

Evolutionary Optimization

Make it work — Make it right — Make it fast



fastly

Peter Bourgon · Metrics, tracing, and logging

Hamburger

Home About Talks Blog

Metrics, tracing, and logging

2017 02 21

Today I had the good fortune of attending the 2017 Distributed Tracing Summit, with lots of rad folks from orgs like AWS/X-Ray, OpenZipkin, OpenTracing, Instana, Datadog, Librato, and many others I regret that I'm forgetting. At one point the discussion took a turn toward project scope and definitions. Should a tracing system also manage logging? What indeed *is* logging, when viewed through the different lenses represented in the room? And where do all of the various concrete systems fit in to the picture?

In short, I felt that we were stumbling a little bit around a shared vocabulary. I thought we could probably map out the domain of instrumentation, or observability, as a sort of Venn diagram. Metrics, tracing, and logging are definitely all parts of a broader picture, and can definitely overlap in some circumstances, but I wanted to try and identify the properties of each that were truly distinct. I had a think over a coffee break and came up with this.

A Venn diagram consisting of three overlapping circles. The top circle is labeled "Metrics" with the subtitle "Aggregatable". The left circle is labeled "Tracing" with the subtitle "Request scoped". The bottom circle is labeled "Logging" with the subtitle "Events". The overlapping areas represent the common properties of the three concepts.

A mediocre white man presents

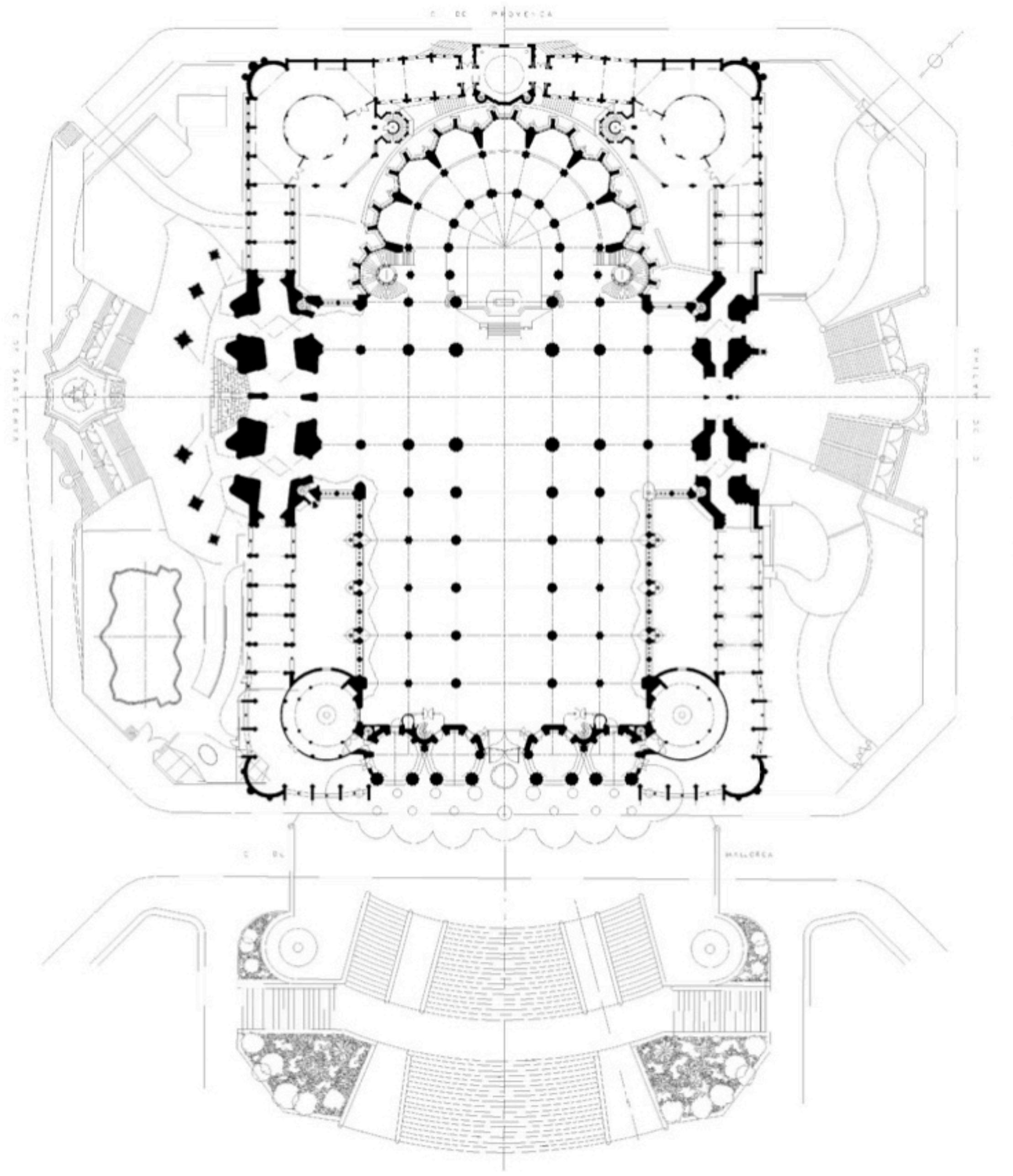
Building OK Log



Design

Make it work





“Weeks of programming can save you hours of planning.”

—A Hackernews

**“Programming is made up of 2 activities:
making decisions (95%) and typing (5%)”**

josephg.com/blog/rplace-in-a-weekend

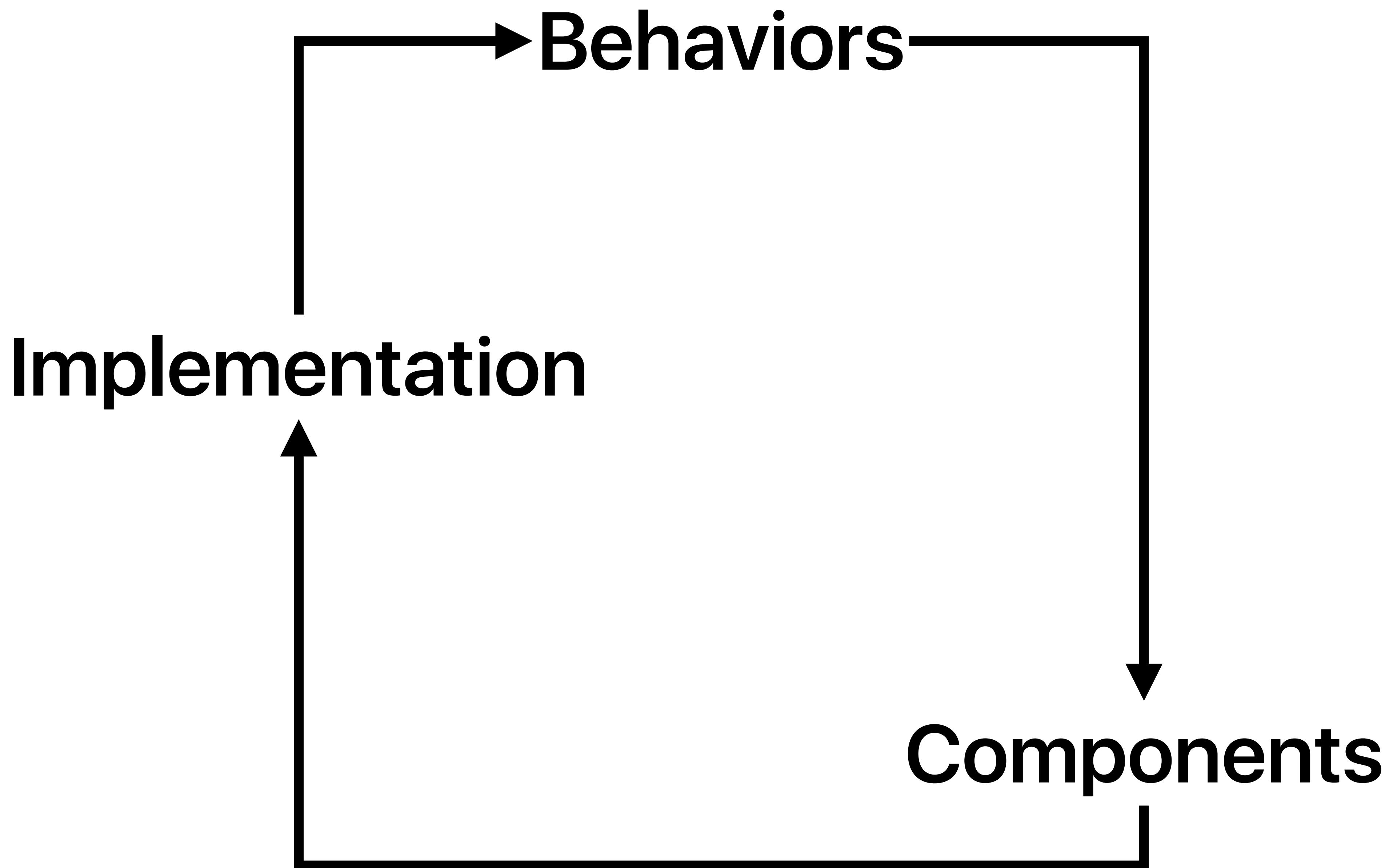
1. Make decisions
2. Type

1. Make decisions Design
2. Type Build

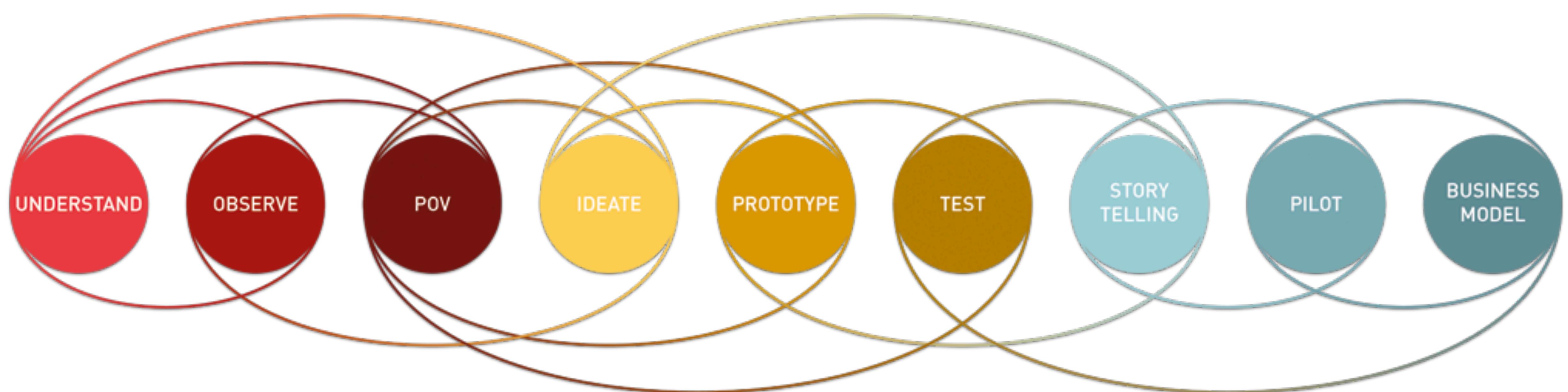
The screenshot shows a web browser window displaying a Dropbox document titled "Prometheus for logs". The document contains a list of steps for segment compaction:

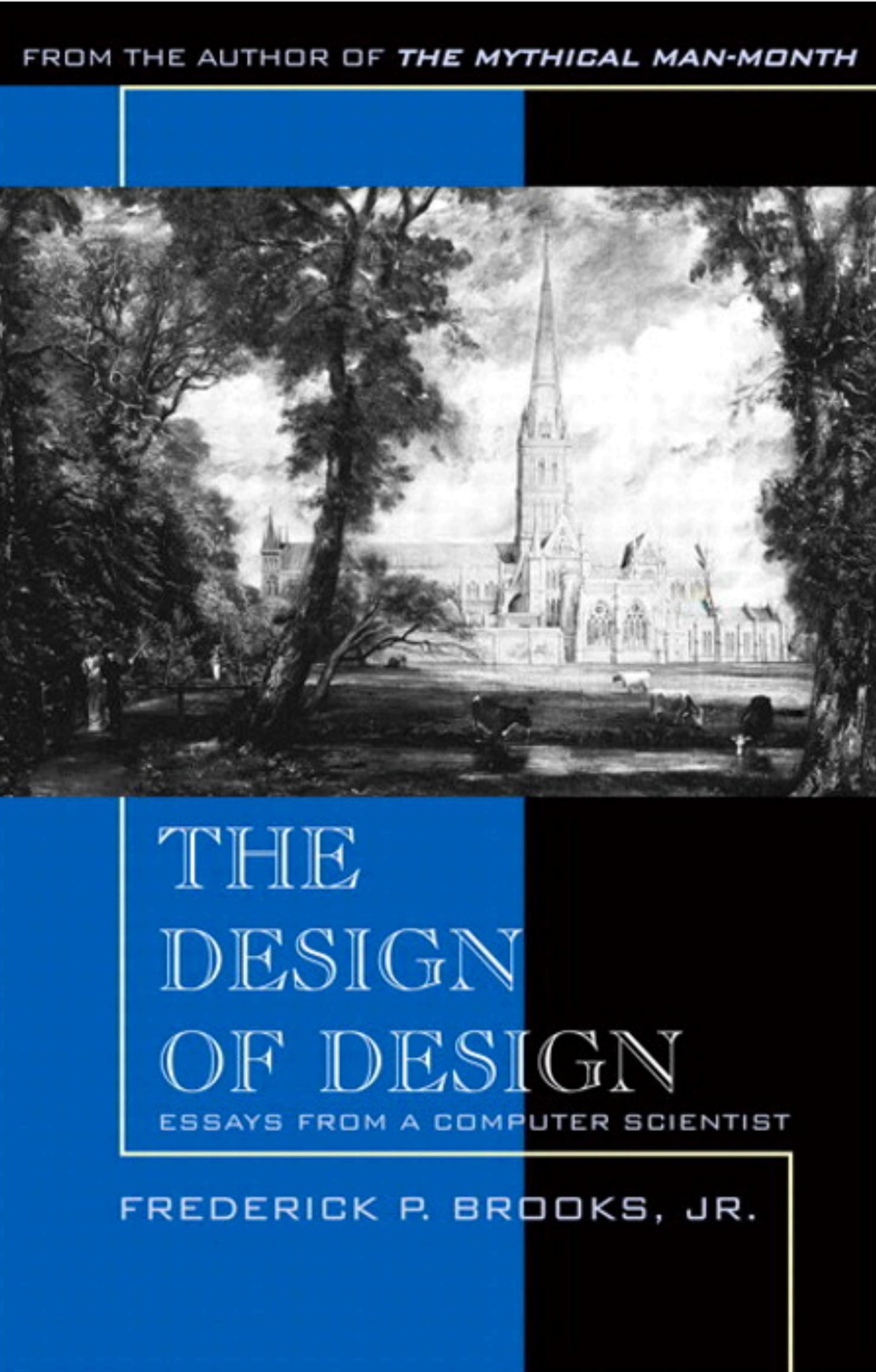
- Every step can fail and will be detected/accommodated with timeouts
- Performs segment compaction
 - Selects a group of local segments with lots of time overlap: A, B, C, D
 - Merges them in time order to produce an übersegment
 - Splits the übersegment to nonoverlapping segments E, F, G
 - All records in {A, B, C, D} are present exactly once in {E, F, G}
 - Note: can't make stronger statements about record location/identity!
 - Mark each new segment as conflicting with (superceding) **every** old segment
 - Segment{ID: E, From: t0, To: t1, Data: d1, Conflicts: {A, B, C, D}}
 - Segment{ID: F, From: t1, To: t2, Data: d2, Conflicts: {A, B, C, D}}
 - Segment{ID: G, From: t2, To: t3, Data: d3, Conflicts: {A, B, C, D}}
 - Replicate each new segment in the cluster, **but don't make them queryable yet**
 - Broadcast a Replace operation to every node in the cluster
 - Replace(Disable: {A, B, C, D}, Enable: {D, E, F})
 - **Atomically** makes D, E, F queryable, and deletes A, B, C, D
 - Note: during the broadcast, queries can see inconsistent state!
 - Note: that inconsistency is detectable!
 - Node 1 returns records from segments A, C – which have no conflicts
 - Node 2 returns records from segment D – which conflicts with A, B, C, D
 - Conflict in result set is detected
 - We can declare at query time how to deal with conflicts
 - FailFast: queries that return conflicting records fail, and may be retried
 - Permissive: merge all results, accepting some duplicate records
 - Duplicate records will likely be easy to spot – UI could help here
 - etc.

A white diagonal banner with torn edges containing the text "This is basically a spec!". The banner is positioned diagonally across the slide, with its top edge pointing towards the top-left and its bottom edge pointing towards the bottom-right. The text is in a large, bold, black sans-serif font. The banner has a slightly irregular shape with visible torn edges along its top and right sides, giving it a hand-drawn or distressed appearance. The background is a solid dark grey.



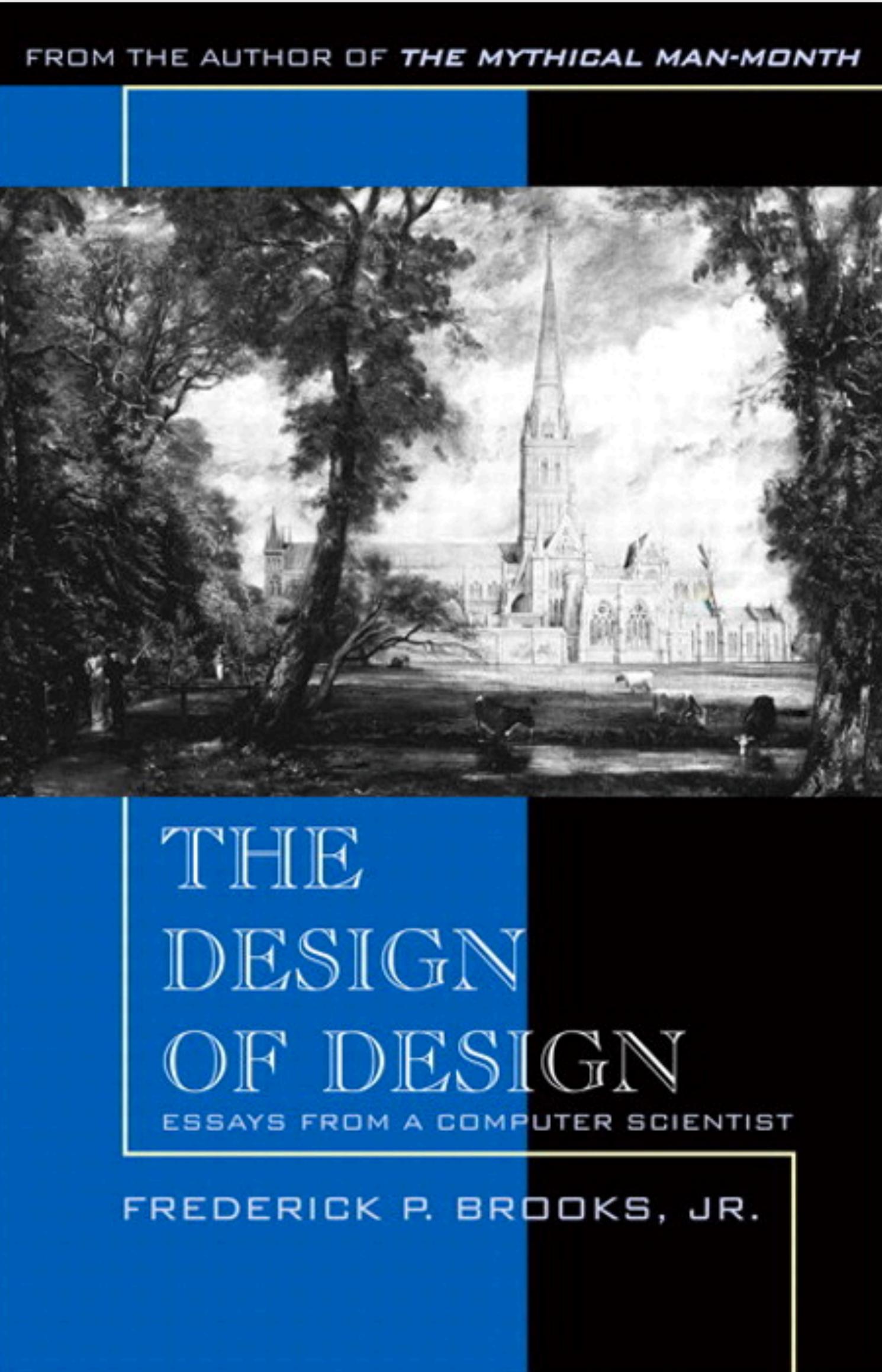
THE DESIGN THINKING PROCESS



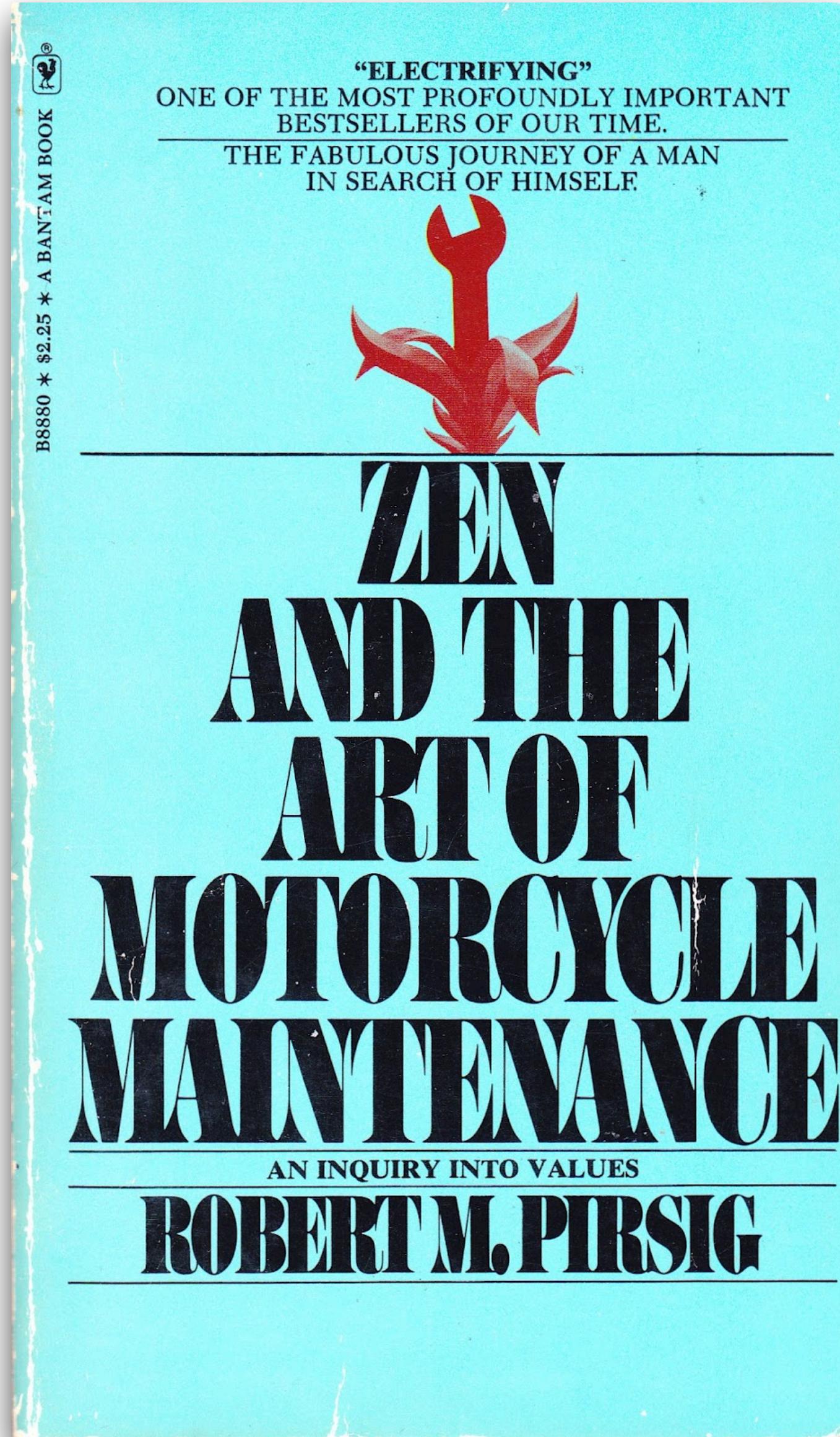


"The whole design problem [is] an intricately interlocked interplay of factors..."

"In practice, designing seems to proceed by oscillating between sub-solution and sub-problem areas..."



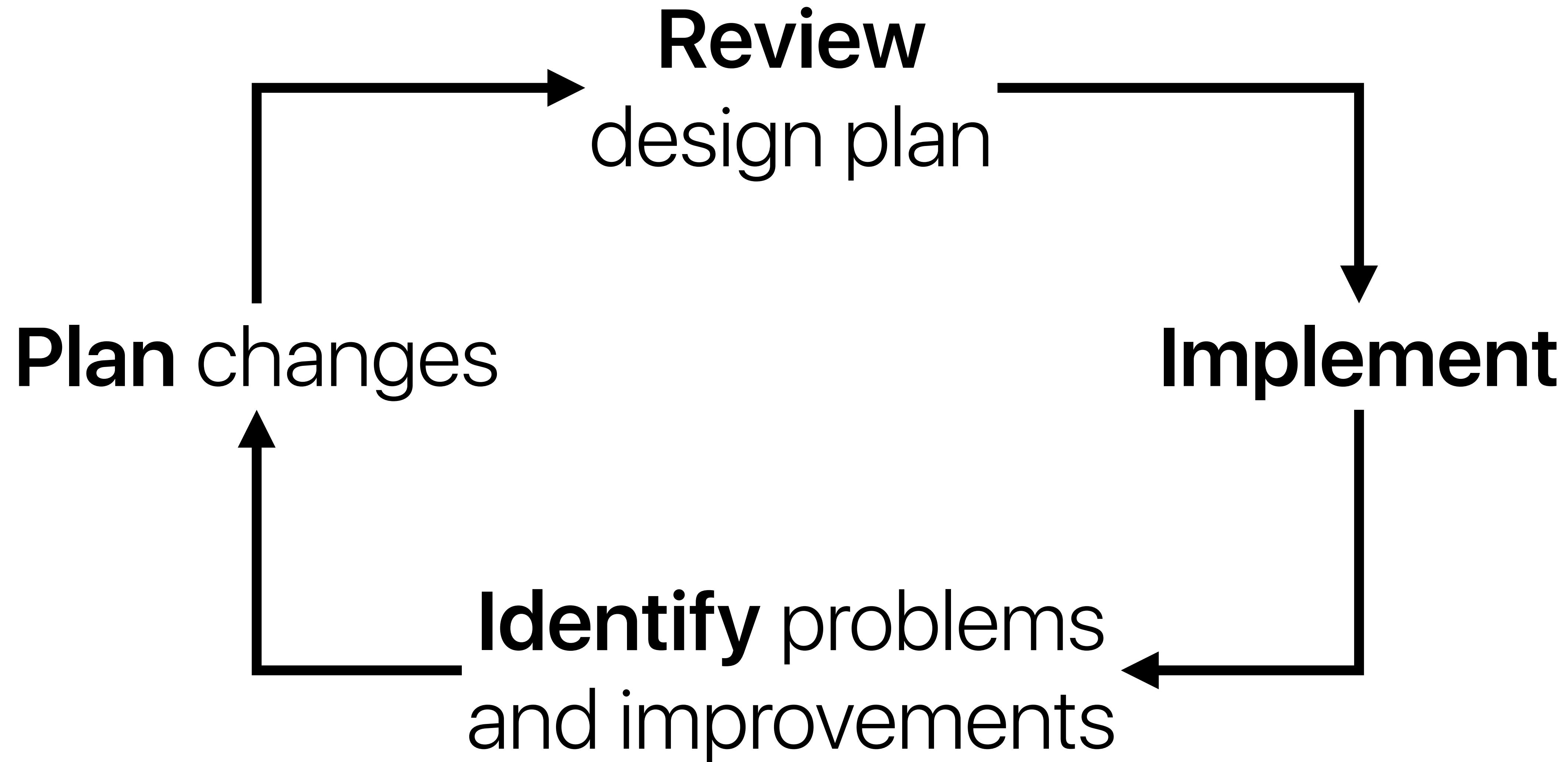
```
LET utility function = ...  
  
UNTIL good enough OR time runs out  
  DO another design, improving utility  
    UNTIL design is complete  
      WHILE design remains feasible  
        Make another design decision  
      END  
      Backtrack up design tree  
      Explore a path not searched before  
    END  
  END  
  Take best design  
END
```



"Phædrus began to pursue ... lateral truths; no longer the frontal truths of science, those toward which the discipline pointed, but the kind of truth you see laterally, out of the corner of your eye. In a laboratory situation, **when your whole procedure goes haywire**, when everything goes wrong or is indeterminate or is so screwed up by unexpected results you can't make head or tail out of anything, **you start looking laterally**. That's a word he later used to describe a growth of knowledge that doesn't move forward like an arrow in flight, but expands sideways, like an arrow enlarging in flight, or like the archer, discovering that although he has hit the bull's eye and won the prize, his head is on a pillow and the sun is coming in the window."

Build

Make it right



The screenshot shows a multi-panel code editor interface with three tabs open:

- fs.go — fs**: Contains the package declaration `package fs` and a single-line comment: `// Package fs abstracts the filesystem.
- real.go — fs**: Contains the package declaration `package fs`.
- virtual.go — fs**: Contains the following Go code:

```
1 package fs
2
3 import (
4     "bytes"
5     "os"
6     "path/filepath"
7     "strings"
8     "sync"
9     "time"
10 )
11
12 // NewVirtualFilesystem yields an in-memory filesystem.
13 func NewVirtualFilesystem() Filesystem {
14     return &virtualFilesystem{
15         files: map[string]*virtualFile{},
16     }
17 }
18
19 type virtualFilesystem struct {
20     mtx sync.RWMutex
21     files map[string]*virtualFile
22 }
23
24 func (fs *virtualFilesystem) Create(path string) (File, error) {
25     fs.mtx.Lock()
26     defer fs.mtx.Unlock()
27     // os.Create truncates any existing file. So we do, too.
28     f := &virtualFile{path, bytes.Buffer{}, time.Now(), time.Now()}
29     fs.files[path] = f
30     return f, nil
31 }
32
33 func (fs *virtualFilesystem) Open(path string) (File, error) {
34     fs.mtx.Lock()
```

The sidebar on the left shows a tree view of the project structure under the 'FS' folder, with 'virtual.go' currently selected. The bottom status bar indicates the current file is 'virtual.go' and shows other details like line count (Ln 1, Col 1), tab size (Tab Size: 4), and encoding (UTF-8).

The screenshot shows a dual-terminal setup within a dark-themed IDE interface. Both terminals are running a Go application.

Terminal 1 (Top):

```
fs.go — fs
import (
    "io"
)

```

Terminal 2 (Bottom):

```
lock_test.go — fs
run test | debug test
func TestLock(t *testing.T) {
    t.Parallel()

    lock := "TESTLOCK"
    for name, filesystem := range map[string]Filesystem{
        "virtual": NewVirtualFilesystem(),
        "real":    NewRealFilesystem(false),
    } {
        t.Run(name, func(t *testing.T) {
            r, existed, err := filesystem.Lock(lock)
            if err != nil {
                t.Fatalf("initial claim: %v", err)
            }
            if existed {
                t.Fatal("initial claim: lock file already exists")
            }
            if _, existed, _ = filesystem.Lock(lock); !existed {
                t.Fatal("second claim: want existed true, have false")
            }
            if err := r.Release(); err != nil {
                t.Fatalf("initial release: %v", err)
            }
            if err := r.Release(); err == nil {
                t.Fatal("second release: want error, have none")
            }
        })
        if filesystem.Exists(lock) {
            t.Errorf("%s: %s still exists", name, lock)
        }
        filesystem.Remove(lock)
    }
}
```

The code in Terminal 2 is a test function named `TestLock` that iterates over two types of filesystems: `virtual` and `real`. It performs a series of operations: attempting to lock a file, checking if it exists, releasing the lock, and then checking again to ensure the lock is released. It uses `testing.T` for parallel execution and `testing.Fatal` and `testing.Fatalf` for assertions.

The image shows two adjacent browser tabs, both titled "fs - The Go Programming Lang". The left tab displays the main package documentation for "fs", while the right tab shows detailed information about the "Filesystem" type.

Left Tab (Main Documentation):

- Header:** The Go Programming Language, Document
- Title:** Package fs
- Code Snippet:** import "github.com/oklog/oklog/pkg/fs"
- Links:** Overview, Index
- Section:** Overview ▾
 - Package fs abstracts the filesystem.
- Section:** Index ▾
 - type File
 - type Filesystem
 - func NewNopFilesystem() Filesystem
 - func NewRealFilesystem(mmap bool) Filesystem
 - func NewVirtualFilesystem() Filesystem
 - type Releaser
- Section:** Package files
 - fs.go
 - nop.go
 - real.go
 - virtual.go

Right Tab (Detailed Type Documentation):

- Header:** Hamburger
- Title:** type Filesystem
- Description:** Filesystem collects the operations we need from a filesystem.
- Code Snippet:**

```
type Filesystem interface {
    Create(path string) (File, error)
    Open(path string) (File, error)
    Remove(path string) error
    Rename(oldname, newname string) error
    Exists(path string) bool
    MkdirAll(path string) error
    Chtimes(path string, atime, mtime time.Time) error
    Walk(root string, walkFn filepath.WalkFunc) error
    Lock(path string) (r Releaser, existed bool, err error)
}
```
- Section:** func NewNopFilesystem
- Code Snippet:** func NewNopFilesystem() Filesystem
- Description:** NewNopFilesystem has methods that always succeed but do nothing.
- Section:** func NewRealFilesystem
- Code Snippet:** func NewRealFilesystem(mmap bool) Filesystem
- Description:** NewRealFilesystem yields a real disk filesystem with optional memory mapping for file reading.
- Section:** func NewVirtualFilesystem ¶
- Code Snippet:** func NewVirtualFilesystem() Filesystem
- Description:** NewVirtualFilesystem yields an in-memory filesystem.
- Section:** type Releaser
- Description:** Releaser is returned by Lock calls.

A **tracer bullet** is "a **spike** with the current architecture,
current technology set, current set of best practices
that results in production quality code."

—Wikipedia on Scrum

main.go — oklog

EXPLORER main.go x

ingest.go — oklog

EXPLORER main.go x ingest.go x store.go

ingest.go — oklog

OPEN EDITORS

- main.go cmd/oklog
- ingest.go cmd/ok...
- store.go cmd/okl...

OKLOG

- cmd
 - forward.go
 - ingest.go
 - ingeststore.go
 - main_test.go
 - main.go
 - query.go
 - store.go
 - stream.go
 - testsrv.go
- data
- dist
- pkg
 - cluster
 - fs
 - fs.go
 - lock_test.go
 - nop.go
 - real.go
 - virtual.go
 - group
 - ingest
 - testdata
 - api.go
 - conn_test.go

```
309     ingestLog,
310     *segmentFlushAge, *segmentFlushSize,
311     connectedClients.WithLabelValues("bulk"),
312     ingestWriterBytes, ingestWriterRecords, ingestWriterSyncs,
313     flushedSegmentAge, flushedSegmentSize,
314   )
315 }, func(error) {
316   bulkListener.Close()
317 })
318 g.Add(func() error {
319   mux := http.NewServeMux()
320   mux.HandleFunc("/ingest/", http.StripPrefix("/ingest", ingest.NewAPI(
321     peer,
322     ingestLog,
323     *segmentPendingTimeout,
324     failedSegments,
325     committedSegments,
326     committedBytes,
327     apiDuration,
328   )))
329   registerMetrics(mux)
330   registerProfile(mux)
331   return http.Serve(apiListener, mux)
332 }, func(error) {
333   apiListener.Close()
334 })
335 }
336 {
337   cancel := make(chan struct{})
338   g.Add(func() error {
339     return interrupt(cancel)
340   }), func(error) {
341     close(cancel)
342   })
343 }
344 return g.Run()
345 }
```

store.go — oklog

EXPLORER main.go ingest.go store.go x

store.go — oklog

EXPLORER main.go ingest.go store.go x

store.go — oklog

OPEN EDITORS main.go cmd/oklog ingest.go cmd/oklog store.go cmd/oklog

OKLOG cmd oklog forward.go ingest.go ingeststore.go main_test.go main.go query.go store.go stream.go testsvc.go

data dist pkg cluster fs fs.go lock_test.go nop.go real.go virtual.go group ingest testdata api.go conn_test.go

```
239 // Execution group.
240 var g group.Group
241 {
242     cancel := make(chan struct{})
243     g.Add(func() error {
244         <-cancel
245         return peer.Leave(time.Second)
246     }, func(error) {
247         close(cancel)
248     })
249 }
250 for i := 0; i < *segmentConsumers; i++ {
251     c := store.NewConsumer(
252         peer,
253         timeoutClient,
254         *segmentTargetSize,
255         *segmentTargetAge,
256         *segmentReplicationFactor,
257         consumedSegments,
258         consumedBytes,
259         replicatedSegments.WithLabelValues("egress"),
260         replicatedBytes.WithLabelValues("egress"),
261         log.With(logger, "component", "Consumer"),
262     )
263     g.Add(func() error {
264         c.Run()
265         return nil
266     }, func(error) {
267         c.Stop()
268     })
269 }
270 {
271     c := store.NewCompacter(
272         storeLog,
273         *segmentTargetSize,
274         *segmentRetain,
275         *segmentPurge,
276         compactDuration
277     )
278 }
```

Ln 187, Col 37 Tab Size: 4 UTF-8 LF Go 😊

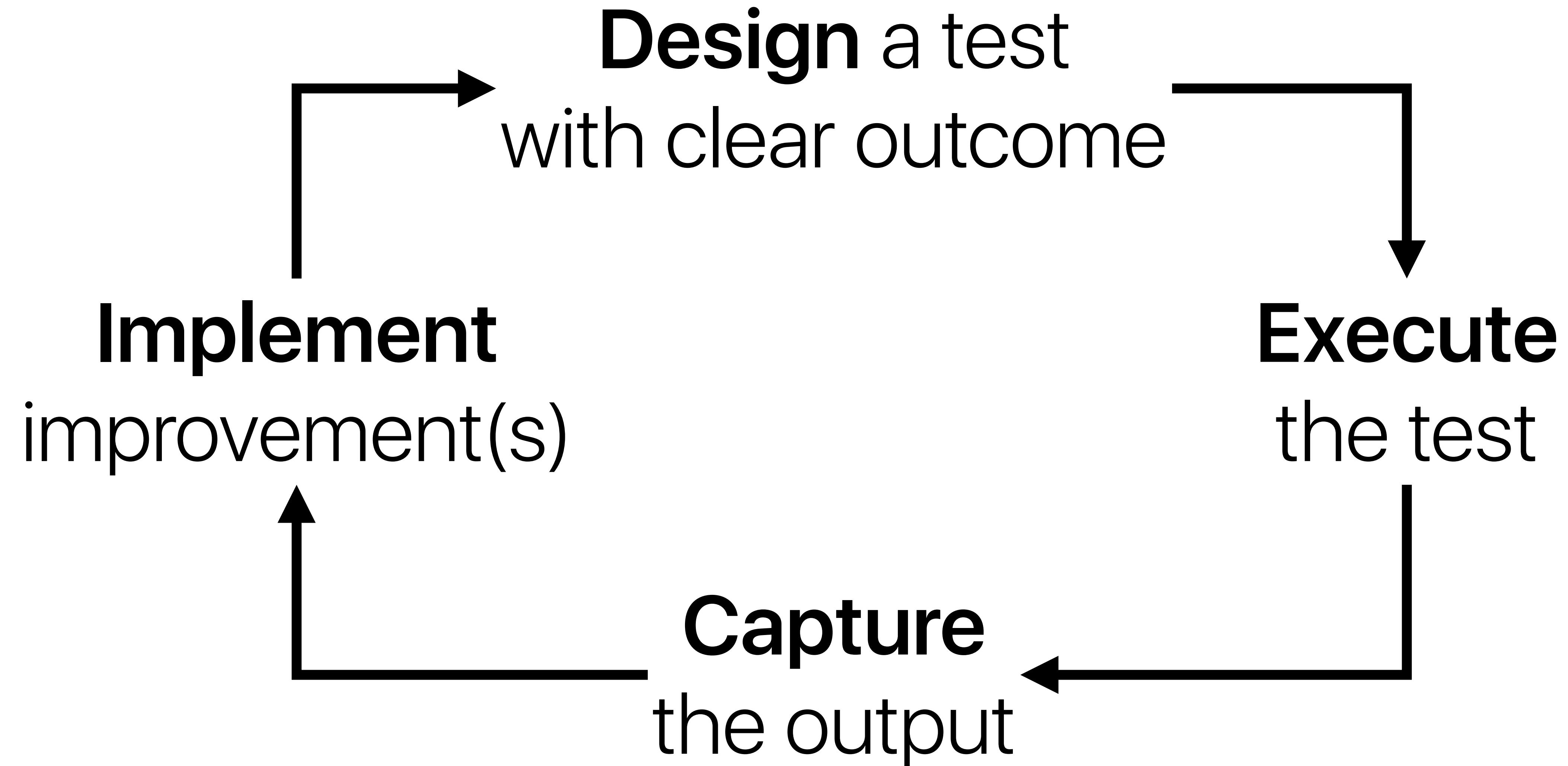
```
    g.Add(func() error {
        return nil
    })
}

// ClusterPeer models cluster.Peer.
type ClusterPeer interface {
    State() map[string]interface{}
}

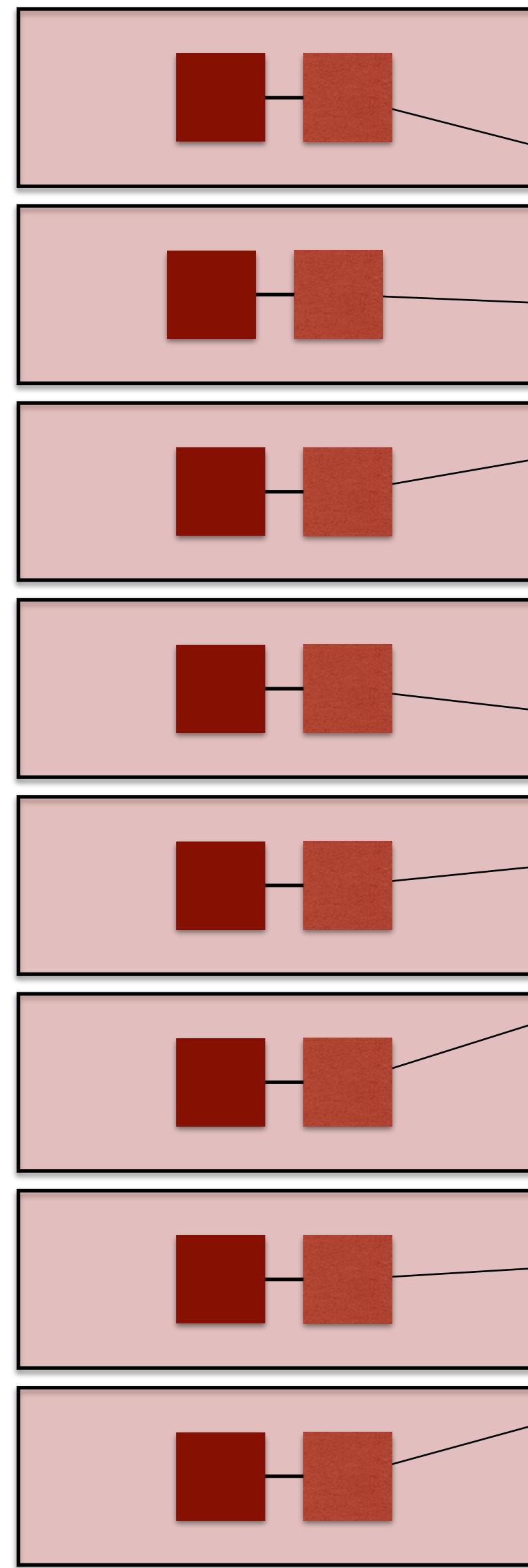
// NewAPI returns a usable ingest API.
func NewAPI(
    peer ClusterPeer,
    log Log,
    pendingSegmentTimeout time.Duration,
    failedSegments, committedSegments, committedBytes prometheus.Counter,
    duration *prometheus.HistogramVec,
) *API {
    a := &API{
        peer:             peer,
        log:              log,
        timeout:         pendingSegmentTimeout,
```

Optimize

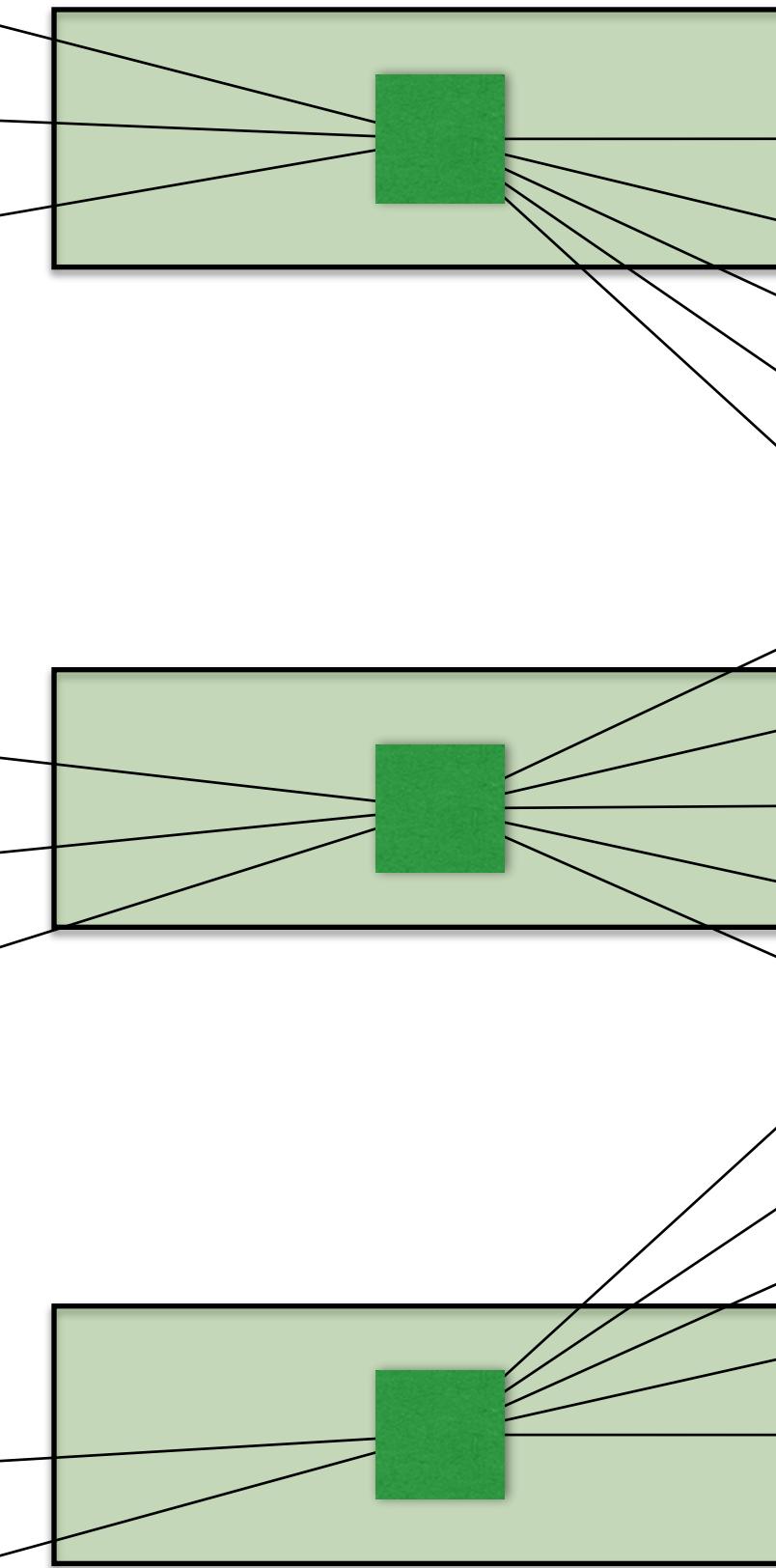
Make it fast



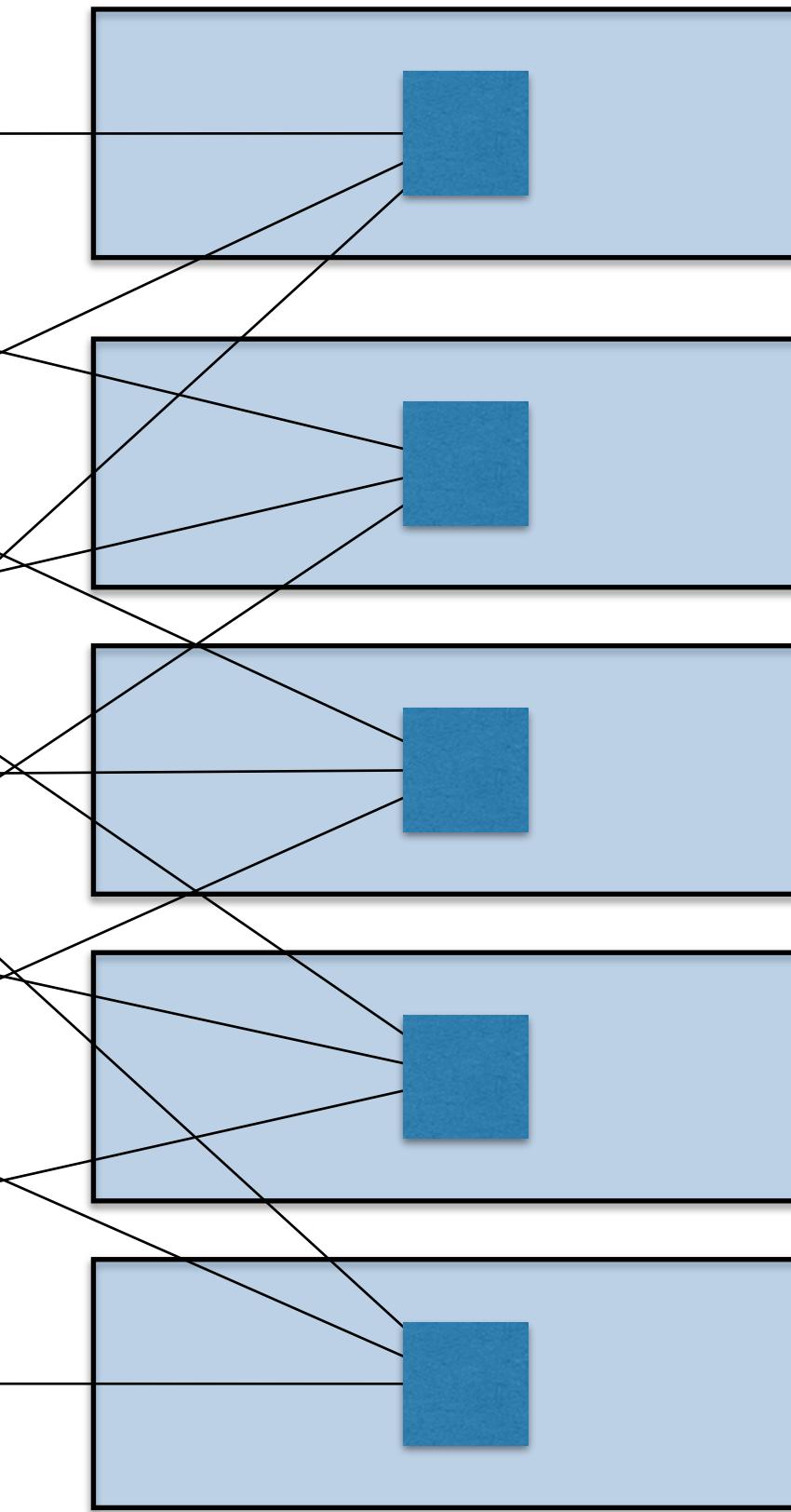
Forward



Ingest



Store



blep ~ _

01:fish*

01:fish*

0831h

0831h

blep ~ _

01:fish*

01:fish*

0834h

0834h

Instrumentation > Logs



USE

Brendan Gregg

Utilization
Saturation
Error rate

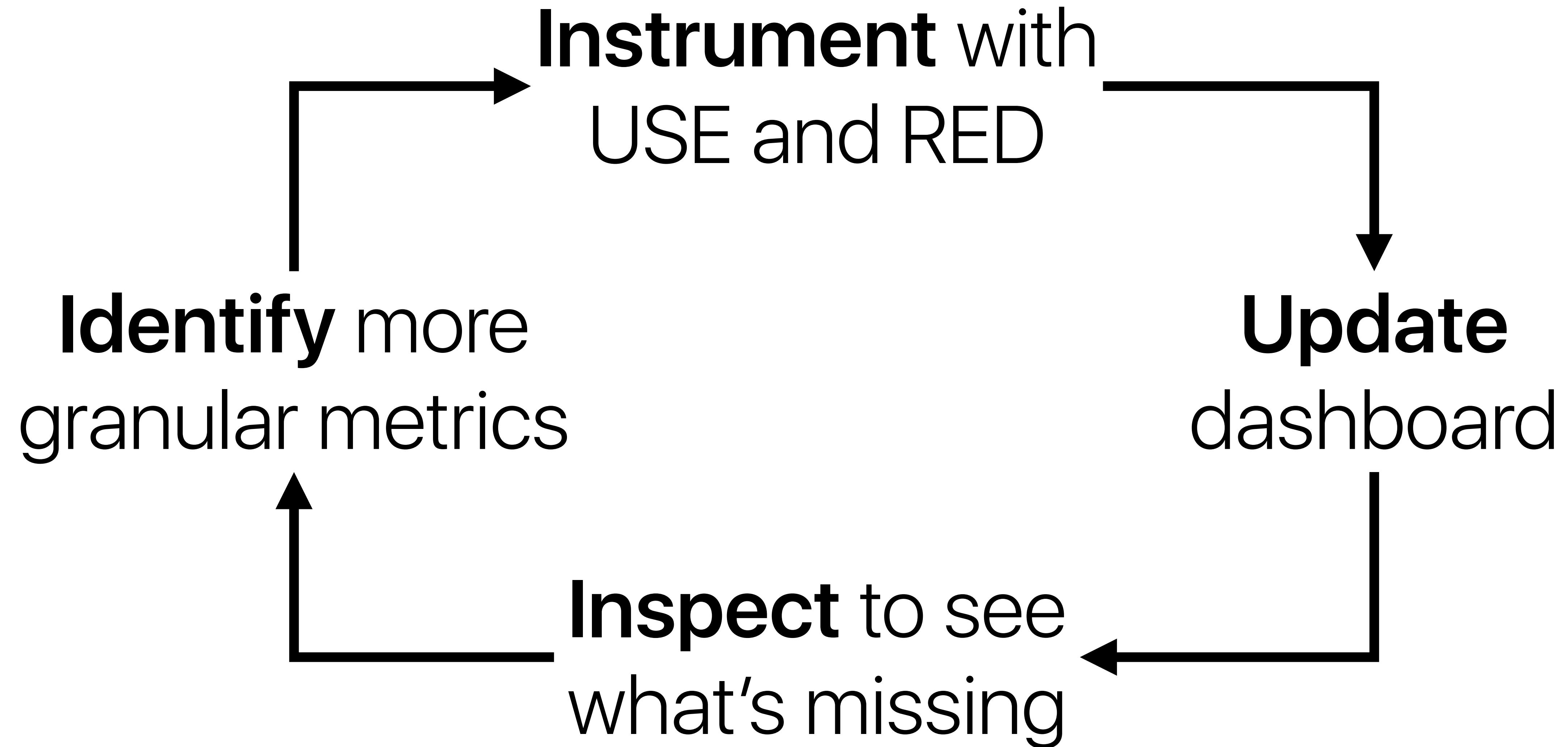
(Resources e.g. queues)

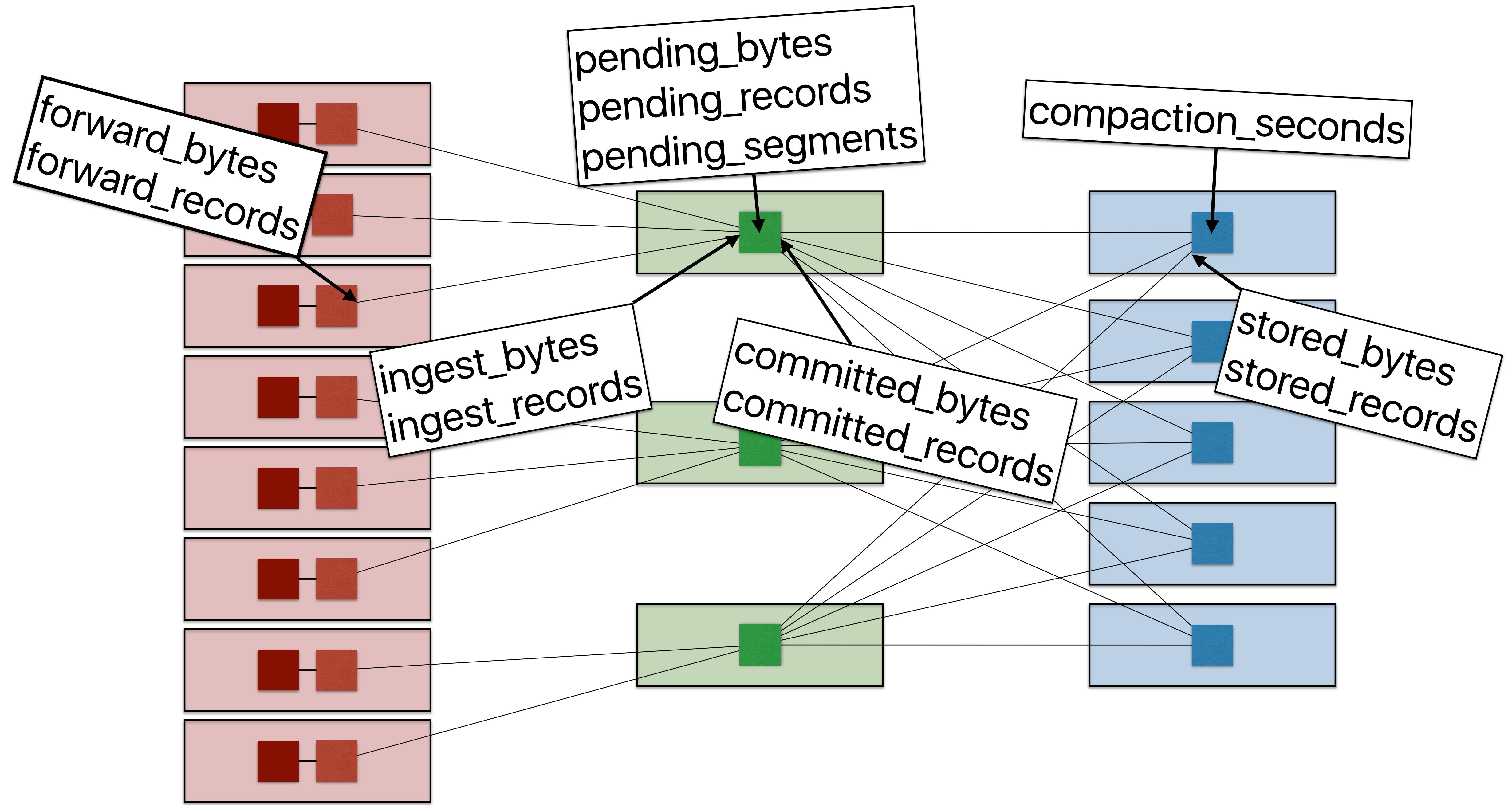
RED

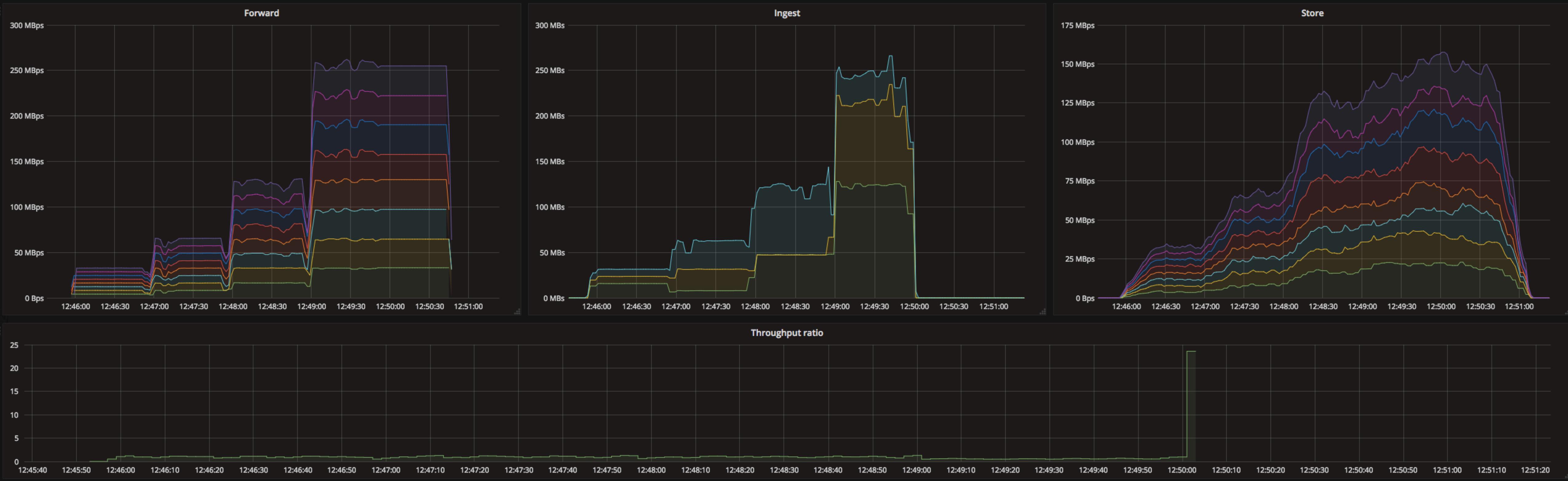
Tom Wilkie

Request rate
Error rate
Duration

(Workers e.g. web servers)







Optimizing ingestion



xxxxx

@peterbourgon

Friends, operators of clusters! I am conducting research! What is the size of your cluster, and your approximate log volume?

6:24 PM - 22 Dec 2016

n = 8

min = 1KBps / node

mean = **~5MBps** / node

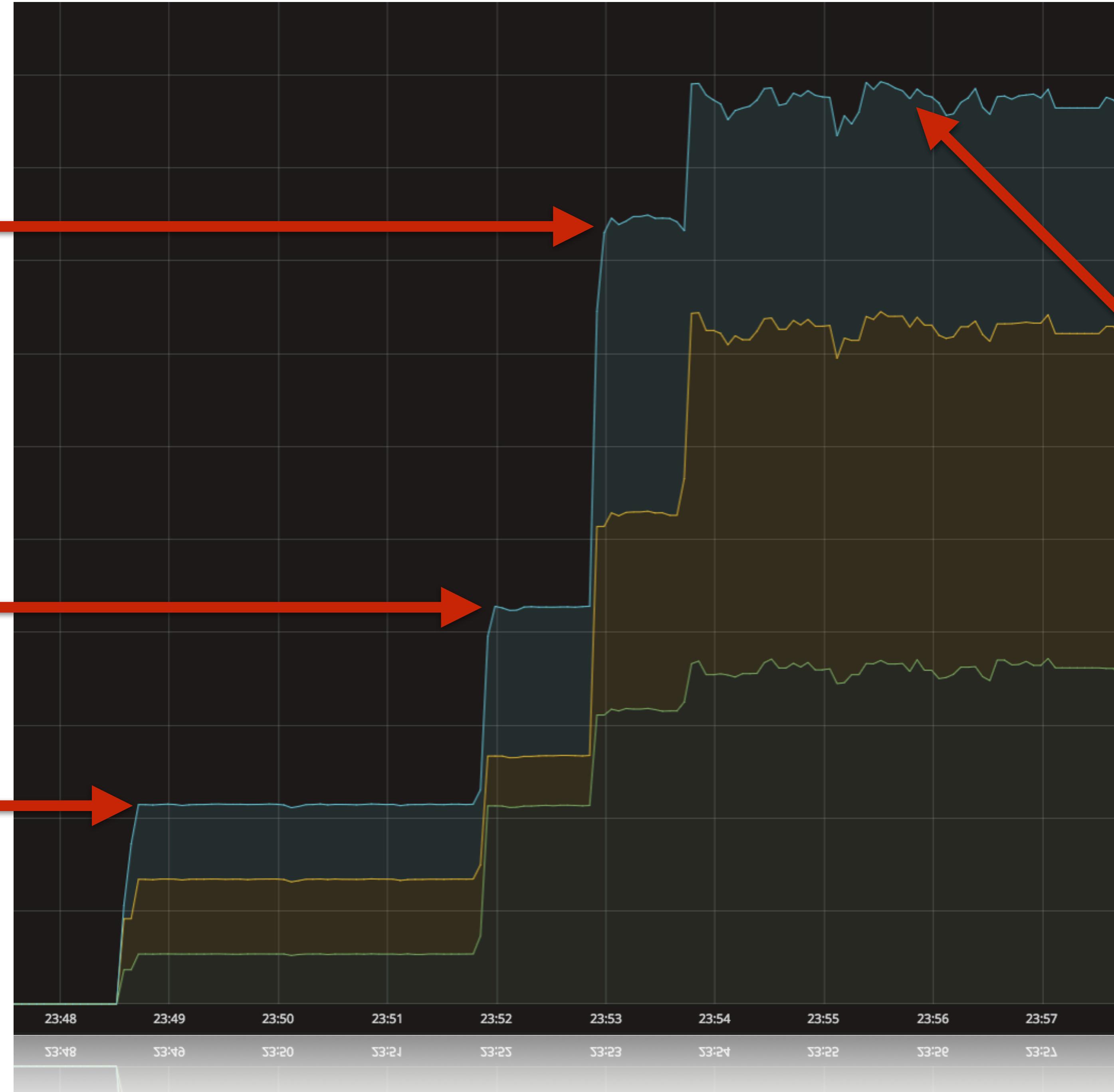
max = 20MBps / node (!)

16 MBps

8 MBps

4 MBps

**not
32 MBps
:(**



CPU [██████████] 62.6% Tasks: 29, 12 thr; 4 running
 Mem [██████████] 111/1000MB Load average: 0.87 0.64 0.41
 Swp [] 0/0MB Uptime: 07:18:34

PID	USER	SHR	S	CPU%	MEM%	TIME+	Command
1	root						
30116	peter	3120	S	0.0	0.5	0:01.62	/sbin/init
30125	peter						g:7651
30123	peter						.org:7651
30122	peter						.org:7651
30119	peter	6316	S	26.1	1.2	0:21.35	.org:7651
30118	peter	6316	S	25.6	1.2	0:07.54	.org:7651
30115	peter	6316	S	25.6	1.2	0:07.54	.org:7651
30153	peter	6316	S	0.0	1.2	0:13.66	.org:7651
30124	peter	6316	S	0.0	1.2	0:00.00	.org:7651
30121	peter	6316	S	0.0	1.2	0:00.00	.org:7651
30120	peter	6316	S	0.0	1.2	0:00.00	.org:7651
19432	peter	6316	S	0.0	1.2	0:00.00	.org:7651
19225	peter	6316	S	0.0	1.2	0:00.00	.org:7651
18947	peter	6316	S	0.0	1.2	0:00.00	.org:7651
17665	peter	6316	S	0.0	1.2	0:00.00	.org:7651
17512	root	6316	S	0.0	1.2	0:00.00	.org:7651
844	root	6316	S	0.0	1.2	0:00.00	.org:7651
12432	root	6316	S	0.0	1.2	0:00.00	.org:7651
12443	peter	6316	S	0.0	1.2	0:00.11	.org:7651
12444	peter	6316	S	0.0	1.2	0:41.02	./oklog testsvc -id for
30847	peter	6956	R	17.8	1.2	0:14.24	./oklog testsvc -id
769	Debian-e	6956	S	0.0	1.2	0:00.38	./oklog testsvc -id
520	root	6956	S	0.0	1.2	0:00.00	./oklog testsvc -id
513	root	6956	S	0.0	1.2	0:11.41	./oklog testsvc -id
512	root						
533	root						
532	root						
531	root						
505	messagebus						
502	root						
500	root						
499	daemon						
498	root						
484	statd						
475	root						
166	root						
161	root						
		20	0	41164	11192	10904	S 0.0 1.1 0:00.73 /lib/systemd/systemd-journald

```
hz := float64(time.Second) / float64(*rate)
for range time.Tick(time.Duration(hz)) {
    fmt.Fprintf(os.Stdout,
        "%s %s\n",
        time.Now().Format(time.RFC3339),
        records[_rant.Intn(len(records))]),
}
}
```

```
var (
    recordsPerCycle = 1
    timePerCycle    = time.Second / *rate
)
for timePerCycle < 50*time.Millisecond {
    recordsPerCycle *= 2
    timePerCycle *= 2
}
```

```
var count int
for range time.Tick(timePerCycle) {
    for i := 0; i < recordsPerCycle; i++ {
        fmt.Fprintf(os.Stdout,
                    "%s %s\n",
                    time.Now().Format(time.RFC3339),
                    records[count%len(records)]),
    }
    count++
}
}
```

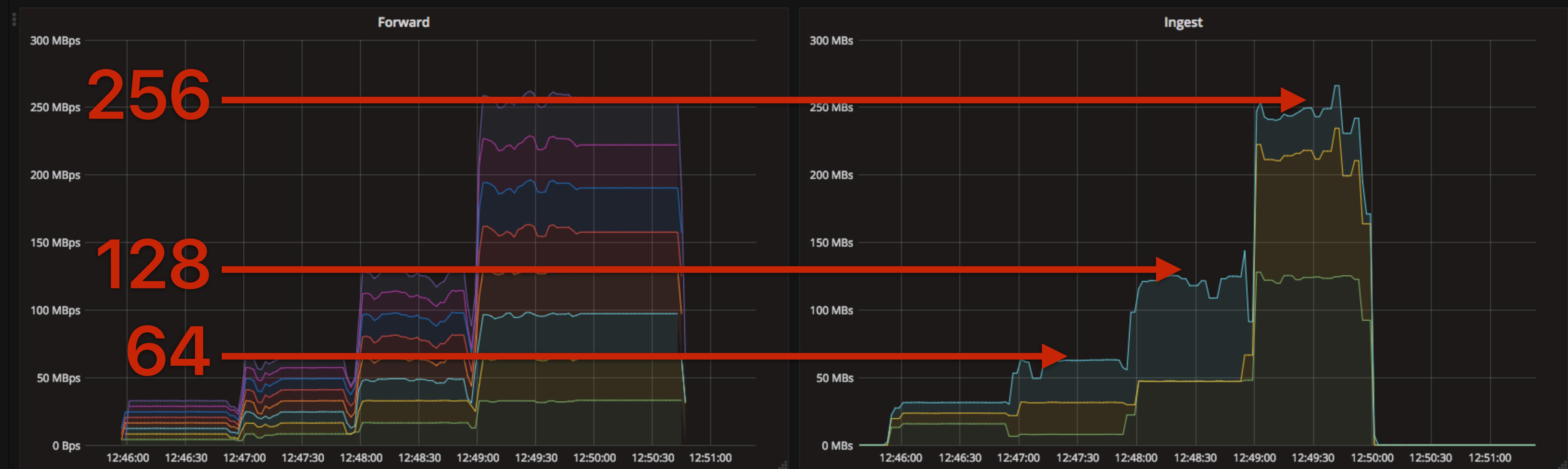
CPU	Mem	Swp		64.6% [110/1000MB]	Tasks: 29, 12 thr; 5 running	Load average: 1.16 0.69 0.55	Uptime: 07:34:49				
PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
1	root	20	0	28572	4744	3120	S	0.0	0.5	0:01.65	/sbin/init
19432	peter	20	0	26712	5132	4512	S	0.4	0.5	0:17.80	- fish ./spam.fish
19225	peter	20	0	26712	5176	4556	S	0.0	0.5	0:18.01	- fish ./spam.fish
18947	peter	20	0	26712	5148	4528	S	0.0	0.5	0:18.01	- fish ./spam.fish
17665	peter	20	0	10704	340	0	S	0.0	0.0	0:00.00	- ssh-agent -c
17512	root	20	0	4256	1112	1012	S	0.0	0.1	0:00.38	- runsvdir -P /etc/service log:
6704	peter	20	0	18840	14764	7848	S	18.5	1.4	0:25.23	- ./oklog forward tcp://ingest1.bourgon.org:7651 tcp://ingest2.bourgon.org:7651 tcp://ingest3.bourgon.org:7651
6713	peter	20	0	18840	14764	7848	S	0.0	1.4	0:04.67	- ./oklog forward tcp://ingest1.bourgon.org:7651 tcp://ingest2.bourgon.org:7651 tcp://ingest3.bourgon.org:7651
6711	peter	20	0	18840	14764	7848	S	18.9	1.4	0:10.57	- ./oklog forward tcp://ingest1.bourgon.org:7651 tcp://ingest2.bourgon.org:7651 tcp://ingest3.bourgon.org:7651
6710	peter	20	0	18840	14764	7848	S	0.0	1.4	0:00.01	- ./oklog forward tcp://ingest1.bourgon.org:7651 tcp://ingest2.bourgon.org:7651 tcp://ingest3.bourgon.org:7651
6707	peter	20	0	18840	14764	7848	S	0.0	1.4	0:00.00	- ./oklog forward tcp://ingest1.bourgon.org:7651 tcp://ingest2.bourgon.org:7651 tcp://ingest3.bourgon.org:7651
6706	peter	20	0	18840	14764	7848	S	0.0	1.4	0:03.82	- ./oklog forward tcp://ingest1.bourgon.org:7651 tcp://ingest2.bourgon.org:7651 tcp://ingest3.bourgon.org:7651
6703	peter	20	0	16728	12004	6892	R	33.9	1.2	0:39.61	- ./oklog testsvc -id forward1 -size 1024 -rate 8000
6714	peter	20	0	16728	12004	6892	S	0.0	1.2	0:03.75	- ./oklog testsvc -id forward1 -size 1024 -rate 8000
6712	peter	20	0	16728	12004	6892	R	12.4	1.2	0:13.87	- ./oklog testsvc -id forward1 -size 1024 -rate 8000
6709	peter	20	0	16728	12004	6892	S	0.0	1.2	0:00.00	- ./oklog testsvc -id forward1 -size 1024 -rate 8000
6708	peter	20	0	16728	12004	6892	R	11.2	1.2	0:12.12	- ./oklog testsvc -id forward1 -size 1024 -rate 8000
844	root	20	0	55184	5400	4728	S	0.0	0.5	0:00.08	- /usr/sbin/sshd -D
12432	root	20	0	82684	5908	5040	S	0.0	0.6	0:00.00	- sshd: peter [priv]
12443	peter	20	0	82684	4048	3192	S	0.0	0.4	0:00.35	- sshd: peter@pts/1
12444	peter	20	0	170M	6380	4696	S	0.0	0.6	0:01.73	- -fish
7624	peter	20	0	25324	4476	2868	R	0.4	0.4	0:00.05	- htop
769	Debian-ex	20	0	53252	3260	2612	S	0.0	0.3	0:00.00	- /usr/sbin/exim4 -bd -q30m
520	root	20	0	12844	1880	1736	S	0.0	0.2	0:00.00	- /sbin/agetty --noclear tty1 linux
513	root	20	0	4256	1672	1528	S	0.0	0.2	0:00.00	- /usr/sbin/acpid
512	root	20	0	256M	3600	2800	S	0.0	0.4	0:00.09	- /usr/sbin/rsyslogd -n
533	root	20	0	256M	3600	2800	S	0.0	0.4	0:00.03	- /usr/sbin/rsyslogd -n
532	root	20	0	256M	3600	2800	S	0.0	0.4	0:00.00	- /usr/sbin/rsyslogd -n
531	root	20	0	256M	3600	2800	S	0.0	0.4	0:00.04	- /usr/sbin/rsyslogd -n
505	messagebu	20	0	42124	3296	2916	S	0.0	0.3	0:00.12	- /usr/bin/dbus-daemon --system --address=systemd: --nofork --nopidfile --systemd-activation
502	root	20	0	19856	2592	2320	S	0.0	0.3	0:00.23	- /lib/systemd/systemd-logind
500	root	20	0	27476	2784	2540	S	0.0	0.3	0:00.12	- /usr/sbin/cron -f
499	daemon	20	0	19024	1752	1584	S	0.0	0.2	0:00.00	- /usr/sbin/atd -f
498	root	20	0	27568	228	4	S	0.0	0.0	0:00.00	- /usr/sbin/rpc.idmapd
484	statd	20	0	37280	2964	2372	S	0.0	0.3	0:00.00	- /sbin/rpc.statd
475	root	20	0	37080	2700	2272	S	0.0	0.3	0:00.13	- /sbin/rpcbind -w
166	root	20	0	40816	3324	2800	S	0.0	0.3	0:00.02	- /lib/systemd/systemd-udevd
161	root	20	0	41164	11196	10908	S	0.0	1.1	0:00.74	- /lib/systemd/systemd-journald

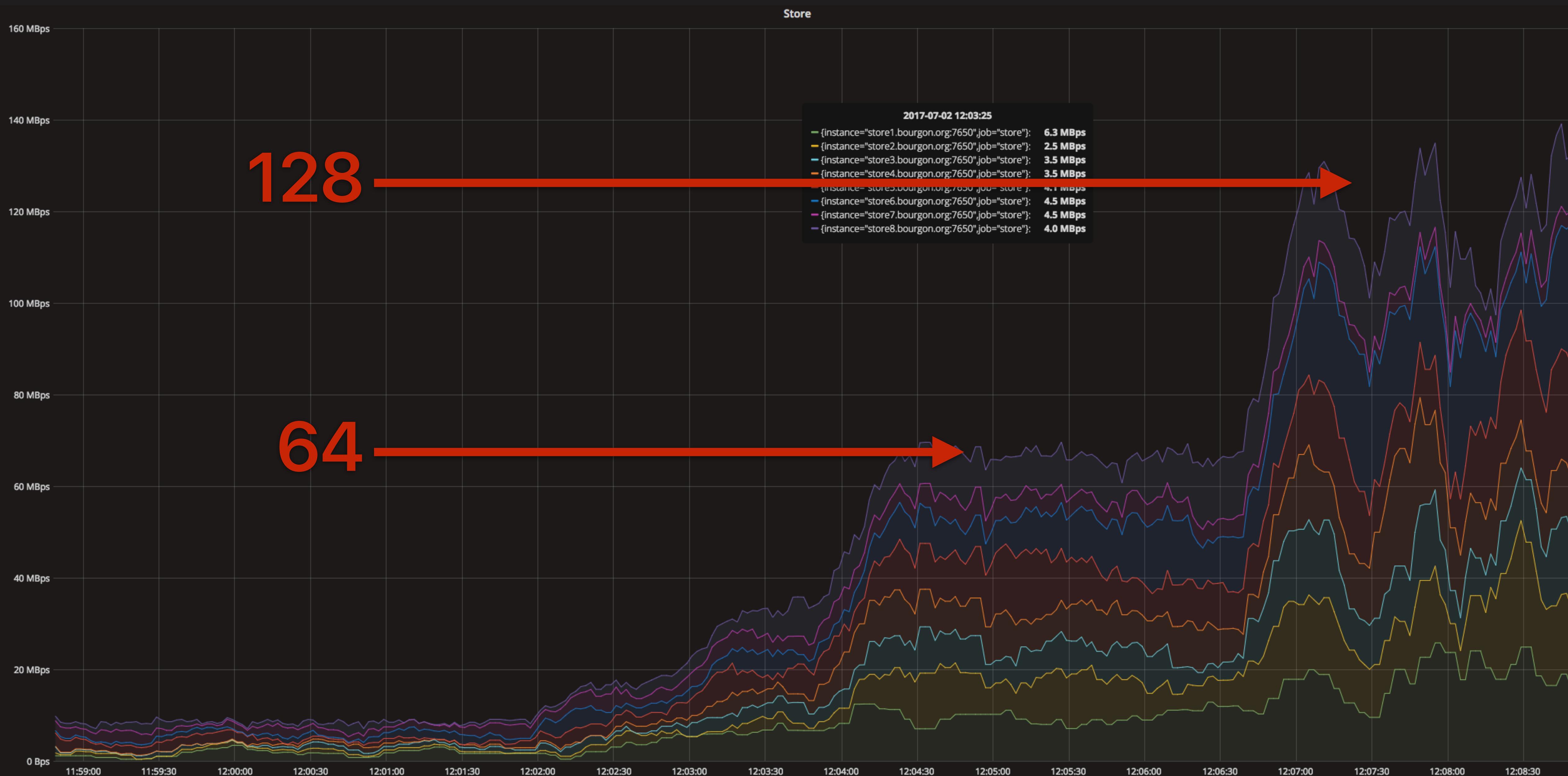
	S	0.0	1.4	0:03.82
92	R	33.9	1.2	0:39.61
92	S	0.0	1.2	0:03.75
92	R	12.4	1.2	0:13.87
92	S	0.0	1.2	0:00.00

```
L ./oklog forward tcp://ing  
./oklog testsvc -id forward1  
E ./oklog testsvc -id forward1  
E ./oklog testsvc -id forward1  
E ./oklog testsvc -id forward1
```

```
./oklog testsvc -id forward1 -size 1024 -rate 8000  
E ./oklog testsvc -id forward1 -size 1024 -rate 8000  
E ./oklog testsvc -id forward1 -size 1024 -rate 8000  
E ./oklog testsvc -id forward1 -size 1024 -rate 8000  
E ./oklog testsvc -id forward1 -size 1024 -rate 8000
```

```
./oklog testsvc -id forward1 -size 256 -rate 16000  
E ./oklog testsvc -id forward1 -size 256 -rate 16000  
E ./oklog testsvc -id forward1 -size 256 -rate 16000  
E ./oklog testsvc -id forward1 -size 256 -rate 16000  
E ./oklog testsvc -id forward1 -size 256 -rate 16000
```





```
1 package main  
2  
3 import (  
4     "flag"  
5     "fmt"  
6     "net"  
7     "net/http"  
8     "net/http/pprof"  
9     "net/url"  
10    "os"  
11    "os/signal"  
12    "runtime"  
13    "strconv"  
14    "strings"  
15    "syscall"  
16    "text/tabwriter"  
17
```

```
148
149 func registerProfile(mux *http.ServeMux) {
150     mux.HandleFunc("/debug/pprof/", pprof.Index)
151     mux.HandleFunc("/debug/pprof/cmdline", pprof.Cmdline)
152     mux.HandleFunc("/debug/pprof/profile", pprof.Profile)
153     mux.HandleFunc("/debug/pprof/symbol", pprof.Symbol)
154     mux.HandleFunc("/debug/pprof/trace", pprof.Trace)
155     mux.Handle("/debug/pprof/block", pprof.Handler("block"))
156     mux.Handle("/debug/pprof/goroutine", pprof.Handler("goroutine"))
157     mux.Handle("/debug/pprof/heap", pprof.Handler("heap"))
158     mux.Handle("/debug/pprof/threadcreate", pprof.Handler("threadcreate"))
159 }
160
161 func usageFor(fs *flag.FlagSet, short string) func() {
162     return func() {
163         fmt.Fprintf(os.Stderr, "USAGE\n")
164         fmt.Fprintf(os.Stderr, "  - %s\n", short)
165         fmt.Fprintf(os.Stderr, "\n")
```

blep ~ _

01:fish*

0957h

*nsit:J6

0957h

blep ~ _

01:fish*

1011h

*nsit:J6

JLJ01

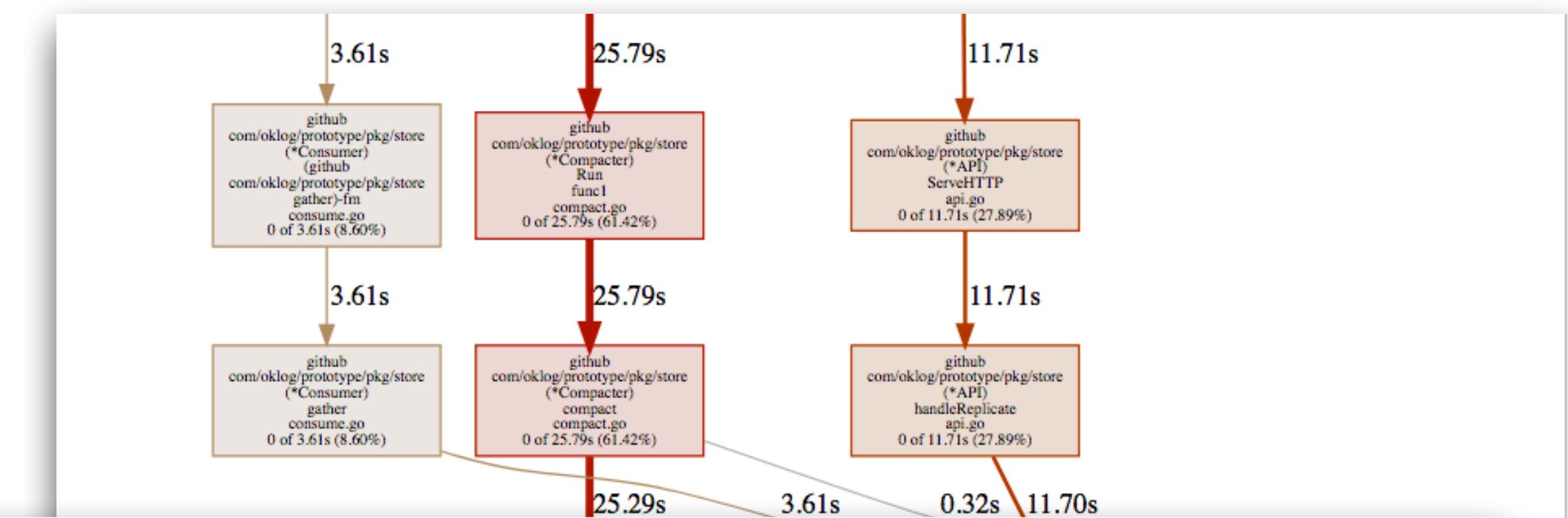
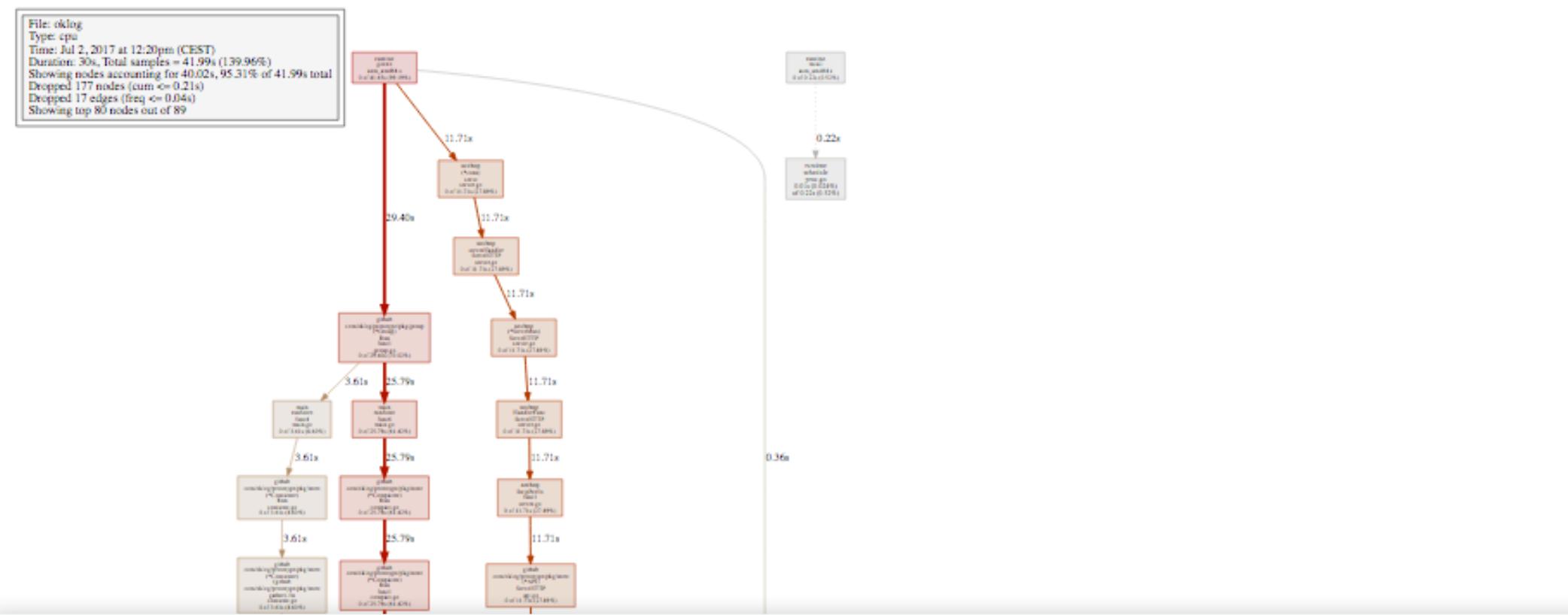
```
ingest1 ~ go tool pprof -alloc_objects bin/oklog http://ingest1.bourgon.org:7650/debug/pprof/heap
Fetching profile from http://ingest1.bourgon.org:7650/debug/pprof/heap
Saved profile in /home/peter/pprof/pprof.oklog.ingest1.bourgon.org:7650.alloc_objects.alloc_space.inuse_objects.inuse_space.007.pb.gz
Entering interactive mode (type "help" for commands)
(pprof) _
```

01:ssh*

1009h

01:ssh*

1009h



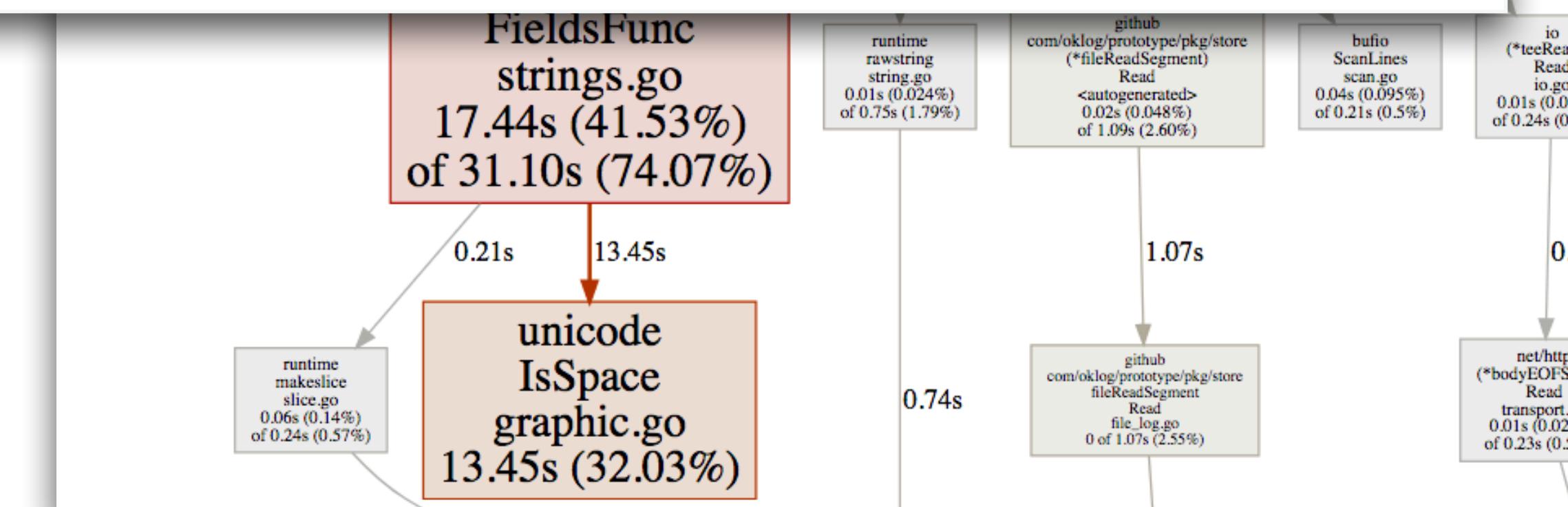
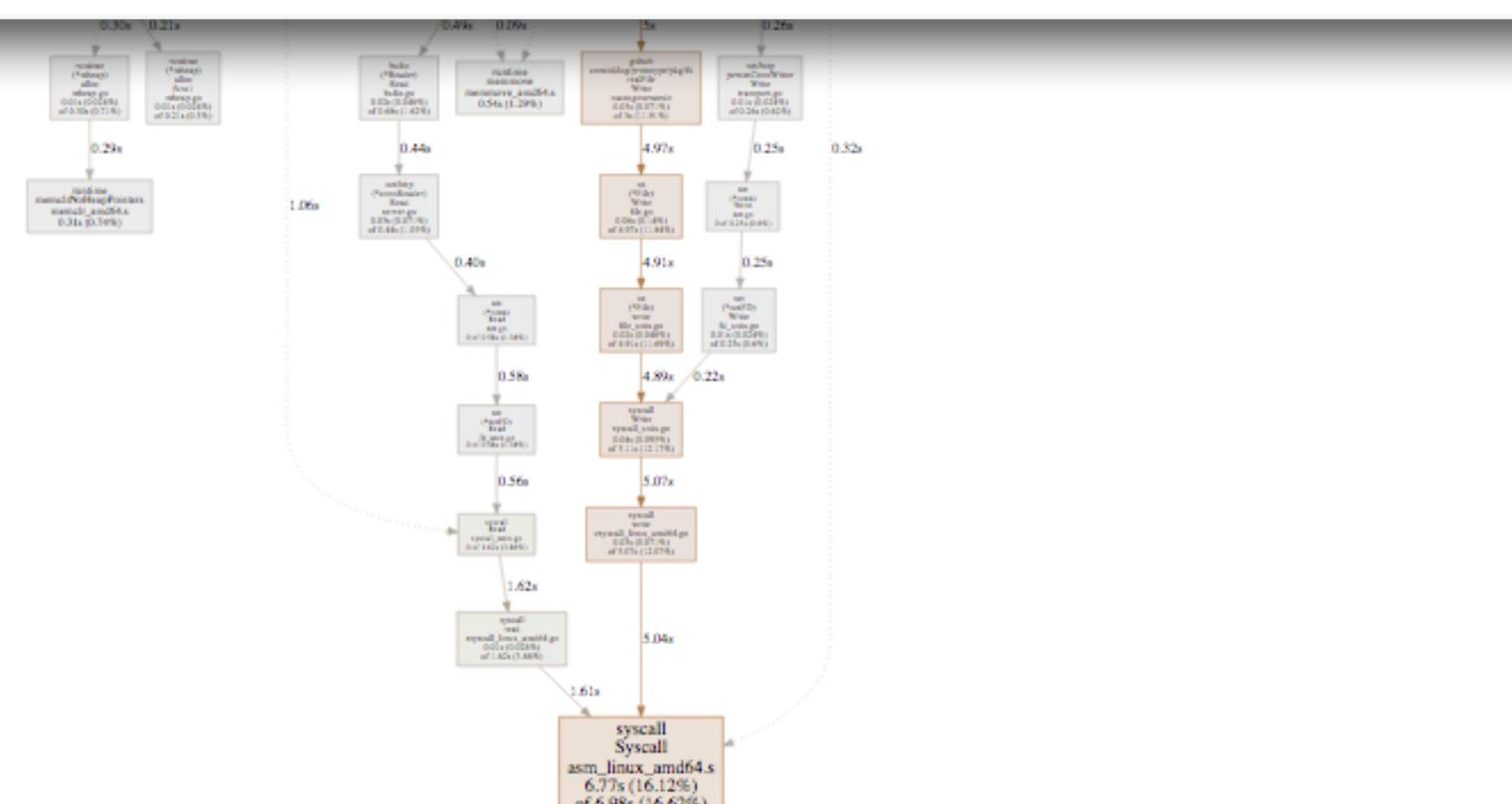
(pprof) top

Showing nodes accounting for 38.51s, 91.71% of 41.99s total

Dropped 177 nodes (cum <= 0.219)

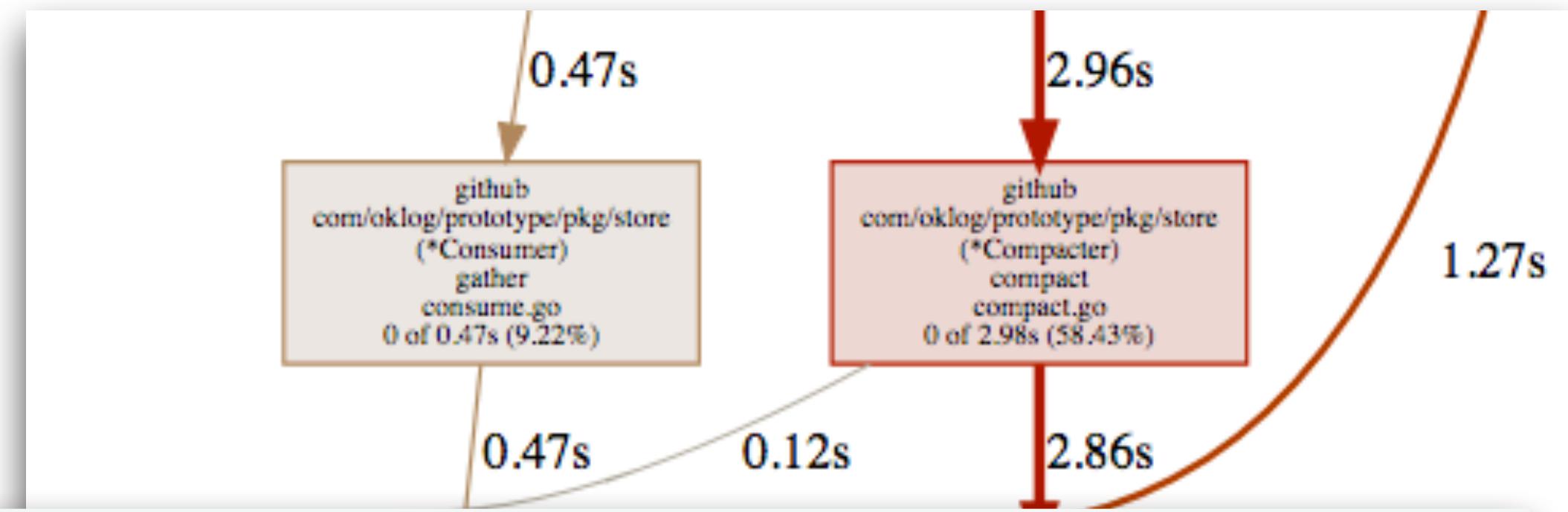
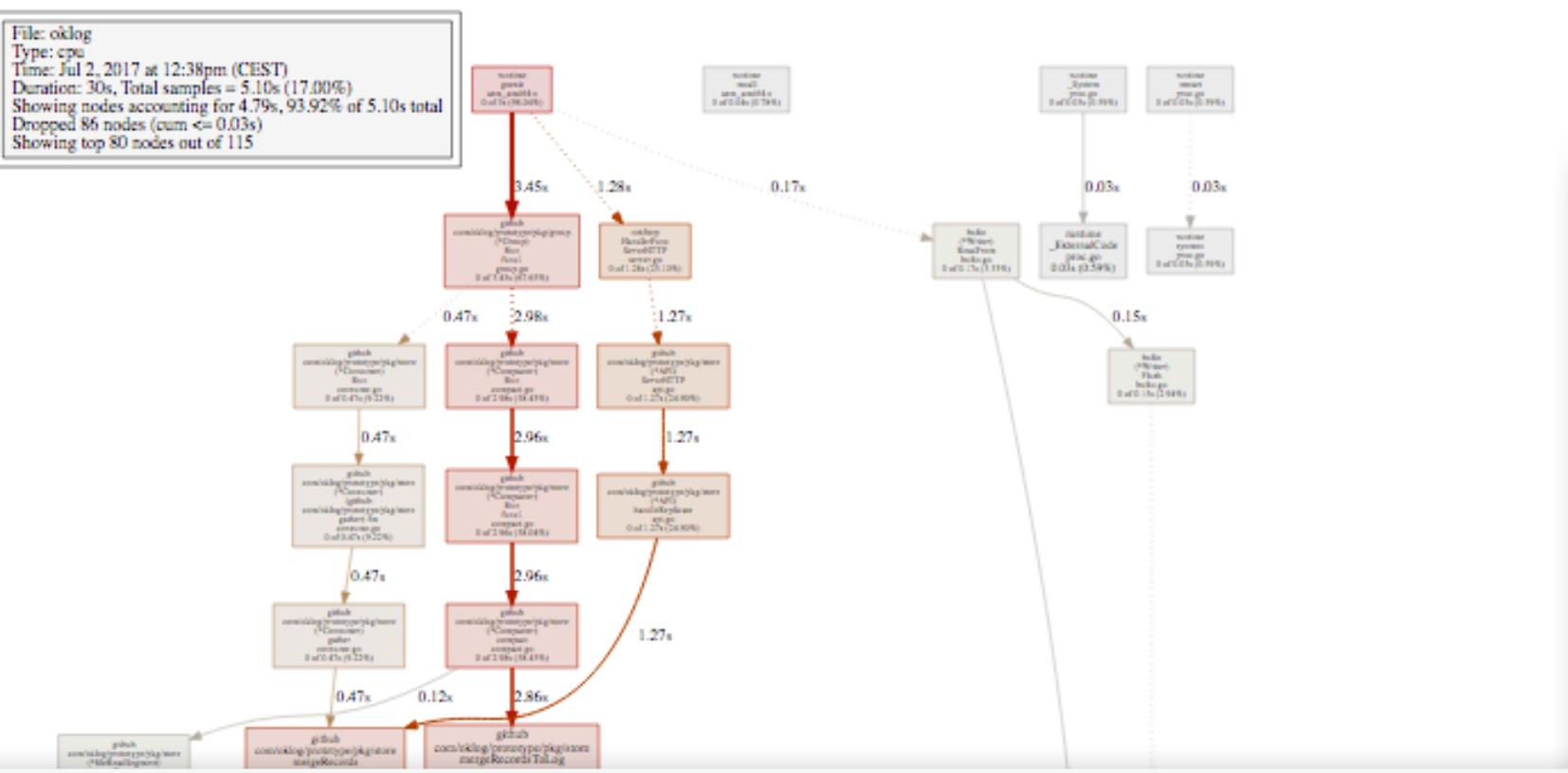
Showing top 5 nodes out of 8

flat	flat%	sum%	cum	cum%	
17.44s	41.53%	41.53%	31.10s	74.07%	strings.FieldsFunc
13.45s	32.03%	73.57%	13.45s	32.03%	unicode.IsSpace
6.77s	16.12%	89.69%	6.98s	16.62%	syscall.Syscall
0.54s	1.29%	90.97%	0.54s	1.29%	runtime.memmove
0.31s	0.74%	91.71%	0.31s	0.74%	runtime.memclrNoHeapPointer



```
// Advance takes the next record from the reader.
advance := func(i int) error {
    if ok[i] = scanner[i].Scan(); ok[i] {
        record[i] = scanner[i].Bytes()
        id[i] = strings.Fields(record[i])[0]
    }
    ...
}
```

```
// Advance takes the next record from the reader.
advance := func(i int) error {
    if ok[i] = scanner[i].Scan(); ok[i] {
        record[i] = scanner[i].Bytes()
        id[i] = record[i][:ulid.EncodedSize]
    }
    ...
}
```



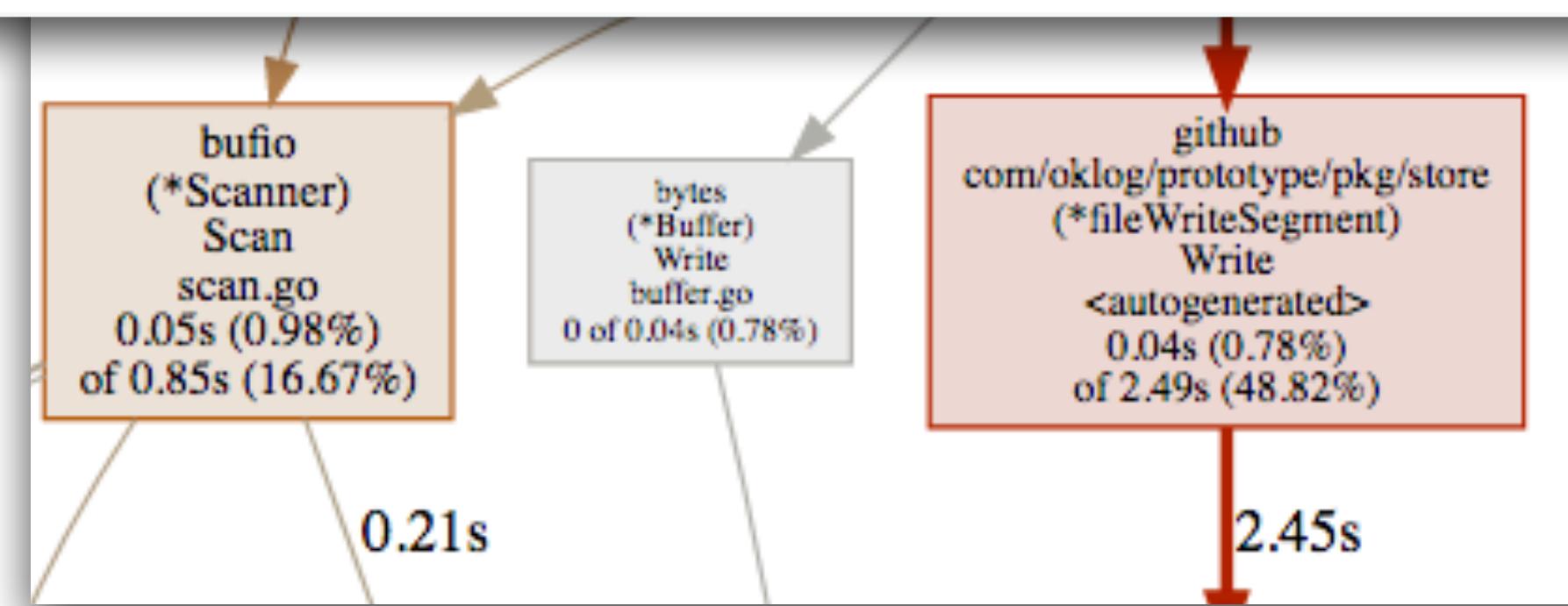
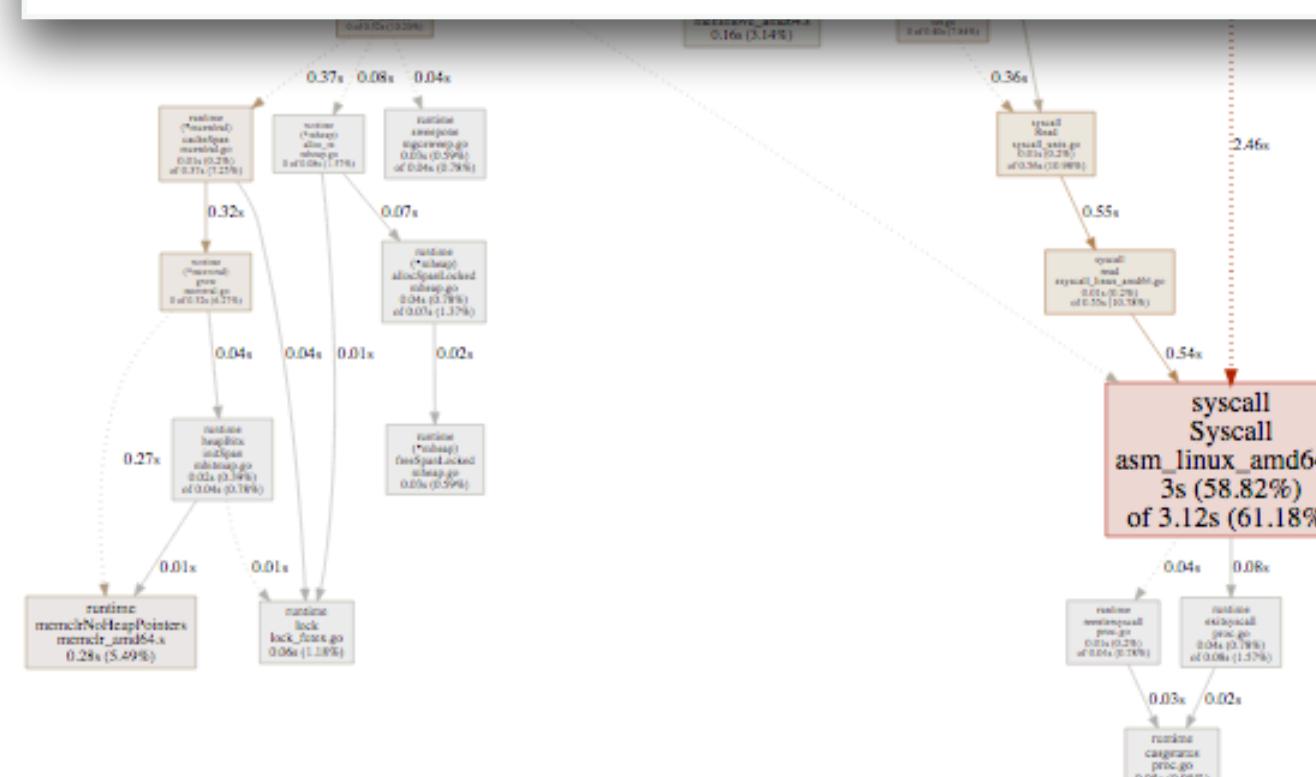
(pprof) top5

Showing nodes accounting for 3700ms, 72.55% of 5100ms total

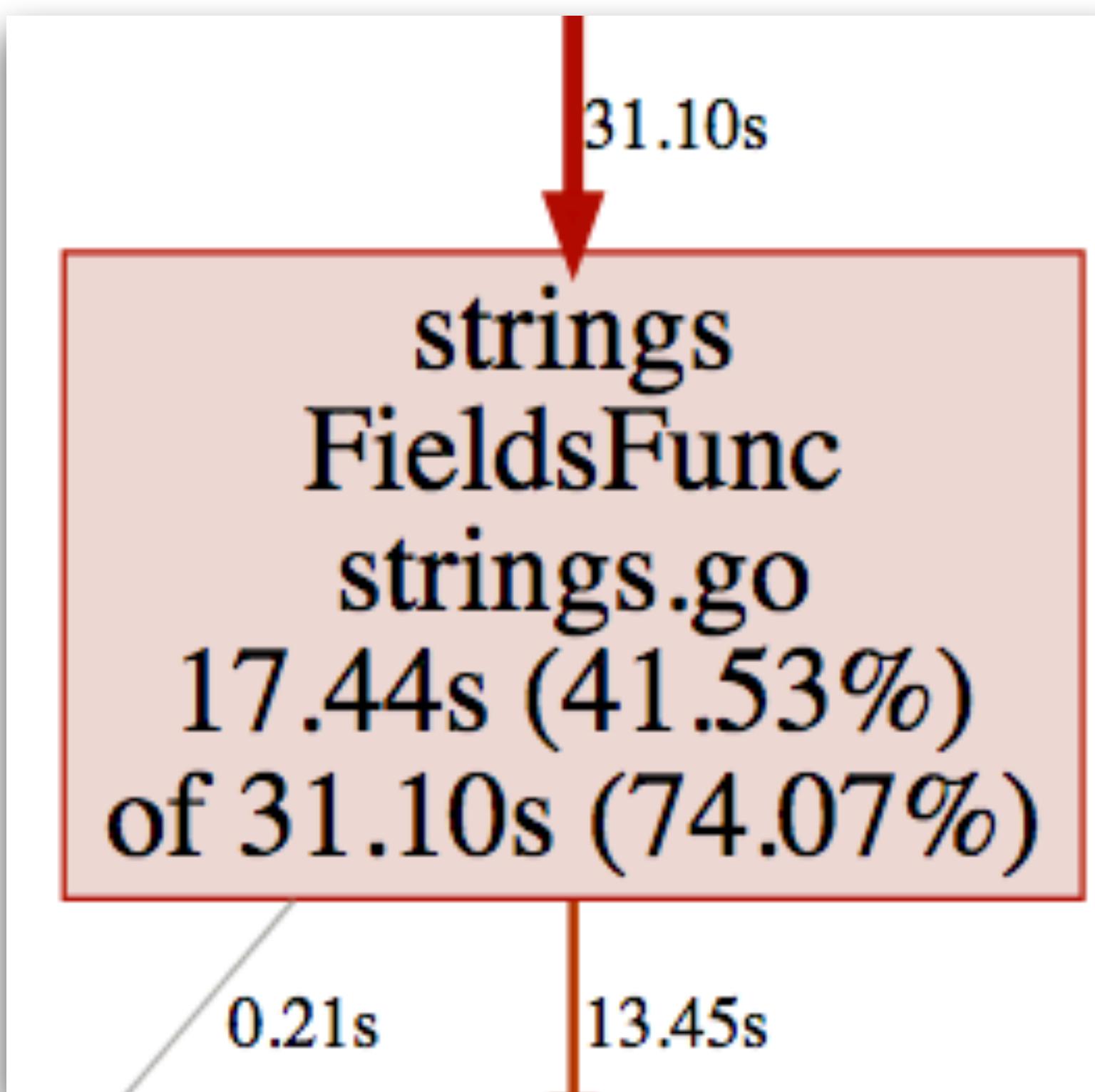
Dropped 86 nodes (cum <= 25.50ms)

Showing top 5 nodes out of 115

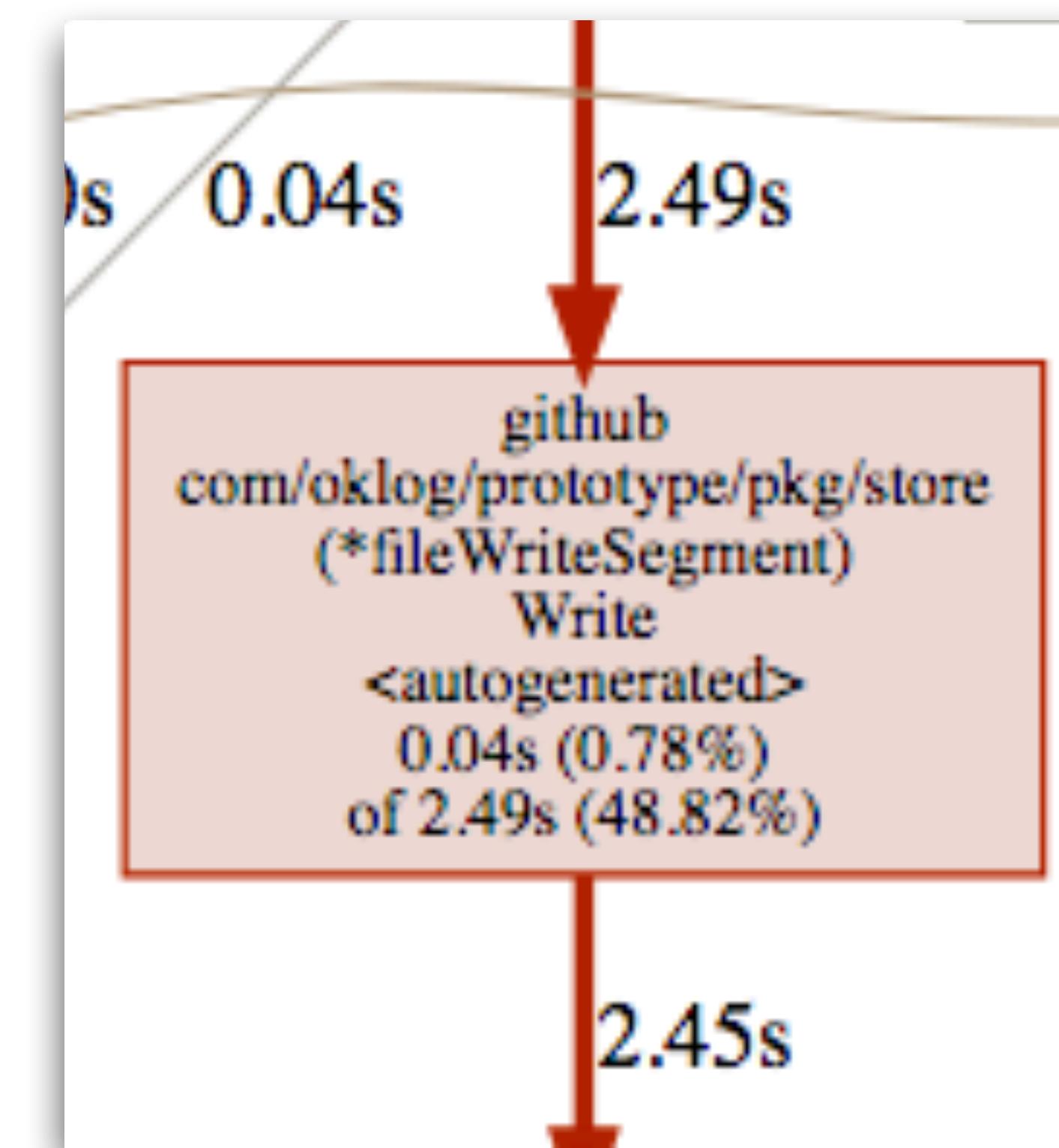
flat	flat%	sum%	cum	cum%	
3000ms	58.82%	58.82%	3120ms	61.18%	syscall.Syscall
280ms	5.49%	64.31%	280ms	5.49%	runtime.memclrNoHeapPointers
200ms	3.92%	68.24%	670ms	13.14%	runtime.mallocgc
160ms	3.14%	71.37%	160ms	3.14%	runtime.memmove
60ms	1.18%	72.55%	60ms	1.18%	runtime.cmpbody

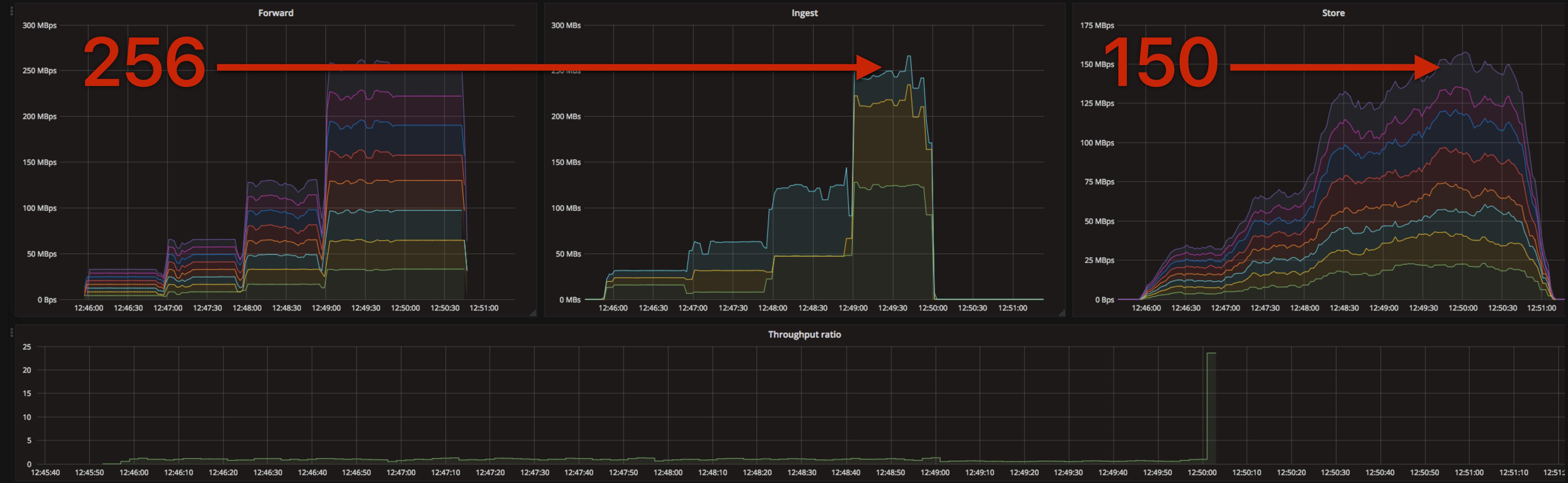


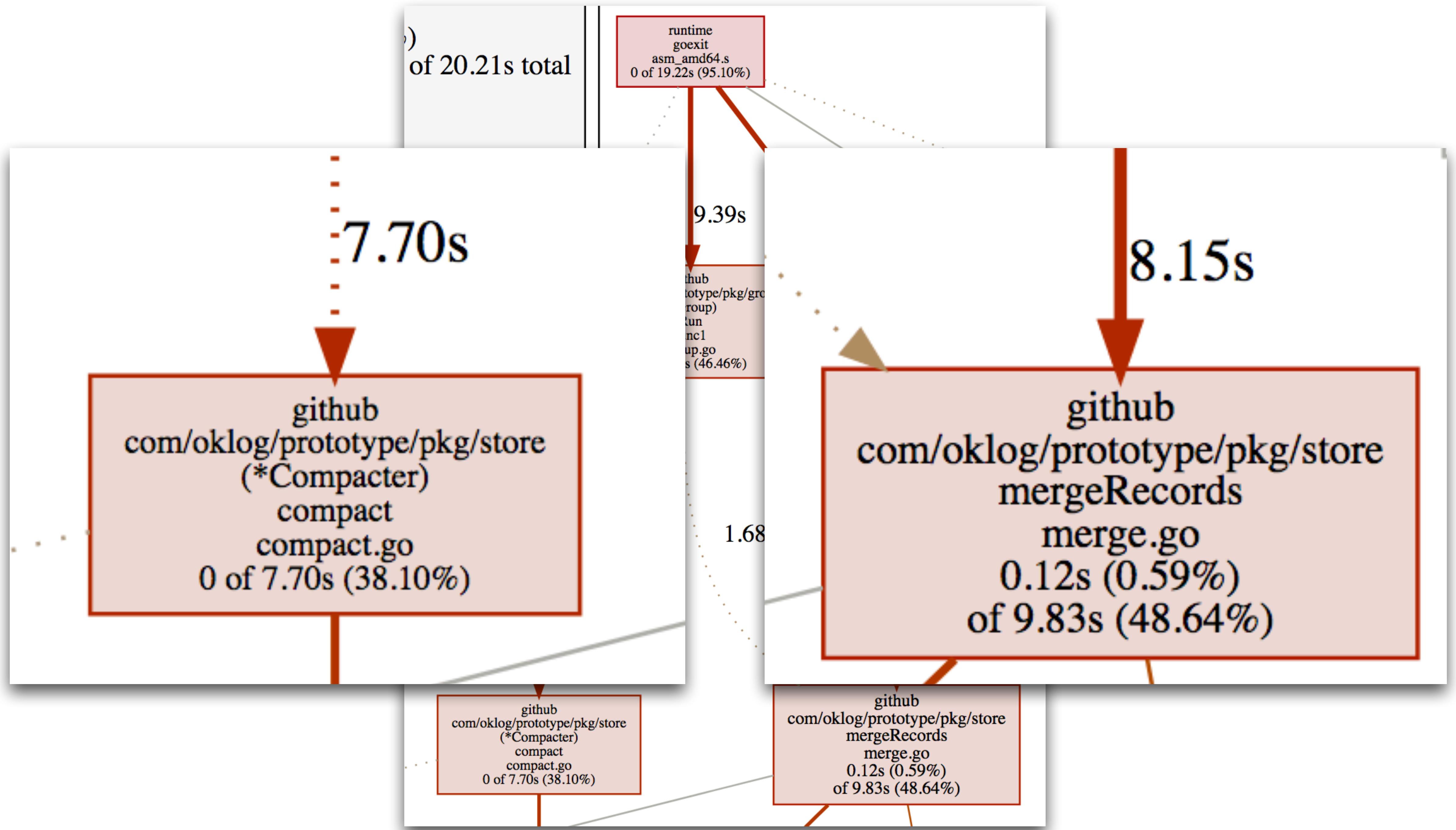
31.10s



2.49s







Optimizing queries

```
store1 ~ du -hs data  
21.2G    data
```

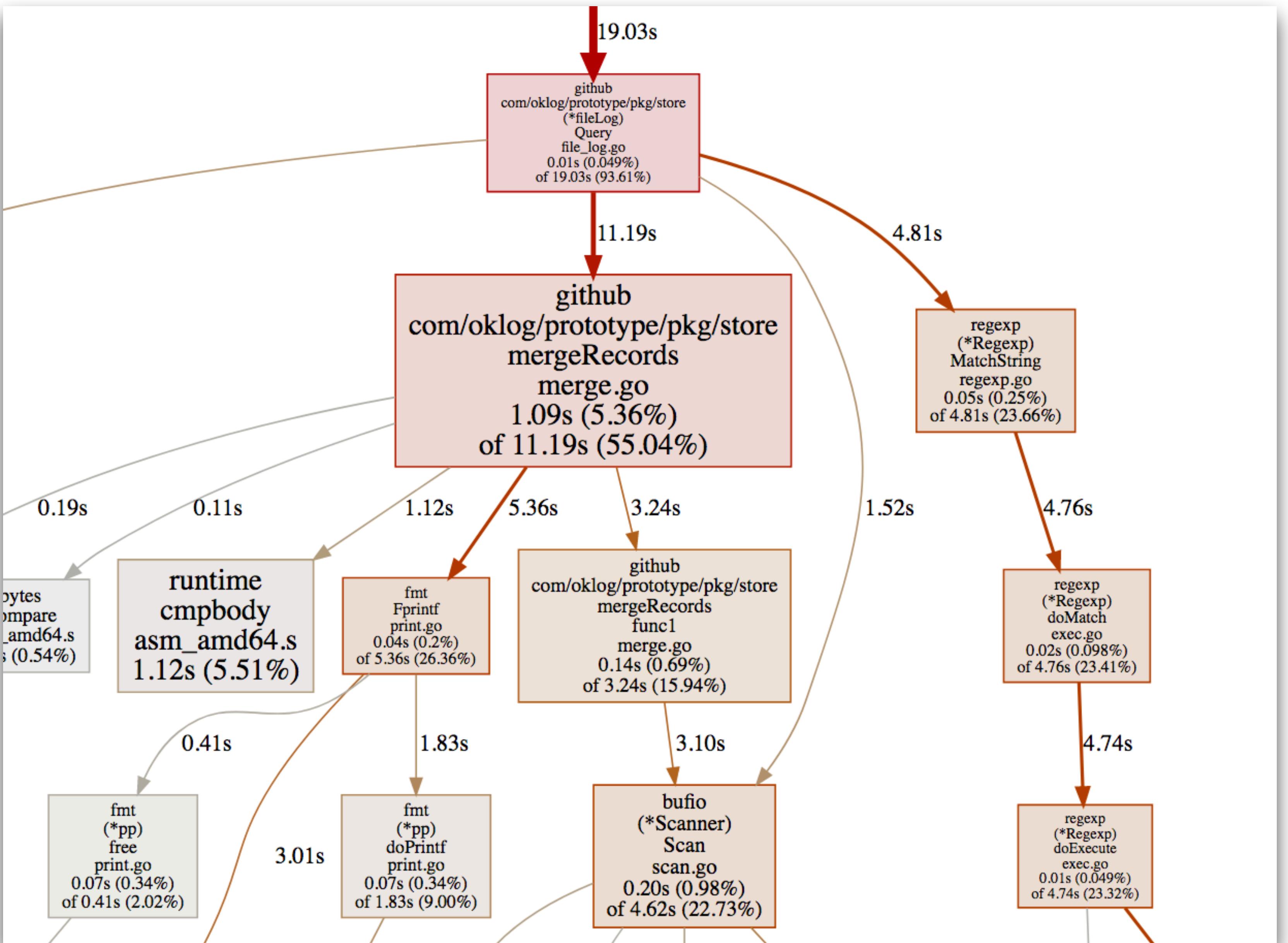
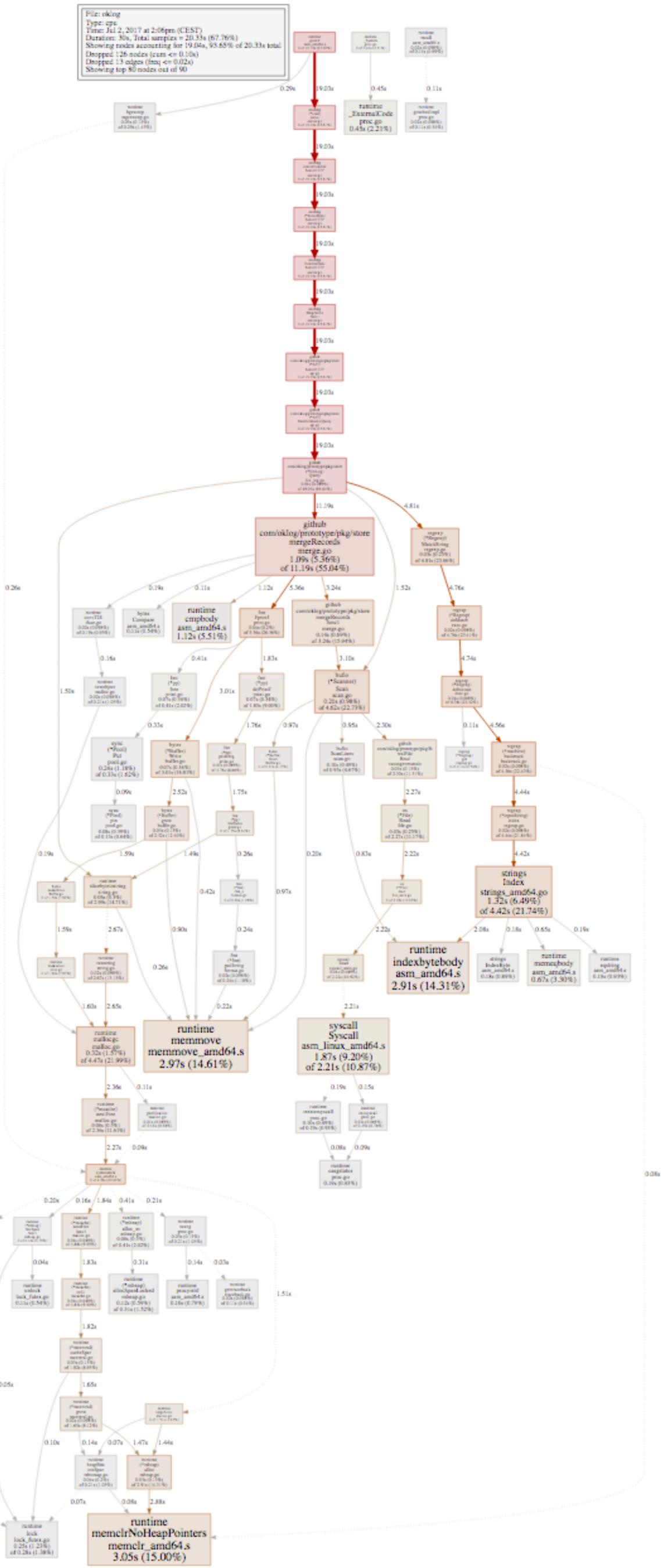
```
store1 ~ for i in (seq 1 5)  
  time rg --count 'ABCD' data >/dev/null  
end  
26.66user ...  
21.40user ...  
18.40user ...  
18.85user ...  
19.01user ...
```

```
store1 ~ rg 'ABCD' data | wc -l  
20287
```

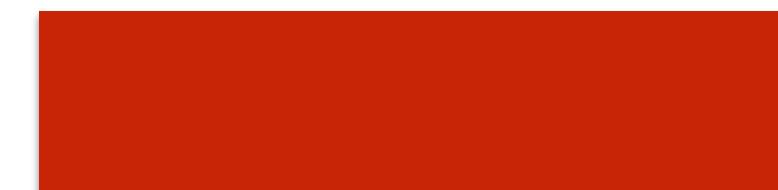
mean = 20.86s
⇒ ~1.01 GBps

```
store1 ~ for i in (seq 1 5)
  time oklog query -q ABCD >/dev/null
end
1:14.01user ...
1:05.51user ...
1:11.87user ...
1:12.07user ...
1:09.98user ...
```

mean = 1m10s
⇒ ~299 MBps



T0-T3



T2-T4



T4-T9



T6-T9

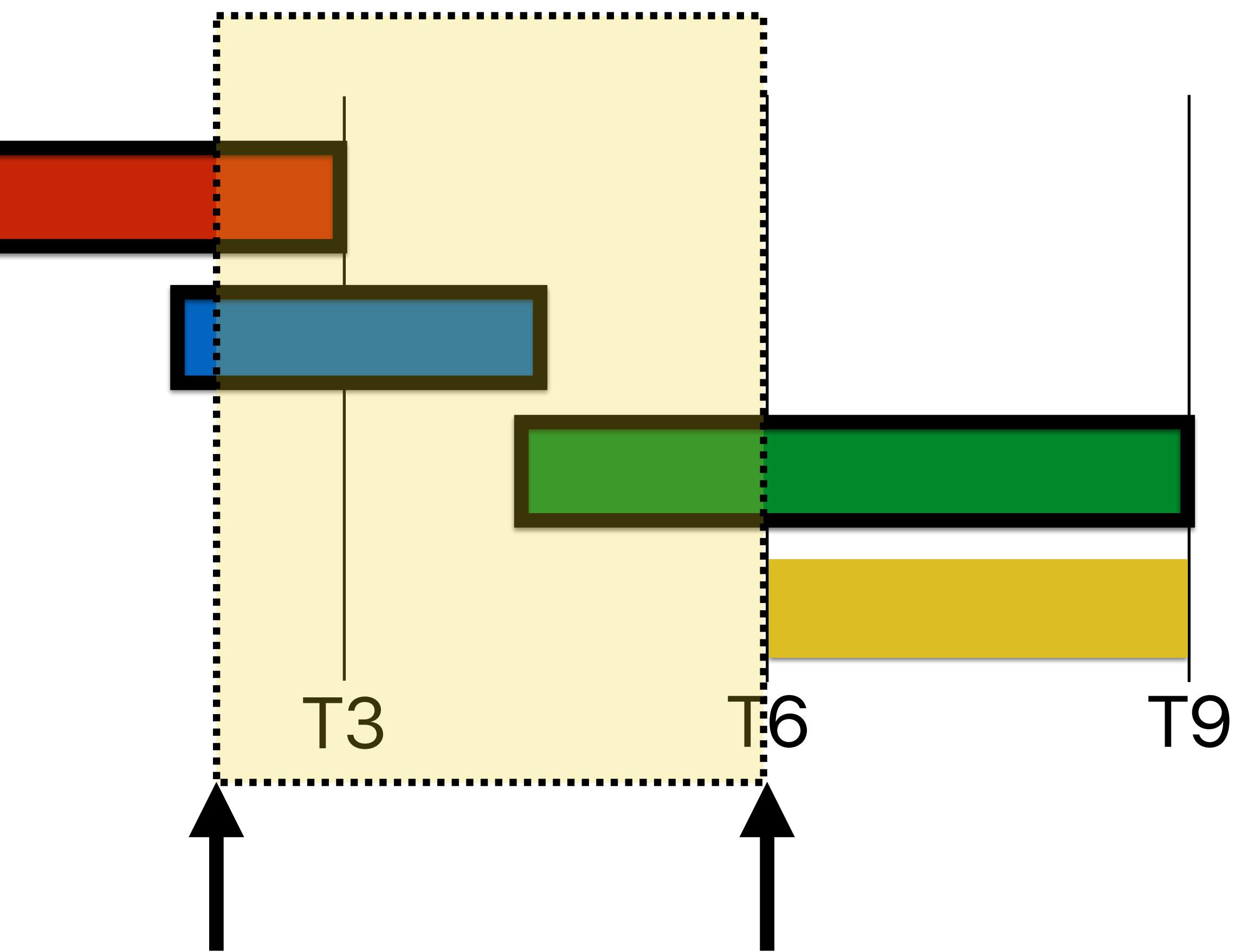


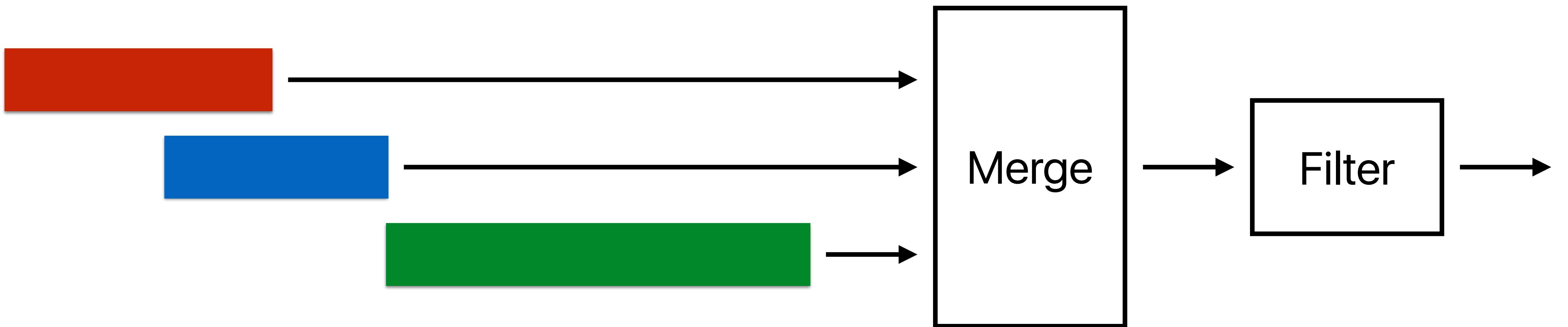
T0

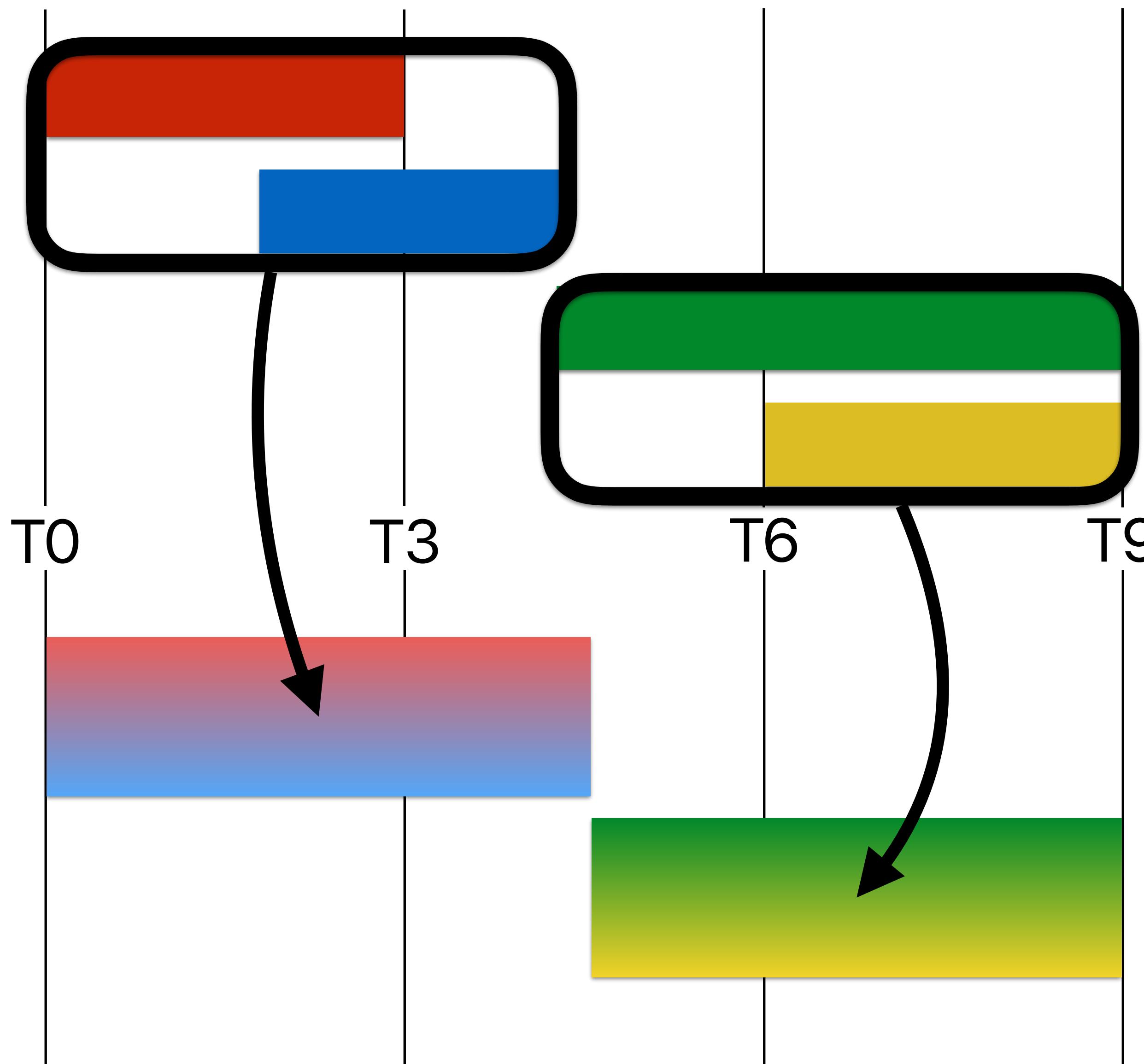
T3

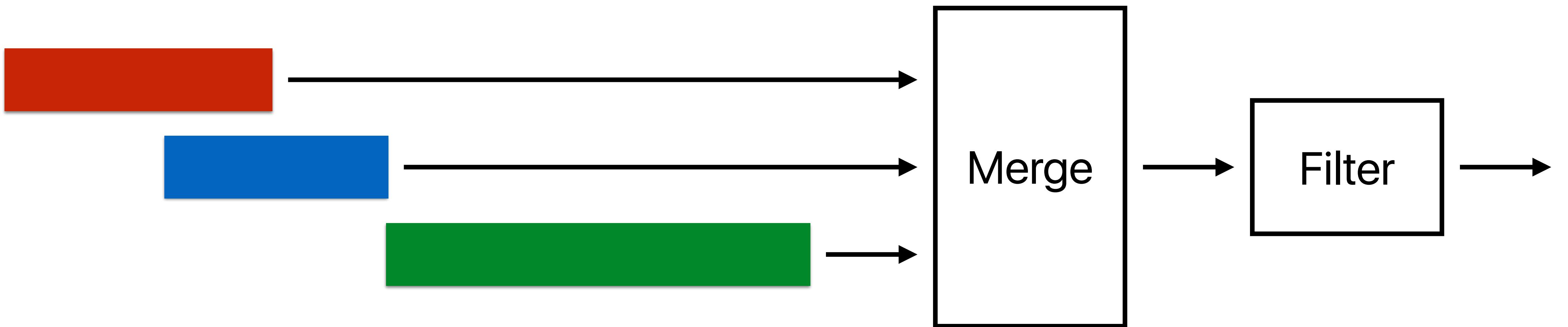
T6

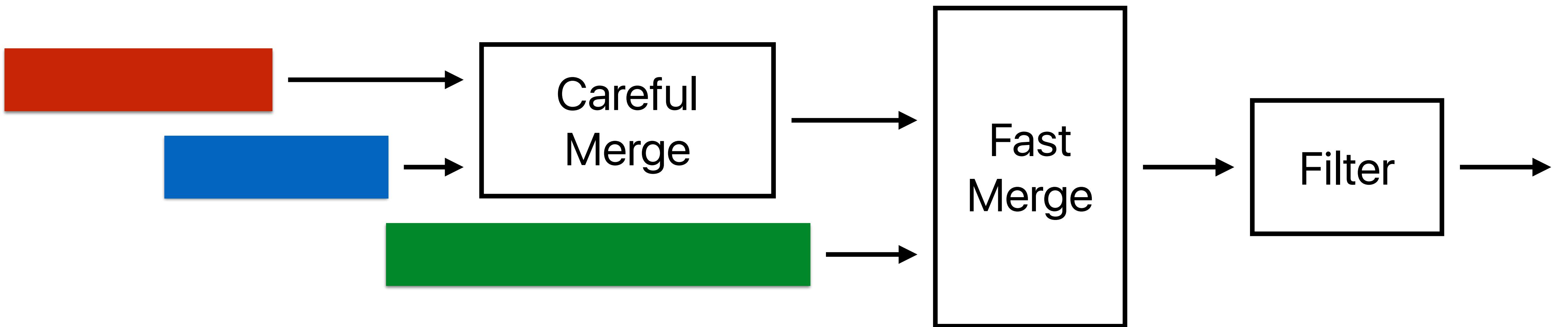
T9

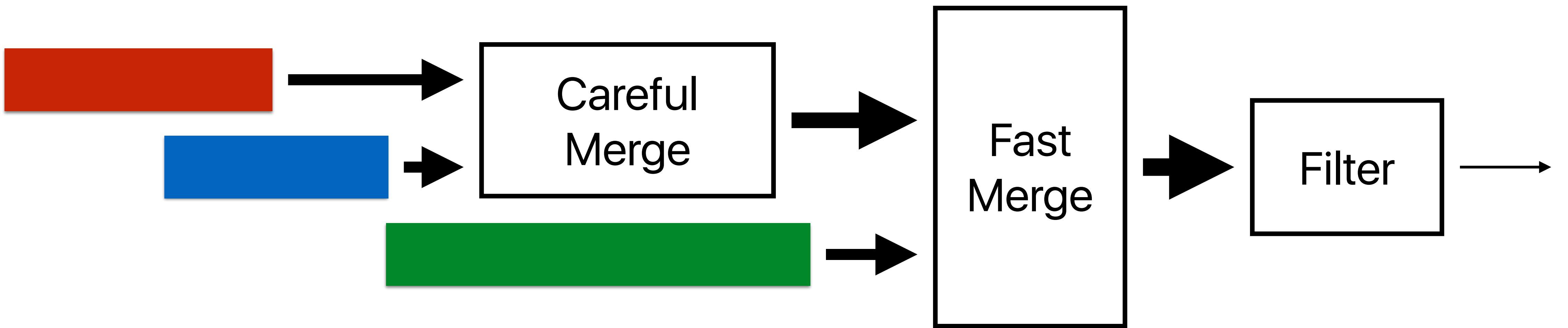


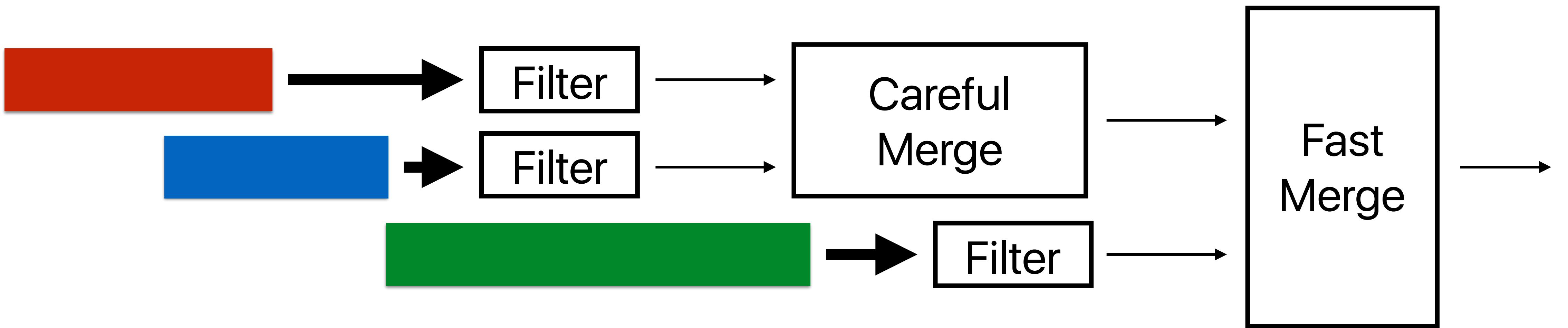


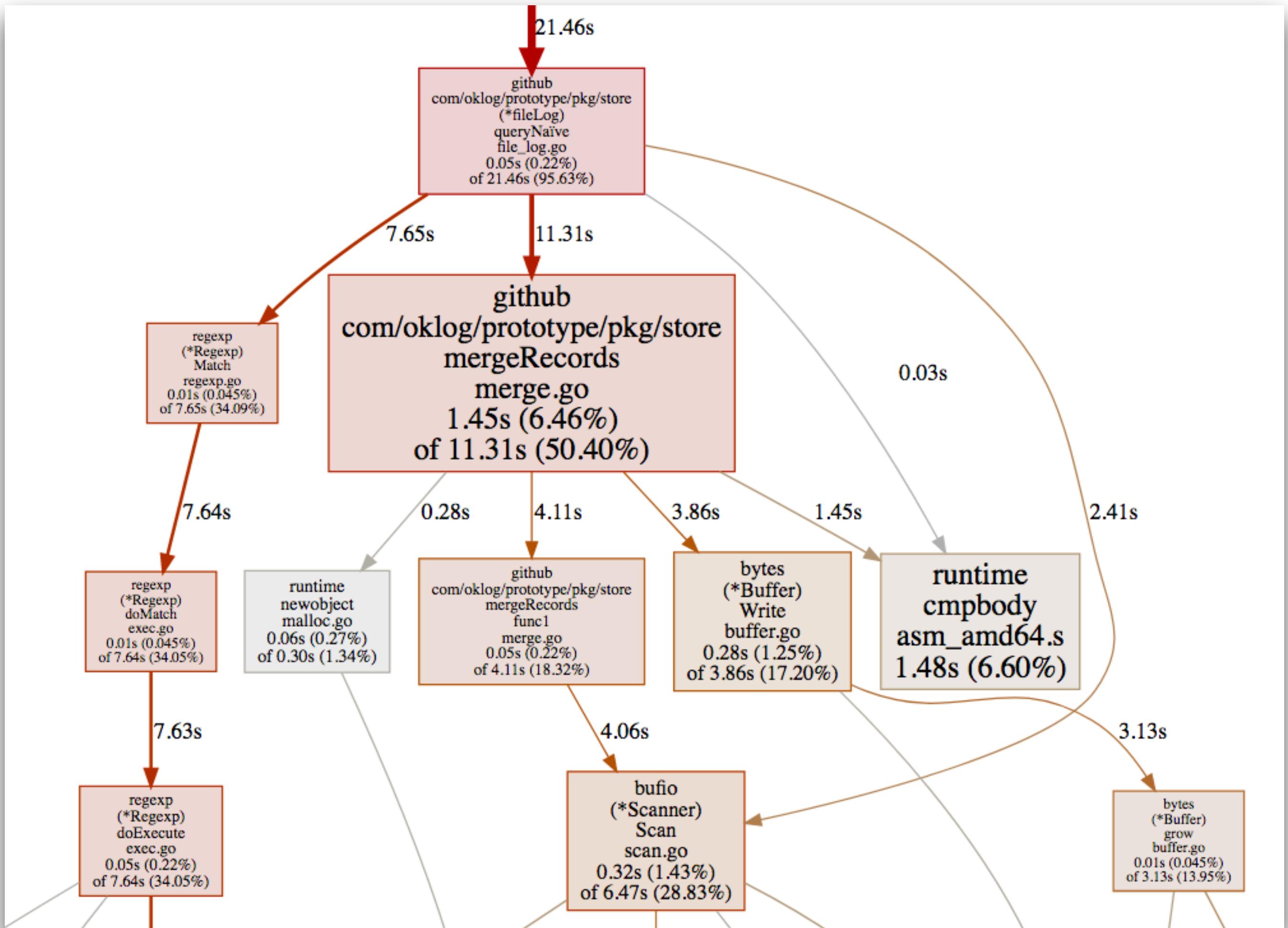
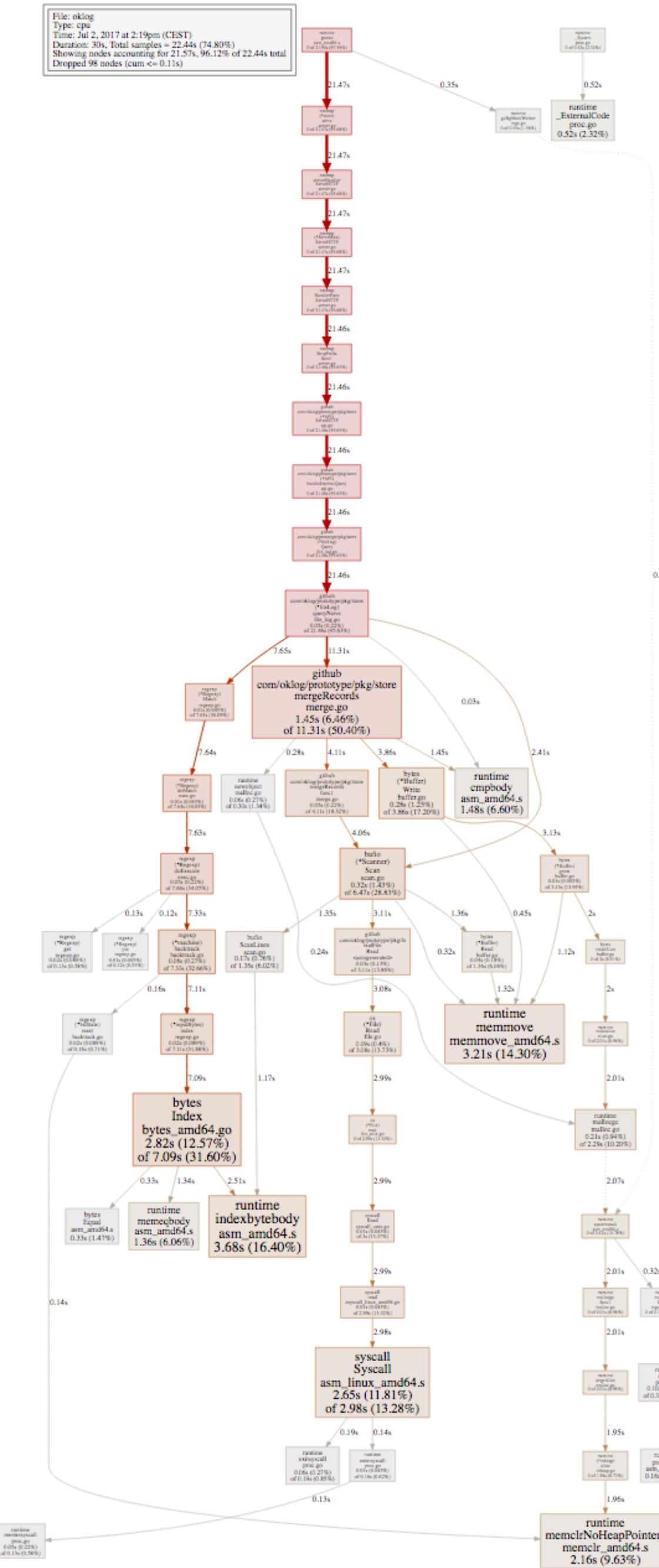












The screenshot shows a GoDoc interface for the `nio` package. The title bar says "GoDoc nio - GoDoc". The URL in the address bar is `godoc.org/github.com/djherbis/nio`. The main content area displays the package's code and documentation.

package nio

```
import "github.com/djherbis/nio"

Package nio provides a few buffered io primitives.
```

Index

```
func Copy(dst io.Writer, src io.Reader, buf Buffer) (n int, err error)
func NewReader(src io.Reader, buf Buffer) io.ReadCloser
type Buffer
type PipeReader
    • func Pipe(buf Buffer) (r *PipeReader, w *PipeWriter)
    • func (r *PipeReader) Close() error
    • func (r *PipeReader) CloseWithError(err error) error
    • func (r *PipeReader) Read(p []byte) (n int, err error)
type PipeWriter
    • func (w *PipeWriter) Close() error
    • func (w *PipeWriter) CloseWithError(err error) error
    • func (w *PipeWriter) Write(p []byte) (int, error)
```

Package Files

- `nio.go`
- `sync.go`

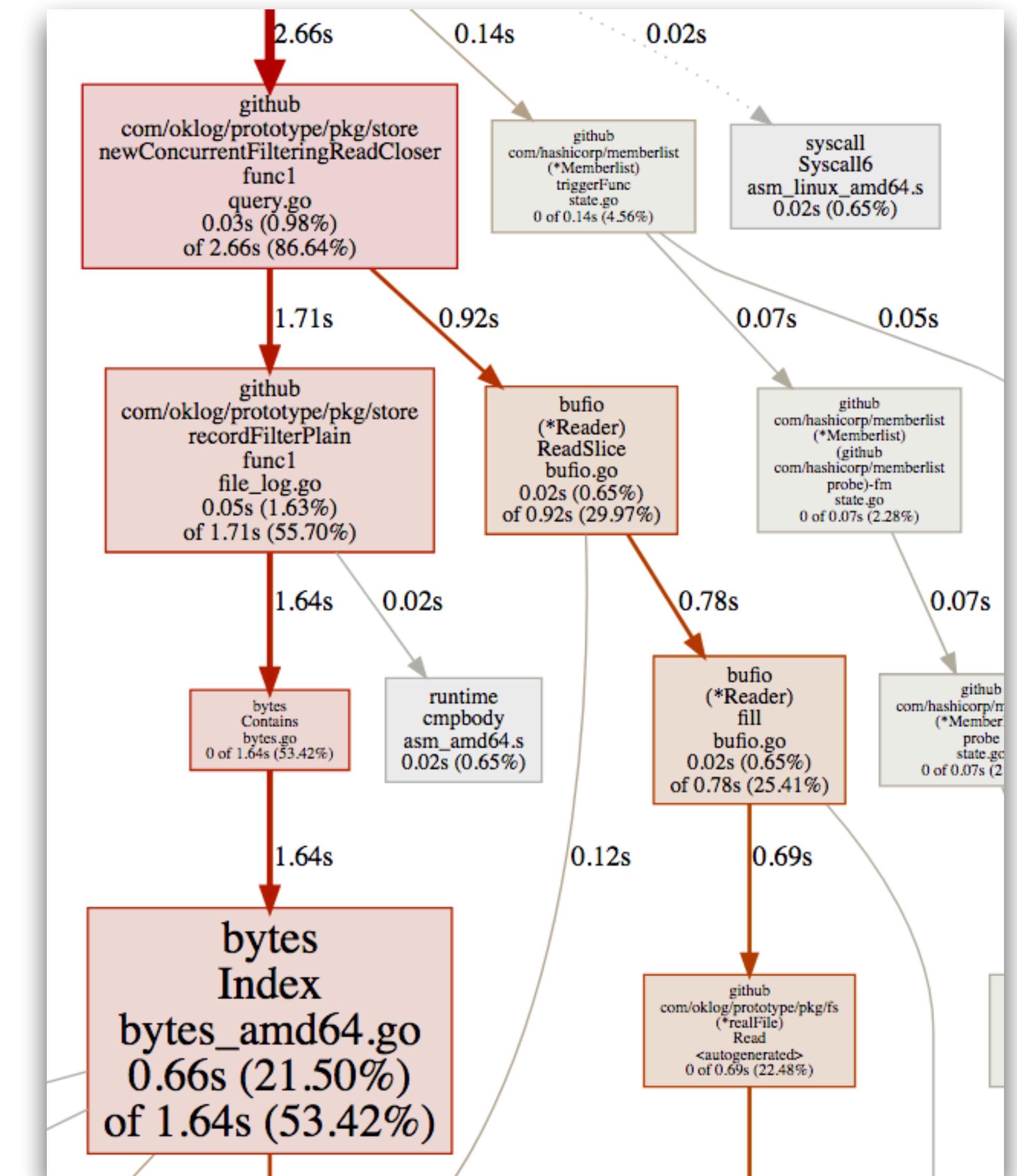
On the right side of the interface, there is a sidebar with the following content:

- func NewReader(src io.Reader, buf Buffer) io.ReadCloser**
- type Buffer**
- type PipeReader**
 - **func Pipe(buf Buffer) (r *PipeReader, w *PipeWriter)**
 - **func (r *PipeReader) Close() error**
 - **func (r *PipeReader) CloseWithError(err error) error**
 - **func (r *PipeReader) Read(p []byte) (n int, err error)**
- type PipeWriter**
 - **func (w *PipeWriter) Close() error**
 - **func (w *PipeWriter) CloseWithError(err error) error**
 - **func (w *PipeWriter) Write(p []byte) (int, error)**

File: oklog
 Type: cpu
 Time: Jul 2, 2017 at 3:15pm (CEST)
 Duration: 30s, Total samples = 3,075 (10.23%)
 Showing 100 samples
 Dropped 0 samples

```

2026 // input:
2027 // SI: data
2028 // BX: data len
2029 // AL: byte sought
2030 // R8: address to put result
2031 TEXT runtime·indexbytebody(SB),NOSPLIT,$0
2032     // Shuffle X0 around so that each byte contains
2033     // the character we're looking for.
2034     MOVD AX, X0
2035     PUNPCKLBW X0, X0
2036     PUNPCKLBW X0, X0
2037     PSHUFL $0, X0, X0
2038
2039     CMPQ BX, $16
2040     JLT small
2041
2042     MOVQ SI, DI
2043
2044     CMPQ BX, $32
2045     JA avx2
2046 sse:
2047     LEAQ -16(SI)(BX*1), AX      // AX = address of last
2048     JMP sseloopentry
2049 sseloop:
2050     // Move the next 16-byte chunk of the data into X1.
2051     MOVOU (DI), X1
2052     // Compare bytes in X0 to X1.
2053     PCMPEQB X0, X1
2054     // Take the top bit of each byte in X1 and put the res
  
```



```
store1 ~ for i in (seq 1 5)
  time oklog query -q ABCD >/dev/null
end
31.04user ...
25.41user ...
27.50user ...
23.01user ...
22.11user ...
```

mean = 25.81s
⇒ ~821 MBps

Conclusions

Think, then act

Work in cycles

NIH isn't a 4-letter word



Thanks y'all

Evolutionary Optimization
GopherCon 2017
@peterbourgon