

elastic{ON}¹⁵

Elasticsearch in Anger

stories from the **GitHub** search clusters

Tim Pease



← Where have we come from

⌚ What we are doing now

→ Where we are going

A Whole New Code Search

 January 23, 2013



TwP

 New Features

 Edit

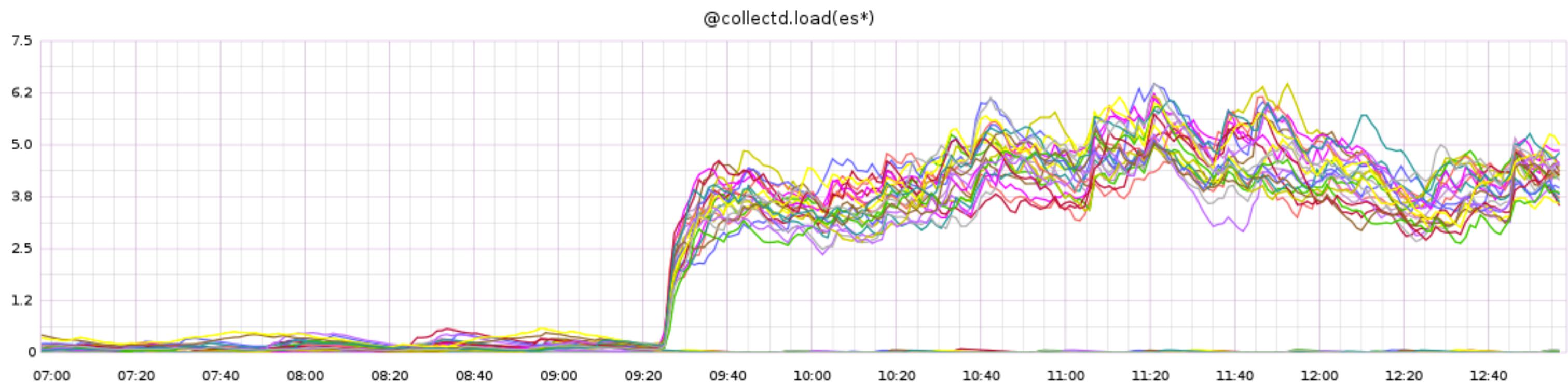
Finding great code on GitHub just got a whole lot easier. Today we're releasing several big improvements to code search.

New Technology

First, we are unveiling our new search infrastructure that will grow to support the immense amount of source code being pushed to GitHub each day.

Under the hood is an [ElasticSearch](#) cluster that live-indexes your code as you push it up to GitHub. Search results will be returned from public and private repositories that you have access to.

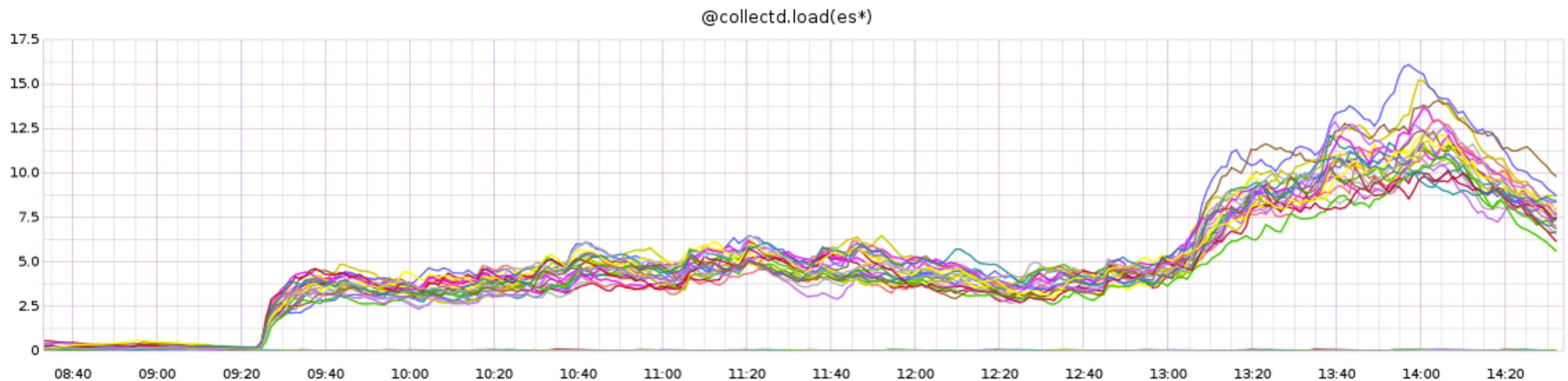
Code Search Load

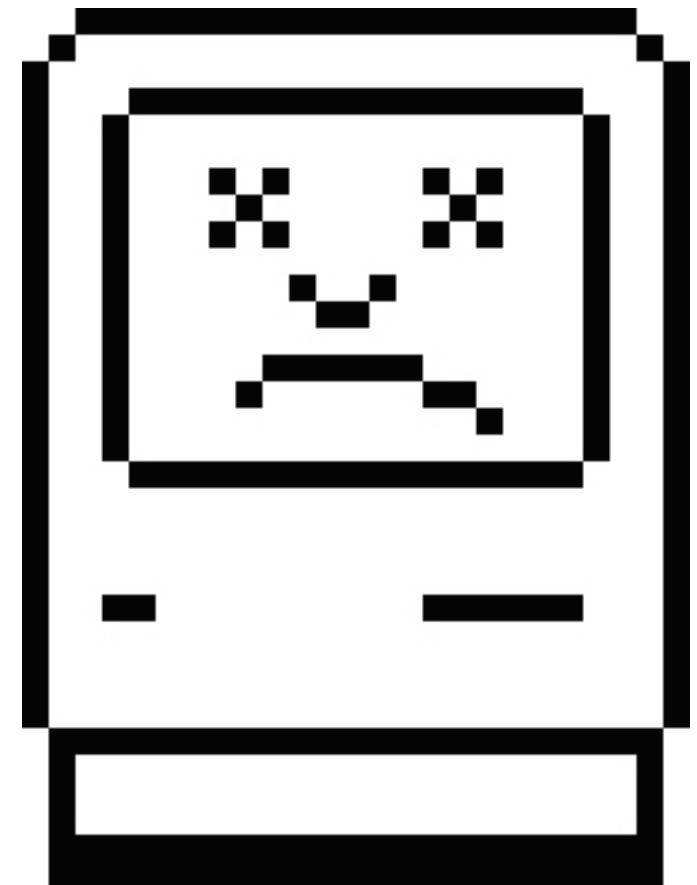


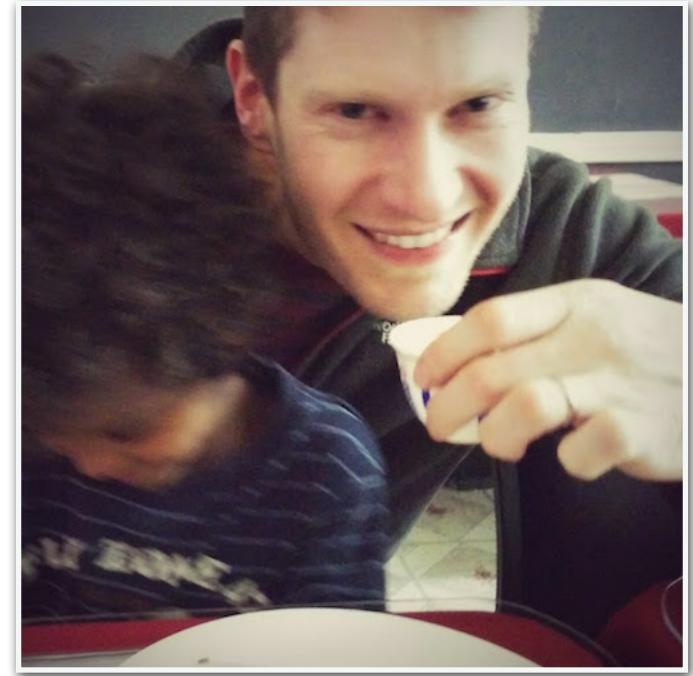


“you can find **secret keys
people have committed to their
public repositories”**

Code Search Load





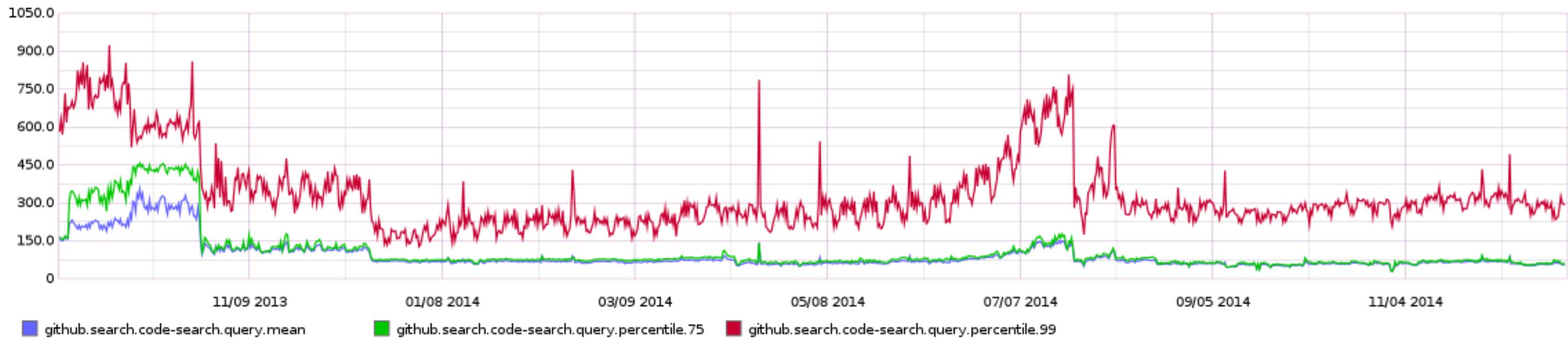


Drew Raines

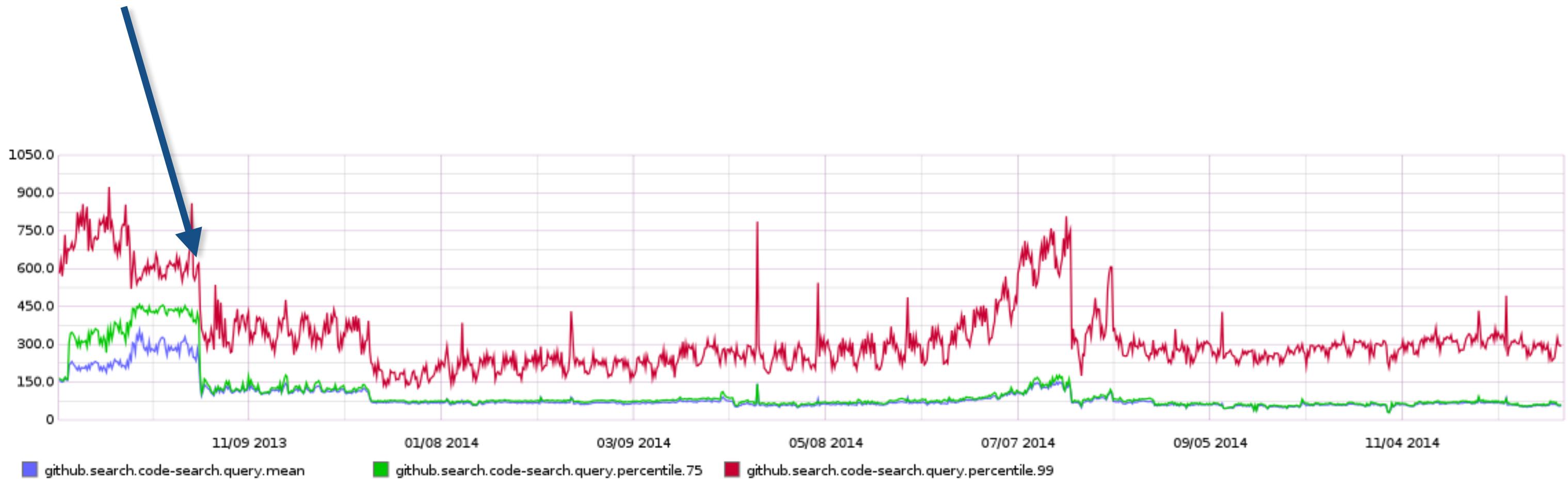
Code Search

- ← We performed inadequate load testing
- ⌚ We had insufficient operations experience
- We need better tools and metrics

Code Search Query Performance



New Cluster

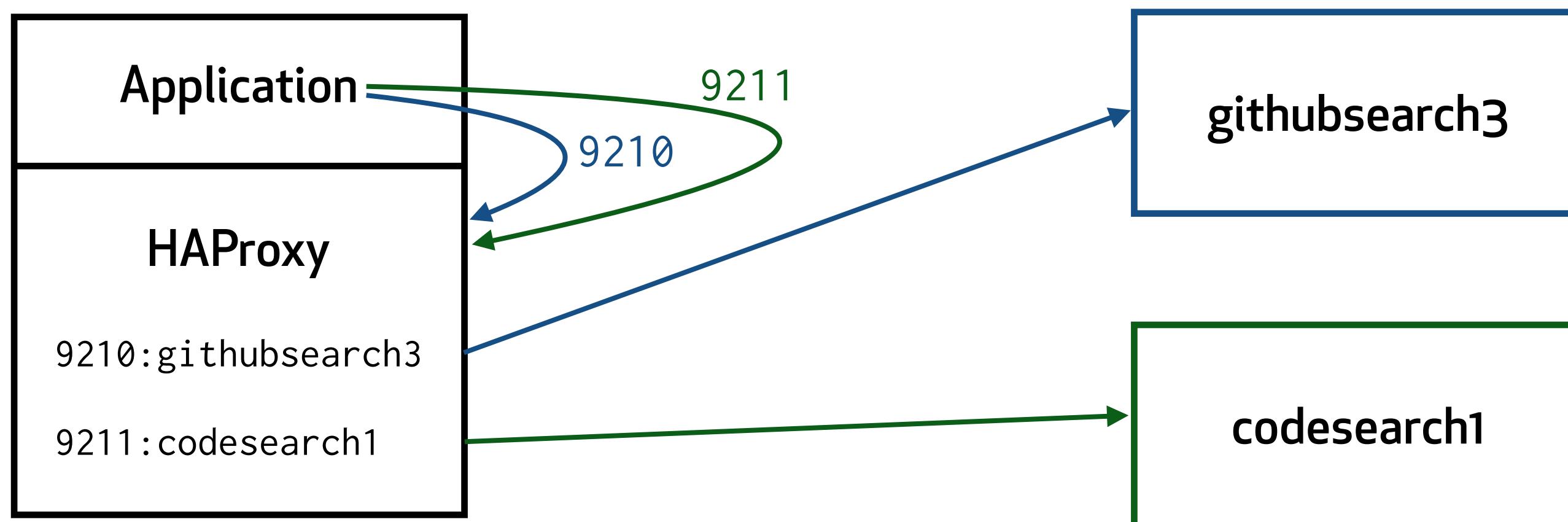


Index Management

[Create new index](#)

Status	Index	Cluster	Repair
	blog-5	githubsearch3	
	blog-6	research1	
	code-search-2	codesearch1	
	issues-5	research1	
	issues-6	githubsearch3	
	pull-requests-6	githubsearch3	
	pull-requests-8	research1	
	repos-6	githubsearch3	
	repos-8	research1	
	repos-9	githubsearch3	

HAProxy for Multiple Clusters



Push Button Index Creation



Index Management				Create new index
Status	Index	Cluster	Repair	
✓	blog-5	githubsearch3		
★ ✓	blog-6	research1		
★ ✗	code-search-2	codesearch1		
✓	issues-5	research1		
★ ✓	issues-6	githubsearch3		
★ ✓	pull-requests-6	githubsearch3		
✓	pull-requests-8	research1		
★ ✓	repos-6	githubsearch3		
✓	repos-8	research1		
★ ✓	searches-2	githubsearch3		

Create a New Index

Create a new index

Select the index to create

blog 

Select the cluster where the index will be created

research1 

 **Make the index searchable**
The index will process search requests when it is the primary index.

 **Make the index writable**
Documents will be added to the index only when it is writable.

Create Index

Backfill Data

Repair the `blog-7` index [Reset](#)

Progress: 72.0%

Processed: 1,300

Added: 1,300

Updated: 0

Removed: 0

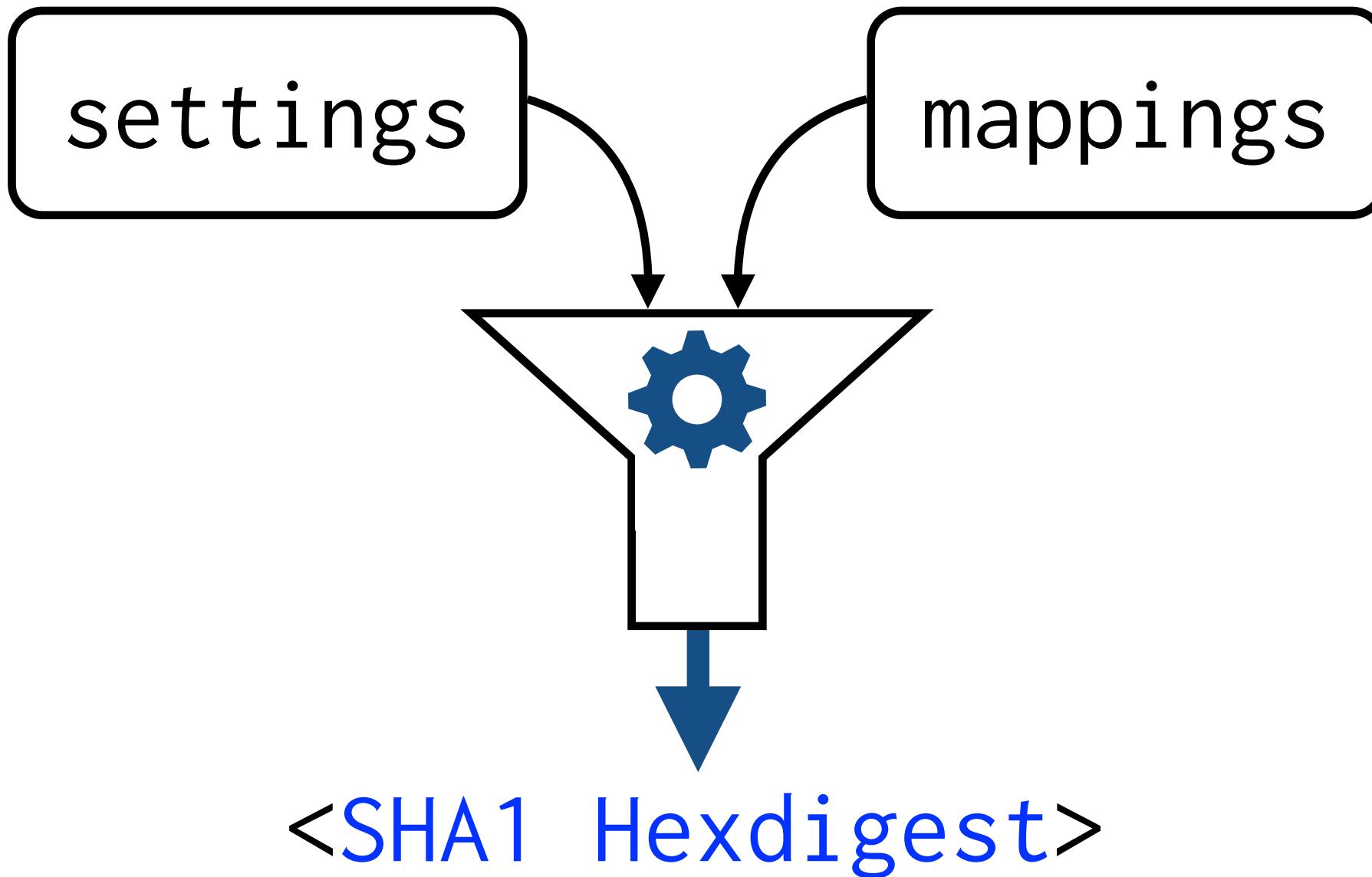
Active workers: 0

Start another worker for this repair job [Start](#)

Pause the repair job [Pause](#)



Index Versions ✘





Index Versions



```
mappings: {  
    index-meta: {  
        _meta: {  
            version: <SHA1 Hexdigest>  
        }  
    }  
}
```

Because ...

GitHub Enterprise



Load Testing

Load Testing

Scientist

github.com/github/scientist

- ▶ control
- ▶ experiment
- ▶ throttling



Load Testing

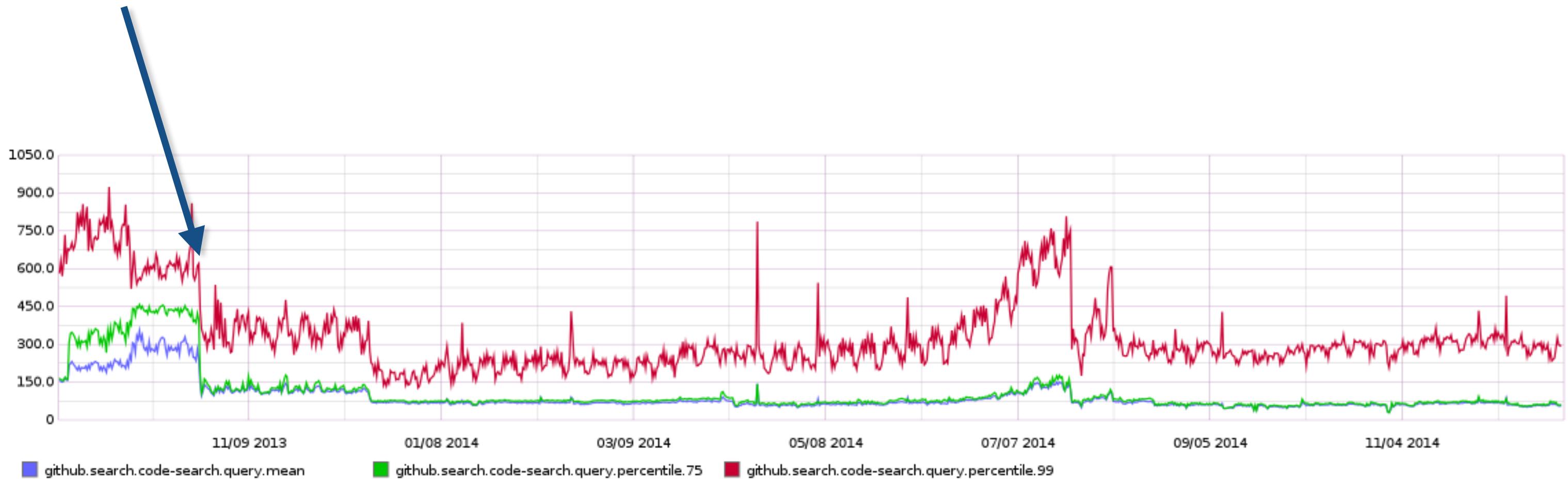
```
require "scientist"

def search(query)
  science "code-search-load-test" do |e|
    e.use { old_index.search(query) }
    e.try { new_index.search(query) }
  end
end
```

New Cluster

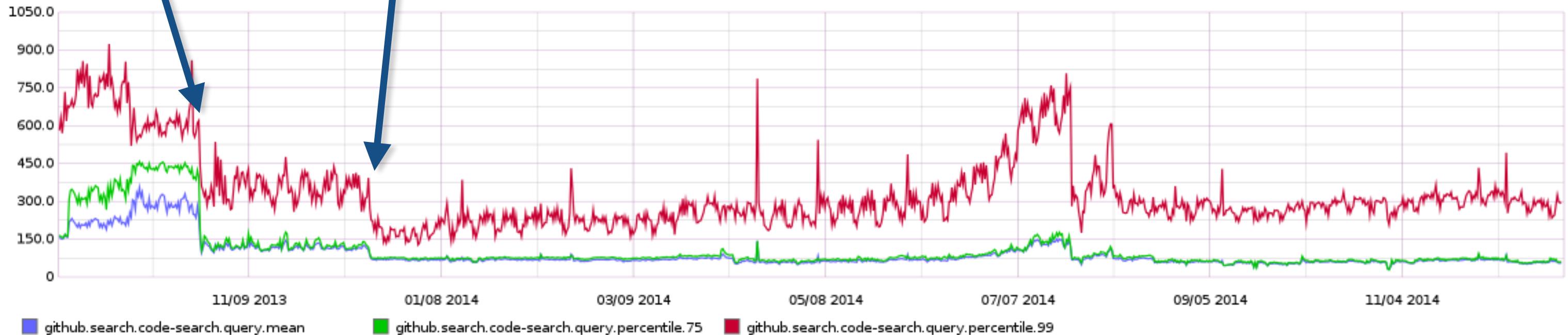
- ← We were outgrowing our old cluster
- ⌚ We created migration tools
- We used production queries for load testing

New Cluster



New Cluster

New Queries

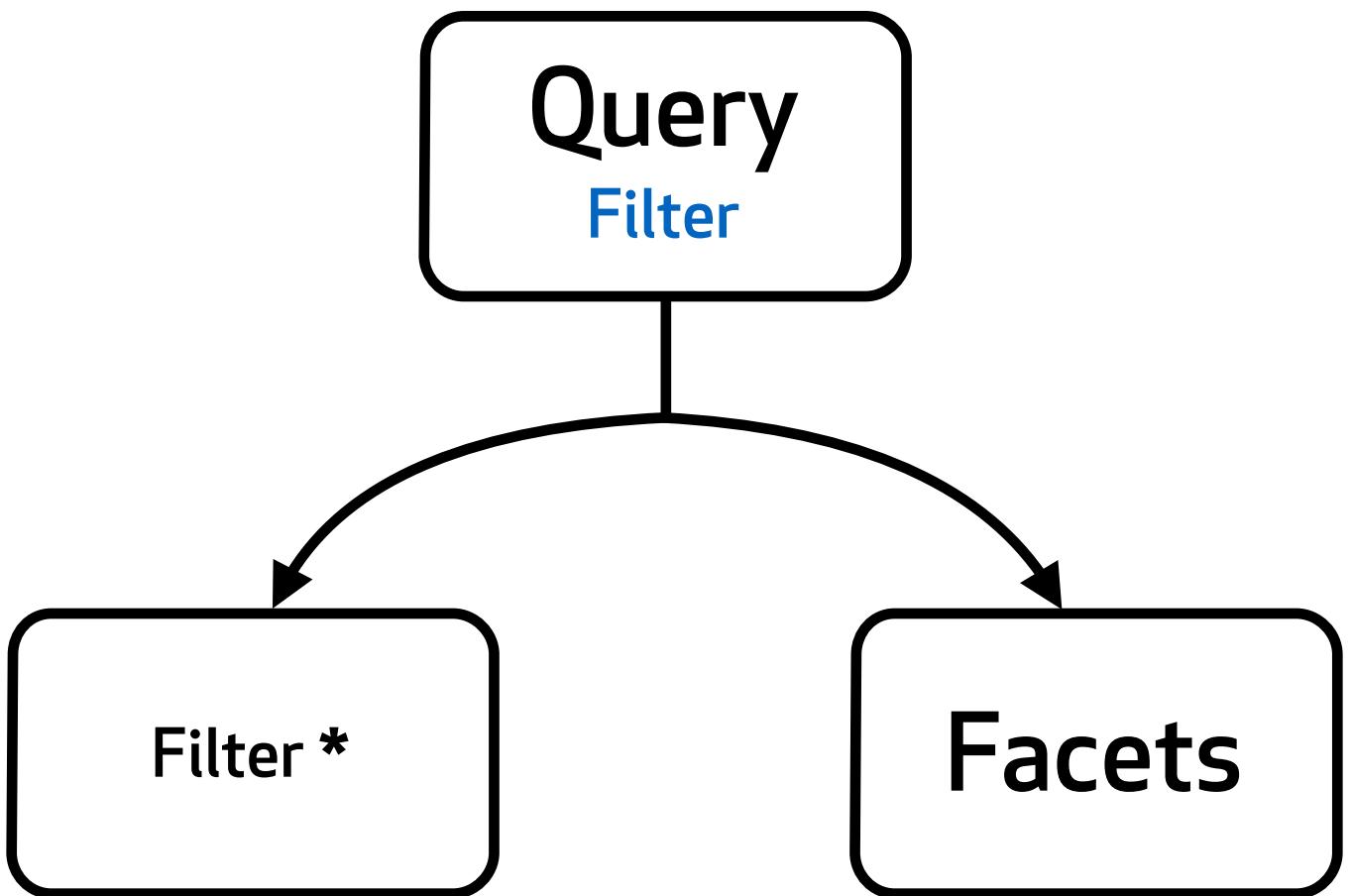
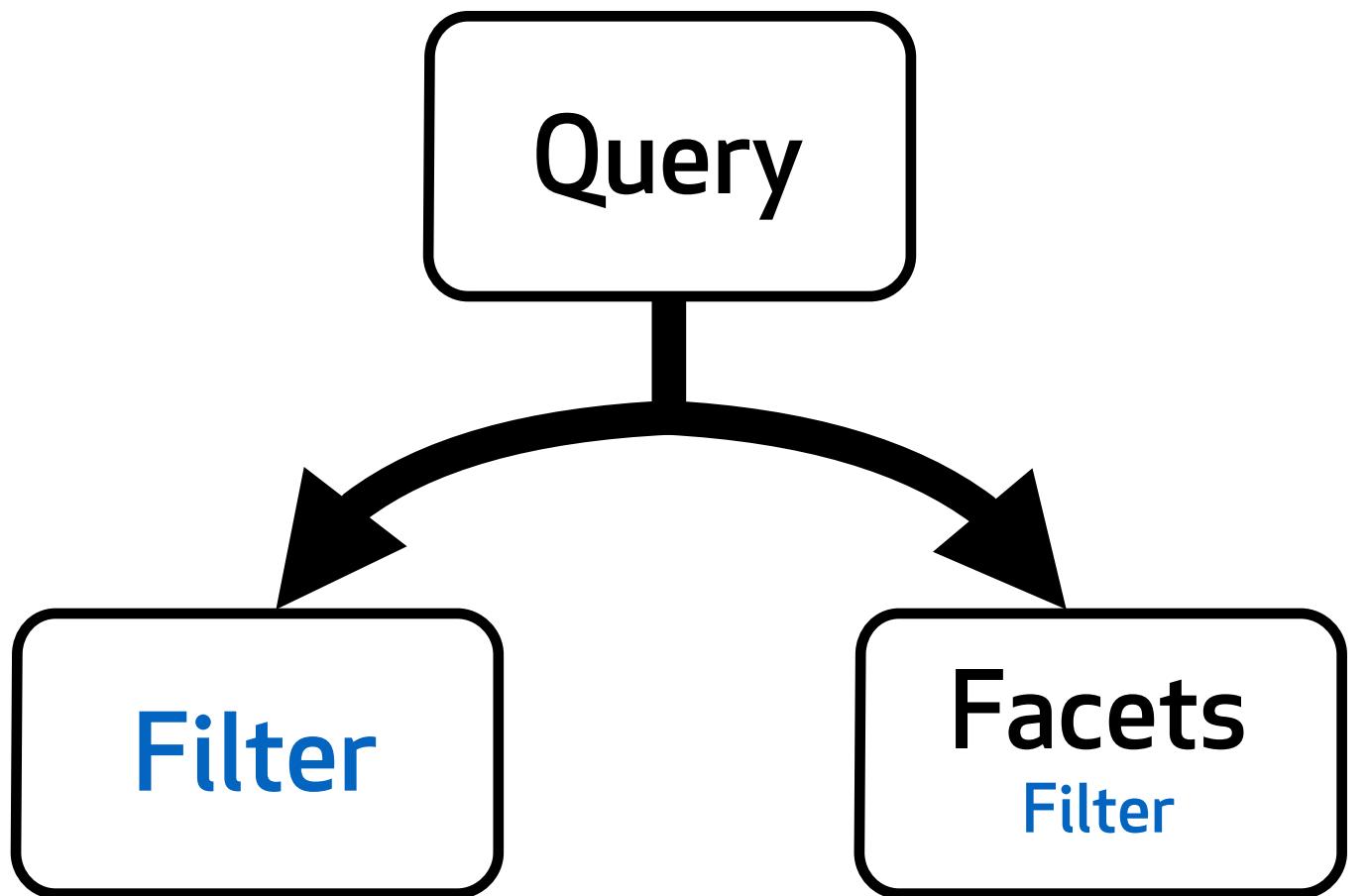


Elasticsearch Training

Tale of Two Queries

```
{  
  query: {match_all: {}},  
  filter: {  
    term: {state: "open"}  
  }  
}
```

```
{  
  query: {constant_score: {  
    filter: {  
      term: {state: "open"}  
    }  
  }}  
}
```



Common Issue Filters

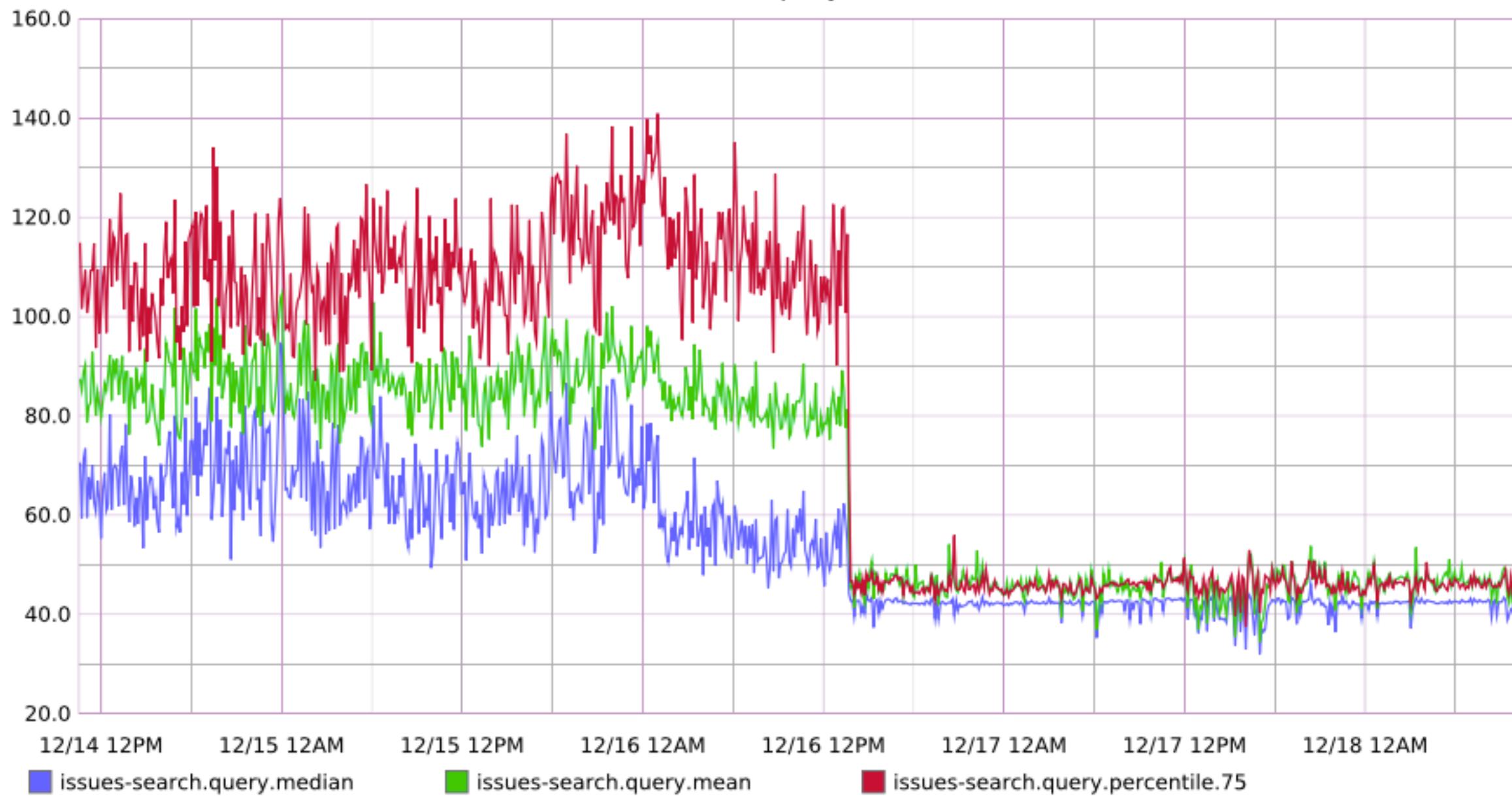
```
author:grantr
assignee:kimchi
mentions:TwP
involves:drewr
label:bug
repo:elasticsearch/elasticsearch
created:<2015-01-01
updated:>2015-03-01
is:issue
is:open
```

Common Issue Filters

author:grantr
assignee:kimchi
mentions:TwP
involves:drewr
label:bug
repo:elasticsearch/elasticsearch
created:<2015-01-01
updated:>2015-03-01
is:issue
is:open



issues-search query metrics

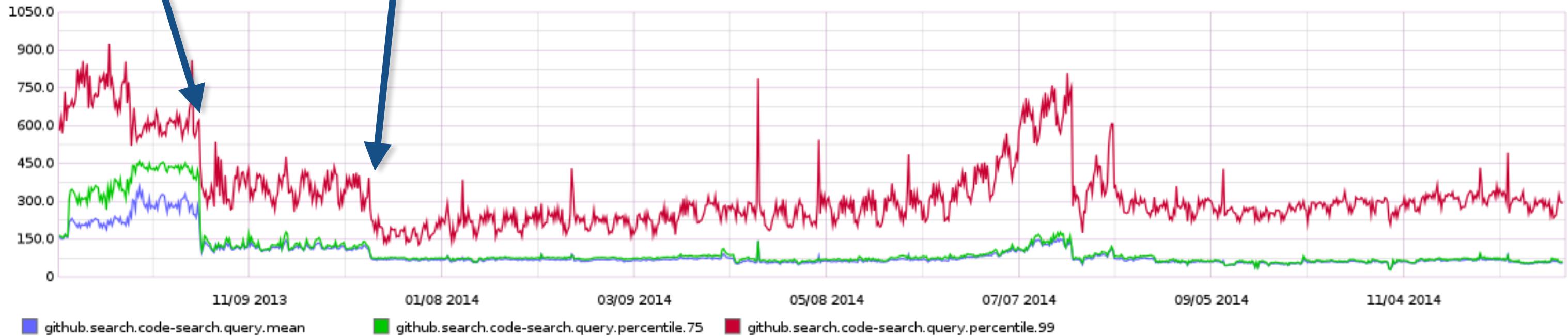


New Queries

- ← We were able to look at query performance
- ⌚ We got some education about filters
- We now enjoy efficient filtered queries

New Cluster

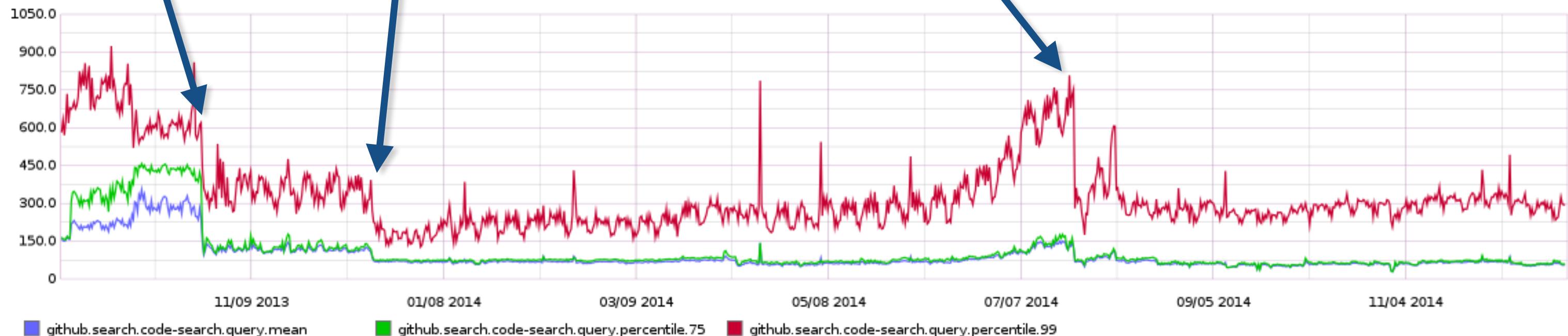
New Queries



New Cluster

New Queries

Heap Exhaustion

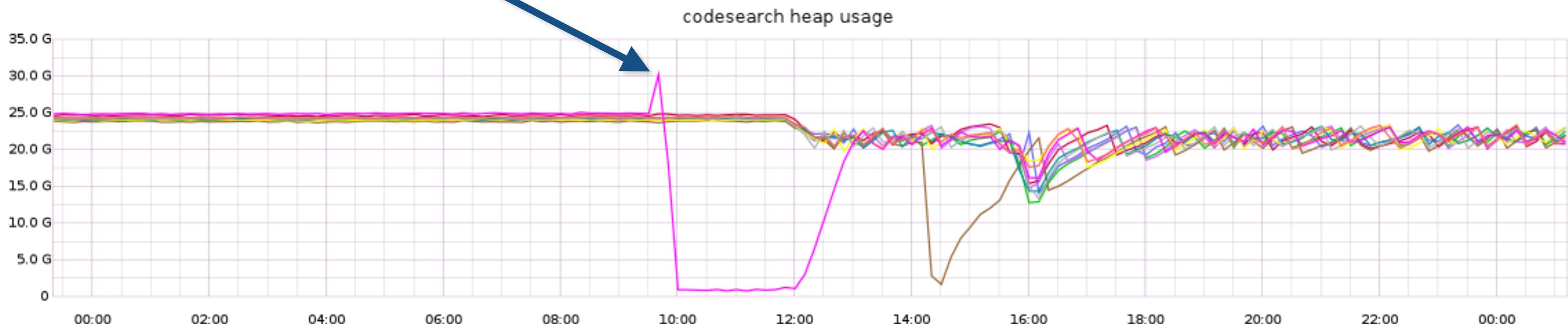


Code Search Heap Usage

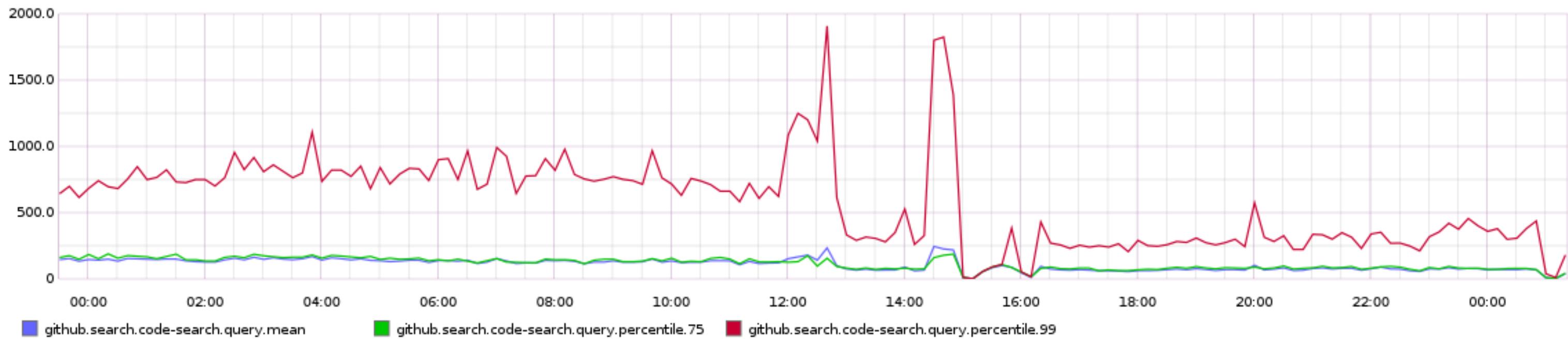
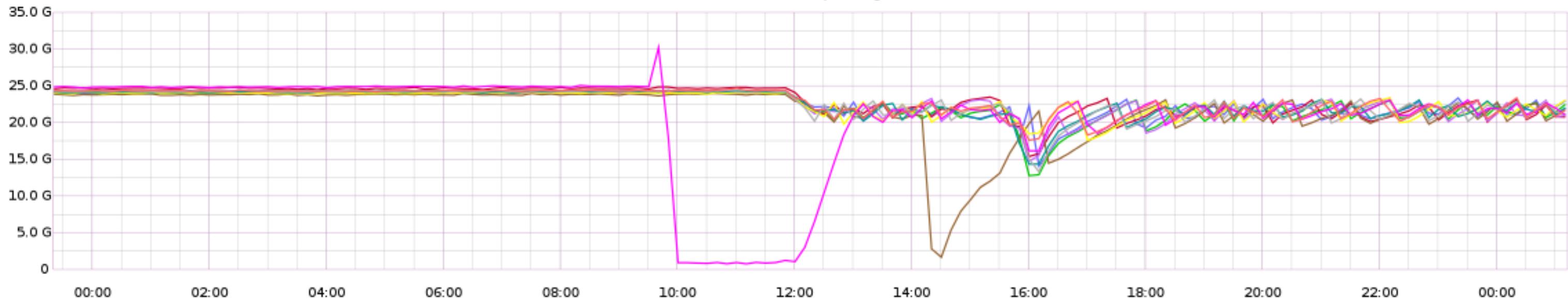


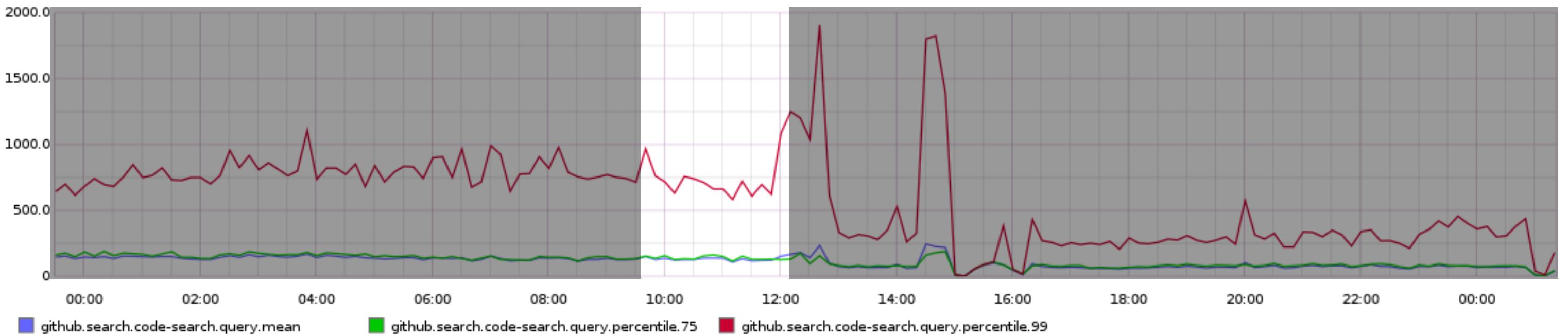
Code Search Heap Usage

Out Of Memory



codesearch heap usage





node name	disk	used	free	percent
codesearch-storage1	6.9T	5.9T	1022G	86%
codesearch-storage2	6.9T	6.2T	699G	91%
codesearch-storage3	6.9T	6.1T	841G	89%
codesearch-storage4	6.9T	6.0T	935G	87%
codesearch-storage5	6.9T	6.3T	630G	92%
codesearch-storage6	6.9T	6.2T	672G	91%
codesearch-storage7	6.9T	6.1T	859G	88%
codesearch-storage8	6.9T	6.1T	843G	88%
codesearch-storage9	6.9T	6.1T	870G	88%
codesearch-storage10	6.9T	6.0T	921G	87%

Add Capacity



Add Capacity

Logical
Volume
Manager

7.6 TB

6.9 TB

7.6 TB

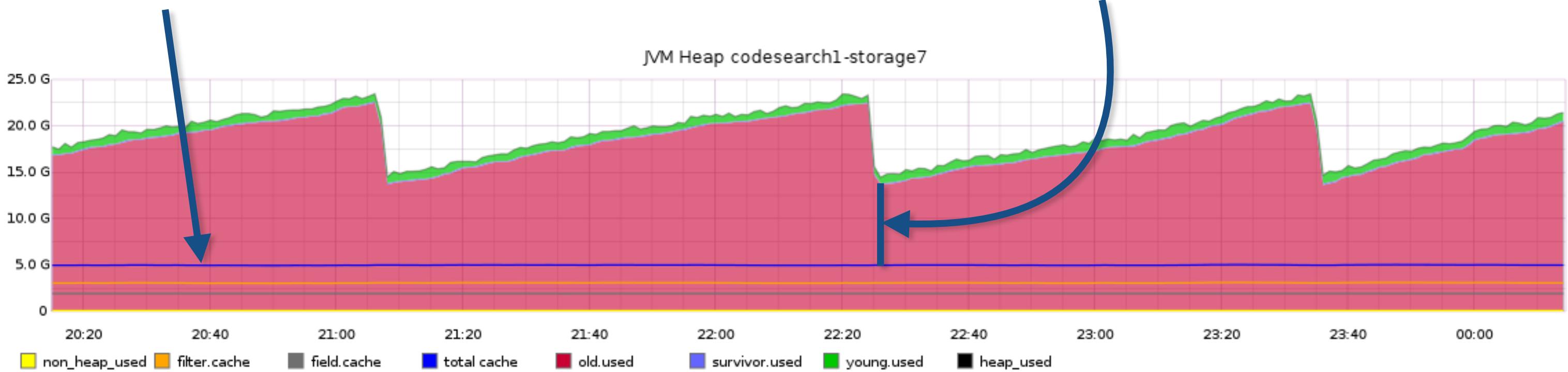
7.6 TB

How Did This Happen?

JVM Heap Usage

Total Cache Size

Lucene Segments



Prevention

Nagios®

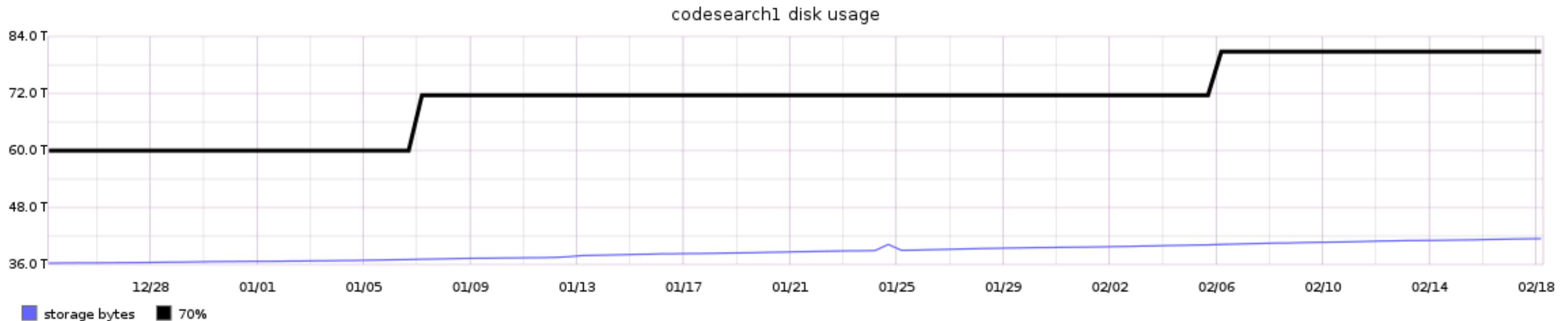
PAGERDUTY



/es forecast disk codesearch1



codesearch1 will reach 70% disk usage
in 302 days (2015-12-31) with 93% confidence



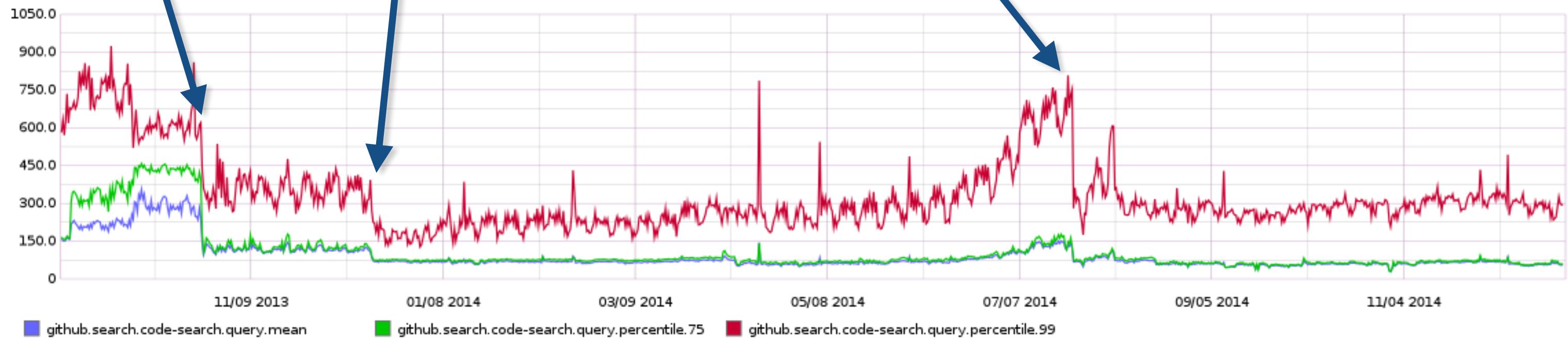
Heap Exhaustion

- ← We were ignoring key metrics
- ⌚ We added alerts for key metrics
- We created tools to forecast growth

New Cluster

New Queries

Heap Exhaustion

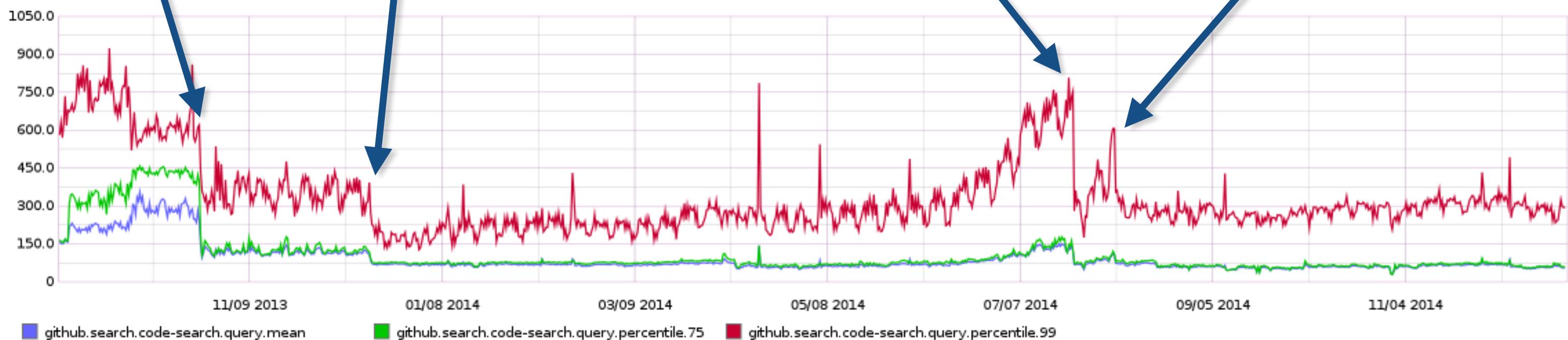


New Cluster

New Queries

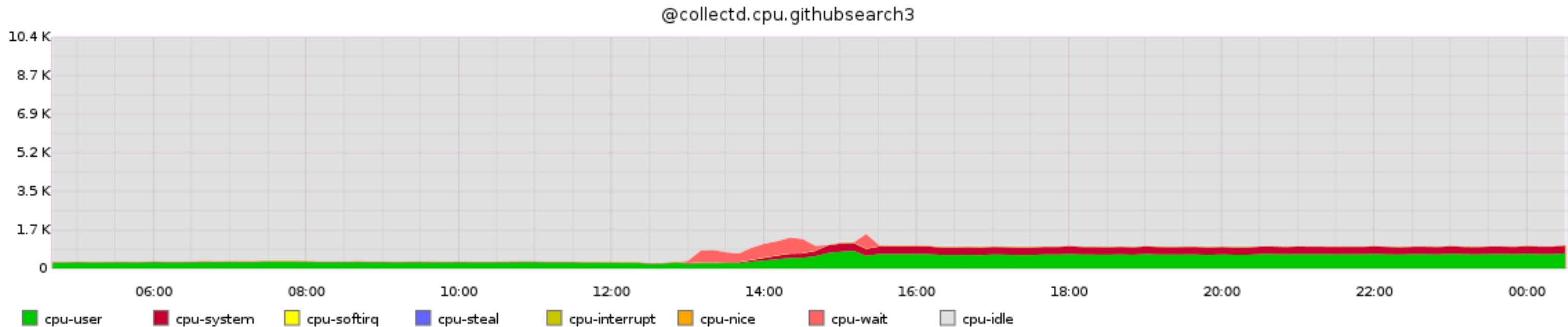
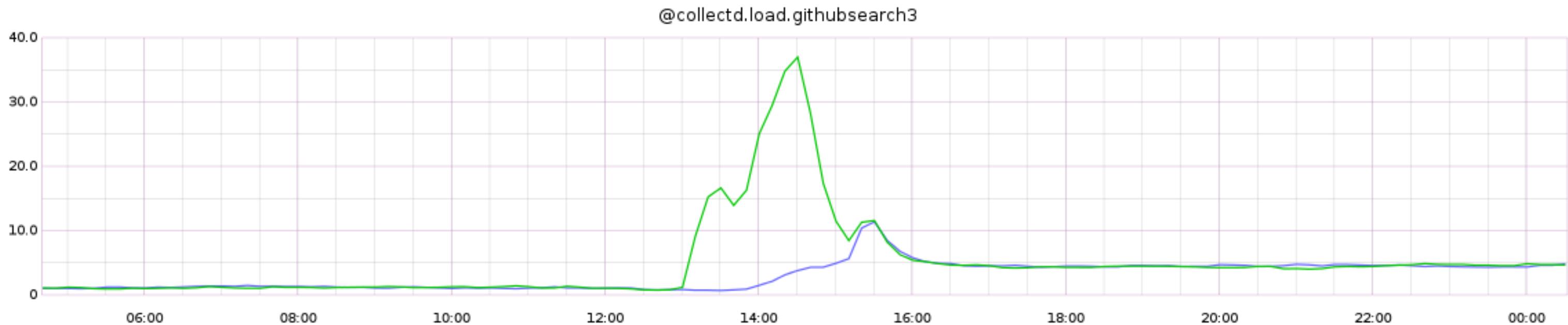
Heap Exhaustion

New Cluster

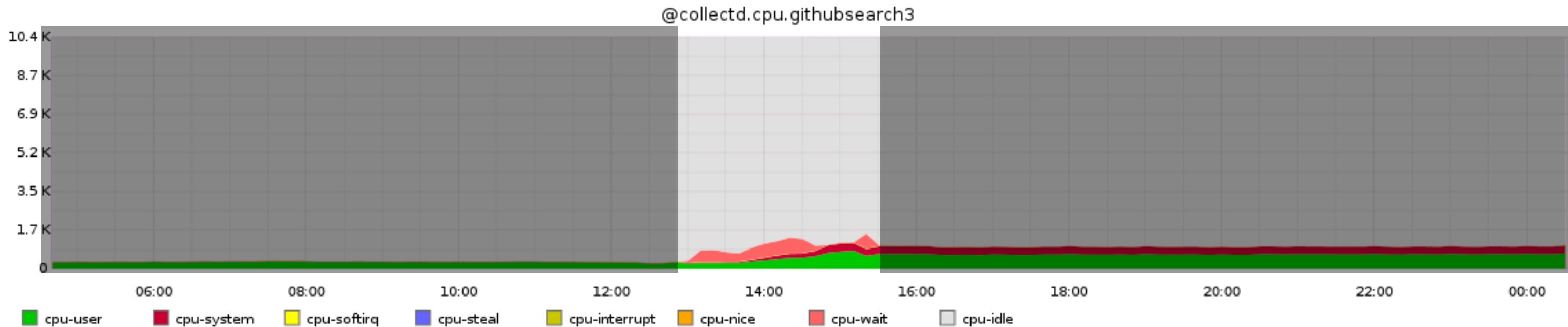
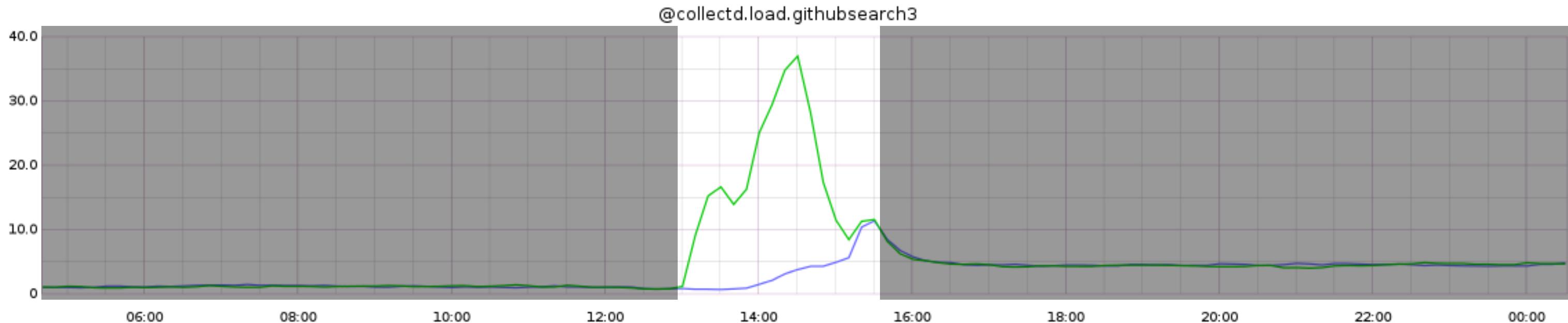


One More Story

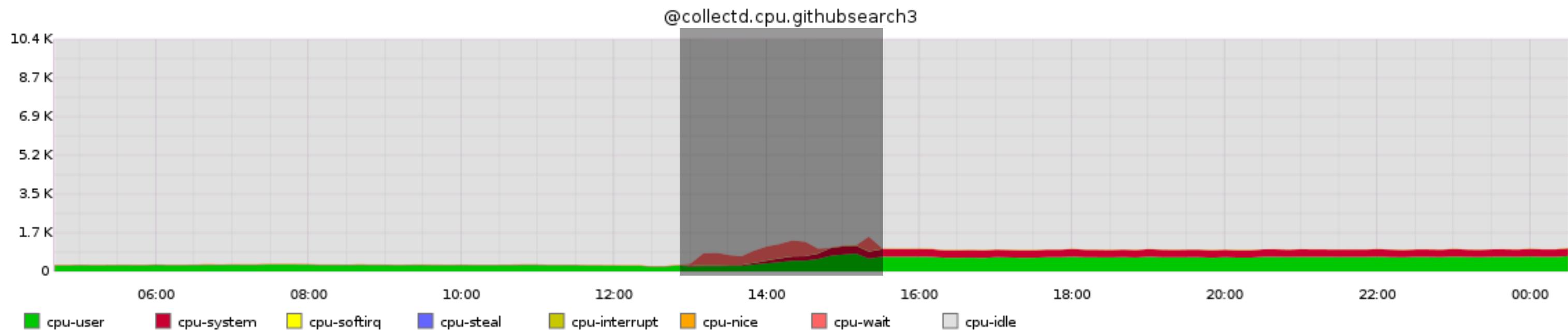
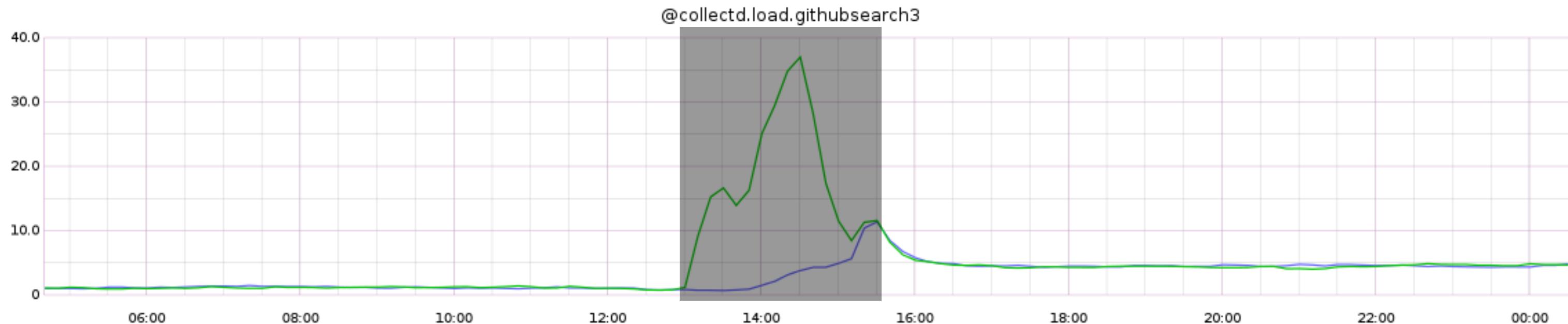
Upgrade to ES 1.4.2



Upgrade to ES 1.4.2



Upgrade to ES 1.4.2



Hot Threads

97.4% (487.1ms out of 500ms) cpu usage by thread 'elasticsearch[githubsearch3-storage1-cp1-prd][management][T#2]'
9/10 snapshots sharing following 9 elements

org.elasticsearch.action.admin.indices.stats.ShardStats.<init>(ShardStats.java:49)

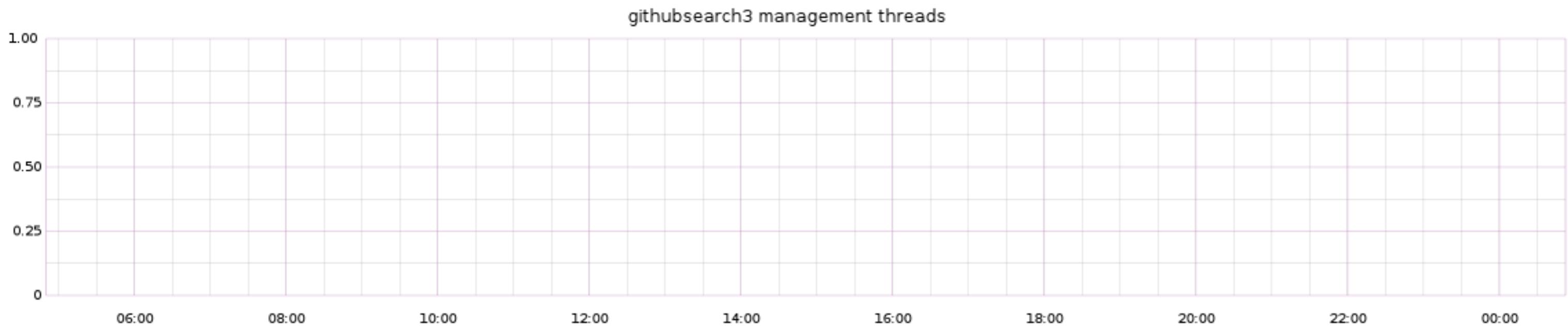
97.3% (486.3ms out of 500ms) cpu usage by thread 'elasticsearch[githubsearch3-storage1-cp1-prd][management][T#3]'
2/10 snapshots sharing following 20 elements

java.io.UnixFileSystem.getLength(Native Method)

96.4% (482.1ms out of 500ms) cpu usage by thread 'elasticsearch[githubsearch3-storage1-cp1-prd][management][T#4]'
2/10 snapshots sharing following 19 elements

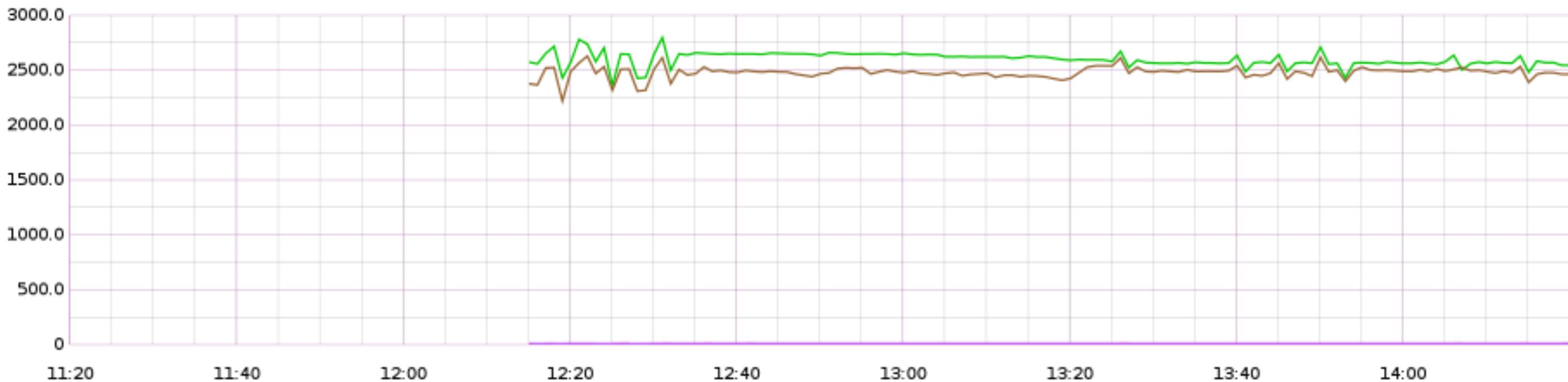
org.apache.lucene.store.FSDirectory.listAll(FSDirectory.java:223)

Management Threads



Amen Sampler

/_nodes/_local/stats



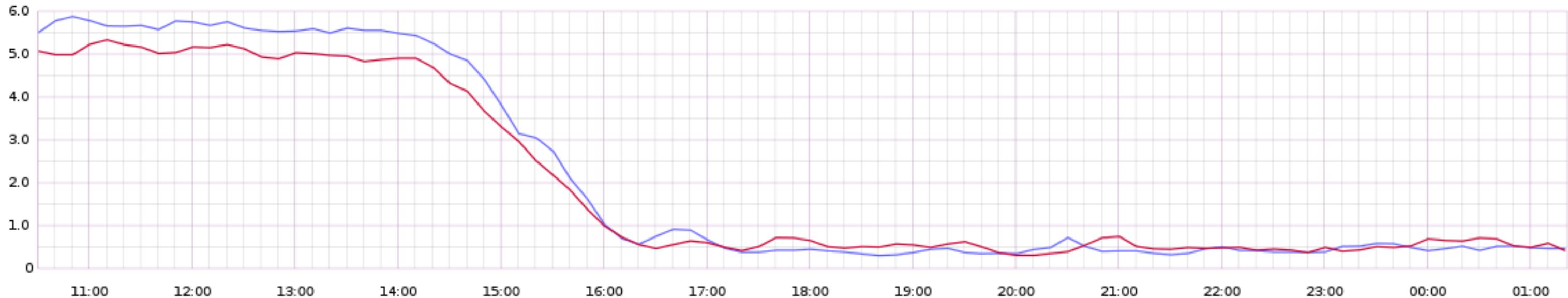
HAProxy Change

~~/_nodes/_local/stats~~ \Rightarrow /

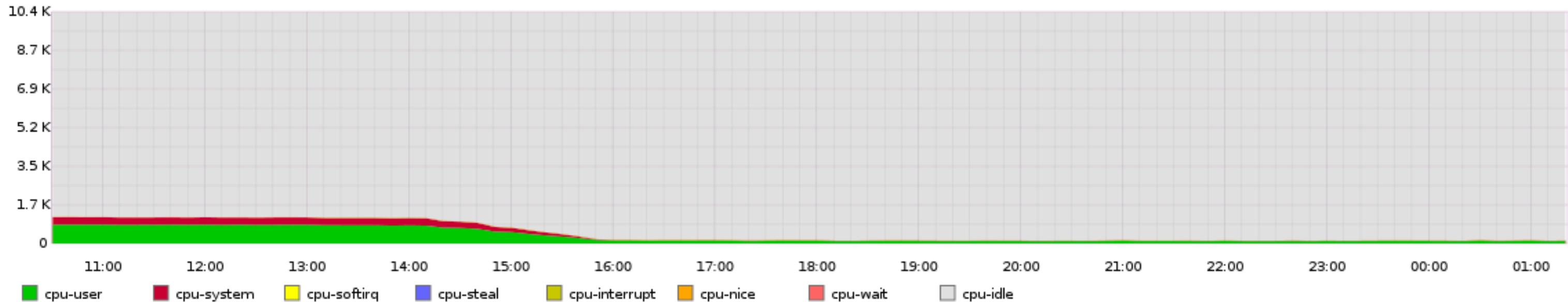


HAProxy Change

@collectd.load.githubsearch3



@collectd.cpu.githubsearch3



Upgrade to ES 1.4.2

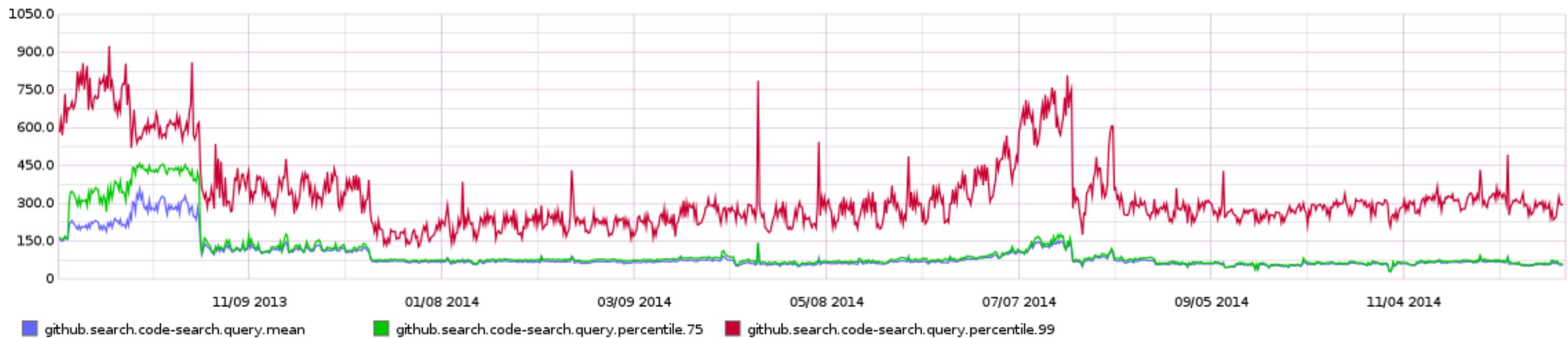
- ← We were missing some important metrics
- ⌚ We have an entire ecosystem
- We have confidence in ES 1.4.2

← Where have we come from

⌚ What we are doing now

→ Where we are going

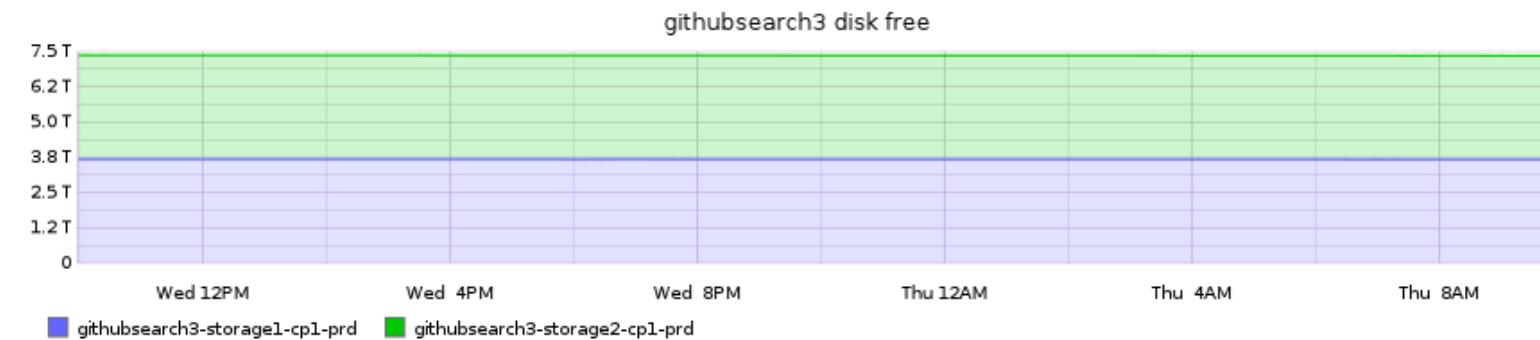
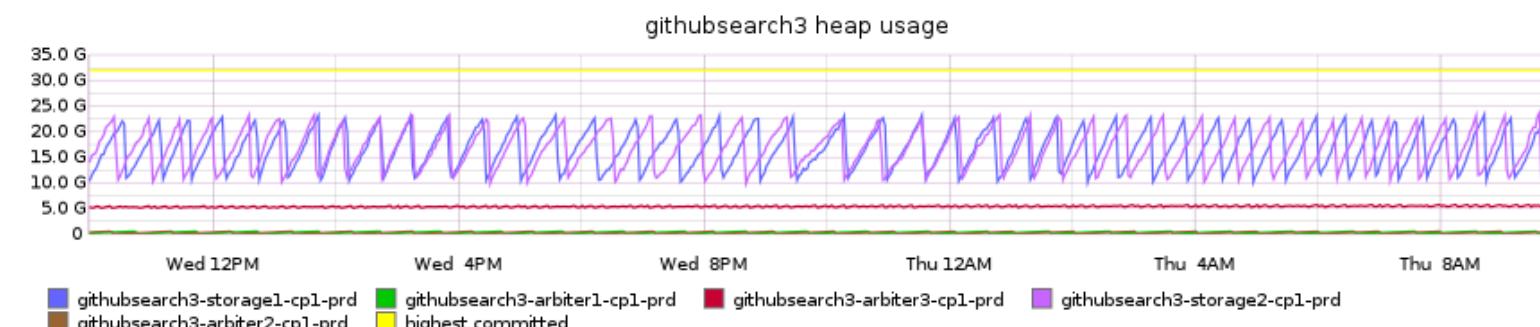
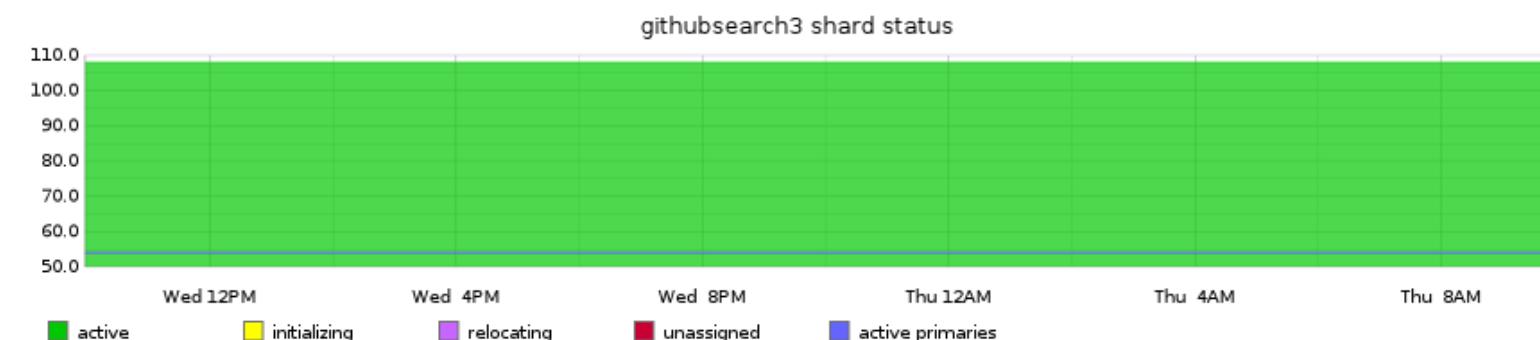
← Where have we come from ←



⌚ What are we doing now ⌚

-30min // -1h // -6h // **-1d** // -1w // -1mon // -3mon // -6mon

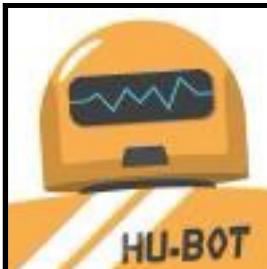
githubsearch3



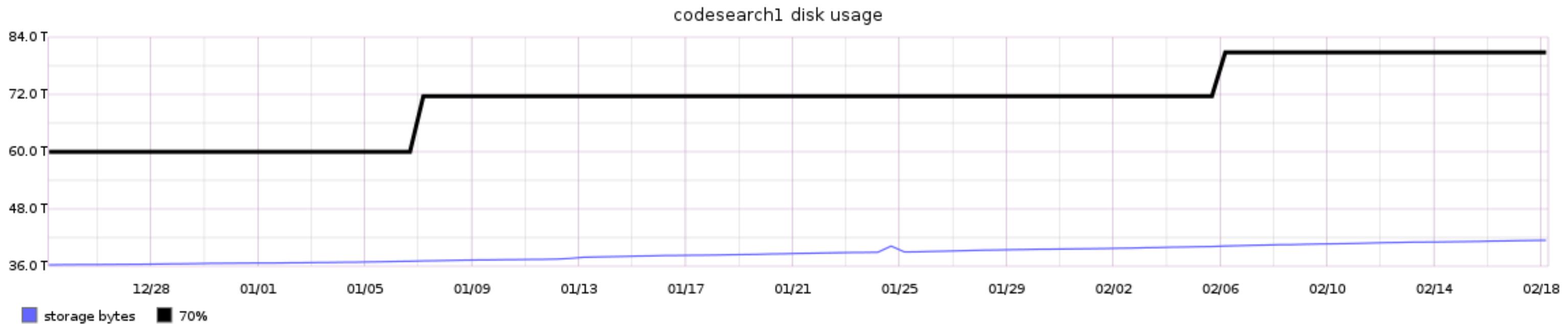
→ Where are we going →



/es forecast disk codesearch1



codesearch1 will reach 70% disk usage
in 302 days (2015-12-31) with 93% confidence



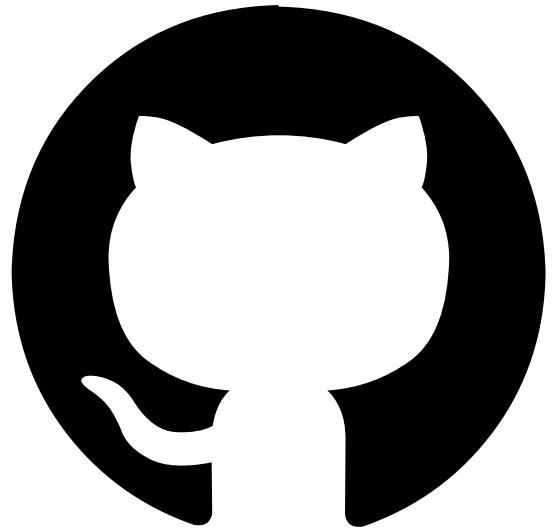
GitHub Search Team



Grant Rodgers
github.com/grantr



Tim Pease
github.com/TwP



The End