

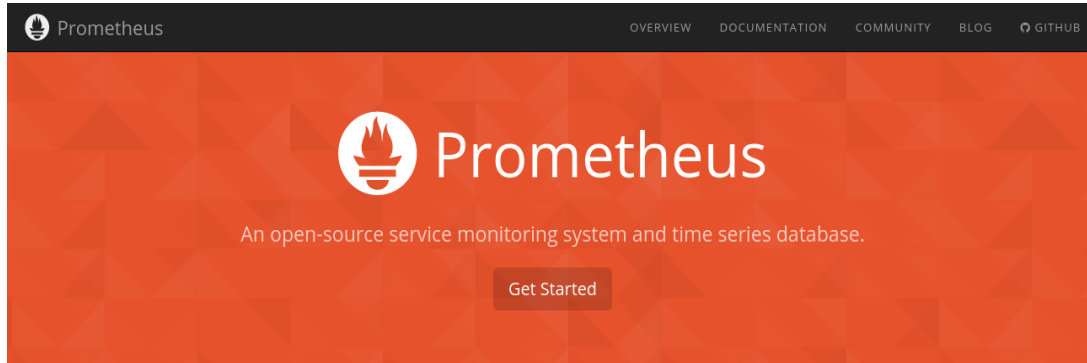


# Prometheus: Designing and Implementing a Modern Monitoring Solution in Go

*Björn "Beorn" Rabenstein, Production Engineer, SoundCloud Ltd.*



# http://prometheus.io



## Data model

Prometheus implements a highly dimensional data model. Time series are identified by a metric name and a set of key-value pairs.

[View details »](#)

## Query language

A flexible query language allows slicing and dicing of collected time series data in order to generate ad-hoc graphs, tables, and alerts.

[View details »](#)

## Visualization

Prometheus has multiple modes for visualizing data: a built-in expression browser, a GUI-based dashboard builder, and a console template language.

[View details »](#)

## Storage

Prometheus stores time series in memory and on local disk in an efficient custom format. Scaling is achieved by functional sharding and federation.

[View details »](#)

## Operation

Each server is independent for reliability, relying only on local storage. Written in Go, all binaries are statically linked and easy to deploy.

[View details »](#)

## Client libraries

Client libraries allow easy instrumentation of services. Currently, Go, Java, and Ruby are supported. Custom libraries are easy to implement.

[View details »](#)

## Alerting

Alerts are defined based on Prometheus's flexible query language and maintain dimensional information. An alertmanager handles notifications and silencing.

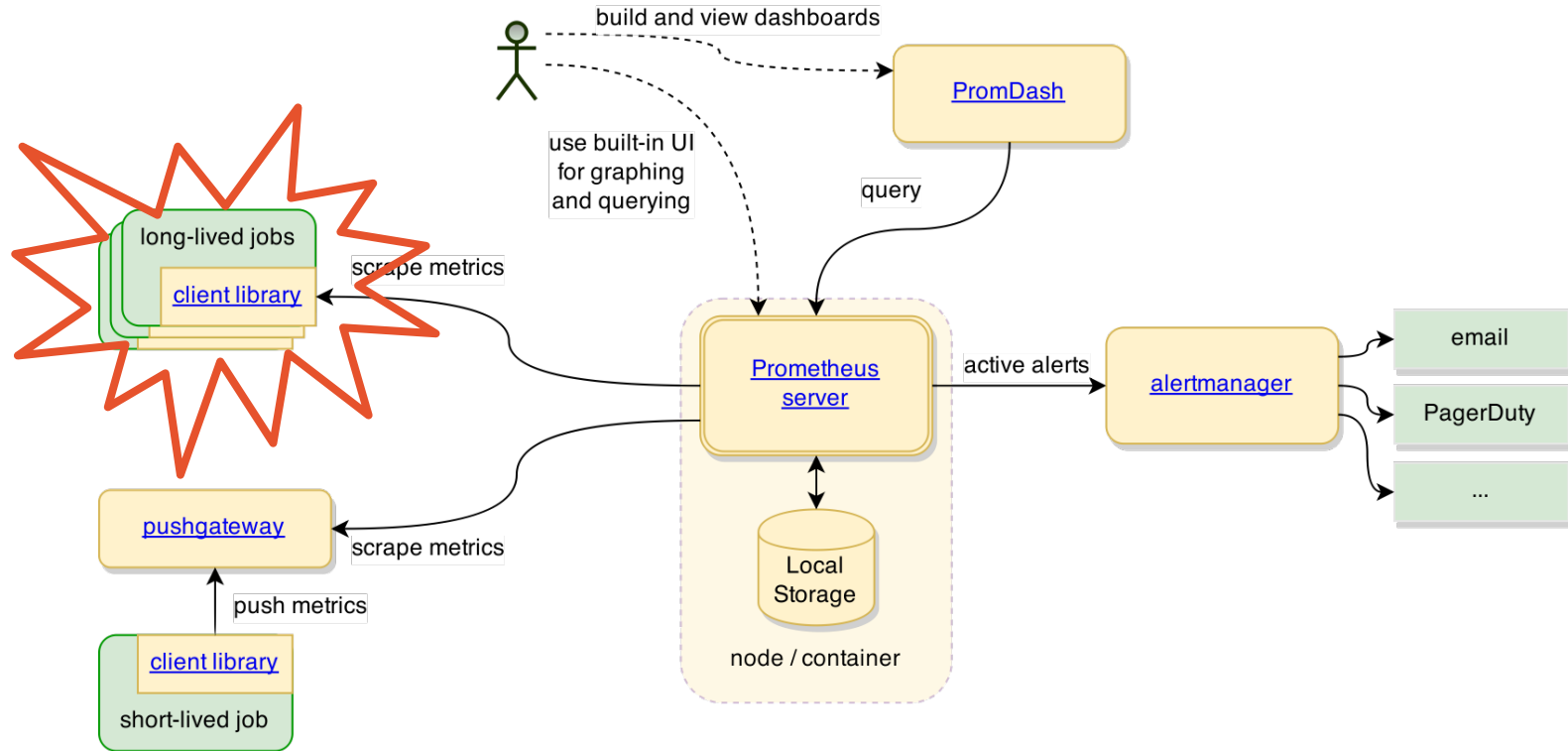
[View details »](#)

## Exporters

Existing exporters allow bridging of third-party data into Prometheus. Examples: system statistics, as well as Docker, HAProxy, StatsD, and JMX metrics.

[View details »](#)

# Architecture



Go client library

# Counter interface

(almost complete)

```
type Counter interface {  
    Metric  
  
    Inc()  
    Add(int)  
}
```

```
type Metric interface {  
    Write(*dto.Metric) error  
}
```

```
type counter struct {  
    value int  
}
```

```
func (c *counter) Add(v int) {  
    if v < 0 {  
        panic(errors.New("counter cannot decrease in value"))  
    }  
    c.value += v  
}
```

```
func (c counter) Write(*dto.Metric) error {  
    // ...  
}
```

```
type counter struct {  
    value int  
    mtx    sync.Mutex  
}  
  
func (c *counter) Add(v int) {  
    c.mtx.Lock()  
    defer c.mtx.Unlock()  
    if v < 0 {  
        panic(errors.New("counter cannot decrease in value"))  
    }  
    c.value += v  
}  
  
func (c *counter) Write(*dto.Metric) error {  
    c.mtx.Lock()  
    defer c.mtx.Unlock()  
    // ...  
}
```

# Performance matters

It's a library, run with a large number of unknown use-cases.

```
func benchmarkAddAndWrite(b *testing.B, c Counter) {  
    for i := 0; i < b.N; i++ {  
        if i%1000 == 0 {  
            c.Write(&dto)  
            continue  
        }  
        c.Add(42)  
    }  
}
```

```
func BenchmarkNaiveCounter(b *testing.B) {  
    benchmarkAddAndWrite(b, NewNaiveCounter())  
}
```

```
func BenchmarkMutexCounter(b *testing.B) {  
    benchmarkAddAndWrite(b, NewMutexCounter())  
}
```

```
$ go test -bench=Counter
```



# Results are in.

Naive counter: 5 ns/op.

*(Probably mostly overhead: function call, for loop...)*

Mutex counter: 150 ns/op.

```

func benchmarkAddAndWrite(b *testing.B, c Counter, concurrency int) {
    b.StopTimer()
    var start, end sync.WaitGroup
    start.Add(1)
    end.Add(concurrency)
    n := b.N / concurrency

    for i := 0; i < concurrency; i++ {
        go func() {
            start.Wait()
            for i := 0; i < n; i++ {
                if i%1000 == 0 {
                    c.Write(&dto)
                    continue
                }
                c.Add(42)
            }
            end.Done()
        }()
    }
    b.StartTimer()
    start.Done()
    end.Wait()
}

func BenchmarkMutexCounter10(b *testing.B) {
    benchmarkAddAndWrite(b, NewMutexCounter(), 10)
}

```

```
$ go test -bench=Counter -cpu=1,4,16 # -race
```

# It's getting worse.

Let's talk about lock contention...

ns/op	1 Goroutine	10 Goroutines	100 Goroutines
GOMAXPROCS=1	150	160	190
GOMAXPROCS=4	150	730	570
GOMAXPROCS=16	150	1100	1100

```
rate(prometheus_local_storage_ingested_samples_total[1m])
```

Execute

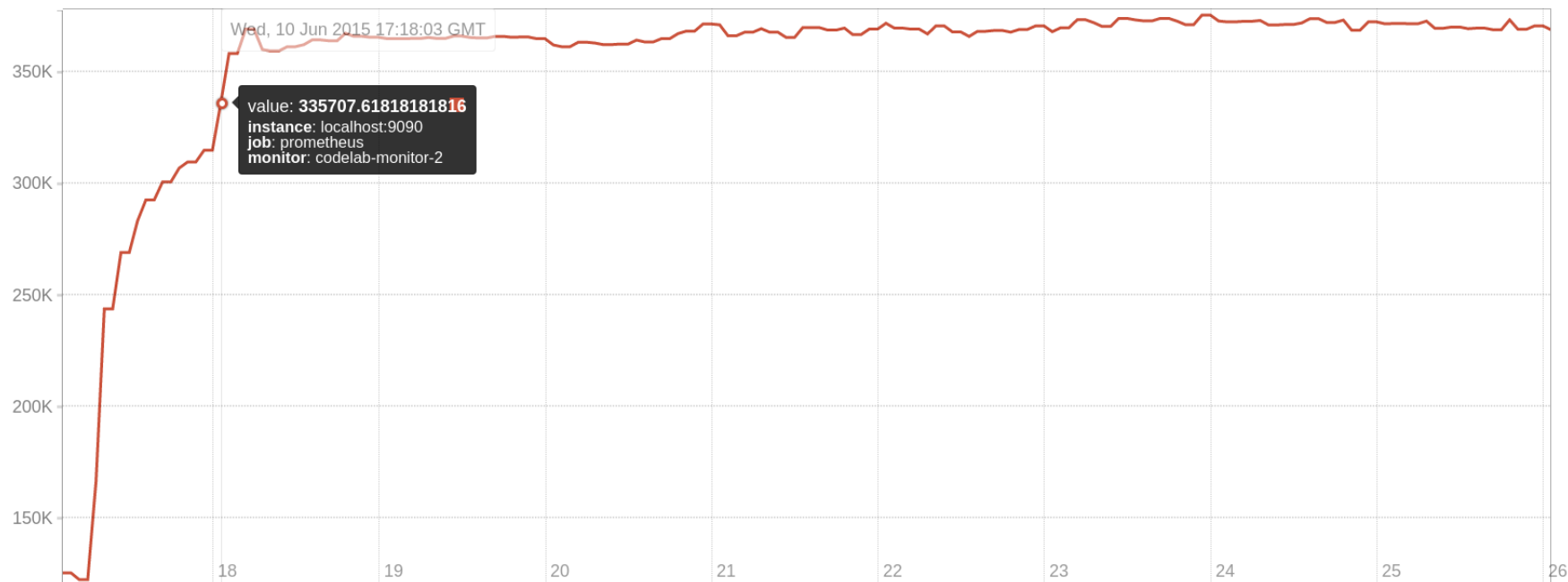
- Insert Metric at Cursor -

Load time: 511ms

Resolution: 3s

Graph Console

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



✓ {instance="localhost:9090",job="prometheus",monitor="codelab-monitor-2"}

Add Graph

Do not communicate by sharing memory;  
share memory by communicating.

*Rob 12:3–4*

```
type counter struct {  
    in chan int // May be buffered.  
    out chan int // Must be synchronous.  
}  
  
func (c *counter) Add(v int) {  
    c.in <- v  
}  
  
func (c *counter) Write(*dto.Metric) error {  
    value <- c.out  
    // ...  
}  
  
func (c *counter) loop() {  
    var value int64  
    for {  
        select {  
        case v := <-c.in:  
            value += v  
        case c.out <- value:  
            // Do nothing.  
        }  
    }  
}
```

# Channel counter.

x / y: Synchronous vs. buffered *in* channel.

ns/op	1 Goroutine	10 Goroutines	100 Goroutines
GOMAXPROCS=1	670 / 310	690 / 320	680 / 360
GOMAXPROCS=4	3600 / 940	2000 / 2000	1600 / 2200
GOMAXPROCS=16	3500 / 850	2300 / 2200	1800 / 2700

```
import "sync/atomic"

type counter struct {
    value int64
}

func (c *counter) Add(v int64) {
    if v < 0 {
        panic(errors.New("counter cannot decrease in value"))
    }
    atomic.AddInt64(&c.value, v)
}

func (c *counter) Write(*dto.Metric) error {
    v := atomic.LoadInt64(&c.value)
    // Process v...
}
```



# Atomic counter.

Yay!

ns/op	1 Goroutine	10 Goroutines	100 Goroutines
GOMAXPROCS=1	15	14	15
GOMAXPROCS=4	14	45	44
GOMAXPROCS=16	14	47	45

# I lied!

Prometheus uses float64 for sample values.

```
type Counter interface {  
    Metric  
  
    Inc()  
    Add(float64)  
}
```

```
type Metric interface {  
    Write(*dto.Metric) error  
}
```

```
type counter struct {
    valueBits uint64
}

func (c *counter) Add(v float64) {
    if v < 0 {
        panic(errors.New("counter cannot decrease in value"))
    }
    for {
        oldBits := atomic.LoadUint64(&c.valueBits)
        newBits := math.Float64bits(math.Float64frombits(oldBits) + v)
        if atomic.CompareAndSwapUint64(&c.valueBits, oldBits, newBits) {
            return
        }
    }
}

func (c *counter) Write(*dto.Metric) error {
    v := math.Float64frombits(atomic.LoadUint64(&c.valueBits))
    // Process v...
}
```

# Atomic “spinning” counter for floats.

Yes, it works...

ns/op	1 Goroutine	10 Goroutines	100 Goroutines
GOMAXPROCS=1	25	23	24
GOMAXPROCS=4	24	97	100
GOMAXPROCS=16	24	120	130

# One last thing.

Read the fine print at the bottom of the page...

## Bugs

☞ On x86-32, the 64-bit functions use instructions unavailable before the Pentium MMX. On non-Linux ARM, the 64-bit functions use instructions unavailable before the ARMv6k core. On both ARM and x86-32, it is the caller's responsibility to arrange for 64-bit alignment of 64-bit words accessed atomically. The first word in a global variable or in an allocated struct or slice can be relied upon to be 64-bit aligned.

Build version go1.4.2.

Except as [noted](#), the content of this page is licensed under the Creative Commons Attribution 3.0 License, and code is licensed under a [BSD license](#).

[Terms of Service](#) | [Privacy Policy](#)

# Prometheus on Raspberry Pi

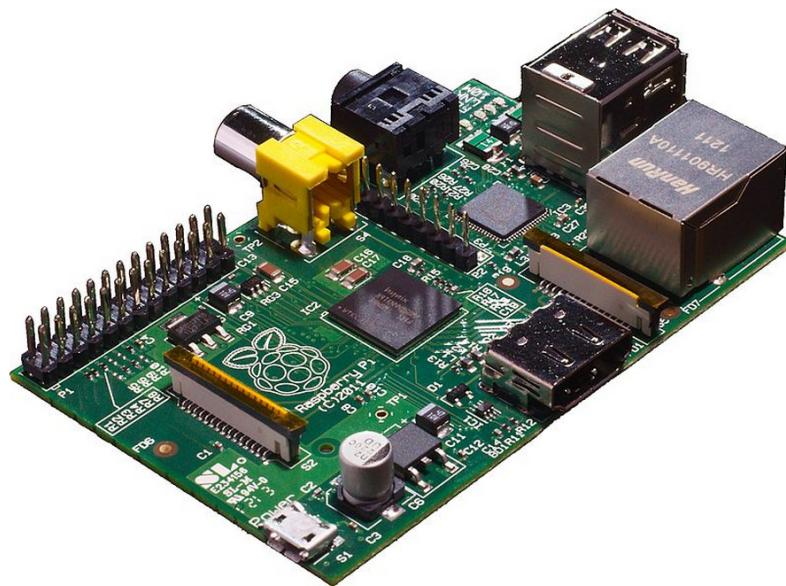
10 FEBRUARY 2015

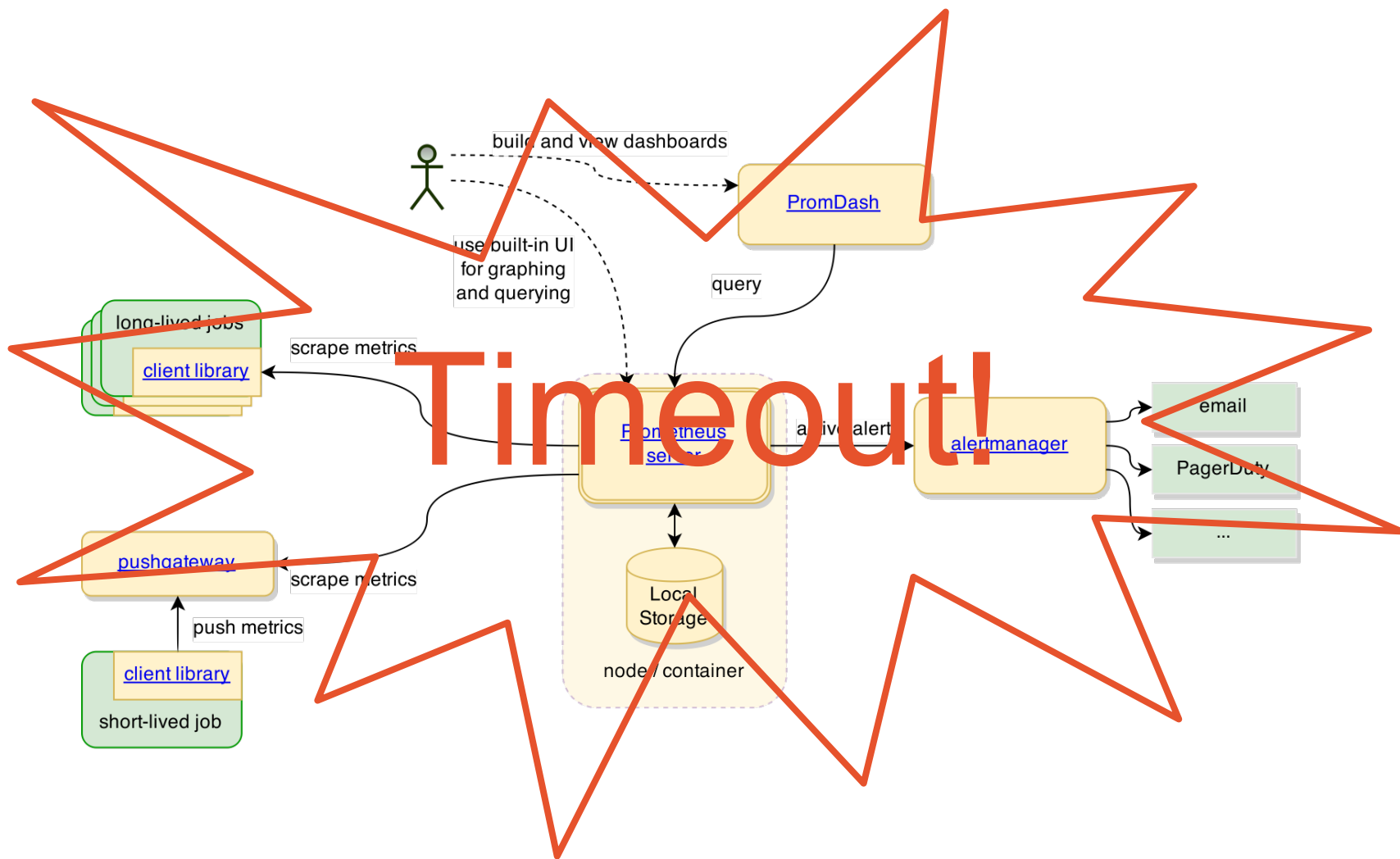
## Prometheus

Prometheus is a new open-source service monitoring system and time series database written in Go.

Check out the [announcement](#) and my article about [monitoring Docker Containers with Prometheus](#) if you don't know what I'm talking about.

## My Stack







# Prometheus: How to increment a numerical value

*Björn "Beorn" Rabenstein, Production Engineer, SoundCloud Ltd.*





# 1. Use -benchmem.

To detect allocation churn.

```
go test -bench=. -cpu=1,4,16 -benchmem
```

*Escape analysis:*

```
go test -gcflags=-m -bench=Something
```

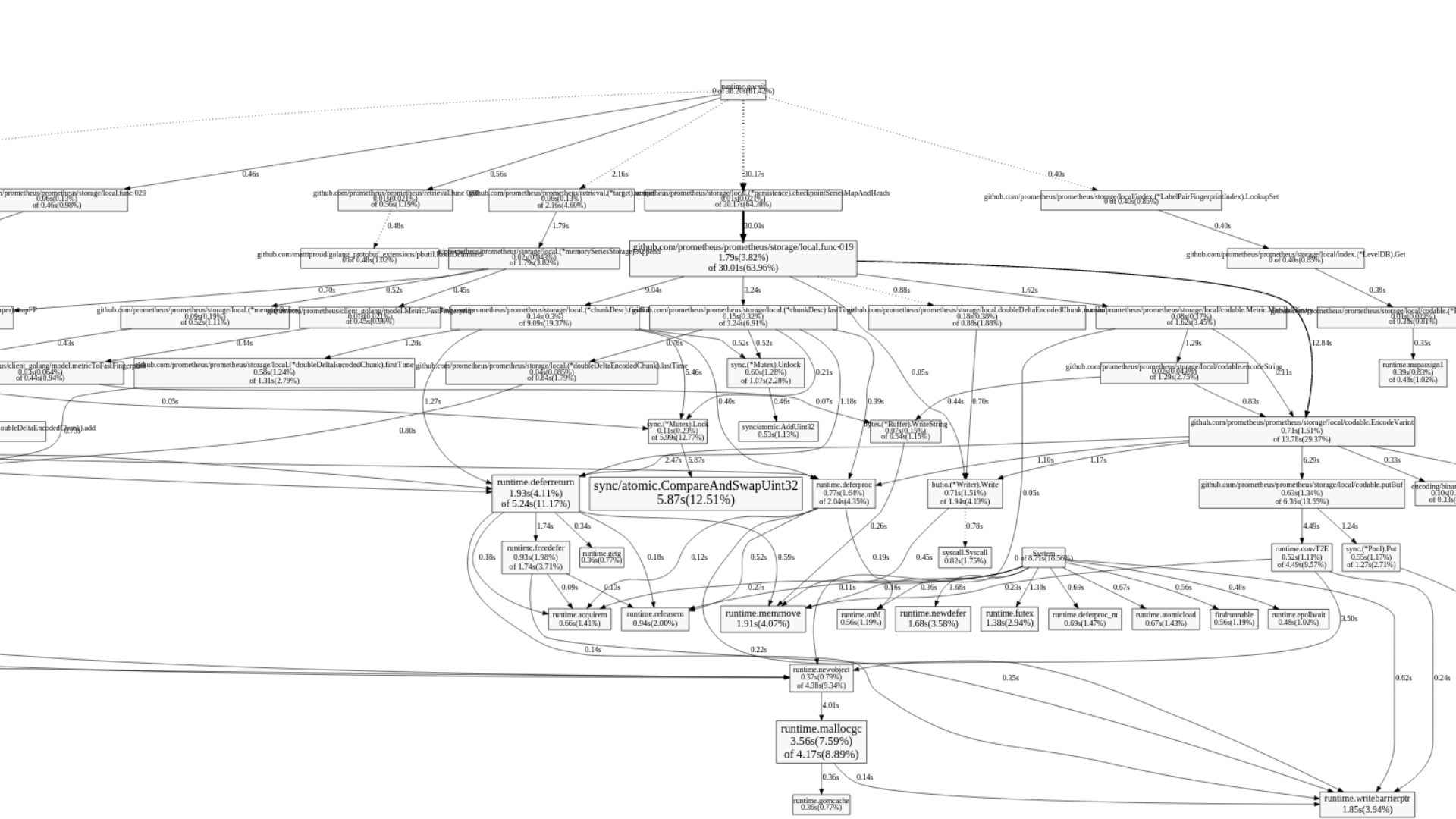
## 2. Use pprof.

For debugging. For runtime and allocation profiling.

```
import _ "net/http/pprof"
```

```
$ go tool pprof http://localhost:9090/debug/pprof/profile  
(pprof) web
```

```
$ go tool pprof http://localhost:9090/debug/pprof/heap  
(pprof) web
```



### 3. Use cgo judiciously.

Highly optimized C libraries can be great. But there is a cost...

- ❑ Loss of certain advantages of the Go build environment.
- ❑ Per-call overhead – dominates run-time if C function runs for  $<1\mu\text{s}$ .
- ❑ Need to shovel input and output data back and forth.

<http://jmoiron.net/blog/go-performance-tales/>

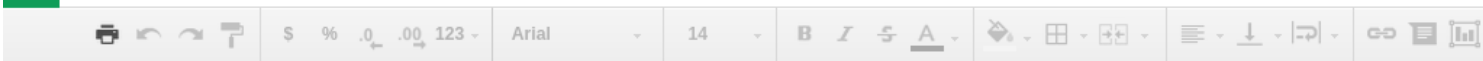
# Special thanks

Matt T. Proud & Julius Volz

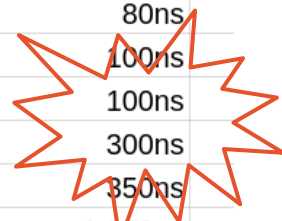
founding fathers of the Prometheus project



Supplementary slides

f<sub>x</sub> | L1 cache reference

	A	B
1	L1 cache reference	1ns
2	Branch mispredict	5ns
3	L2 cache reference	7ns
4	Mutex lock/unlock	25ns
5	Linear seek in 100 long slice @ F1	80ns
6	Main memory reference	100ns
7	Go sync.Mutex lock/unlock @ F1	100ns
8	Go channel send/receive @ F1	300ns
9	Go native map retrieval @ F1	350ns
10	Go native map insertion @ F1	1.000ns
11	Call method reflectively @ F1	1.000ns
12	Compress 1kb /w cheap compression algorithm	3.000ns
13	Send 2kb over 1GBPS network	20.000ns
14	Read 1Mb sequentially from memory	250.000ns
15	Roundtrip within datacenter	500.000ns
16	Memcache access @ F1	3.000.000ns
17	Datastore put @ F1	6.000.000ns
18	Datastore get @ F1	6.000.000ns
19	Datastore query @ F1	6.500.000ns
20	Disk seek	10.000.000ns
21	Read 1MB sequentially from disk	20.000.000ns
22	Roundtrip over the atlantic	150.000.000ns
23		





```
type counter struct {  
    value int  
    mtx    sync.RWMutex  
}  
  
func (c *counter) Add(v int) {  
    c.mtx.Lock()  
    defer c.mtx.Unlock()  
    if v < 0 {  
        panic(errors.New("counter cannot decrease in value"))  
    }  
    c.value += v  
}  
  
func (c *counter) Inc() {  
    c.Add(1)  
}  
  
func (c *counter) Write(*dto.Metric) error {  
    c.mtx.RLock()  
    defer c.mtx.RUnlock()  
    // ...  
}
```

# RWMutex

ns/op	1 Goroutine	10 Goroutines	100 Goroutines
GOMAXPROCS=1	170	180	210
GOMAXPROCS=4	170	820	680
GOMAXPROCS=16	170	1300	1200

```
func (c *counter) loop() {  
    var value float64  
    for {  
        select {  
        case v := <-c.write:  
            value += v  
        default:  
            select {  
            case v := <-c.write:  
                value += v  
            case c.read <- value:  
                // Do nothing.  
            }  
        }  
    }  
}
```

# Tricky channel counter.

ns/op	1 Goroutine	10 Goroutines	100 Goroutines
GOMAXPROCS=1	117 ↓	130	164
GOMAXPROCS=4	389 ↑↑	707	1044 ↑↑
GOMAXPROCS=16	388 ↑↑	1297	1707 ↑

# Channel counter without Write.

ns/op	1 Goroutine	10 Goroutines	100 Goroutines
GOMAXPROCS=1	240 / 73	254 / 75	260 / 82
GOMAXPROCS=4	1040 / 150	760 / 290	500 / 630
GOMAXPROCS=16	1040 / 150	700 / 360	510 / 460