

海量日志实时分析

日志搜索引擎

DTCC

2015中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2015

大数据技术探索和价值发现



提纲

- 日志的应用场景
- 过去的做法
- 现在的做法
- 日志搜索引擎
- 日志易产品架构



一条 Apache Access 日志

- 180.150.189.243 - - [15/Apr/2015:00:27:19 +0800] "POST /report HTTP/1.1" 200 21 "https://rizhiyi.com/search/" "Mozilla/5.0 (Windows NT 6.1; WOW64; rv:37.0) Gecko/20100101 Firefox/37.0" "10.10.33.174" 0.005 0.001
- 字段：
 - Client IP: 180.150.189.243
 - Timestamp: 15/Apr/2015:00:27:19 +0800
 - Method: POST
 - URI: /report
 - Version: HTTP/1.1
 - Status: 200
 - Bytes: 21
 - Referrer: <https://rizhiyi.com/search/>
 - User Agent: Mozilla/5.0 (Windows NT 6.1; WOW64; rv:37.0) Gecko/20100101 Firefox/37.0
 - X-Forward: 10.10.33.174
 - Request_time: 0.005
 - Upstream_request_time:0.001



日志：时间序列机器数据

- IT 系统信息
 - 操作系统
 - 应用软件
- 用户信息
 - 用户行为
- 各种传感器信息
- 日志反映的是事实数据
 - 深度解析LinkedIn大数据平台（<http://www.csdn.net/article/2014-07-23/2820811/1>）



应用场景

- 运维可用性监控
- 应用性能监控
 - Application Performance Monitoring (APM)
- 安全审计
 - Security Information Event Management (SIEM)
- 用户数据统计分析
- 物联网
 - 智能家电
 - Nest Lab 采集的智能恒温器数据使用 Splunk 来分析
 - 车联网



过去的做法

- 日志没有集中处理
 - 登陆每一台服务器，使用脚本命令或程序查看
- 日志被删除
 - 磁盘满了删日志
 - 黑客删除日志，抹除入侵痕迹
- 日志只做事后追查
 - 没有实时监控、分析
- 使用数据库存储日志
 - 无法适应TB级海量日志
 - 数据库的schema无法适应千变万化的日志格式
 - 无法提供全文检索
- Complex Event Processing (CEP)
 - 难以处理大数据量



现在的做法

- Hadoop
 - 批处理，不够及时
 - 查询慢
 - 可作基于日志的用户数据离线挖掘，无法做 OLAP (On Line Analytic Processing)
- Storm
 - 历史久，停止开发
 - 任务调度差
- Spark
 - 生态圈完整
 - DataBricks 专门支持
- Storm vs. Spark Streaming
 - Storm 是真正的流式处理，Spark Streaming 是 mini-batch
 - Exactly Once vs. At Least Once
 - 延时与吞吐率的取舍
- Hadoop/Storm/Spark都只是一个开发框架，不是拿来即用的产品



对日志准实时搜索、分析

- 日志准实时搜索引擎
- Splunk
- ELK (Elasticsearch/Logstash/Kibana)

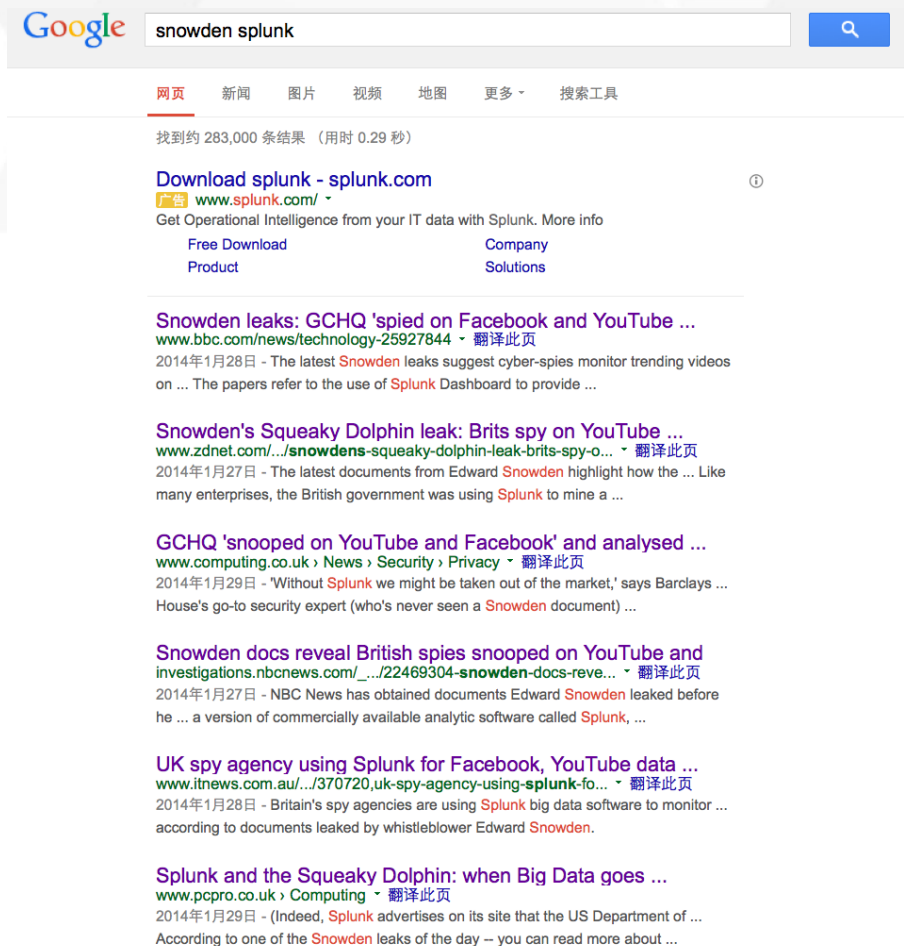


Splunk

- 首创用准实时搜索的方法来分析日志
- 功能非常丰富
 - Search Processing Language
 - 类似 Linux 命令，支持管道，子查询等功能
- Gartner report (2014/7/18)
 - Splunk在日志检索时抽取日志的关键字段，检索速度慢
 - Splunk按每天处理的日志量收费，价格较贵



Splunk 与 棱镜门



Google search results for "snowden splunk".

找到约 283,000 条结果 (用时 0.29 秒)

Download splunk - splunk.com
www.splunk.com/
Get Operational Intelligence from your IT data with Splunk. More info
Free Download Company
Product Solutions

Snowden leaks: GCHQ 'spied on Facebook and YouTube ...
www.bbc.com/news/technology-25927844 · 翻译此页
2014年1月28日 - The latest **Snowden** leaks suggest cyber-spies monitor trending videos on ... The papers refer to the use of **Splunk** Dashboard to provide ...

Snowden's Squeaky Dolphin leak: Brits spy on YouTube ...
www.zdnet.com/.../snowdens-squeaky-dolphin-leak-brits-spy-o-... · 翻译此页
2014年1月27日 - The latest documents from Edward **Snowden** highlight how the ... Like many enterprises, the British government was using **Splunk** to mine a ...

GCHQ 'snooped on YouTube and Facebook' and analysed ...
www.computing.co.uk/News/Security/Privacy · 翻译此页
2014年1月29日 - 'Without **Splunk** we might be taken out of the market,' says Barclays ... House's go-to security expert (who's never seen a **Snowden** document) ...

Snowden docs reveal British spies snooped on YouTube and investigations.nbcnews.com/.../22469304-snowden-docs-reve... · 翻译此页
2014年1月27日 - NBC News has obtained documents Edward **Snowden** leaked before he ... a version of commercially available analytic software called **Splunk**, ...

UK spy agency using Splunk for Facebook, YouTube data ...
www.itnews.com.au/.../370720,uk-spy-agency-using-splunk-fo-... · 翻译此页
2014年1月28日 - Britain's spy agencies are using **Splunk** big data software to monitor ... according to documents leaked by whistleblower Edward **Snowden**.

Splunk and the Squeaky Dolphin: when Big Data goes ...
www.pcpro.co.uk/Computing · 翻译此页
2014年1月29日 - (Indeed, **Splunk** advertises on its site that the US Department of ... According to one of the **Snowden** leaks of the day – you can read more about ...



ELK

- 三个独立的开源套件
 - 一些著名互联网公司对 ELK 做二次开发，使用过百台的 ELK 集群分析日志
- 存在问题
 - 运维管理不方便，三个独立的系统，没有统一的部署、管理工具
 - 没有告警功能
 - 没有用户认证及权限管理
 - 统计、分析功能有限

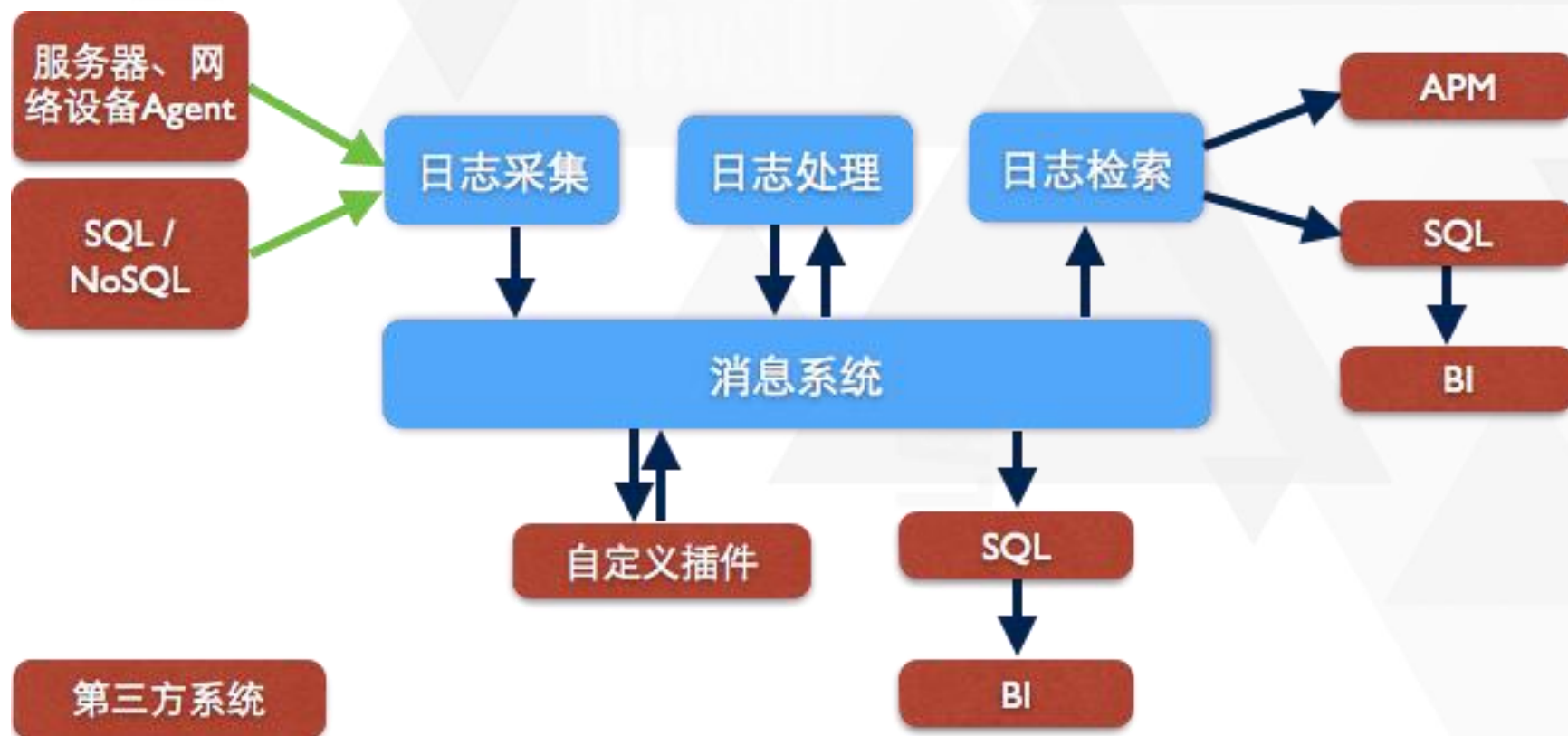


日志易

- 日志搜索分析平台
 - rizhiyi.com
- 企业部署版
- SaaS 版
 - 每天500MB日志处理免费



日志易架构



日志易功能

- 搜索
- 告警
- 统计
 - 事务关联
- 配置解析规则，识别任何日志
- 安全攻击自动识别
- 开放API，对接第三方系统
- 高性能、可扩展分布式架构
 - 10万 EPS (Event Per Second), 每天TB级日志

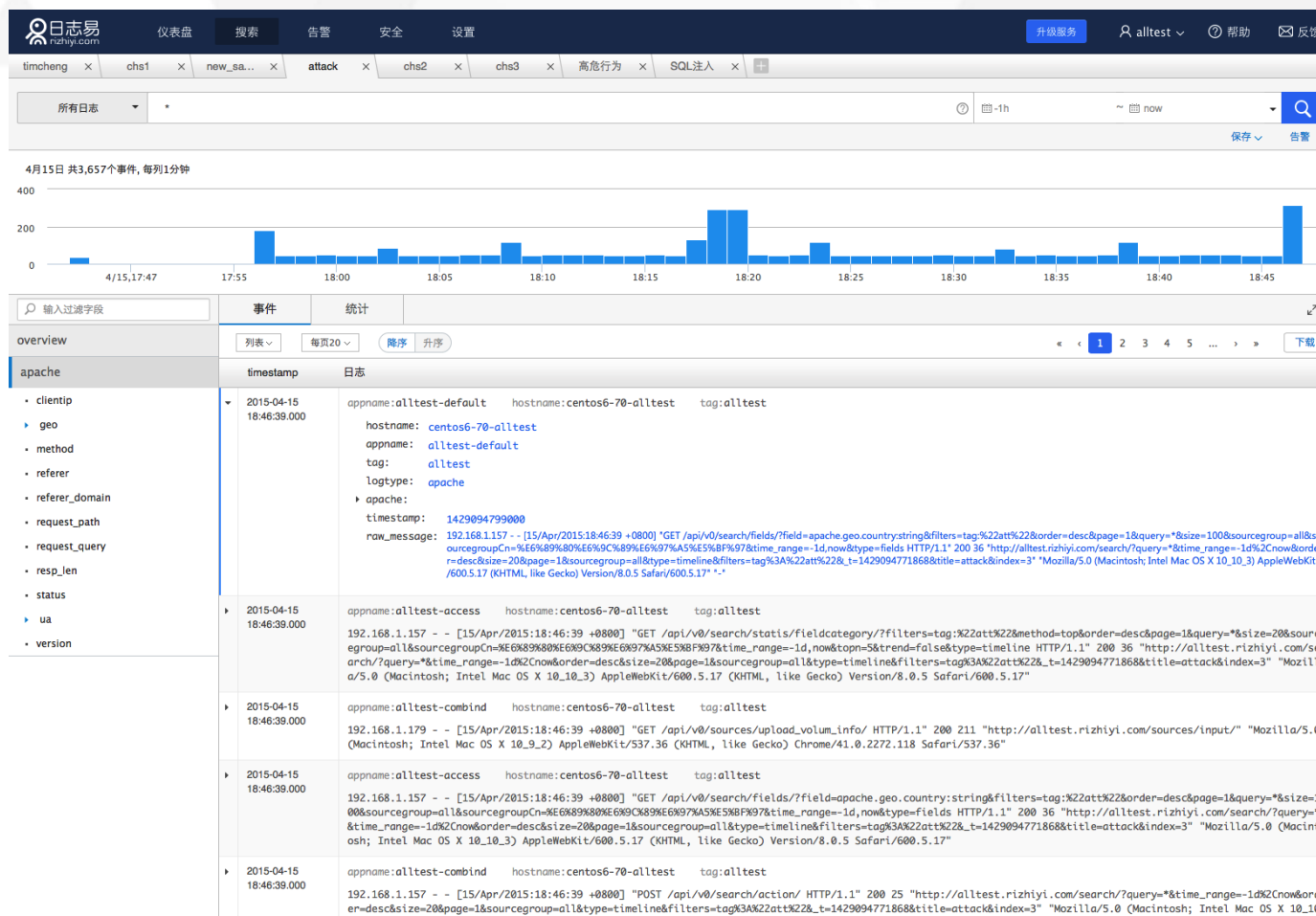


日志易 vs. Splunk

- Splunk
 - 日志进入系统时不抽取关键字段，直接做索引，在检索时抽取关键字段
 - 灵活、索引文件小，但检索延时大
- 日志易
 - 日志进入系统时就抽取关键字段，然后做索引
 - 索引文件大，但检索延时小
 - 可在索引前配置解析规则，抽取日志里的任何关键字段



日志易介绍：总览

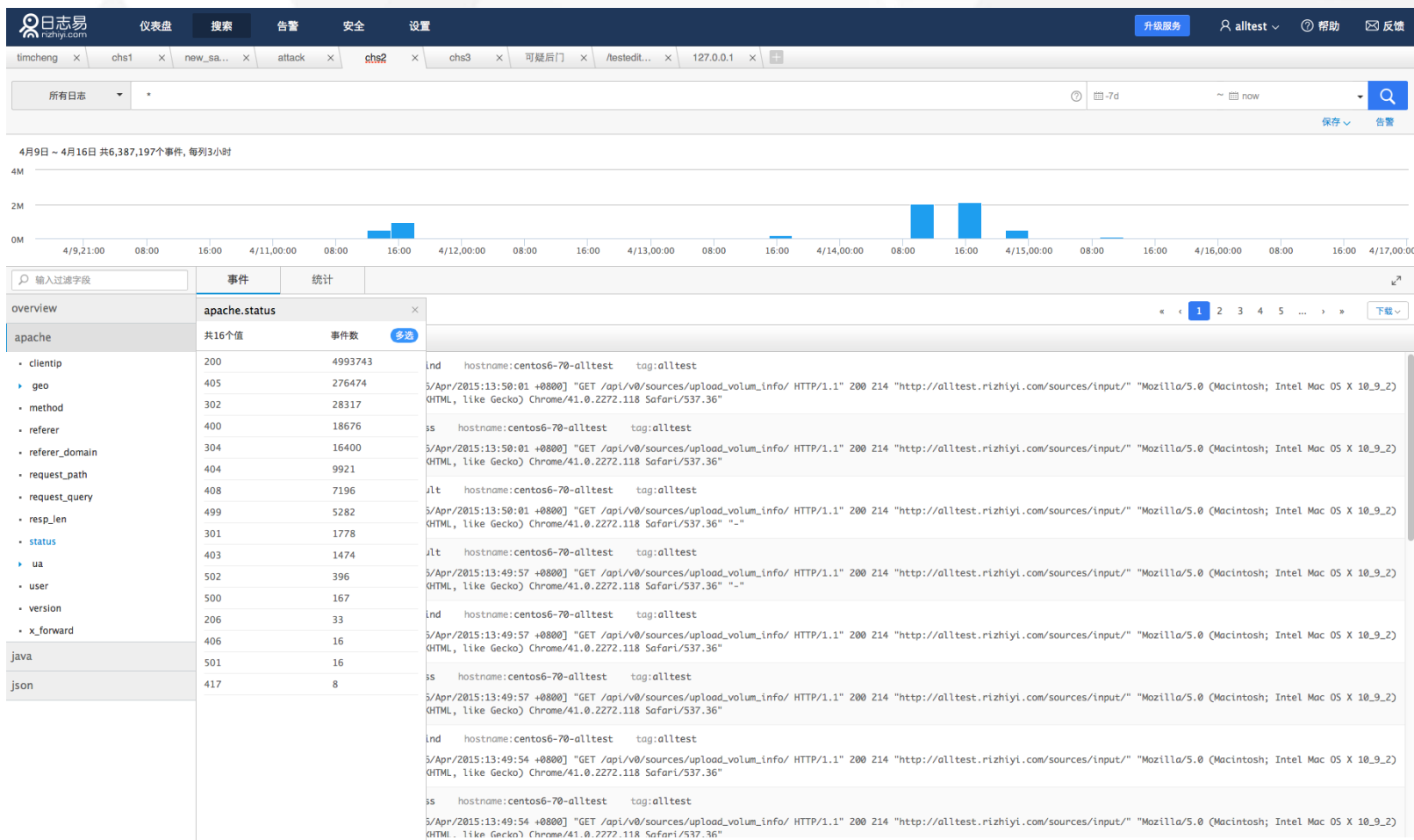


日志易介绍：日志结构化

事件	统计
列表 ▾	每页20 ▾
降序	升序
timestamp	日志
2015-04-16 13:50:01.000	<div>appname:alltest-combind hostname:centos6-70-alltest tag:alltest</div> <div>hostname: centos6-70-alltest</div> <div>appname: alltest-combind</div> <div>tag: alltest</div> <div>logtype: apache</div> <div>▼ apache:</div> <div>clientip: 192.168.1.179</div> <div>▶ geo:</div> <div>method: GET</div> <div>referer: http://alltest.rizhiyi.com/sources/input/</div> <div>referer_domain: alltest.rizhiyi.com</div> <div>request_path: /api/v0/sources/upload_volum_info/</div> <div>resp_len: 214</div> <div>status: 200</div> <div>▼ ua:</div> <div>browser: Chrome</div> <div>browser_v: Chrome 41.0.2272</div> <div>device: Other</div> <div>os: Mac OS X</div> <div>os_v: Mac OS X 10.9.2</div> <div>version: 1.1</div> <div>timestamp: 1429163401000</div> <div>raw_message: 192.168.1.179 - - [16/Apr/2015:13:50:01 +0800] "GET /api/v0/sources/upload_volum_info/ HTTP/1.1" 200 214 "http://alltest.rizhiyi.com/sources/input/" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_2) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/41.0.2272.118 Safari/537.36"</div>



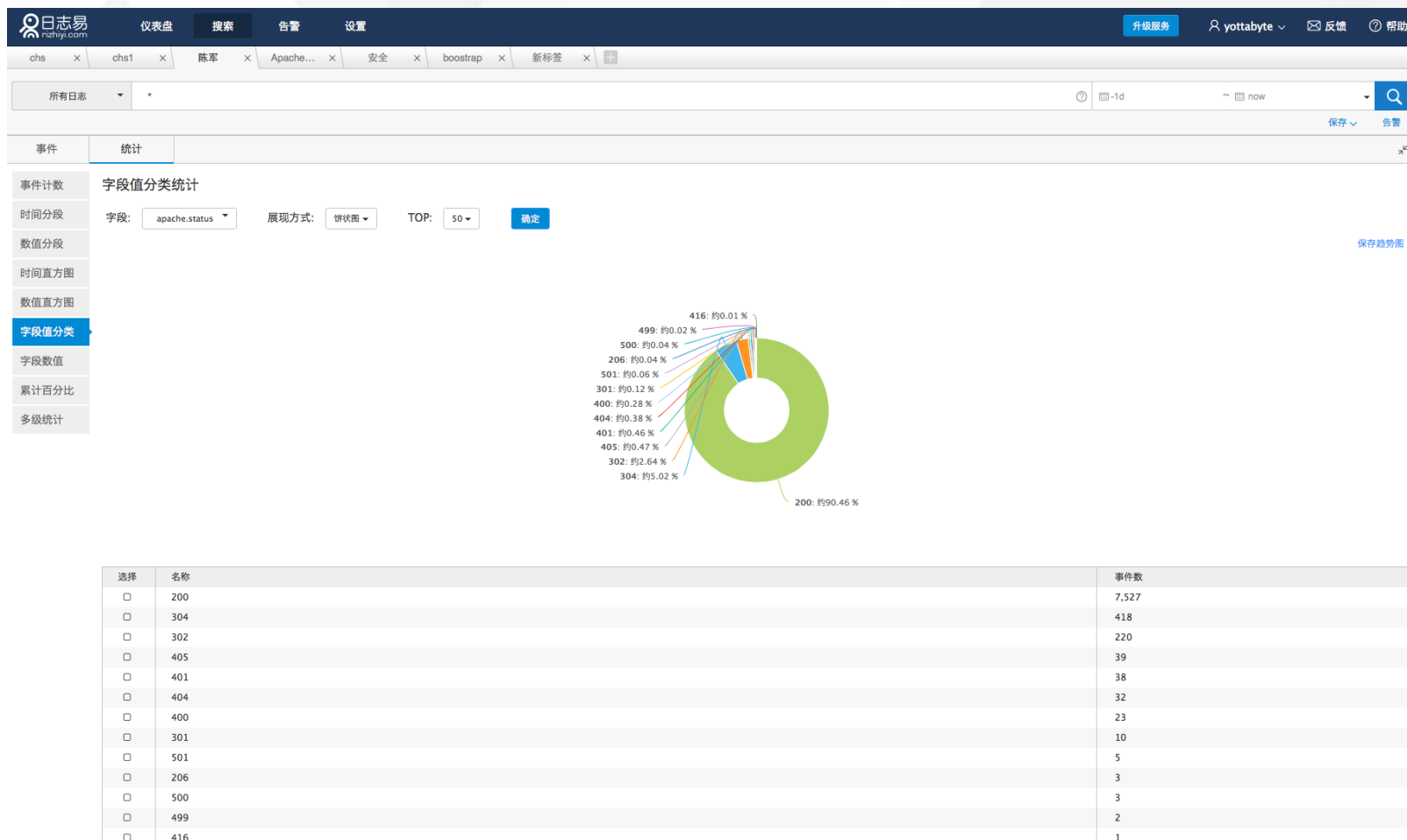
日志易介绍：字段抽取、统计



日志易介绍：搜索



日志易介绍：统计



日志易介绍：告警

 日志易
rzhily.com

仪表盘 搜索 告警 设置

升级服务

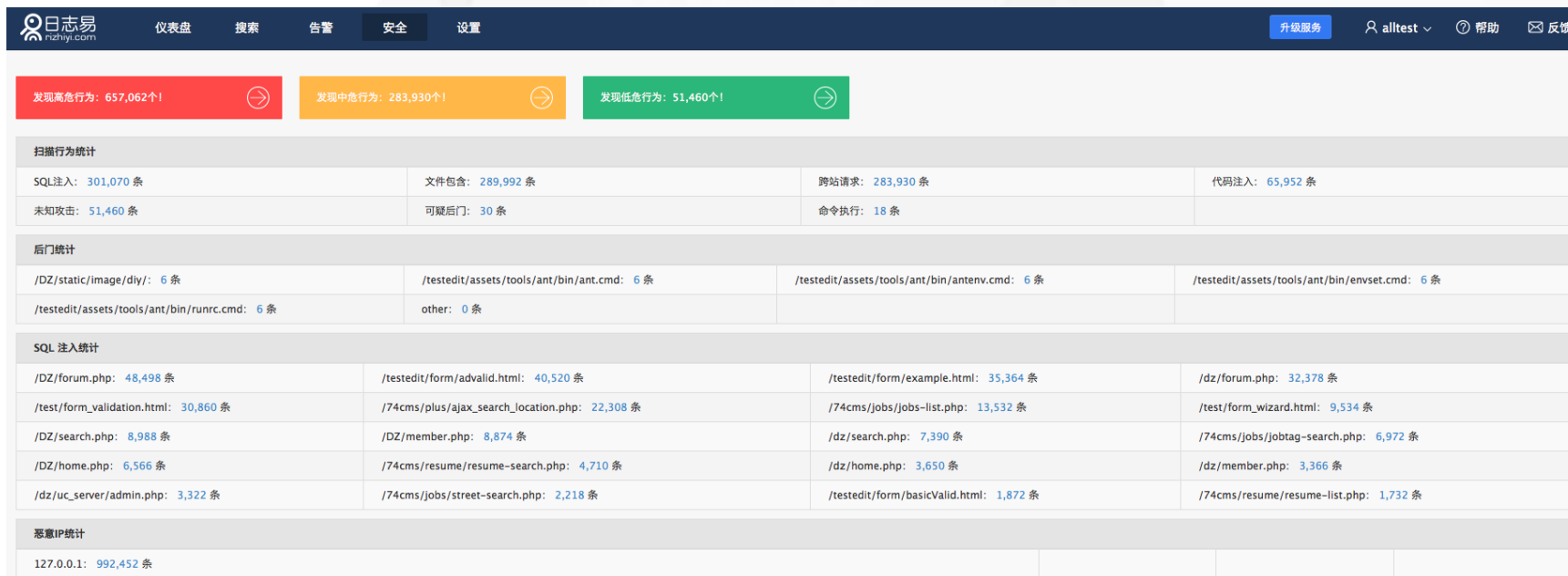
yottabyte

反馈

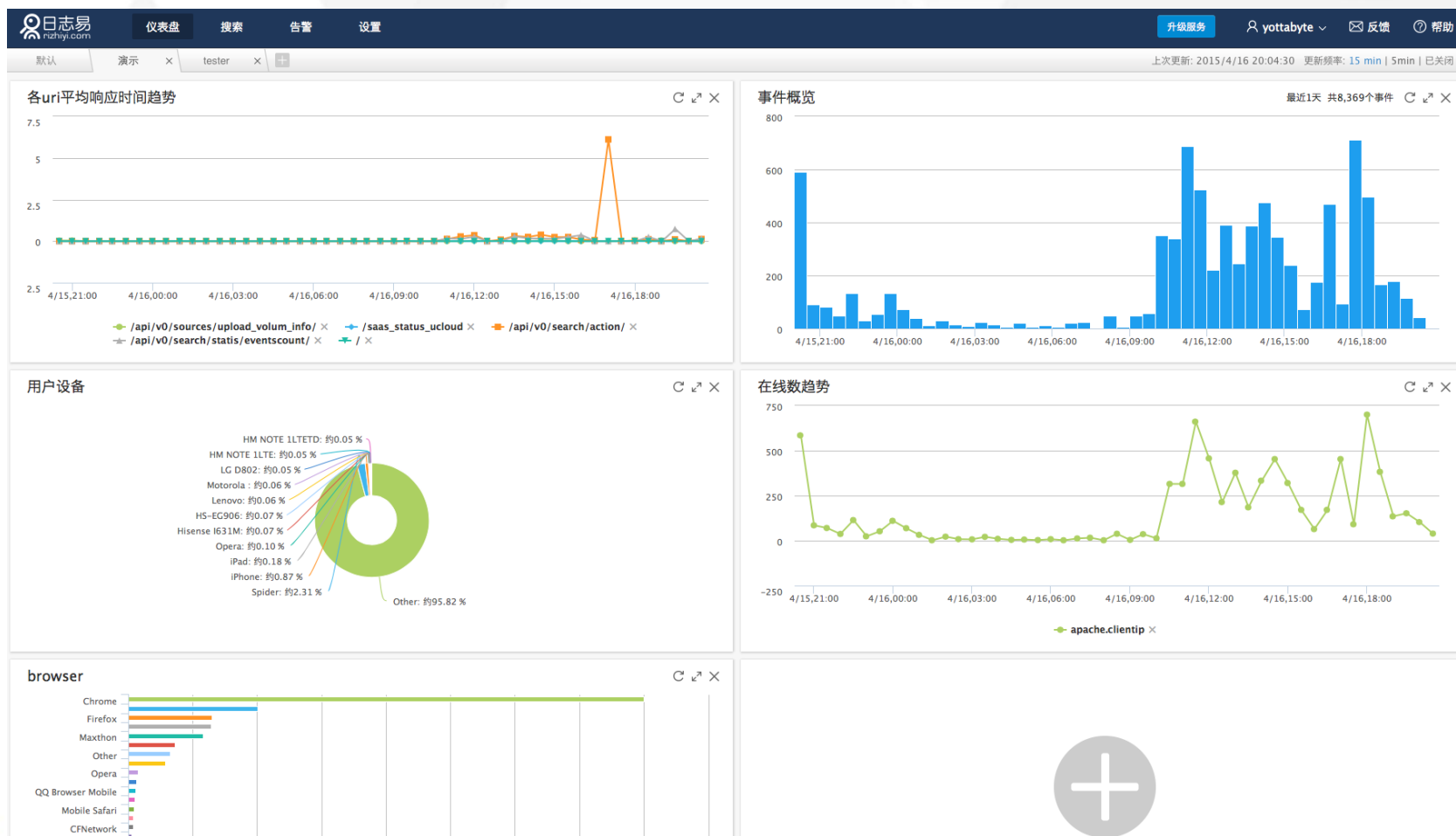
帮助

用户告警(7) [+ 新建](#)输入关键字 

日志易介绍：安全审计



日志易介绍：仪表盘



日志易，日志分析更容易

rizhiyi.com

