

如何打造大规模互联网企业的 监控告警平台

-- 以携程hickwall为例



唐锐华
rhtang@ctrip.com



促进软件开发领域知识与创新的传播



关注InfoQ官方信息
及时获取QCon软件开发者
大会演讲视频信息



[北京站] 2016年12月2日-3日

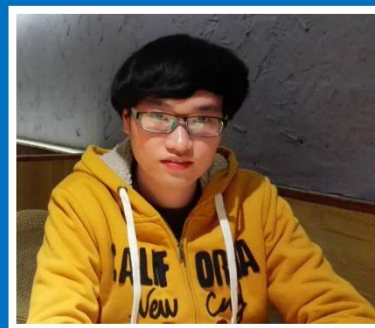
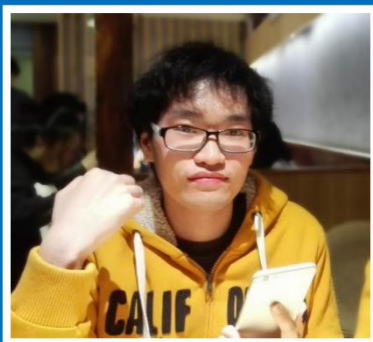
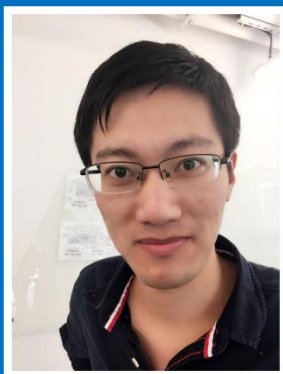
咨询热线: 010-89880682



[北京站] 2017年4月16日-18日

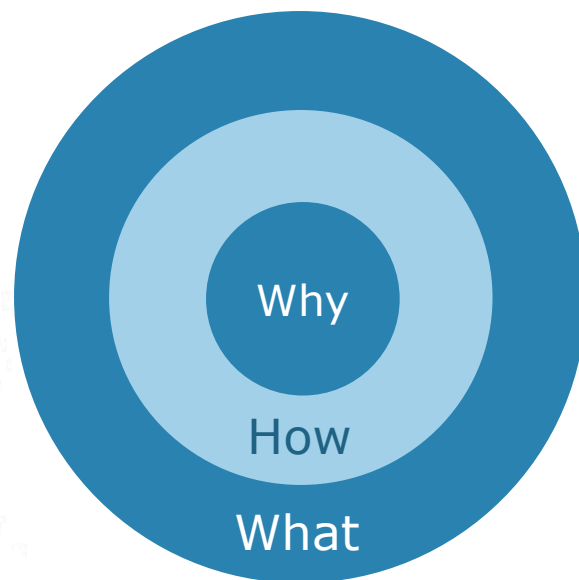
咨询热线: 010-64738142

自我介绍



目录

- 1 原动力 和 架构设计
- 2 数据采集 和 加工
- 3 告警模块设计
- 4 数据展现 以及 配置中心
- 5 可靠性 与 吞吐量



- 1 原动力 和 架构设计
- 2 数据采集 和 加工
- 3 告警模块设计
- 4 数据展现 以及 配置中心
- 5 可靠性 与 吞吐量

原动力 和 架构设计

- 规模
- 数据入口
- 集成

千万级监控点

数据被转手多次

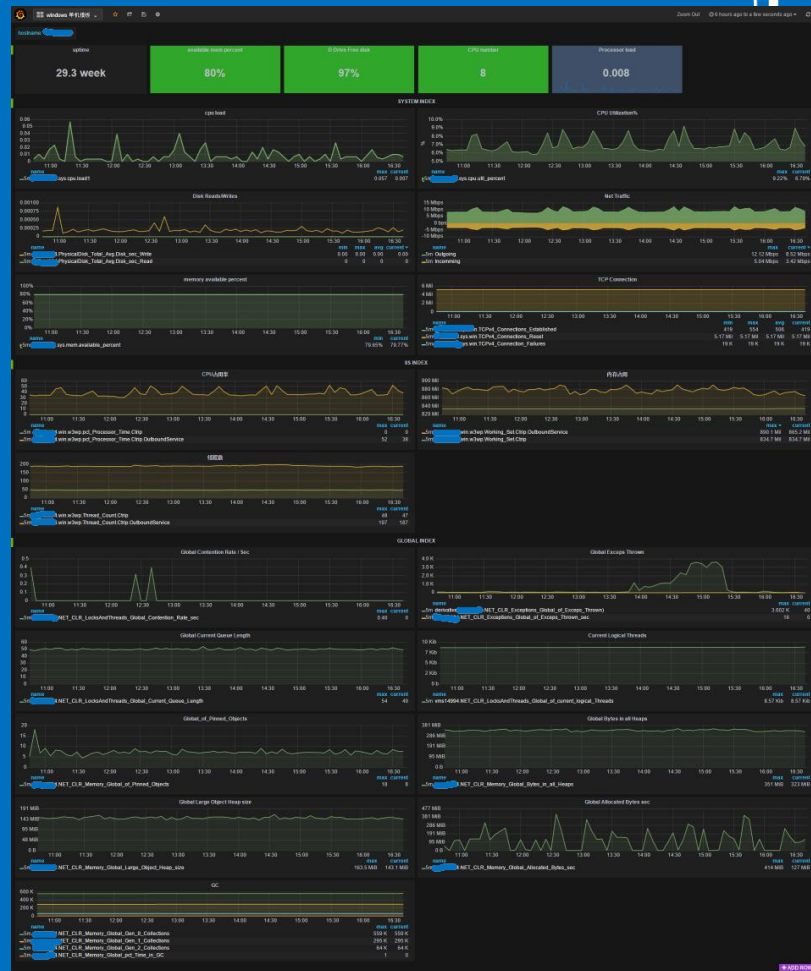
工具多，使用方法迥异



原动力 和 架构设计

目标

- 10 Million+ 监控指标
- 200K+/s 数据流量
- 强大的数据展现
- 集中配置管理
- 灵活的告警规则



原动力 和 架构设计

选择

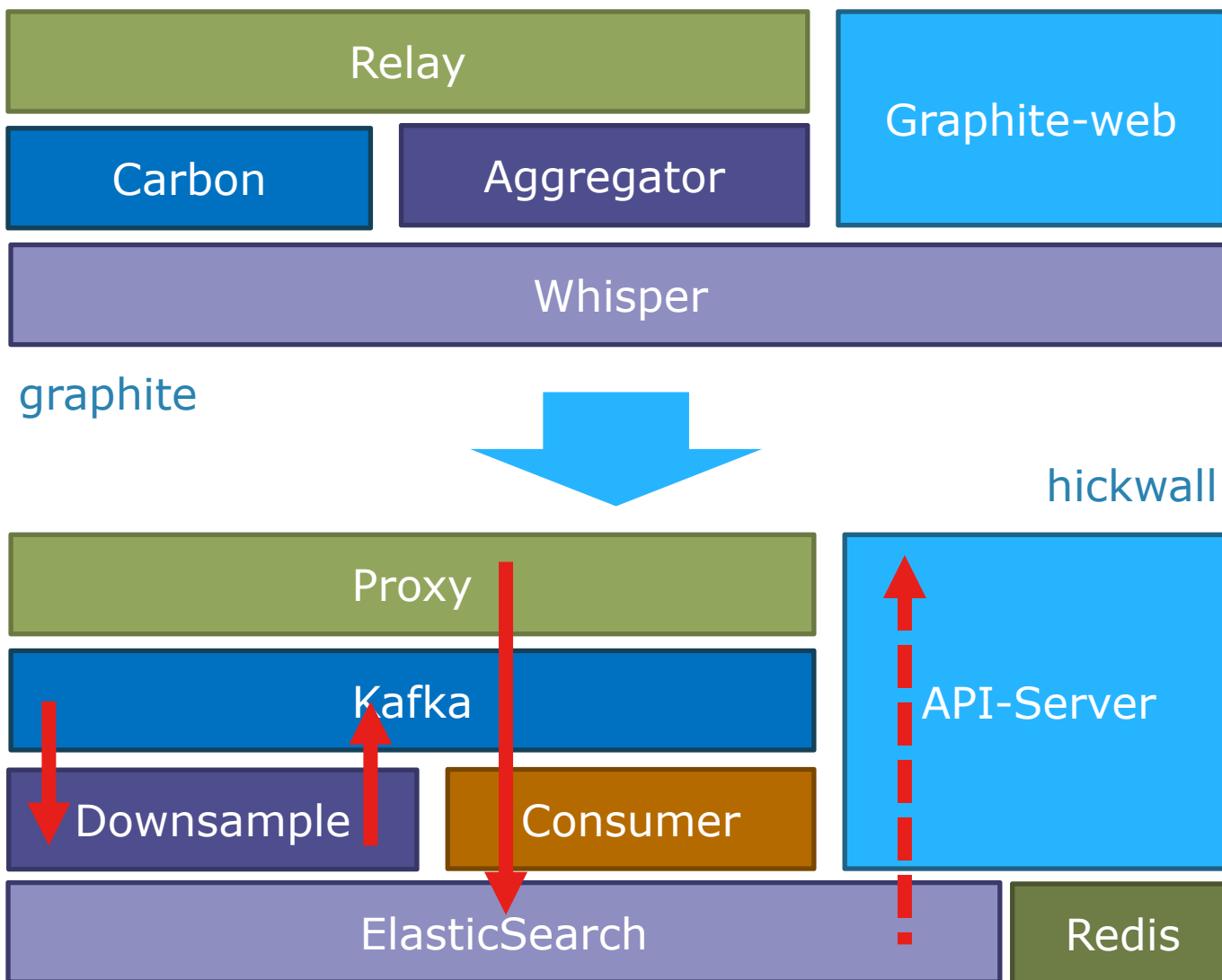
- Graphite生态 vs RRD
- RRD + Consistent Hashing
- Influxdb
- OpenTSDB



Graphite Over ElasticSearch

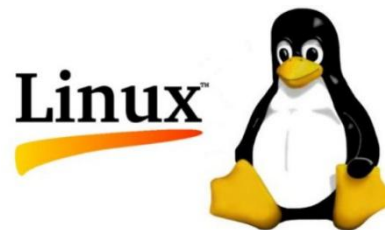
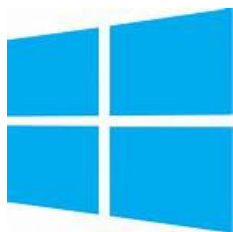
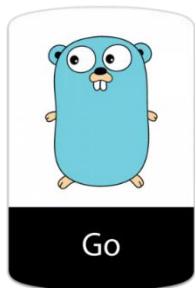
收益

- 享受整个 Graphite 生态系统的红利
 - Grafana
 - 多语言的数据上报类库
 - 现有工具迁移成本低
- 享受 ES 的红利
 - PB级稳定, 强大的存储和搜索
 - 水平扩展, 便于管理



- 1 原动力 和 架构设计
- 2 数据采集 和 加工
- 3 告警模块设计
- 4 数据展现 以及 配置中心
- 5 可靠性 与 吞吐量

* hickwall agent



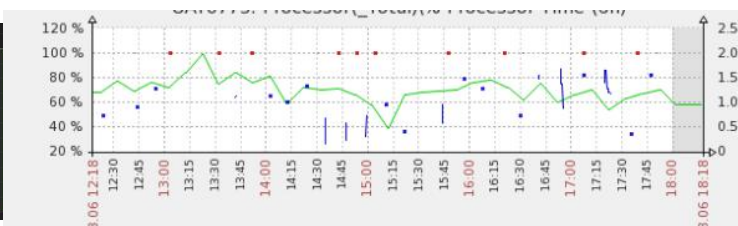
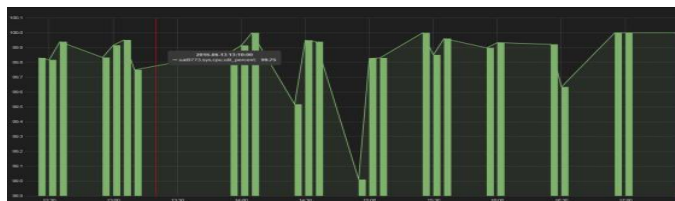
* hickwall remote collector



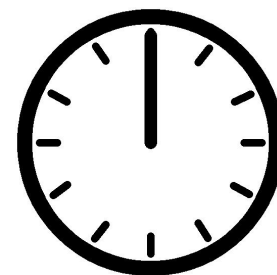
Scala

Or 任意常用语言

与各种诡异情况做斗争



变态时钟



数据采集 和 加工

Proxy

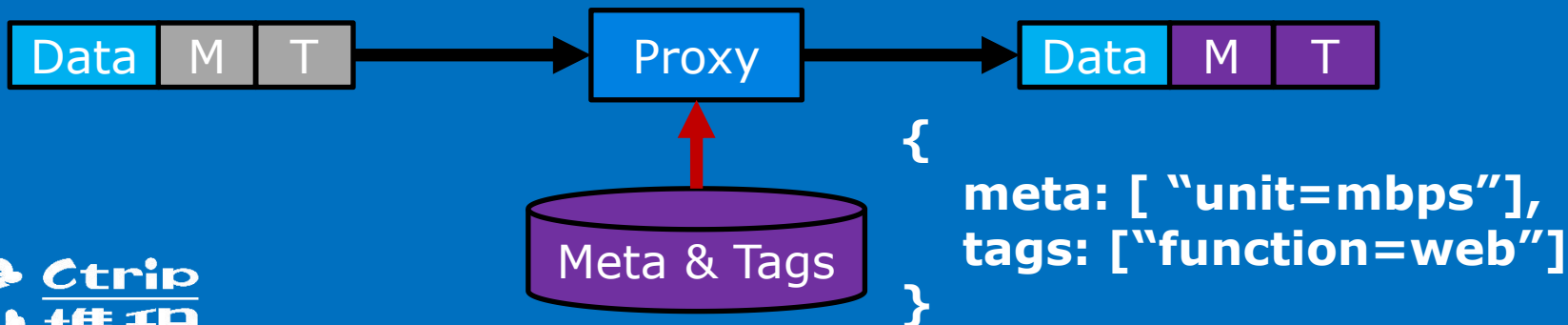
```
PORT=2003
SERVER=graphite.your.org
echo "local.random.diceroll 4 `date +%s`" | nc -q0 ${SERVER} ${PORT}
```

(Note: In the original image, 'local.random.diceroll' is highlighted with a red box and an arrow points to the word 'Metric' in the header above.)



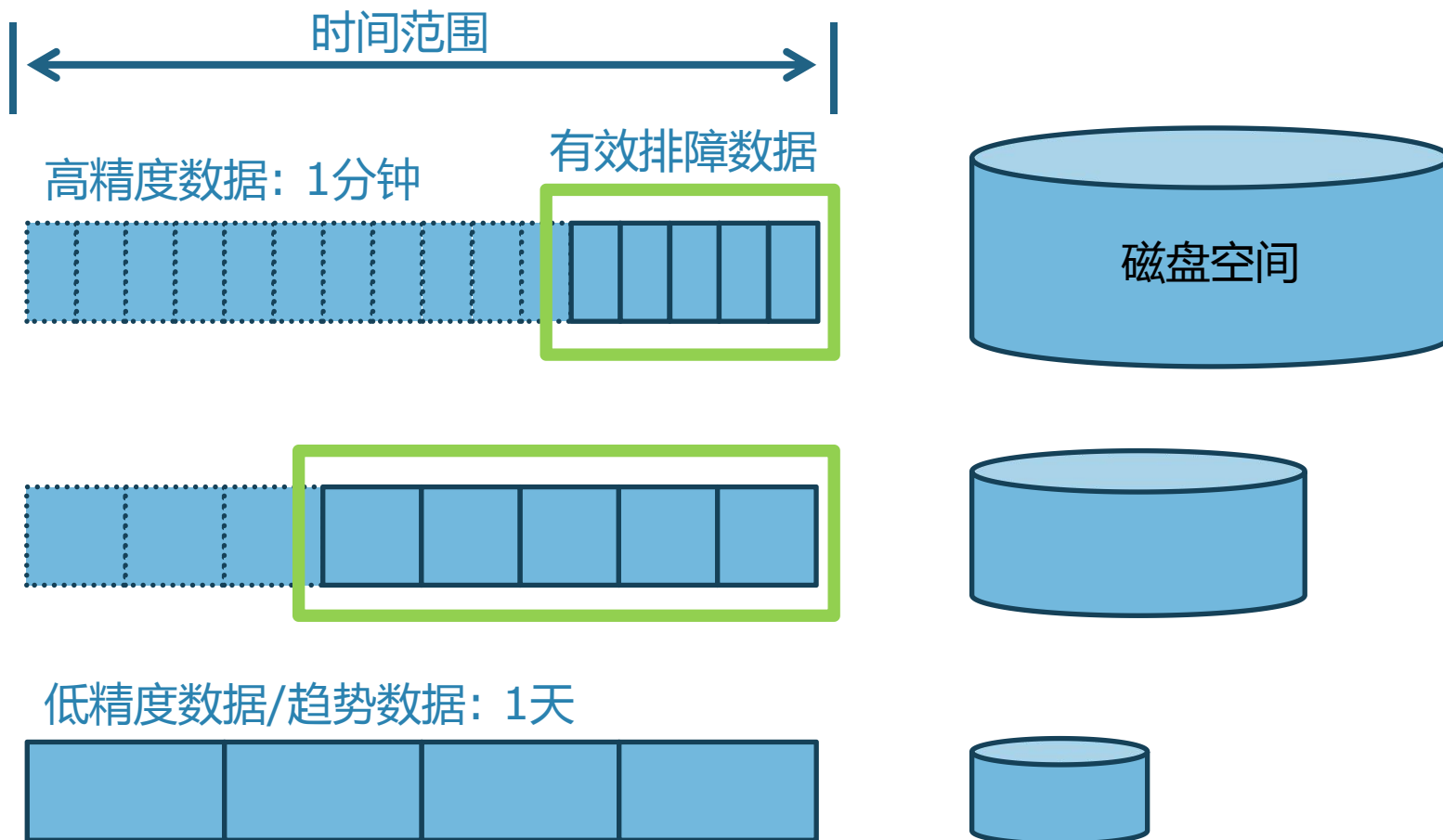
```
struct DataPoint
{
    1: required string metric,
    2: required i64 timestamp,
    3: optional double value,
    4: optional string str_value,
    5: optional map<string, string> tags,
    6: optional map<string, string> meta,
}
```

(Note: In the original image, lines 5 and 6 are highlighted with a red box.)



数据采集 和 加工

Downsample - 数据精度 vs 存储资源



- 1 原动力 和 架构设计
- 2 数据采集 和 加工
- 3 告警模块的设计
- 4 数据展现 以及 配置中心
- 5 可靠性 与 吞吐量

Local(Agent)

Server

本地监控指标

全局监控指标

元信息单一

补全后元信息丰富

数据保留时间短

数据保留时间长

服务端压力小

服务端压力大

受宕机影响, 范围小

易维护但宕机影响范围广

多监控指标告警, 同比/环比告警等

host.fs.free < 10%

cluster.net.out_going.wow < 30%

告警模块设计

- DSL - Domain Specific Language
- JavaScript

告警规则:

连续5次可用内存百分比低于10%, 告警

```
Check(  
  Get("sys.mem.available_percent", 5).count("lt", 10) >= 5,  
  CRITICAL ,  
  "可用内存低于10 %"  
)
```

```
Critical(Get("sys.mem.available_percent", 5).count("lt", 10) >= 5)
```

简约 vs 灵活

Init DSL

```
trigger.requireMetric('m1', 'sys.mem.available_percent', "6m");
```

run DSL

```
if ( V.m1.length <= 4 ) {  
    return UNKNOWN("data missing!")  
}
```

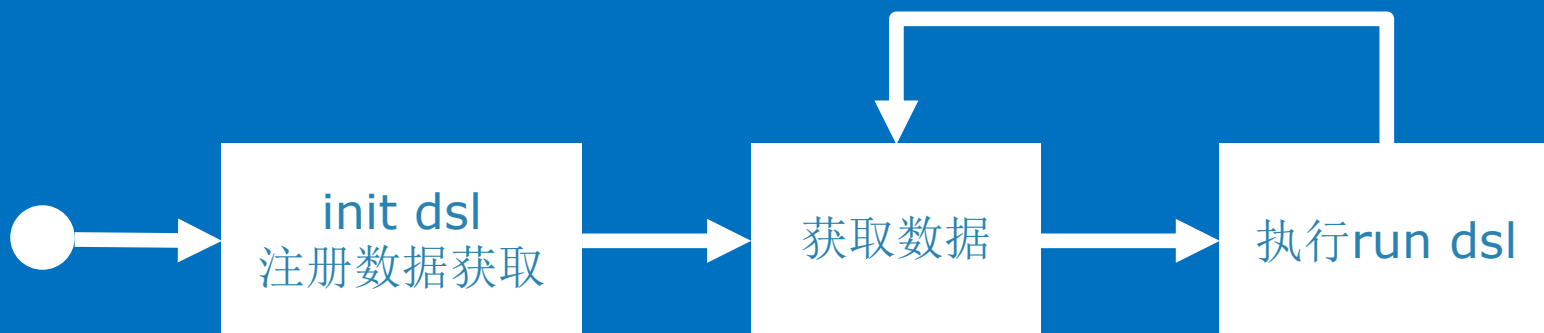
```
if ( V.m1.count('<', 10 ) >= 5 ) {  
    return CRITICAL("sys.mem.available_percent lower than 10");  
}
```

```
return OK('length:' + V.m1.length);
```

故意写这么复杂

告警模块设计

服务器端告警需要优化性能



提高后端提供数据的效率

- 1 原动力 和 架构设计
- 2 数据采集 和 加工
- 3 告警模块设计
- 4 数据展现 和 配置中心
- 5 可靠性 与 吞吐量

数据展现 以及 配置中心

API

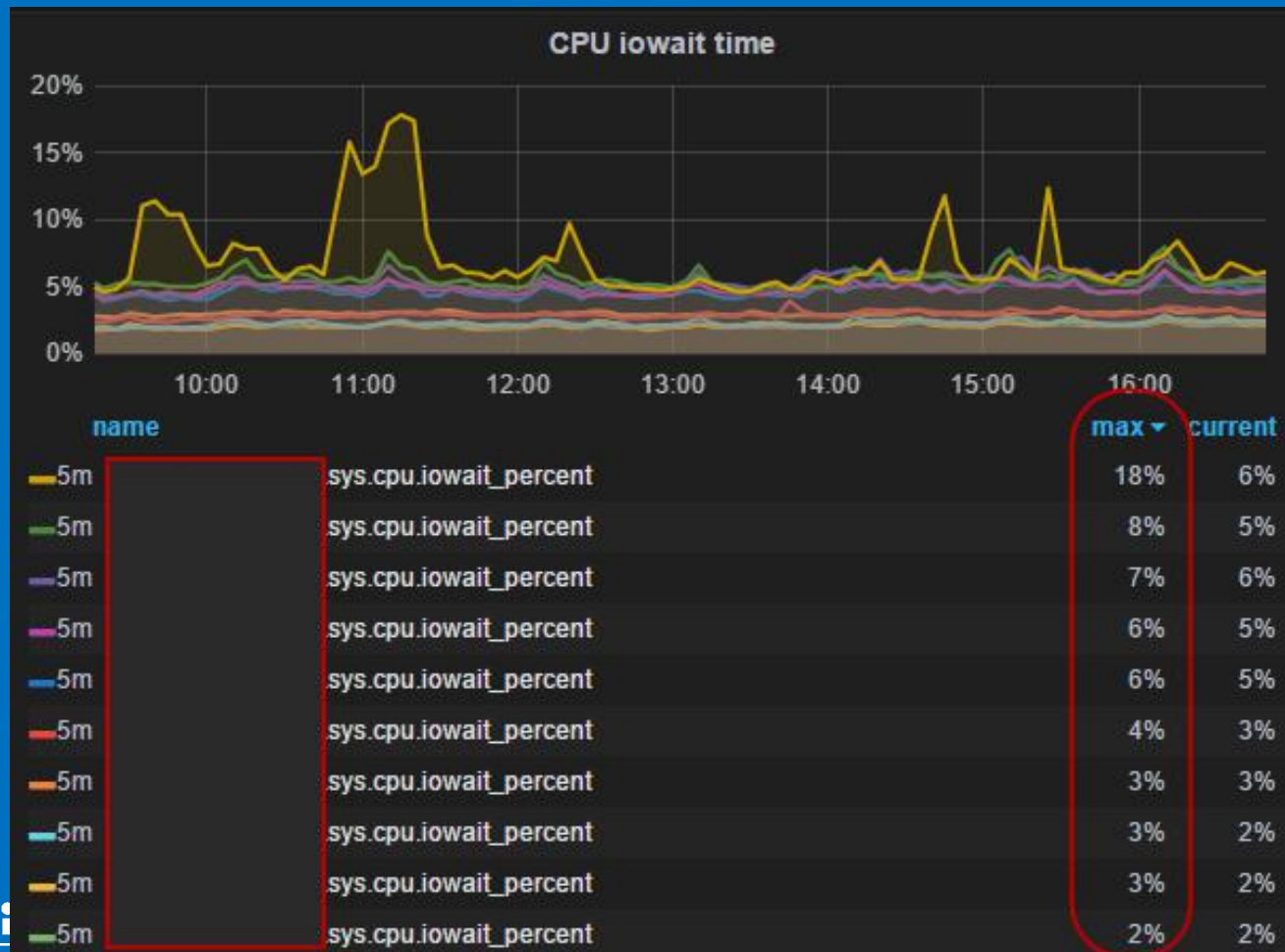
- 实现常用Graphite API
- 实现大量常用Graphite Functions
 - E.g.: alias, cumulative, currentAbove, derivative, mostDeviant, scaleToSeconds, timeShift, ...
- 增强: 返回 标签 以及 元信息
 - E.g: Meta[Unit], Meta[Ip]
- 变态查询过滤

数据展现 以及 配置中心

Grafana 不用改造直接可用

Max 排序

Grafana 功能

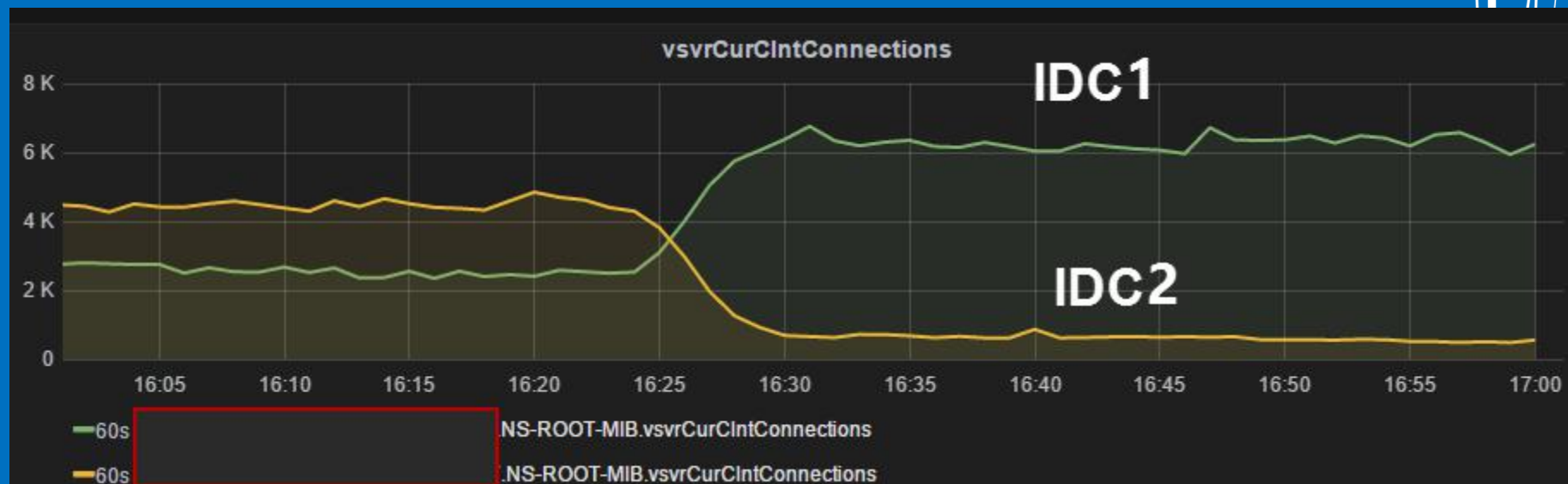


数据展现 以及 配置中心

Grafana 不用改造直接可用

Grafana 功能

IDC 切换

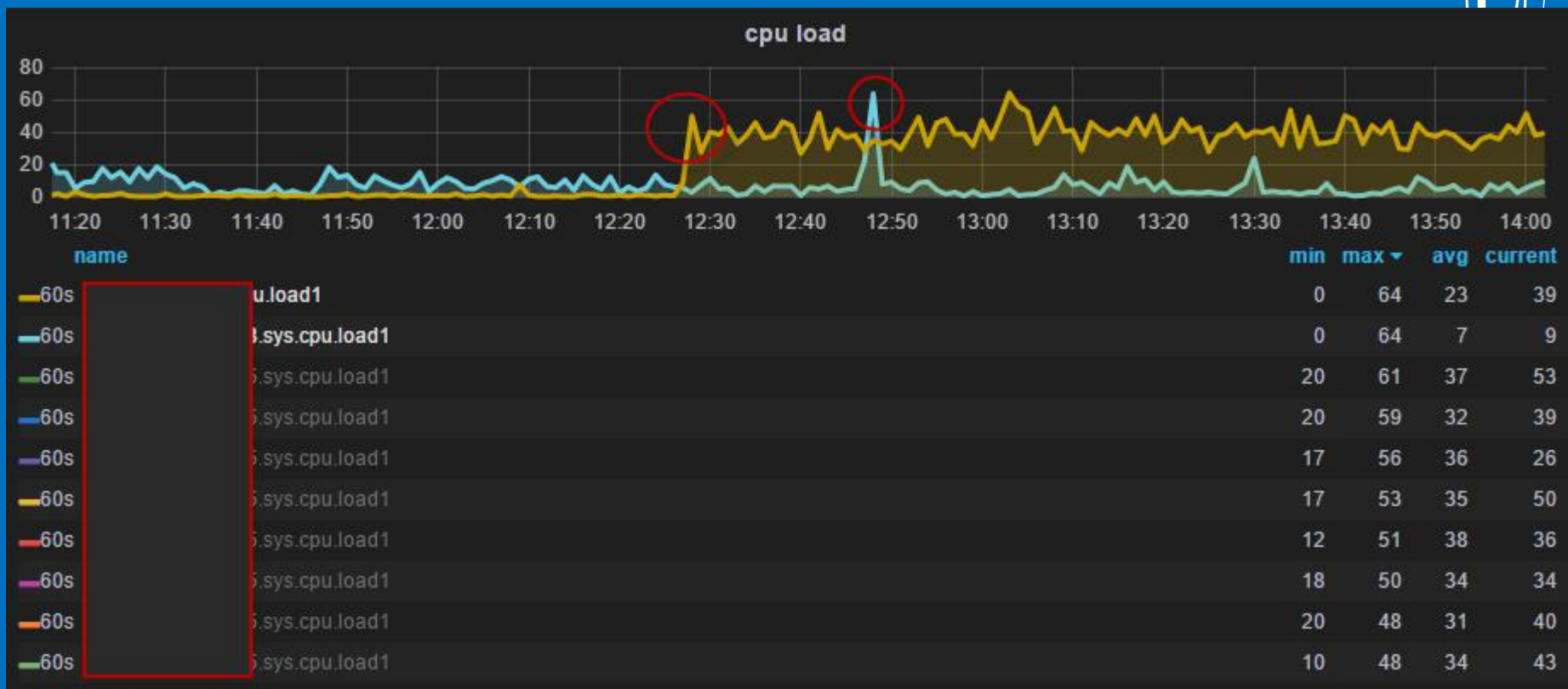


数据展现 以及 配置中心

Grafana 不用改造直接可用

Graphite API

mostDeviant 一段时间内变化最大的

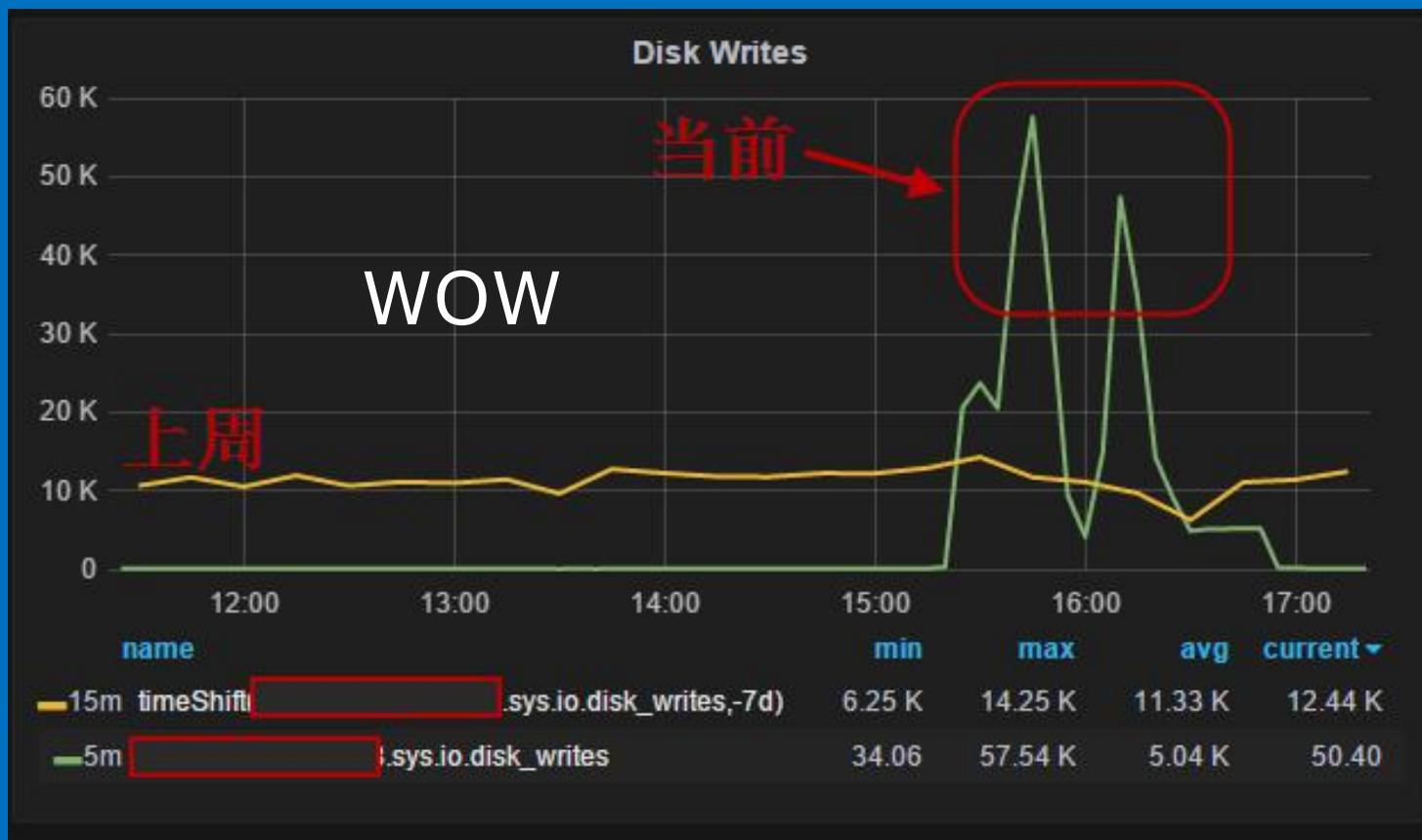


数据展现 以及 配置中心

Grafana 不用改造直接可用

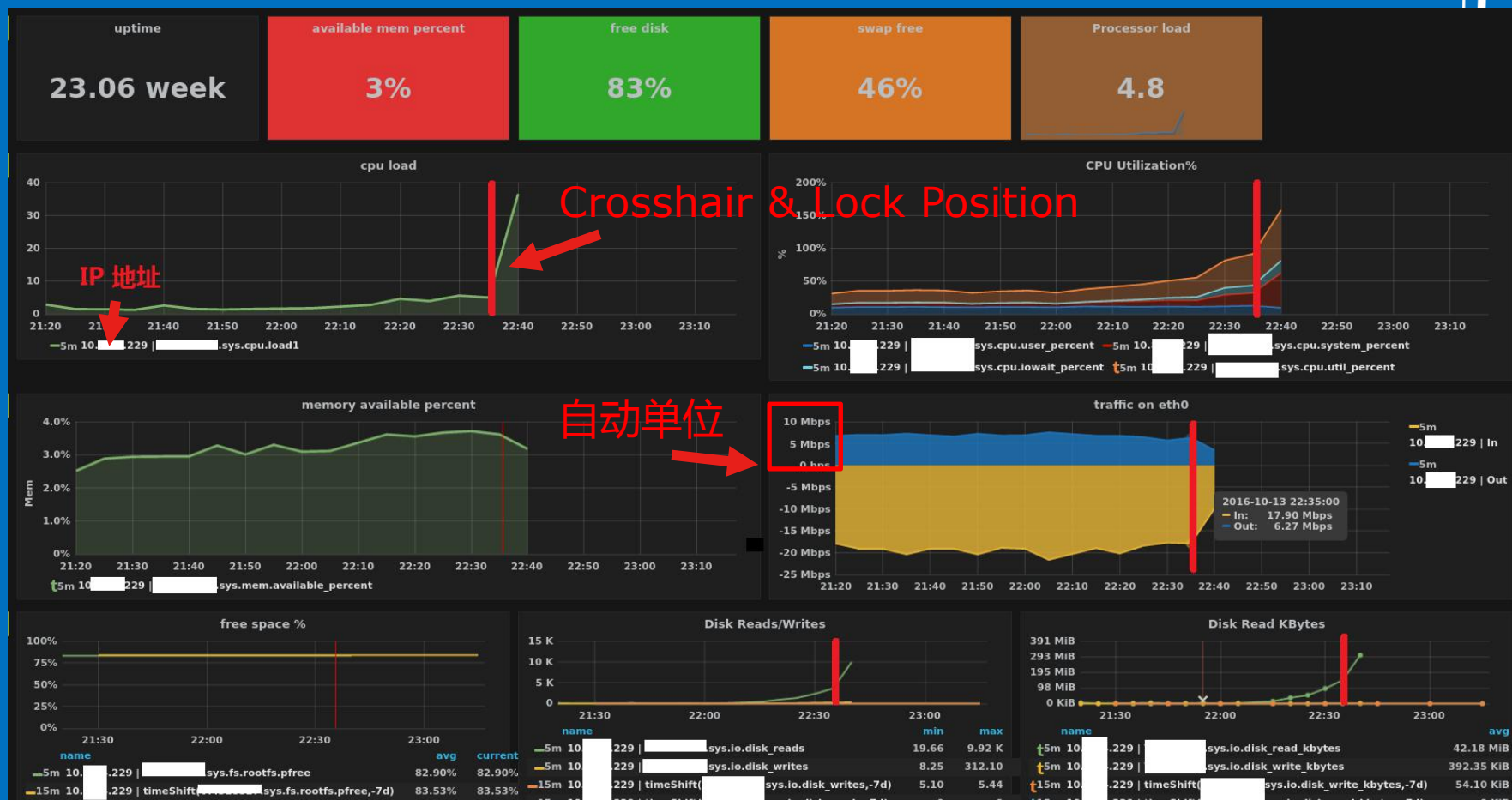
Graphite API

wow 周环比



数据展现 以及 配置中心

定制开发锦上添花



数据展现 以及 配置中心

进程级Top5



请告诉我这条线背后的故事，而不是一堆难懂的日志。

23431 57.541 MB
23431 57.541 MB
23431 57.541 MB
23431 57.541 MB
23431 57.541 MB

Top 1		Top 2		Top 3		Top 4		Top 5	
Value	Pid	Value	Pid	Value	Pid	Value	Pid	Value	Pid
14.325 GB	5559	239.37 MB	23431	57.541 MB	21388	20.242 MB	1	19.562 MB	1
14.325 GB	5559	239.37 MB	23431	57.541 MB	21388	20.242 MB	1	19.562 MB	1
14.325 GB	5559	239.37 MB	23431	57.541 MB	21388	20.242 MB	1	19.562 MB	1
14.325 GB	5559	239.37 MB	23431	57.541 MB	21388	20.242 MB	1	19.562 MB	1
14.325 GB	5559	239.37 MB	23431	57.541 MB	21388	20.242 MB	1	19.562 MB	1
14.406 GB	27606	9.004 GB	5559	235.532 MB	27727	204.628 MB	6135	28.021 MB	6135
14.406 GB	27606	9.004 GB	5559	235.532 MB	27727	204.632 MB	6135	28.164 MB	6135
14.406 GB	27606	9.004 GB	5559	235.532 MB	27727	177.492 MB	6135	28.172 MB	6135
14.406 GB	27606	9.004 GB	5559	235.532 MB	27727	177.5 MB	6135	28.172 MB	6135
14.406 GB	27606	9.004 GB	5559	235.532 MB	27727	146.817 MB	6135	28.172 MB	6135

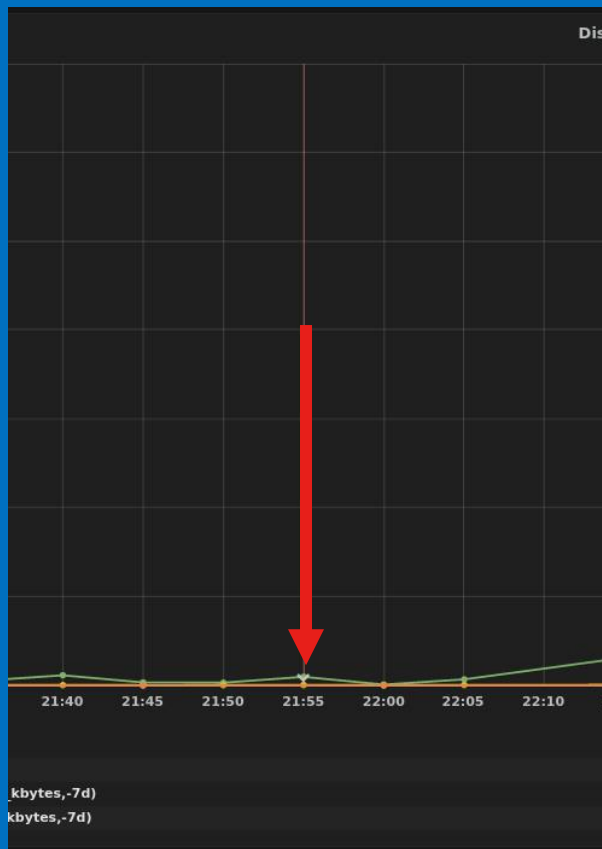
27606 9.004 GB
27606 9.004 GB
27606 9.004 GB
27606 9.004 GB
27606 9.004 GB

Pid And Commands

Pid	Commands
1	systemd
5559	java
6135	salt-minion
16271	java
21388	hickwall
23431	salt-minion
27606	salt-minion
27727	logagent-linux-

数据展现 以及 配置中心

另一个例子



TopN

Metric Compare ☒ Compare Current ☐

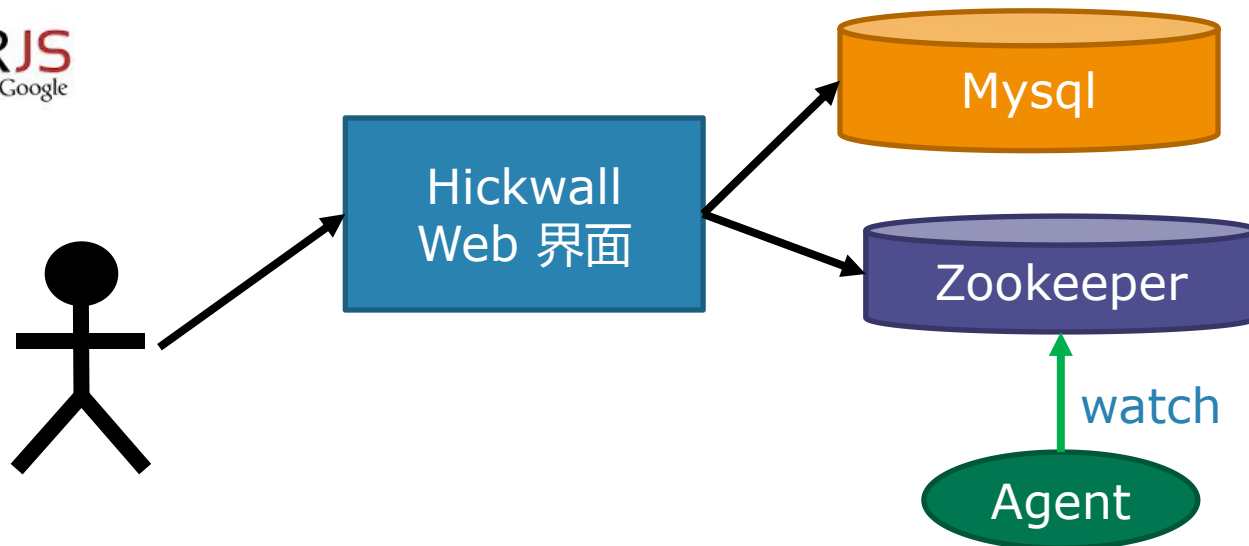
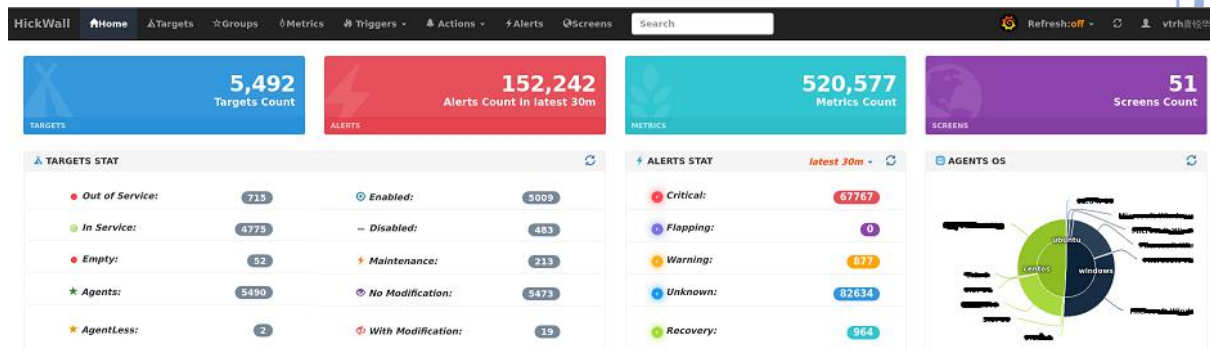
Time	Value	Top 1		Top 2		Top 3		Top 4		Top 5	
		Pid	Value	Pid	Value	Pid	Value	Pid	Value	Pid	Value
1476369307066	271176.08777	18572	5177	28118	3659.84	1003	2835.93	1265	1641.68	14900	1273.72
1476369229066	203242.84252	18572	7265.42	28118	4242.98	28821	3238.02	14900	2078.58	1003	1387.32
1476369121066	134394.48292	18572	2810.4	28118	2085.55	28821	1430.73	14900	868.22	1003	746.02
1476369050066	112528.33054	18572	3749.96	28118	2334.31	28821	1640.52	14900	1096.09	1003	649.76
1476368967066	16520.41537	18572	525.04	28118	331.16	1003	276	1565	203.91	28821	200.84
1476366915070	330.37823	31690	20.96	31693	11.2	29738	7.65	14924	0.41	30049	0.27
1476366855070	386.37884	18572	17.61	1	10.17	30049	0.07	787	0	122	0
1476366796070	1972.39423	1565	48.33	1265	32.9	1003	23.07	1008	21.03	998	18.02
1476366738070	9698.19592	18572	284.4	28118	118.72	28821	89.57	29736	36.45	14900	36.18
1476366677070	9337.60597	18572	309.32	28118	101.03	28821	87.04	1265	70.38	1565	53.66

Pid And Commands

Pid	Commands
1	
122	
787	
998	
1003	
1008	
1265	
1565	
14900	
14924	
18572	salt-minion
28118	supervisord
28821	
29736	

数据展现 以及 配置中心

配置中心



数据展现 以及 配置中心

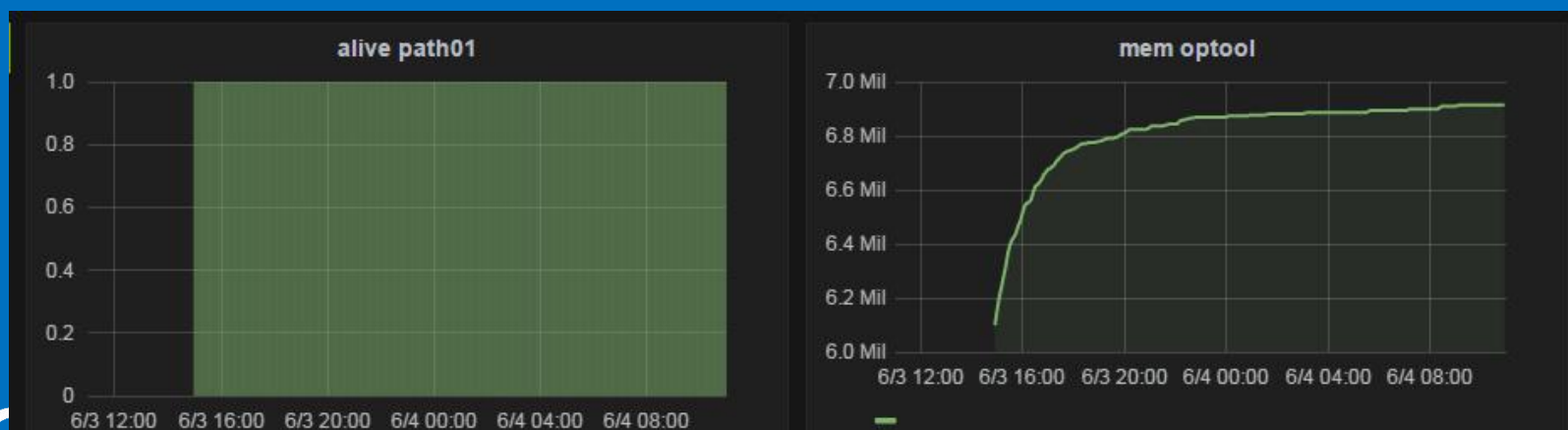
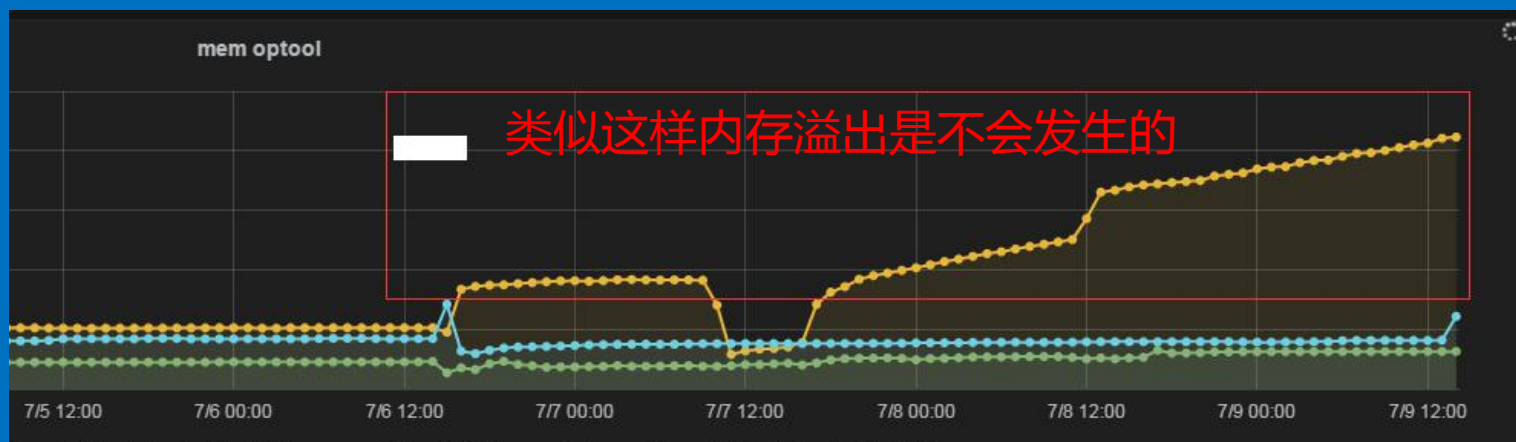
目标

- 100K+ 机器
- 配置更新实时生效
- 模板化配置
- 告警实时开关
- 标签管理，用模板自动化生成图表等

- 1 原动力 和 架构设计
- 2 数据采集 和 加工
- 3 告警模块设计
- 4 数据展现 以及 配置中心
- 5 可靠性 与 吞吐量

可靠性 与 吞吐量

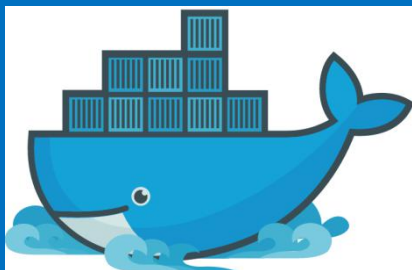
Agent端: 资源限制(内存, IO, 连接数, ...), 自杀



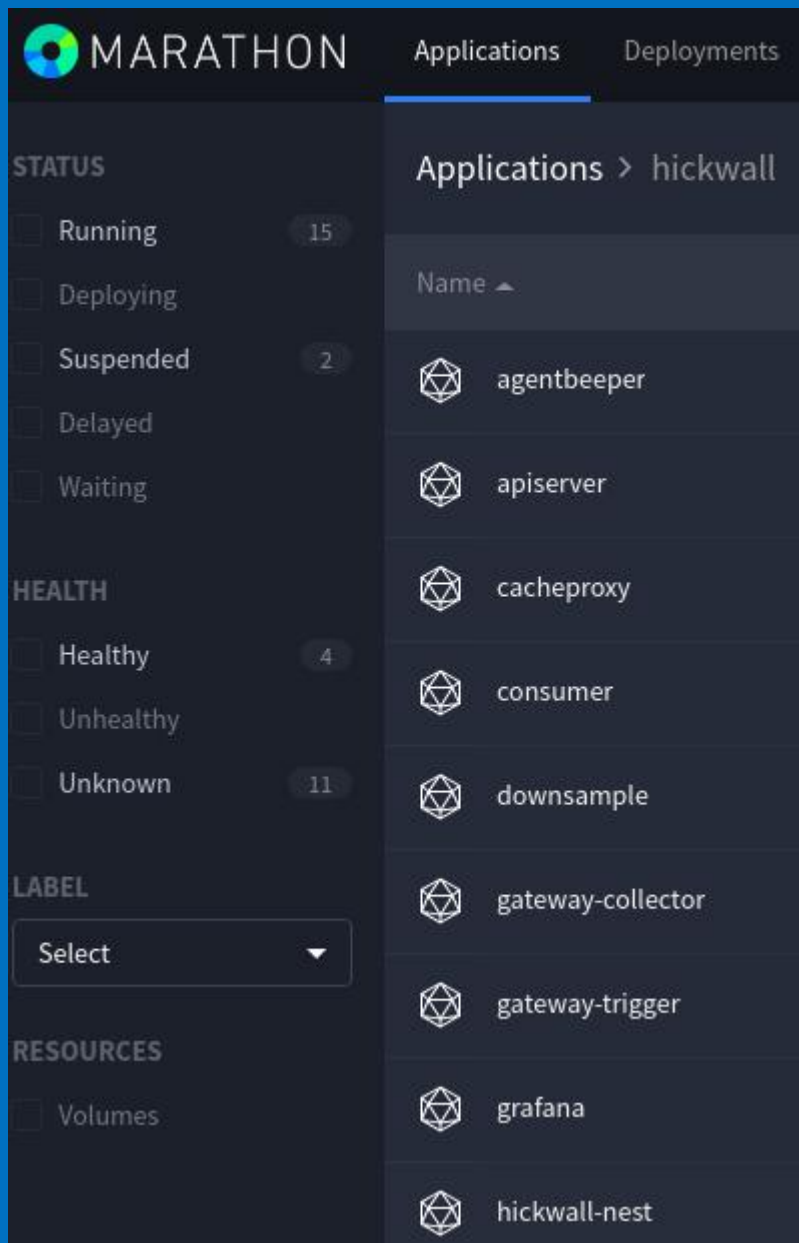
可靠性 与 吞吐量

服务端

- 享受docker红利
- 容器化发布
- 自恢复



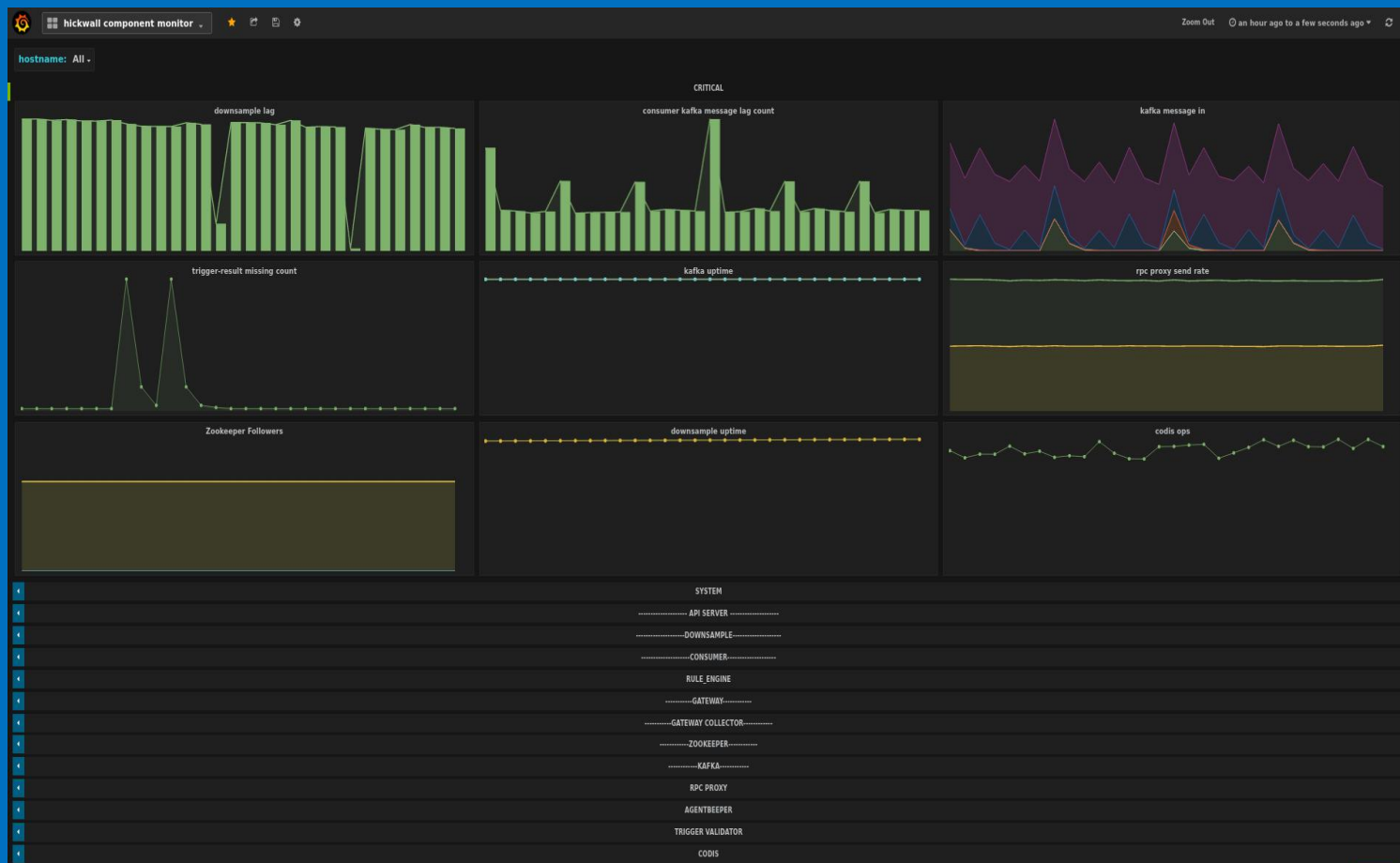
Apache
MESOS



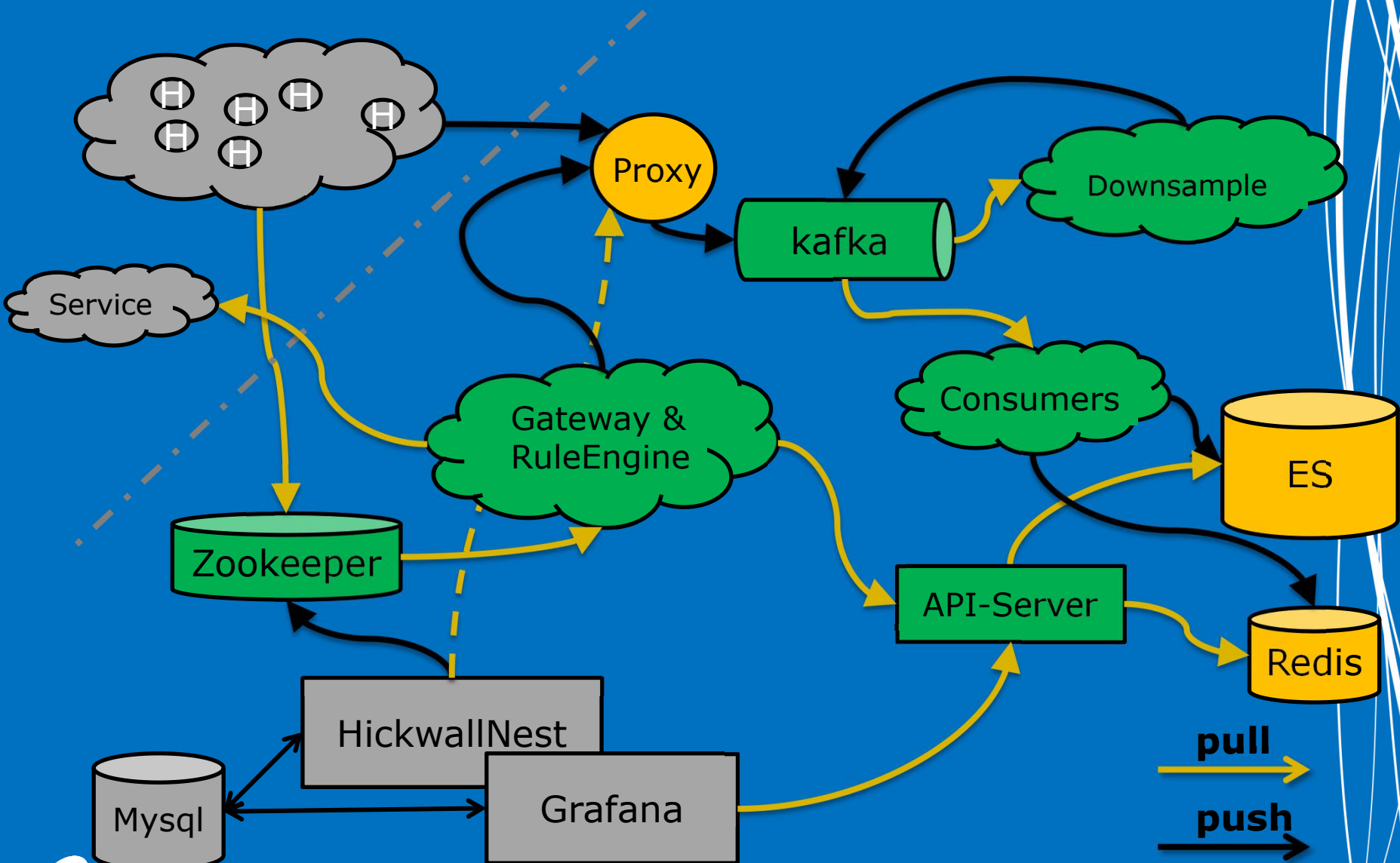
可靠性 与 吞吐量

服务端

组件第三方监控

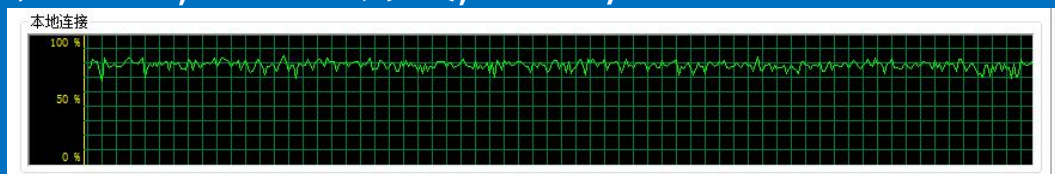


可靠性与吞吐量

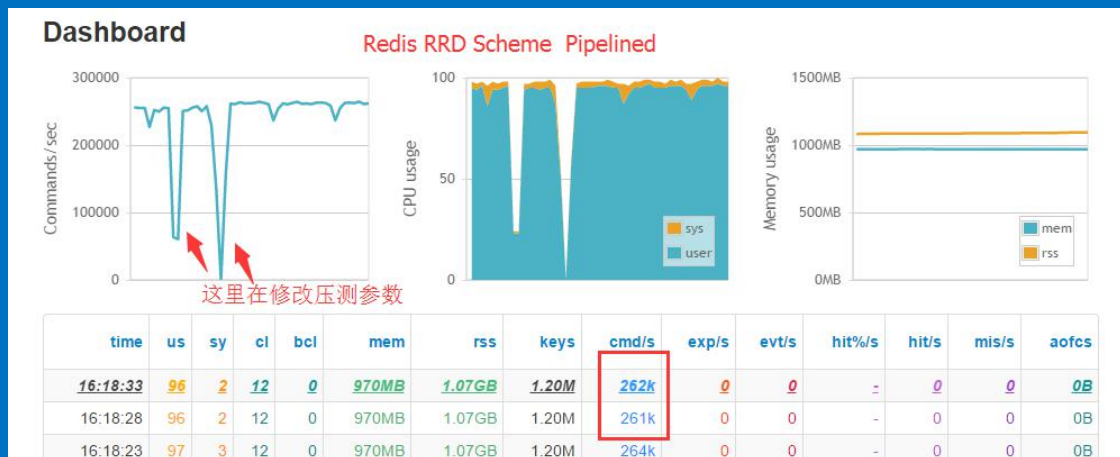


可靠性 与 吞吐量

- proxy 不解包：
 - 100M网卡跑到 75%, batch 方式, 17.5w / s
 - 10G 网卡跑到 4.5%, batch 方式, 90w / s



- redis rrd scheme pipelined
 - 1.2 Million keys, 每次1hour 数据, 压测26w cmd/s



现状 和 展望

- 生产全面部署
- 不断改进，未来开源

“

Q & A

”



谢谢！



微信公众号
携程技术保障中心