

# SIG Instrumentation

Best practices for cluster observability through metrics and logging across all Kubernetes components

# Optimizing Metric Rendering in *kube-state-metrics*

@mxinden



**Max Inden**

IndenML@gmail.com  
@mxinden

Red Hat

# Optimizing Metric Rendering in *kube-state-metrics*

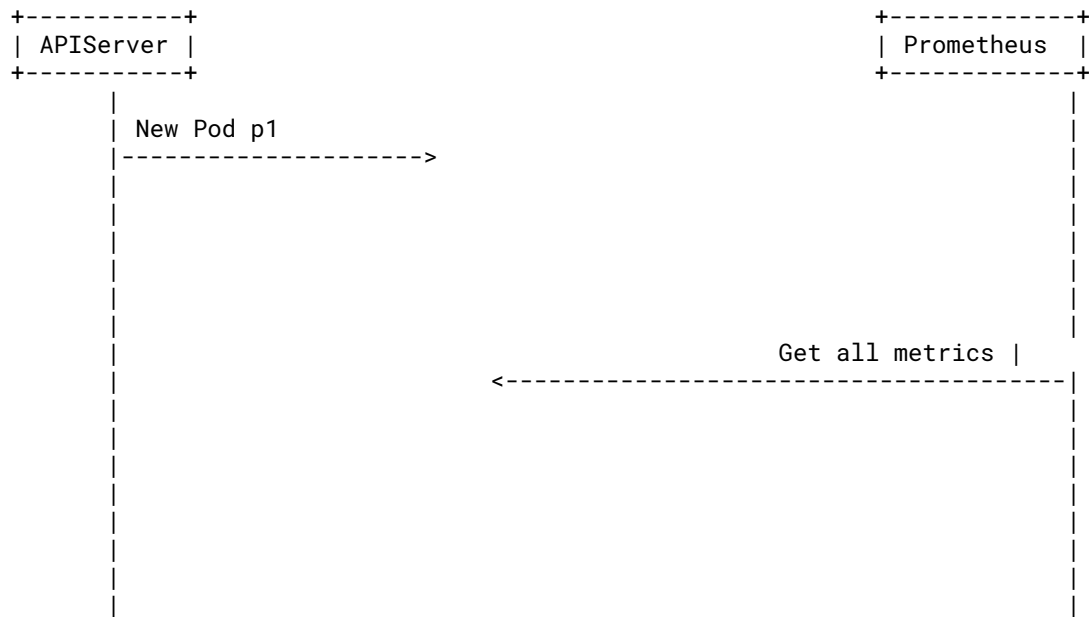
@mxinden

# Performance Optimizing with *Metrics* on *Kubernetes*

@mxinden

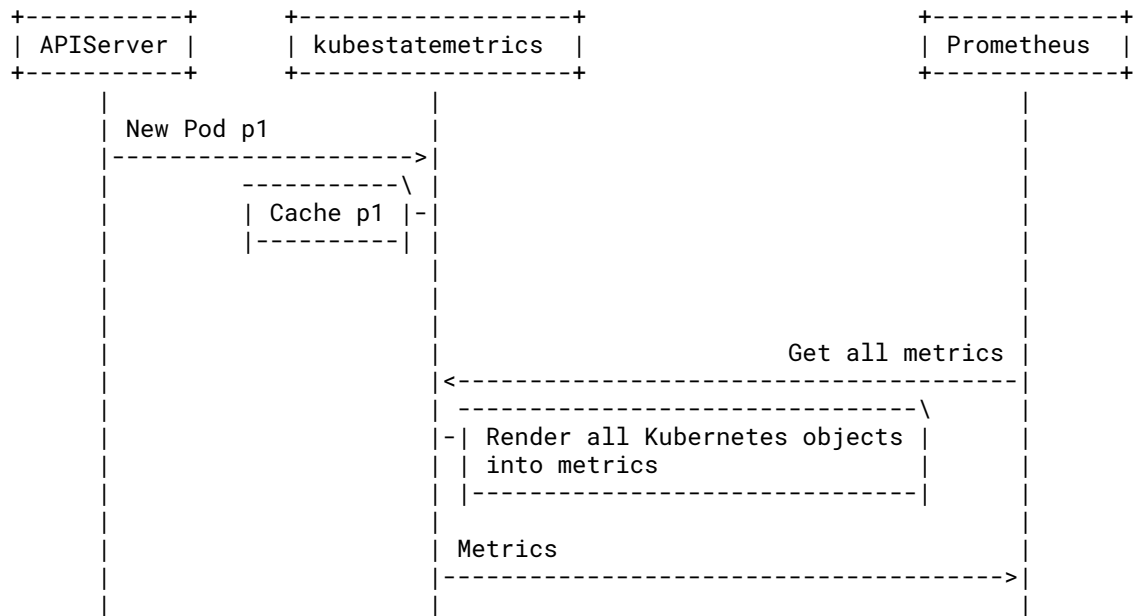
# What is kube-state-metrics? Old version.

Exposes the state of a Kubernetes cluster in Prometheus metrics.



# What is kube-state-metrics? Old version.

Exposes the state of a Kubernetes cluster in Prometheus metrics.



## Kubernetes Object:

```
- apiVersion: v1
  kind: Pod
  metadata:
    labels:
      app: kube-state-metrics
      pod-template-hash: 5fc64f676f
  name: kube-state-metrics-5fc64f676f-gl6v6
  namespace: monitoring
```

## Prometheus Metric:

```
kube_pod_container_info{container="kube-state-metrics",namespace="monitoring",pod="kube-state-metrics-5fc64f676f-gl6v6"} 1
```

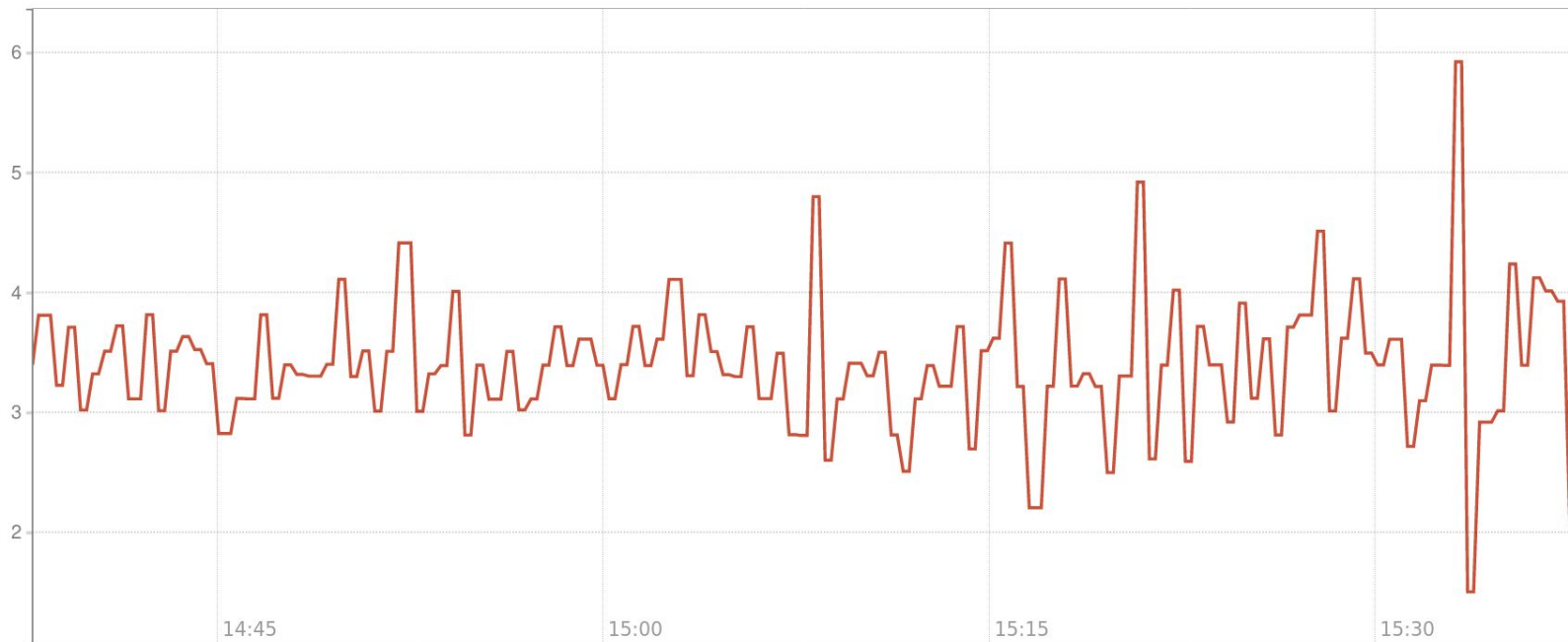
```
kube_pod_labels{label_app="kube-state-metrics",label_pod_template_hash="5fc64f676f",namespace="monitoring",pod="kube-state-metrics-5fc64f676f-gl6v6"} 1
```



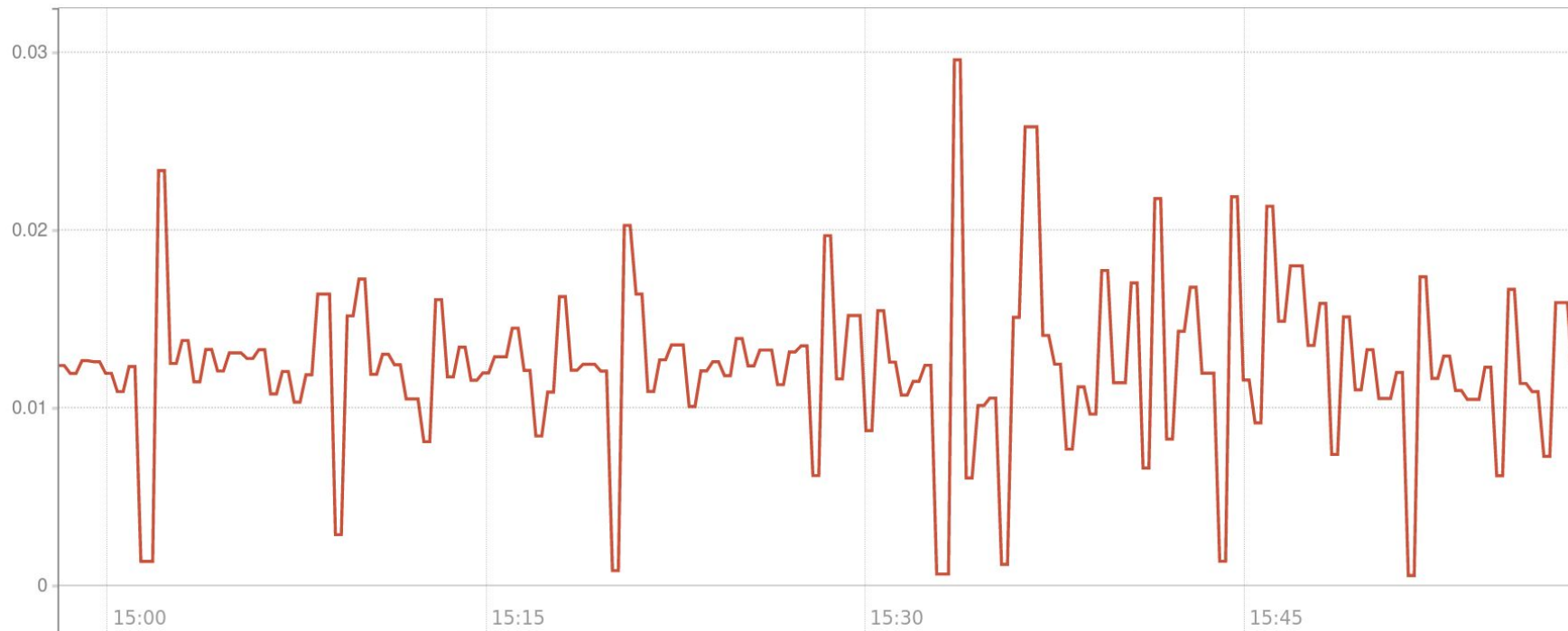
# Problem

- High response times
  - 10s - 20s on big production clusters with ~50 mb of metrics
- High & unstable resource usage
  - Difficult to predict resource limits
  - ...

# scrape\_duration\_seconds



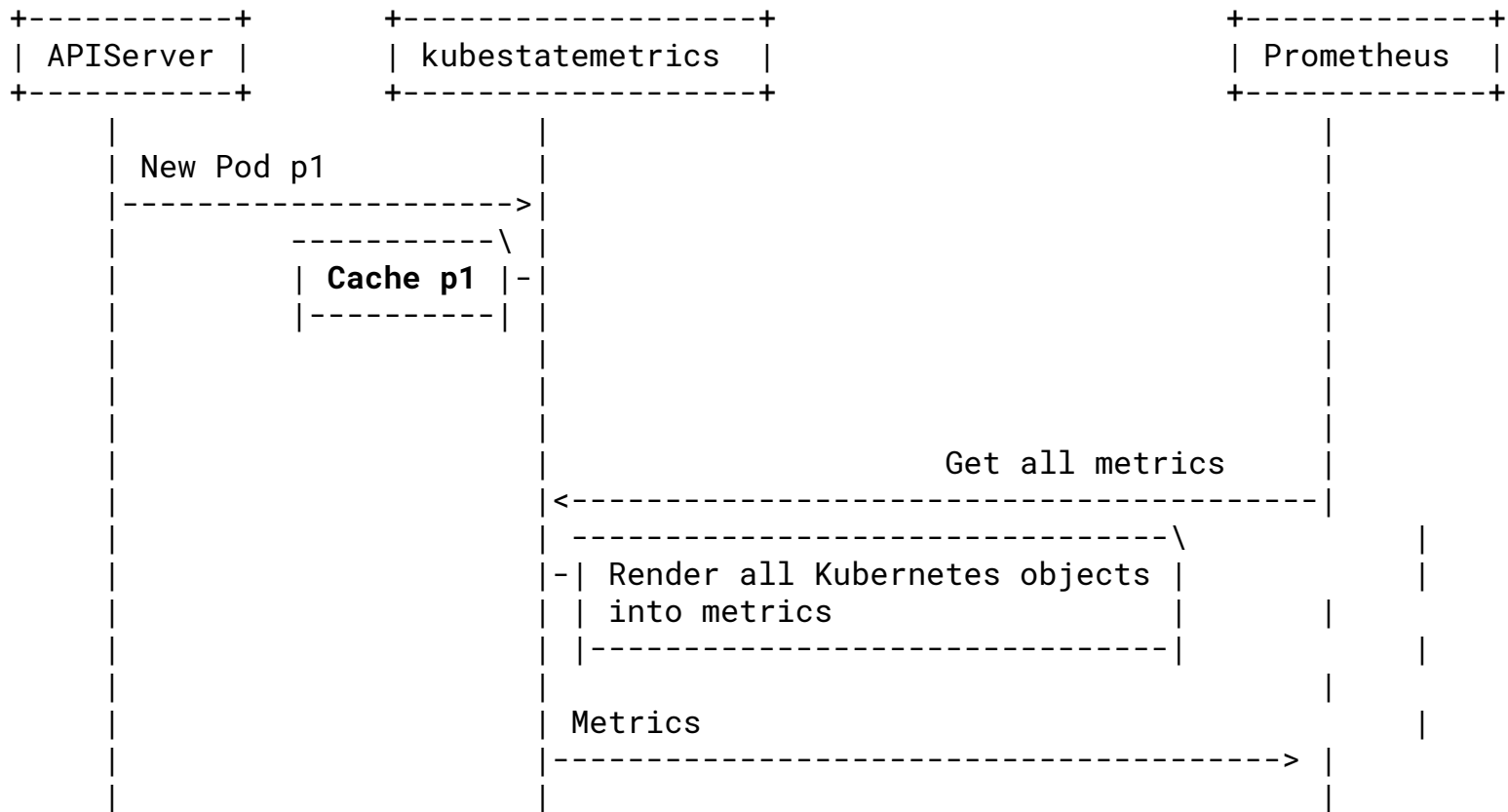
# container\_cpu\_usage\_seconds



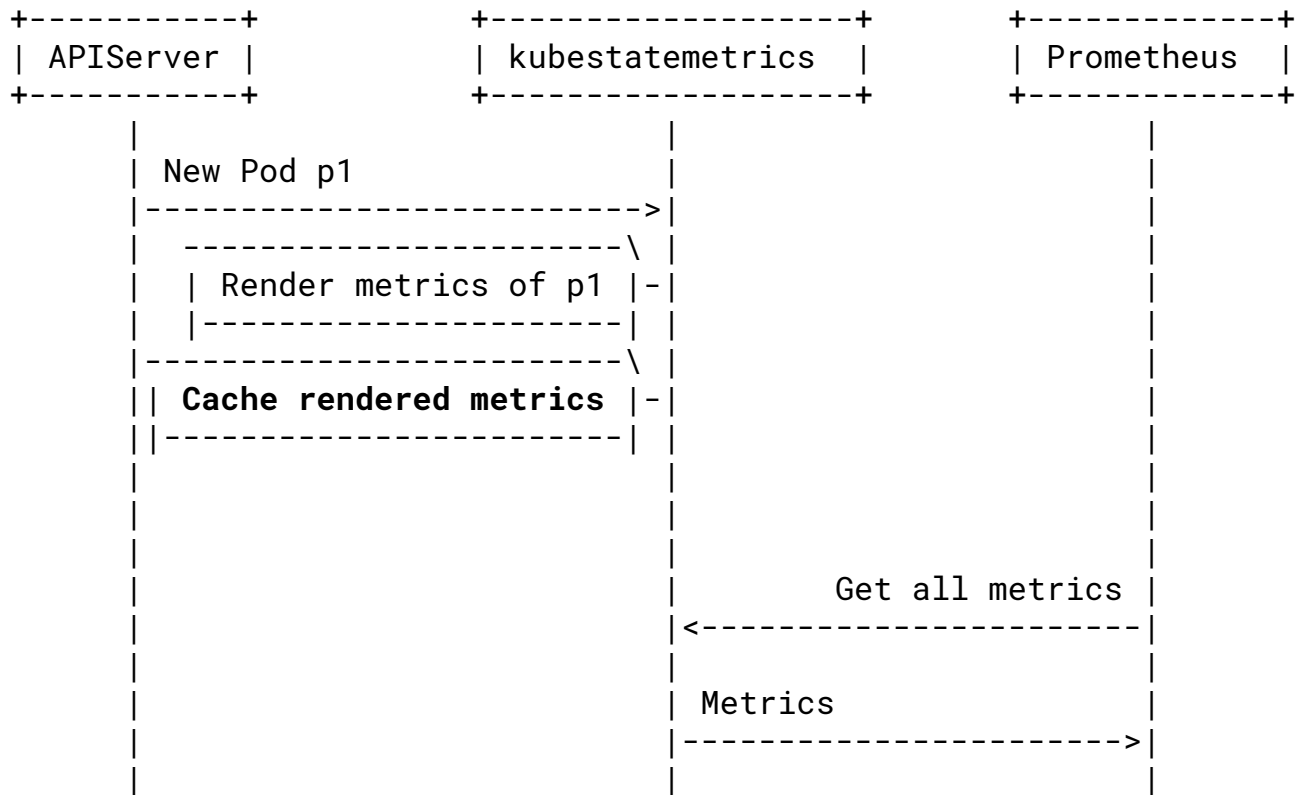
# container\_memory\_usage\_bytes



# Caching



# Caching



# What to keep in cache?

## Old: Kubernetes Object

```
- apiVersion: v1
  kind: Pod
  metadata:
    labels:
      app: kube-state-metrics
      pod-template-hash: 5fc64f676f
      name: kube-state-metrics-5fc64f676f-g16v6
      namespace: monitoring
```

## New: Prometheus Metric

```
kube_pod_container_info{container="kube-state-metrics",namespace="monitoring",pod="kube-state-metrics-5fc64f676f-g16v6"} 1
```

```
kube_pod_labels{label_app="kube-state-metrics",label_pod_template_hash="5fc64f676f",namespace="monitoring",pod="kube-state-metrics-5fc64f676f-g16v6"} 1
```

# Performance improvement through caching

scrape\_duration\_seconds(endpoint="http-metrics")

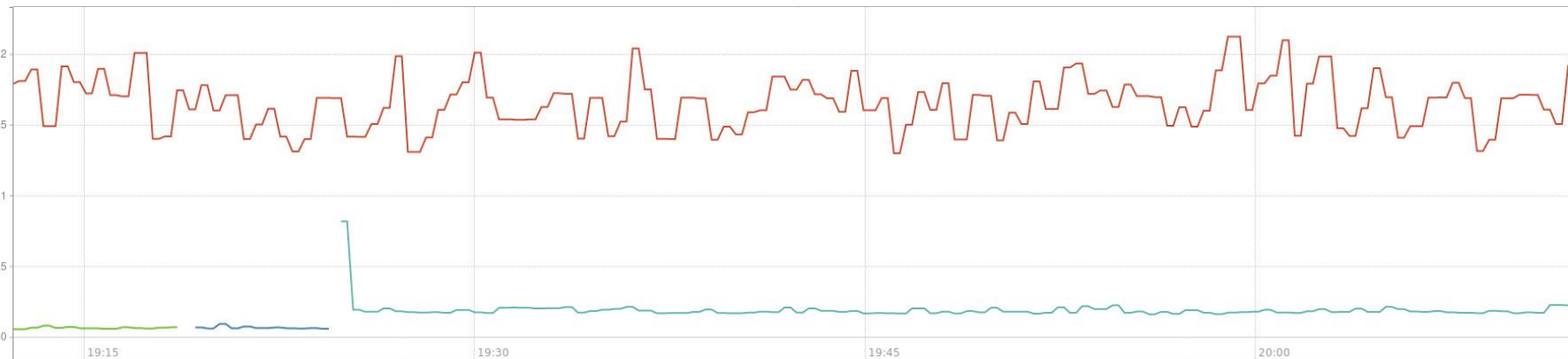
Load time: 256ms  
Resolution: 14s  
Total time series: 4

Execute

- insert metric at cursor -

Graph Console

- 1h + << Until >> Res. (s)  stacked





# Use just Reflector not whole Informers

```
// Reflector watches a specified resource
// and causes all changes to be reflected
// in the given store.
```

```
type Reflector struct {
    // [...]
```

```
    // The type of object we expect
    // to place in the store.
```

```
    expectedType reflect.Type
```

```
    // The destination to sync up
    // with the watch source
```

```
    store Store
```

```
    // listerWatcher is used to perform
    // lists and watches.
```

```
    listerWatcher ListerWatcher
```

```
    // [...]
```

```
}
```

```
type Store interface {
    Add(obj interface{}) error
    Update(obj interface{}) error
    Delete(obj interface{}) error
    List() []interface{}
    ListKeys() []string
    Get(obj interface{}) (item interface{}, exists bool, err error)
    GetByKey(key string) (item interface{}, exists bool, err error)

    // [...]
}
```

# Compression

```
# HELP kube_secret_type Type about secret.  
# TYPE kube_secret_type gauge  
kube_secret_type{namespace="default",secret="test-0",type="Opaque"} 1  
kube_secret_type{namespace="default",secret="test-1",type="Opaque"} 1  
kube_secret_type{namespace="default",secret="test-2",type="Opaque"} 1
```

<https://golang.org/pkg/compress/gzip/>

<https://github.com/NYTimes/gziphandler>

Load time: 242ms  
Resolution: 3s  
Total time series: 4

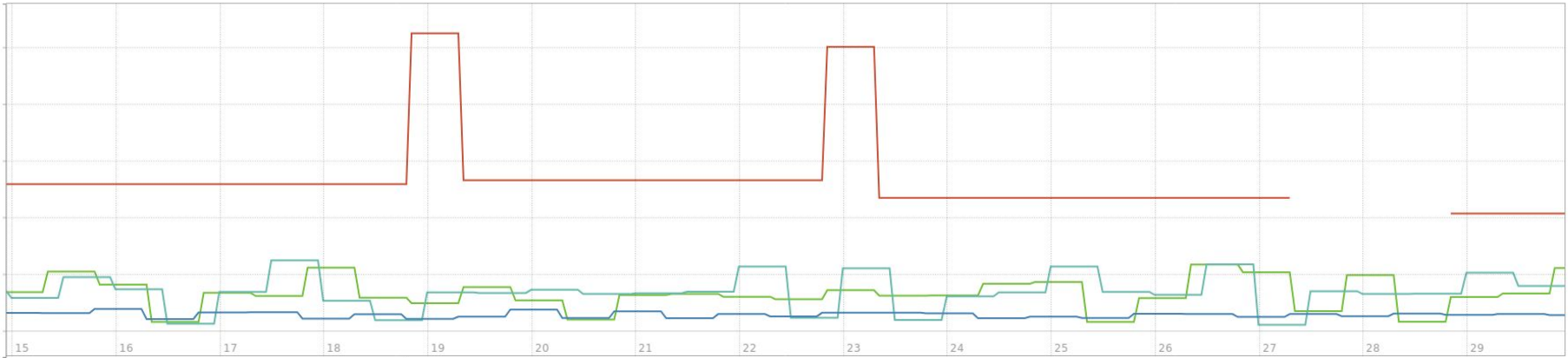
```
sum((irate(container_cpu_usage_seconds_total(container_name="kube-state-metrics",namespace="mxinden")[5m]))*100) by (pod_name)
```

Execute - insert metric at cursor -

Graph

Console

- 15m + << Until >> Res. (s)  stacked



- pod\_name="kube-state-metrics-perf-5dc6b4b799-stzbd"
- pod\_name="kube-state-metrics-nytgzip-76d6c4d64-hnb65"
- pod\_name="kube-state-metrics-gzip-767dfbd40fn7j4"
- pod\_name="kube-state-metrics-b4db9b4f7-5fmxp"

Remove Graph

# metrics + gzip(metrics) > metrics

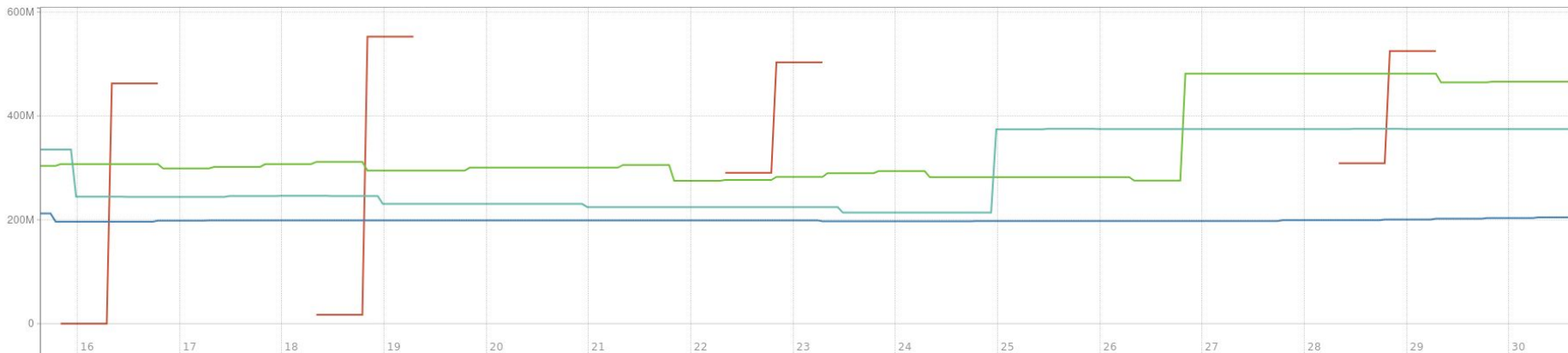
sum(container\_memory\_usage\_bytes(container\_name="kube-state-metrics", namespace="mxinden")) by (pod\_name)

Load time: 697ms  
Resolution: 3s  
Total time series: 4

Execute - insert metric at cursor -

Graph Console

- 15m + << Until >> Res. (s) stacked



■ {pod\_name="kube-state-metrics-perf-5dc6b4b799-stzd"}  
■ {pod\_name="kube-state-metrics-nytgzip-76d6c44d64-hnb65"}  
■ {pod\_name="kube-state-metrics-gzip-767dfbd48fn7j4"}  
■ {pod\_name="kube-state-metrics-b4db9b417-5fmxp"}

Remove Graph

Load time: 823ms  
Resolution: 3s  
Total time series: 4

sum(scrape\_duration\_seconds(endpoint="http-metrics")) by (pod)

Execute - insert metric at cursor -

Graph

Console

- 15m + ◀ Until ▶ Res. (s)  stacked



■ {pod="kube-state-metrics-perf-5dc6b4b799-stzbd"}  
■ {pod="kube-state-metrics-nytgzip-76d6c44d64-hnb65"}  
■ {pod="kube-state-metrics-gzip-767dfbd48En77j4"}  
■ {pod="kube-state-metrics-b4db9b4f7-5fmxp"}

[Remove Graph](#)

# Compression

`size(metrics + gzip(metrics)) > size(metrics)`

Improved network throughput < Lost CPU cycles compressing

Higher CPU utilization

# Hard code float-to-string common cases

```
func writeFloat(w *strings.Builder, f float64) {
    switch {
    case f == 1:
        w.WriteByte('1')
    case f == 0:
        w.WriteByte('0')
    case f == -1:
        w.WriteString("-1")
    case math.IsNaN(f):
        w.WriteString("NaN")
    case math.IsInf(f, +1):
        w.WriteString("+Inf")
    case math.IsInf(f, -1):
        w.WriteString("-Inf")
    default:
        // [...]
    }
}
```

# Golang strings.Builder

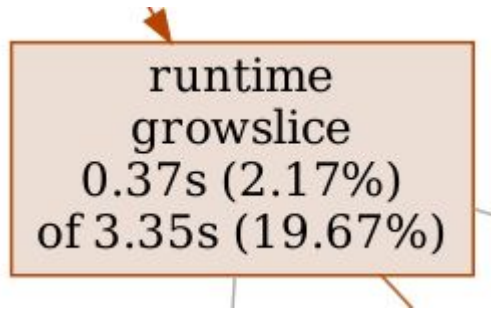
```
// A Builder is used to efficiently build
// a string using Write methods.
// It minimizes memory copying.
type Builder struct {
    // [...]
    buf []byte
}
```



<b>Version</b>	<b>Scrape duration</b>	<b>CPU usage</b>	<b>RSS peak</b>	<b>% in GC</b>
v1.4	22s	1.0	2.5GiB	0.015
newprom_gzip	16s	0.8	2.2 GiB	0.010
newprom_nogzip	13s	0.6	2.4 GiB	0.014
mxinden_gzip	11s	0.3	0.81 GiB	0.006
mxinden_nogzip	7s	0.16	0.73 GiB	0.004

v1.5.0

# Memory pre-allocation



runtime  
growslice  
0.37s (2.17%)  
of 3.35s (19.67%)

```
// growslice handles slice growth during append.  
// It is passed the slice element type, the old  
// slice, and the desired new minimum capacity,  
// and it returns a new slice with at least that  
// capacity, with the old data copied into it.  
// [...]  
func growslice(et *_type, old slice, cap int) slice {
```

```
{
  Name: "kube_pod_container_info",
  // [...]
  GenerateFunc: wrapPodFunc(func(p *v1.Pod) *metric.Family {
    ms := []*metric.Metric{}
    // [...]

    for _, cs := range p.Status.ContainerStatuses {
      ms = append(ms, &metric.Metric{
        // [...]
      })
    }

    // [...]
  })),
},
```

```
{
    Name: "kube_pod_container_info",
    // [...]
    GenerateFunc: wrapPodFunc(func(p *v1.Pod) *metric.Family {
        ms := make([]*metric.Metric, len(p.Status.ContainerStatuses))
        // [...]

        for i, cs := range p.Status.ContainerStatuses {
            ms[i] = &metric.Metric{
                // [...]
            }
        }

        // [...]
    })),
},
```

```
sum(rate(container_cpu_usage_seconds_total(container_name="kube-state-metrics", pod_name=~".v1-6-0-rc-0-.*")[5m])) by (pod_name)
```

Load time: 312ms  
Resolution: 172s  
Total time series: 3

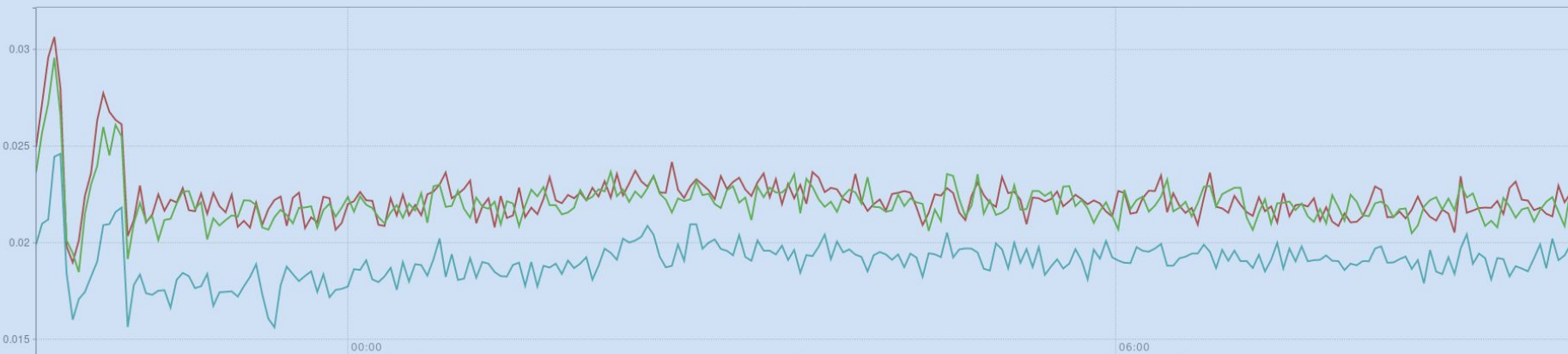
Execute

- insert metric at cursor -

Graph

Console

- 12h + << Until >> Res. (s)  stacked



■ {pod\_name="kube-state-metrics-v1-6-0-rc-0-5-7cf69f8d48-5wn74"}  
■ {pod\_name="kube-state-metrics-v1-6-0-rc-0-4-758fb8d5ff-jg9ts"}  
■ {pod\_name="kube-state-metrics-v1-6-0-rc-0-0-756b46798c-zhgft"}

Remove Graph

```
func growslice(et *_type, old slice, cap int) slice {
    // [...]

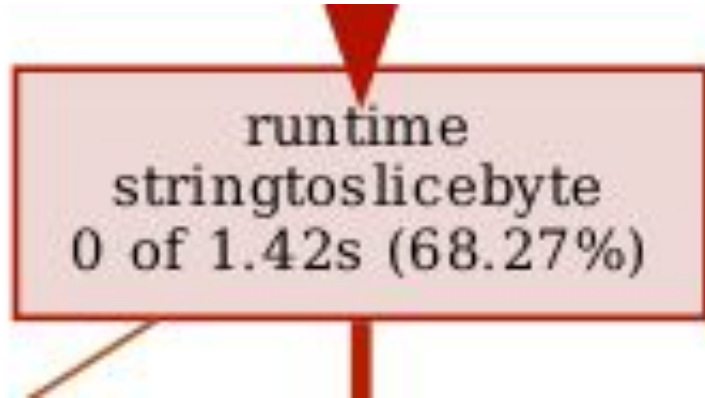
    newcap := old.cap
    doublecap := newcap + newcap

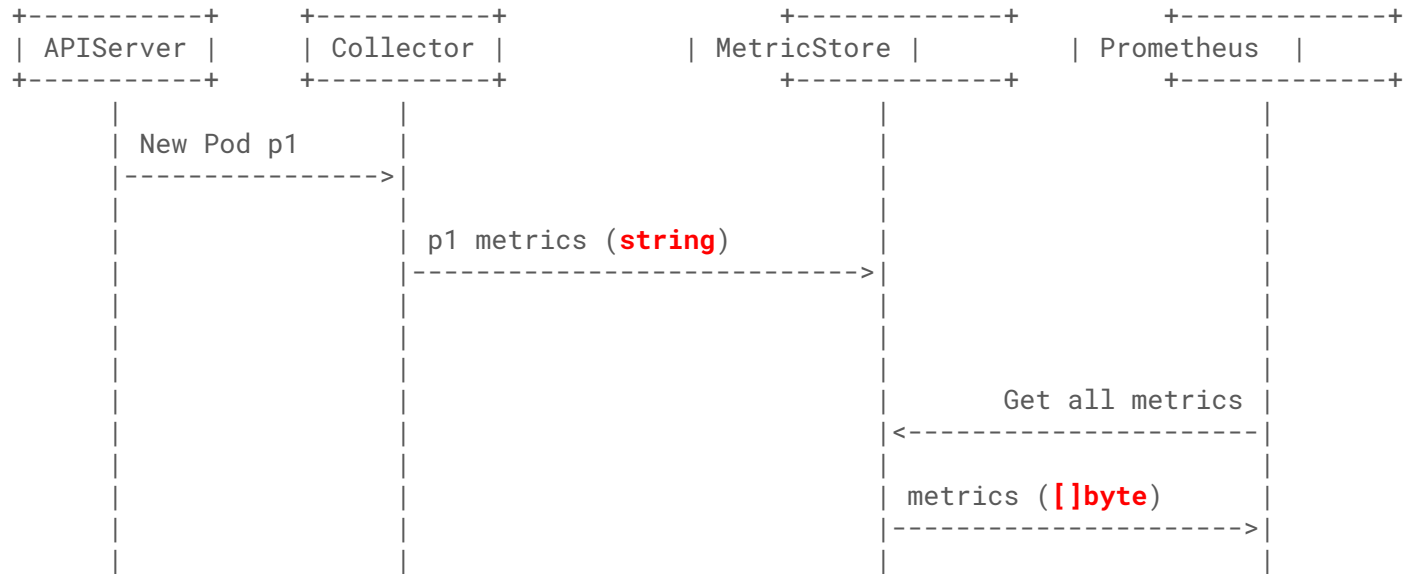
    if cap > doublecap {
        newcap = cap
    } else {
        if old.len < 1024 {
            newcap = doublecap
        } else {
            // [...]
        }
    }
}
```



# []byte != string

- Strings are immutable





Load time: 312ms  
Resolution: 172s  
Total time series: 3

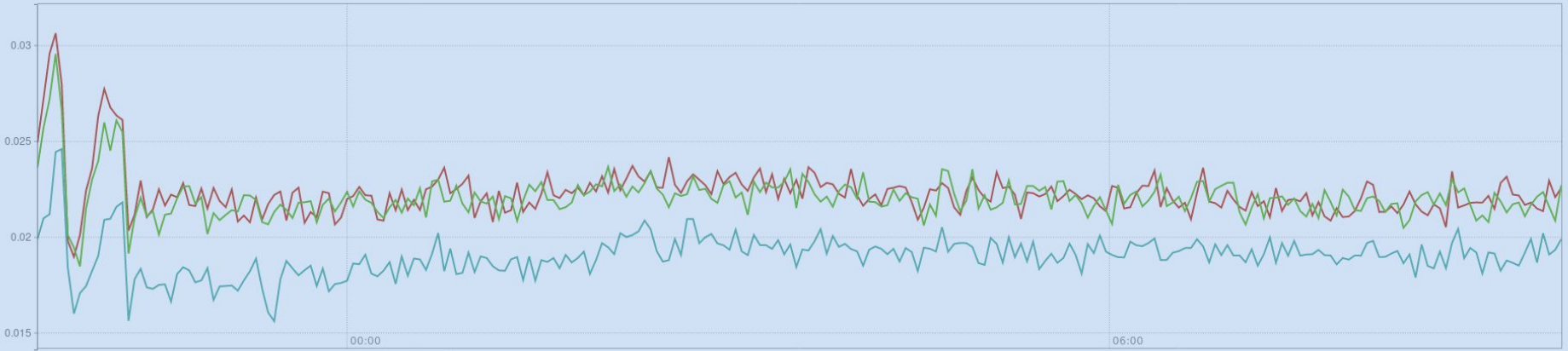
```
sum(rate(container_cpu_usage_seconds_total{container_name="kube-state-metrics", pod_name=~".v1-6-0-rc-0-.*"}[5m])) by (pod_name)
```

Execute

- insert metric at cursor -

Graph Console

- 12h + << Until >> Res. (s)  stacked



■ {pod\_name="kube-state-metrics-v1-6-0-rc-0-5-7cf69f8d48-5wn74"}  
■ {pod\_name="kube-state-metrics-v1-6-0-rc-0-4-758fb8d5ff-jg9ts"}  
■ {pod\_name="kube-state-metrics-v1-6-0-rc-0-0-756b46798c-zhgft"}

Remove Graph

v1.6.0

# Future

- Introducing sharding based on object id instead of Kubernetes type
- Reduce heap escapes
- Optimize memory alignment on structs in hot path
- ...
- (Commutative compression)

# Take aways

- Monitor your cluster and application!
- Use **Monitoring** for alerting **and performance analysis**

Golang memory model documentation:

- > If you must read the rest of this document to understand the behavior of your program, you are being too clever.
- > Don't be clever.