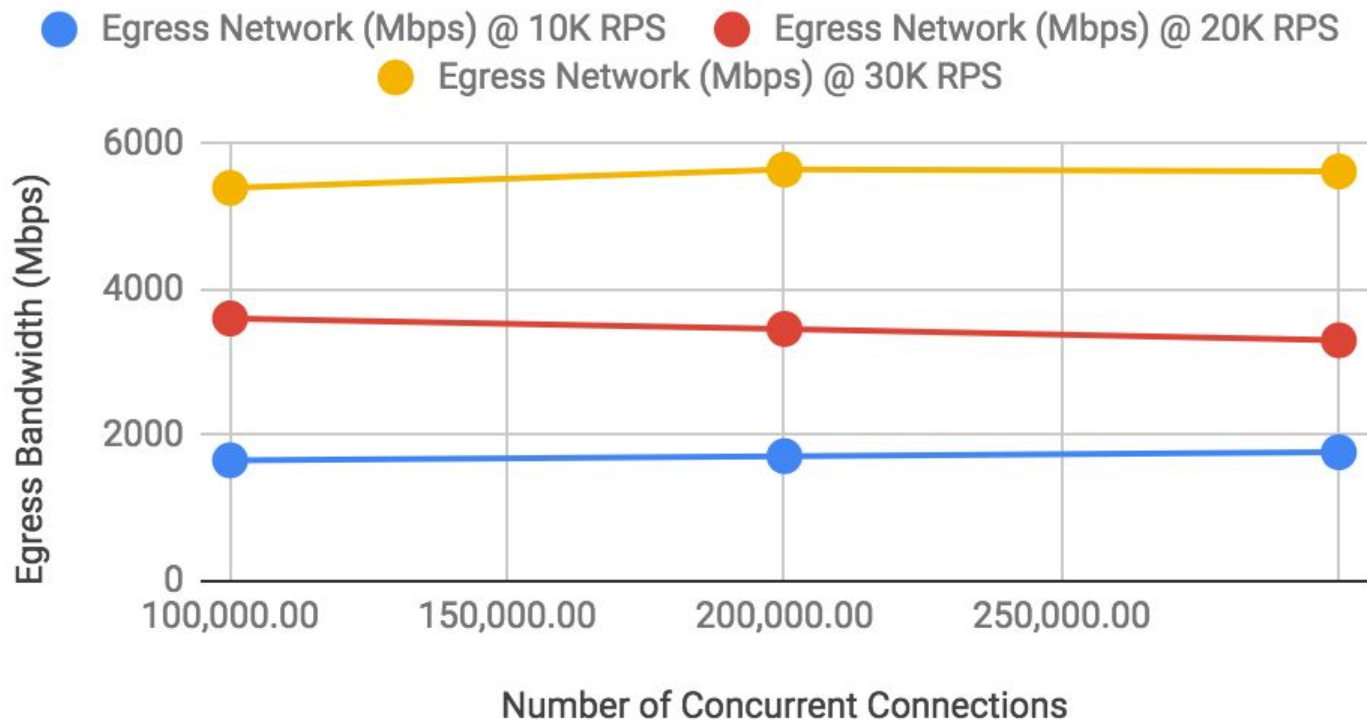
● CPU Cores @ 10K RPS ● CPU Cores @ 20K RPS ● CPU Cores @ 30K RPS



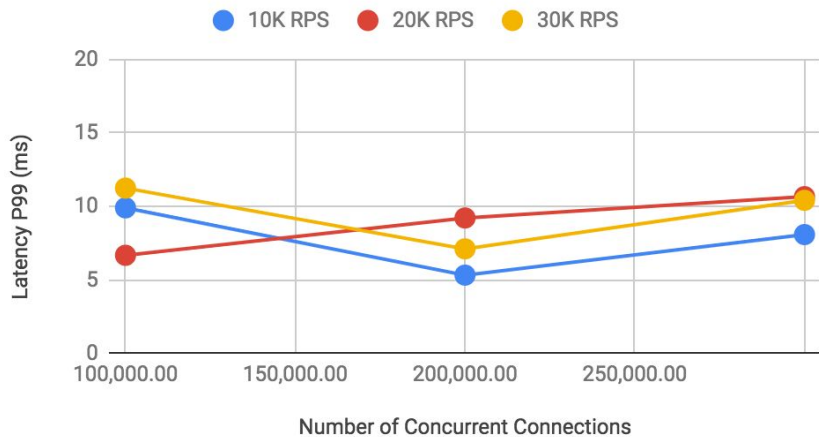
Gimbal: Impact of Concurrency & RPS on Memory



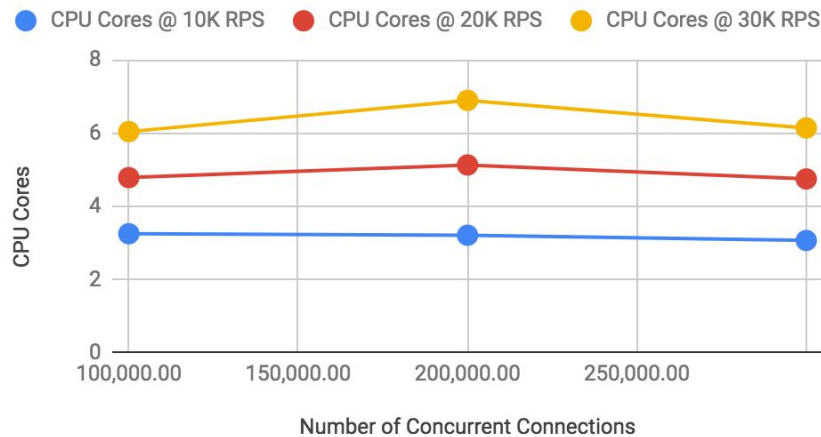
Gimbal: Impact of Concurrency & RPS on Network



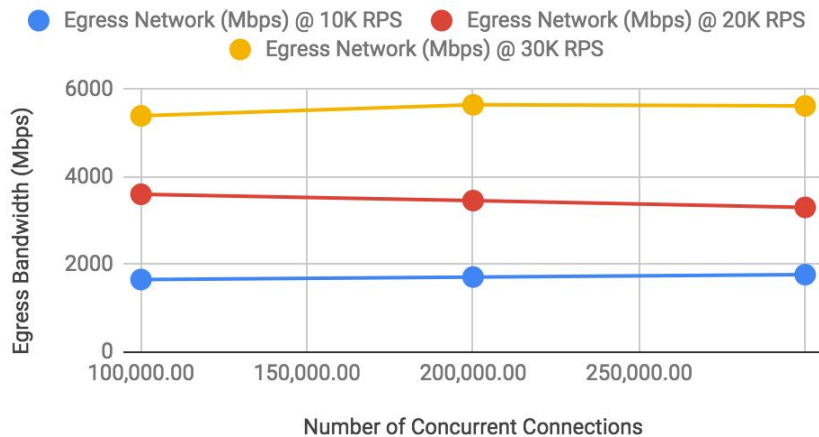
Gimbal: Impact of Concurrency & RPS on Latency



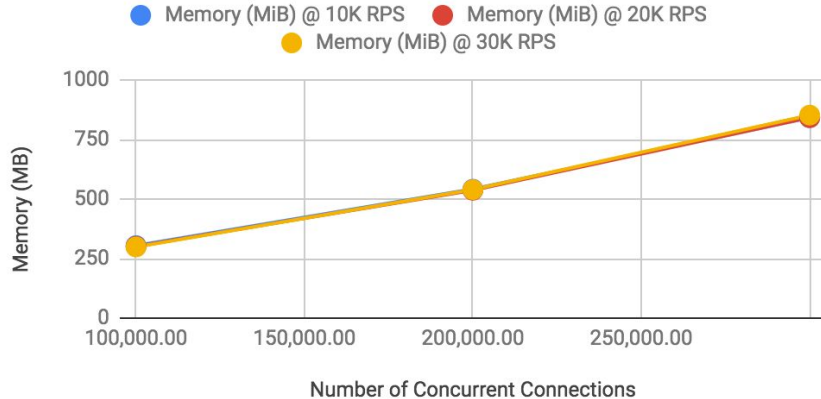
Gimbal: Impact of Concurrency & RPS on CPU

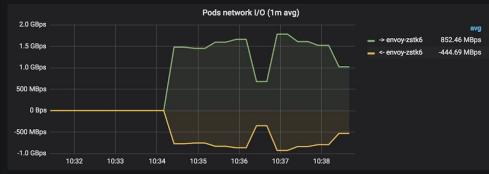
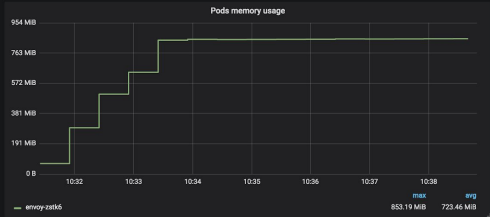
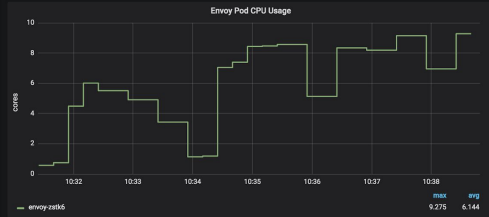
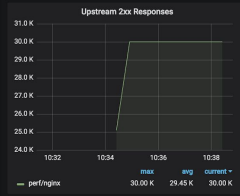
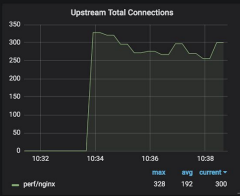
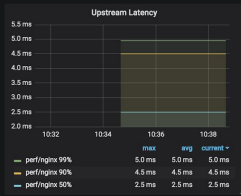
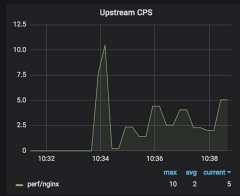
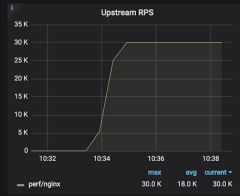
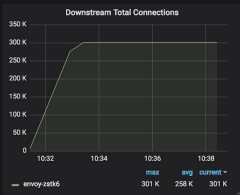
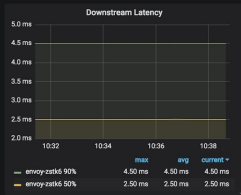
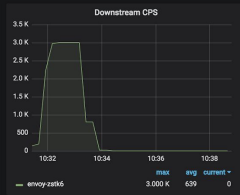
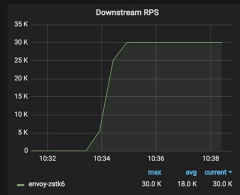


Gimbal: Impact of Concurrency & RPS on Network



Gimbal: Impact of Concurrency & RPS on Memory

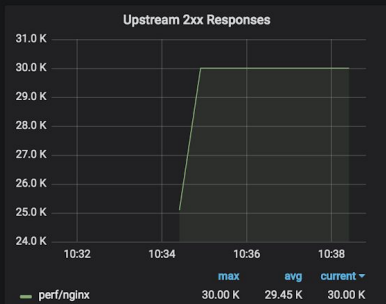
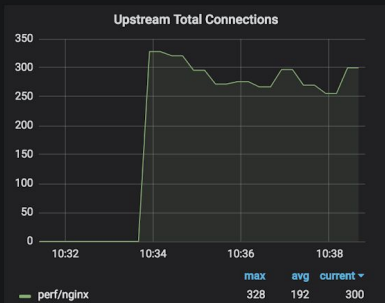
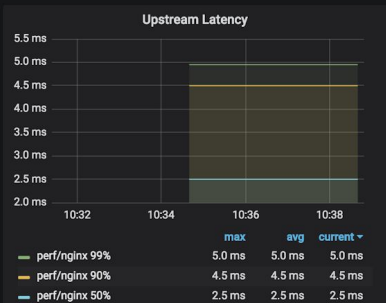
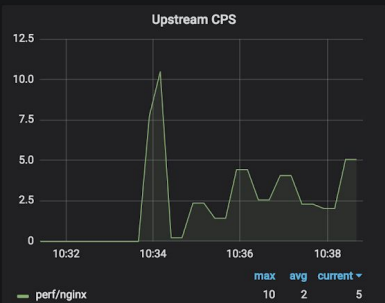
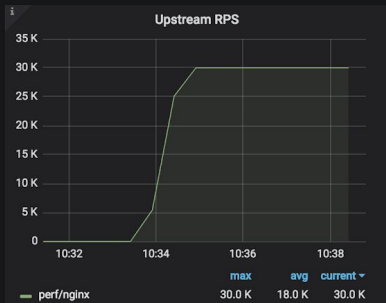
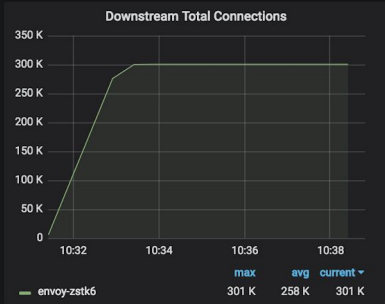
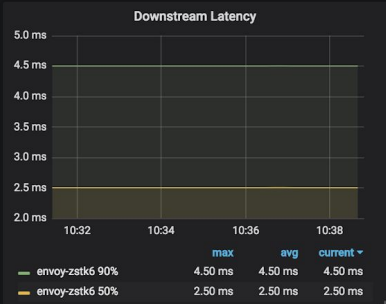
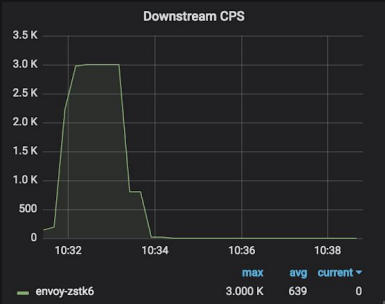
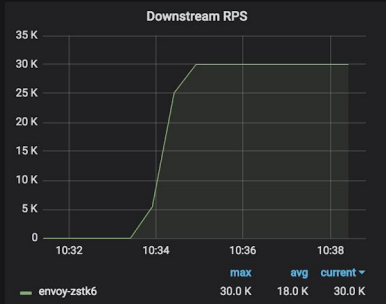


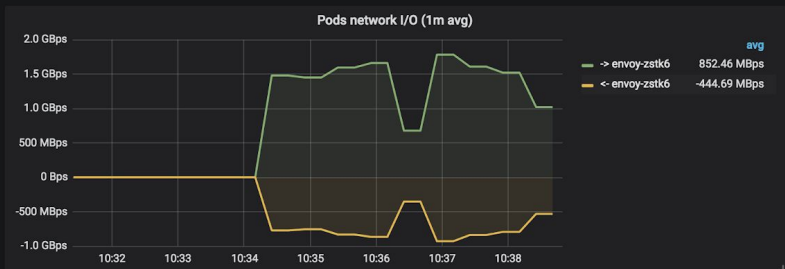
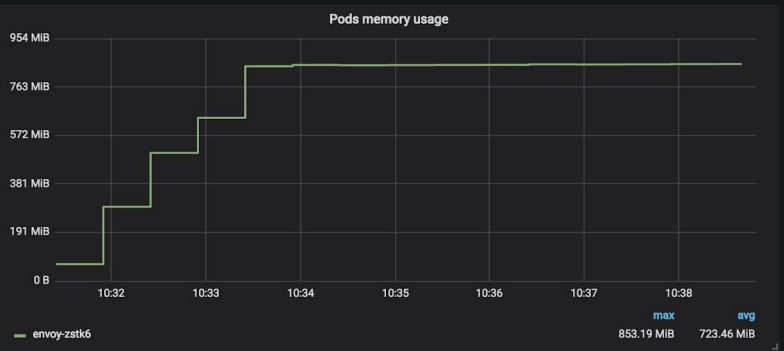
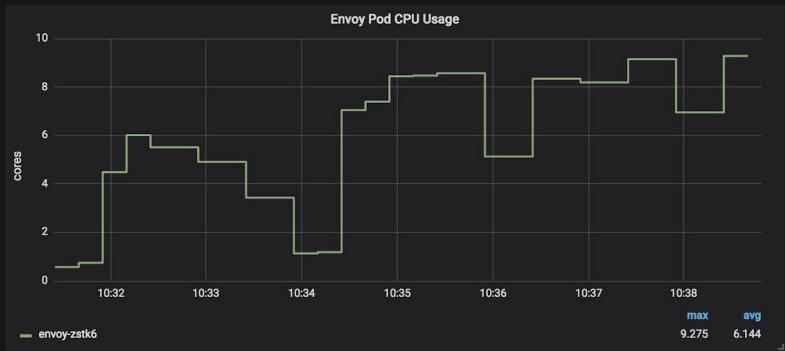
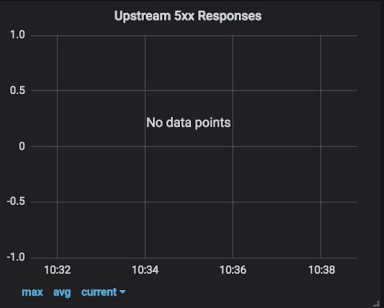
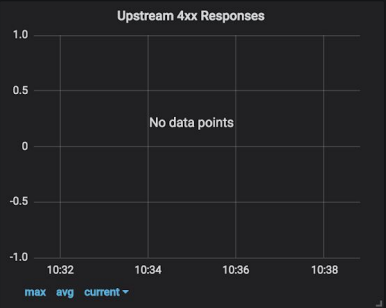
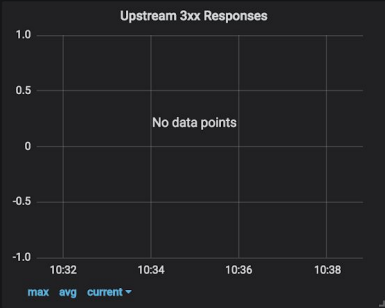
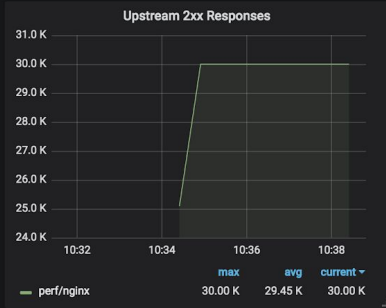
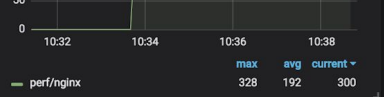
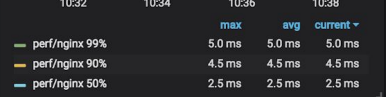
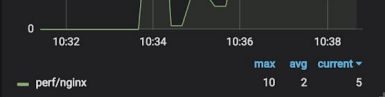
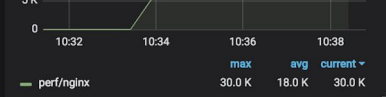




Namespace All

Service All





Lessons Learned

Document everything



- Create a plan that outlines what and how you are going to test
- Create a results table (or document) to capture the numbers you care about
- Document the environment's characteristics and specifications
- Keep a journal or scratchpad while you are running tests

Observability is paramount



- Prometheus and Grafana proved to be indispensable
- Envoy, Contour, Ginkgo discoverers all produce useful metrics
- Node-level visibility via Prometheus node_exporter
- Create test-specific dashboards
- Metrics are key to understand the system under test
- Don't fly blind

Short tests can be deceptive



- Test should be in the order of minutes instead of seconds
- Allows all components in the test path to warm up
- Reduces network jitter over a long measurement period
- Prometheus can obtain a larger set of data points

Check the network



- Understand your network's capacity before running any tests
- Keep the capacity in mind when designing test cases
- The network will limit your testing if the pipes are overfilled
- For example, we observed very different performance when running on AWS vs bare-metal lab
- Use `iperf3` to measure network bandwidth

Tweak the kernel



- The kernel can get in your way
- System and kernel logs can be helpful
- Can use init containers in Kubernetes

Timebox rabbit holes



- Weird things will happen at scale
- Some might be one-offs, some might be actual issues
- Take a note of what happened
- Investigate, but make sure to set a timebox on it

Understand first, automate later

- Resist temptation to automate everything from the get-go
- The test plan or strategy might change along the way
- Once the strategy is solid, document it
- Creating end-to-end automation might not be worth it



Happy testing!



heptio/gimbal | heptio/contour



@alexbrand