

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



Mobike大数据基础平台建设

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



曹永鹏

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



目录

- Mobike 介绍
- 团队介绍
- 平台演变
- 平台架构
- 平台建设
- 未来展望

Mobike 介绍

- 全球第一大互联网出行服务平台
- 10+个国家200个城市
- 日订单量3000万+
- 700万+ 摩拜单车
- 2亿+ 用户

团队介绍

- 职责

- 提供一站式大数据服务平台
- 提供大数据存储计算平台
- 提供实时搜索服务平台
- 标准化数据计算逻辑和管理

- 用户

数据团队、研发工程师

平台历程

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

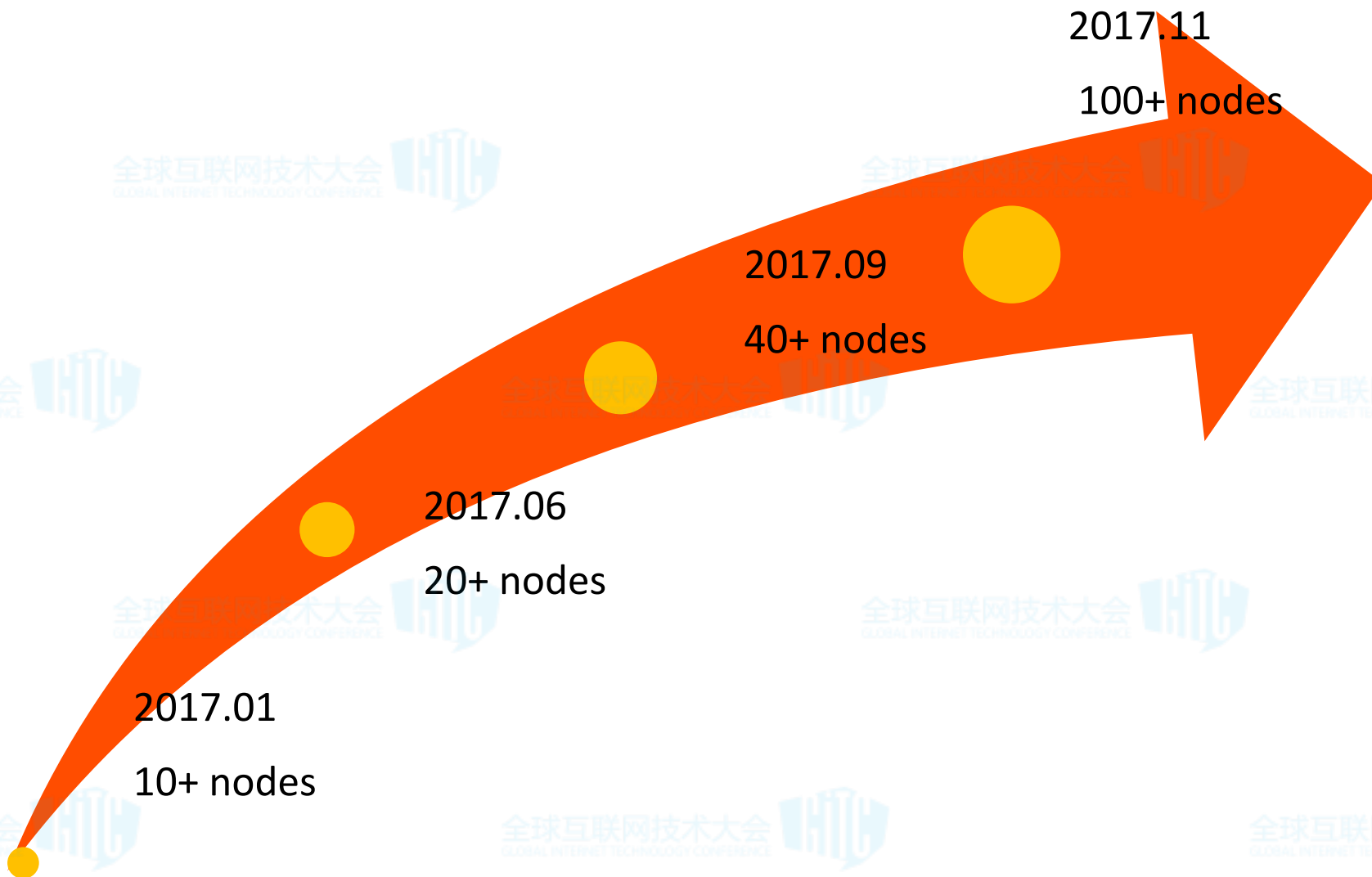
全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



业务

- 日志处理
- 数据仓库
- 报表
- 红包车
- 反作弊
- 实时计算
- 车辆调度
-

数据量

- 规模

- 2 Hadoop cluster
- 100 + nodes

- 存储

- ~1PB
- ~日增6TB

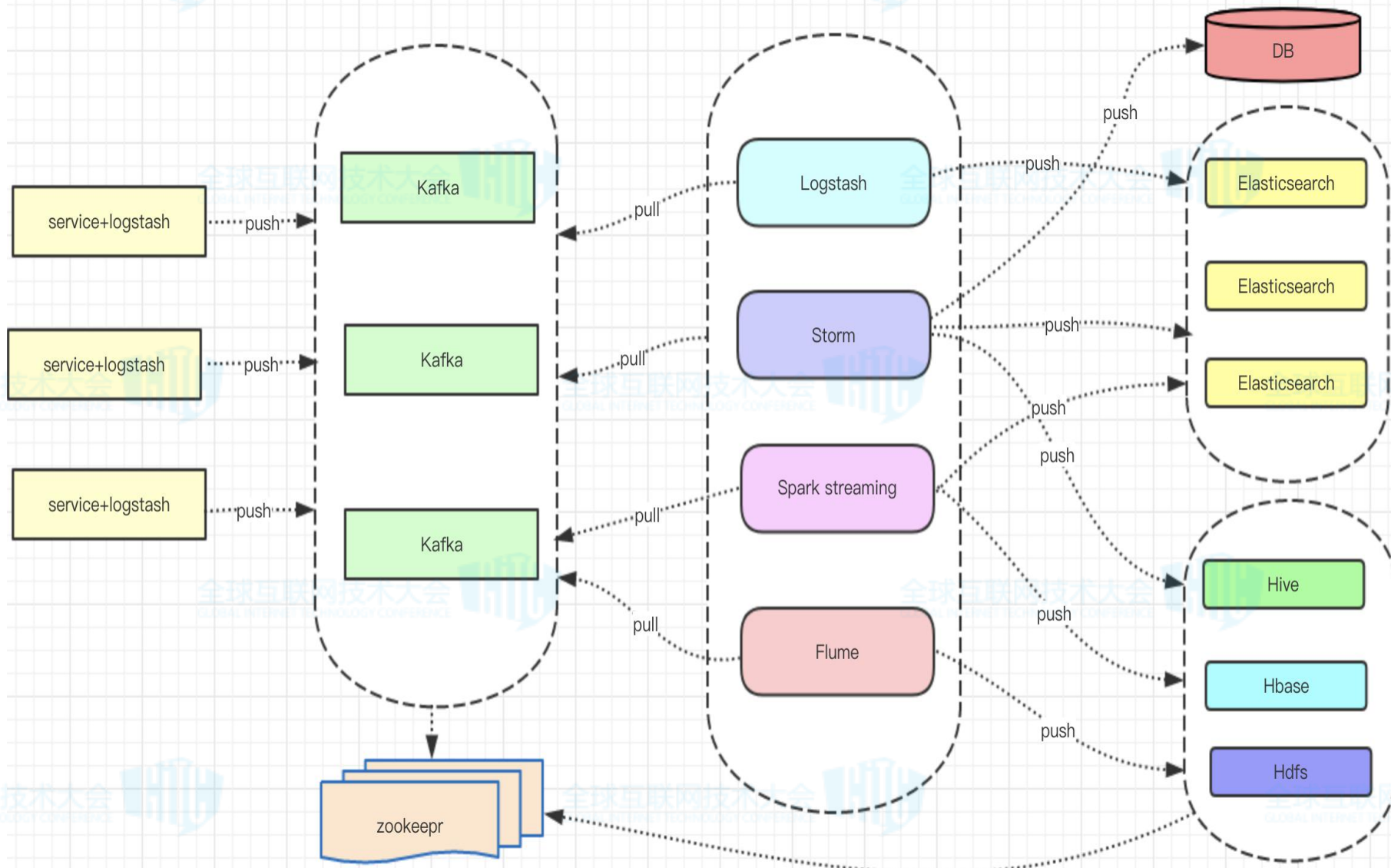
- 计算

- 800+ jobs/day
- 400000+ Tasks/day

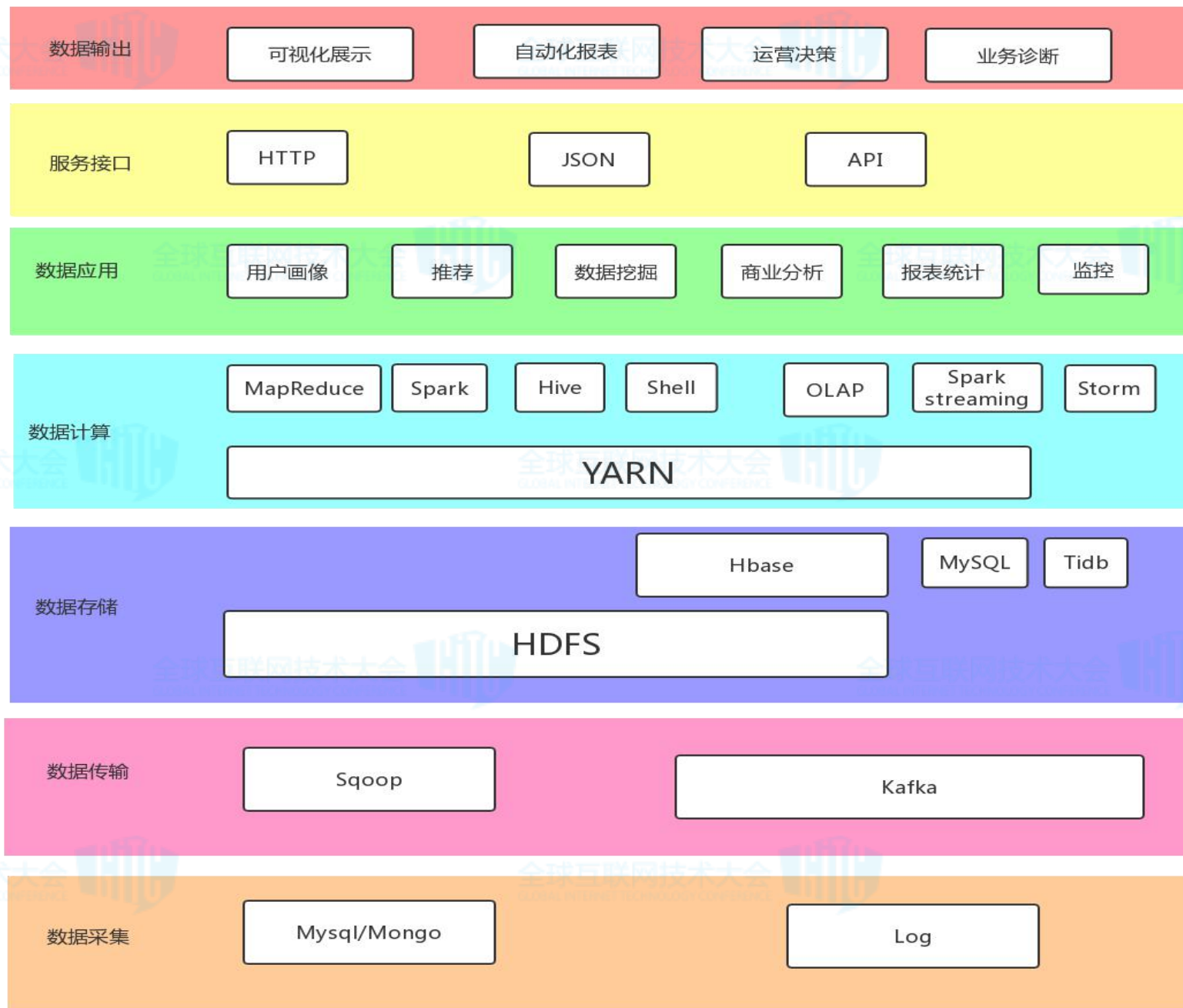
- 工具

- 以Hadoop为核心，Kafka、HBase、Hive、Spark、Storm、ELK

平台架构



平台架构



集群portal

数据管理

系统监控

用户管理

资源调度

集群状态

机器管理

工具管理

权限管理

流程管理

配额管理

平台建设

- 日志收集

- Logstash + Kafka + Flume-ng

- 离线处理

- HDFS HA + RM HA / All on Yarn

- Hive 数仓

- Spark Mllib 模型训练

- 实时处理

- Storm

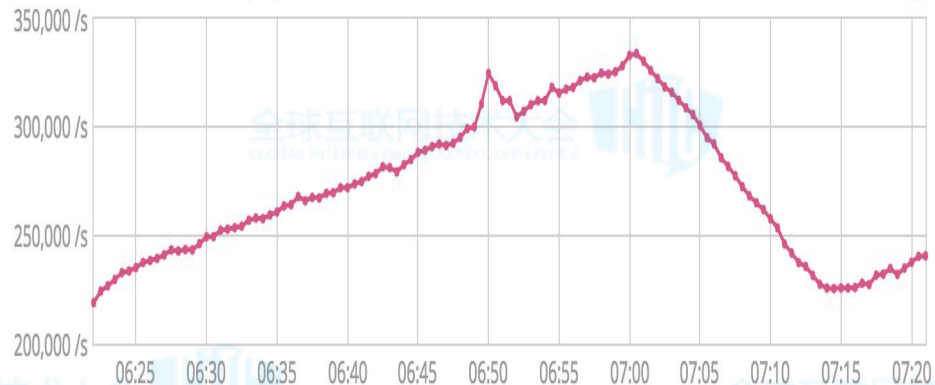
- Spark streaming

- Es实时搜索服务

- 全链路实时监控

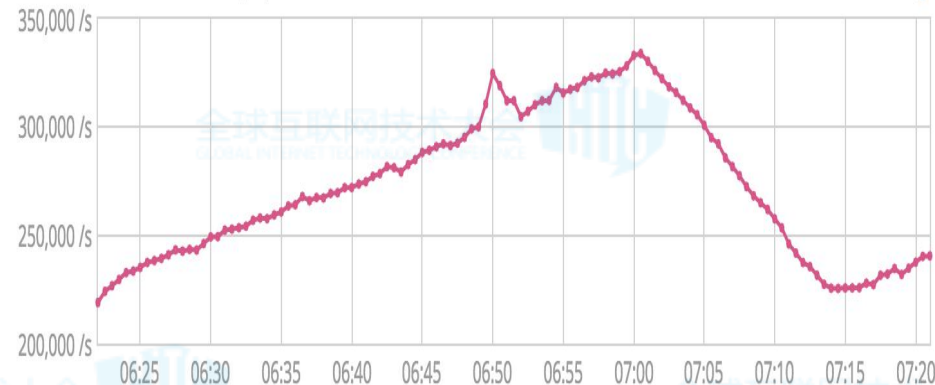
全链路监控

Events Received Rate (/s)



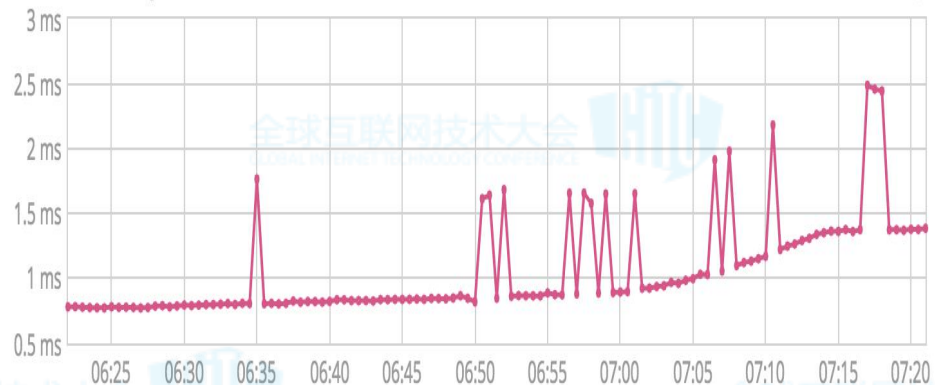
● Events Received Rate 240,522.87 /s

Events Emitted Rate (/s)



● Events Emitted Rate 240,505.1 /s

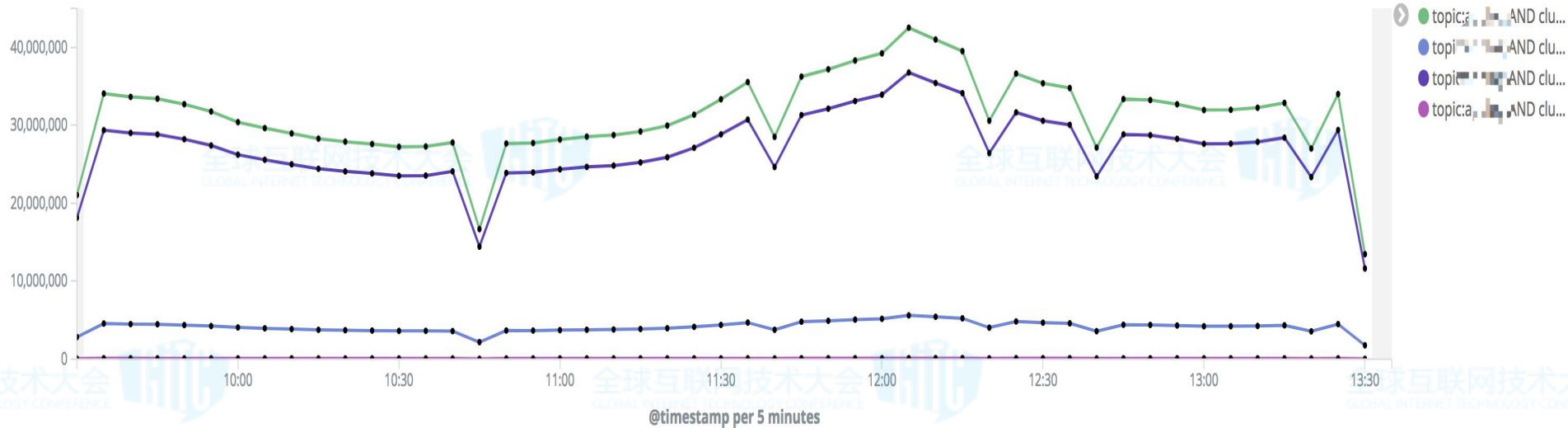
Event Latency (ms)



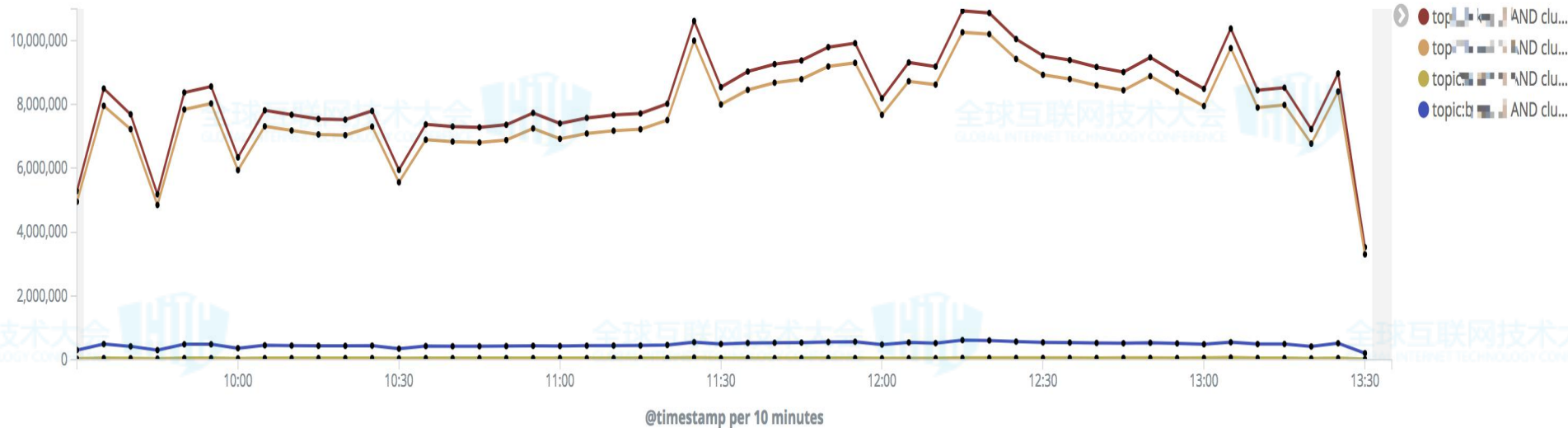
● Event Latency 1.39 ms

全链路监控

kafka sum



kafka num



全链路监控

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

Partition	LogSize	Offset	Lag	Owner	Created	Modify
88	3418519	3418519	0	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:21	2017-11-22 07:03:21
89	3418751	3418748	3	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:19	2017-11-22 07:03:19
110	3418496	3418495	1	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:21	2017-11-22 07:03:21
111	3414934	3414932	2	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:21	2017-11-22 07:03:21
112	3418558	3418557	1	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:20	2017-11-22 07:03:20
113	3416078	3416076	2	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:20	2017-11-22 07:03:20
114	3415324	3415323	1	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:21	2017-11-22 07:03:21
115	3418705	3418704	1	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:22	2017-11-22 07:03:22
116	3418757	3418757	0	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:22	2017-11-22 07:03:22
117	3418755	3418755	0	10.1.110.189-consumer-6-6be88375-db7b-48c8-8fd2-620a99260531	2017-11-22 07:03:21	2017-11-22 07:03:21

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

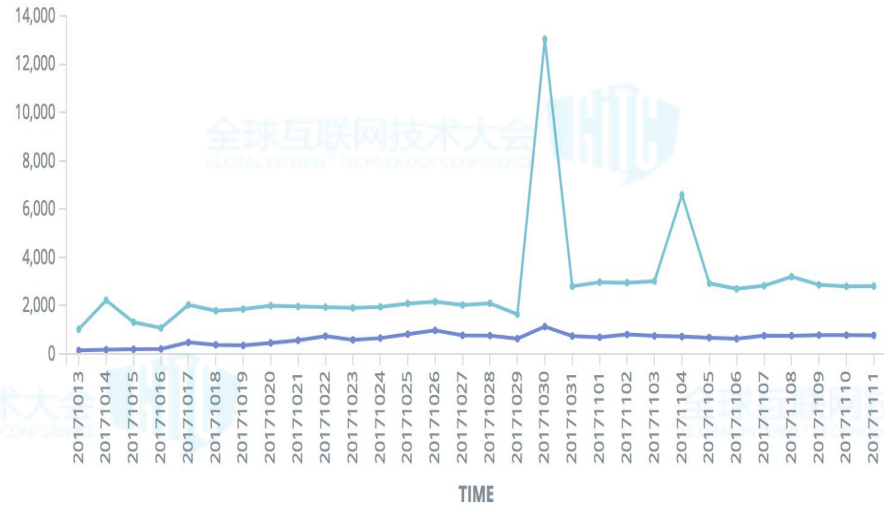
全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE

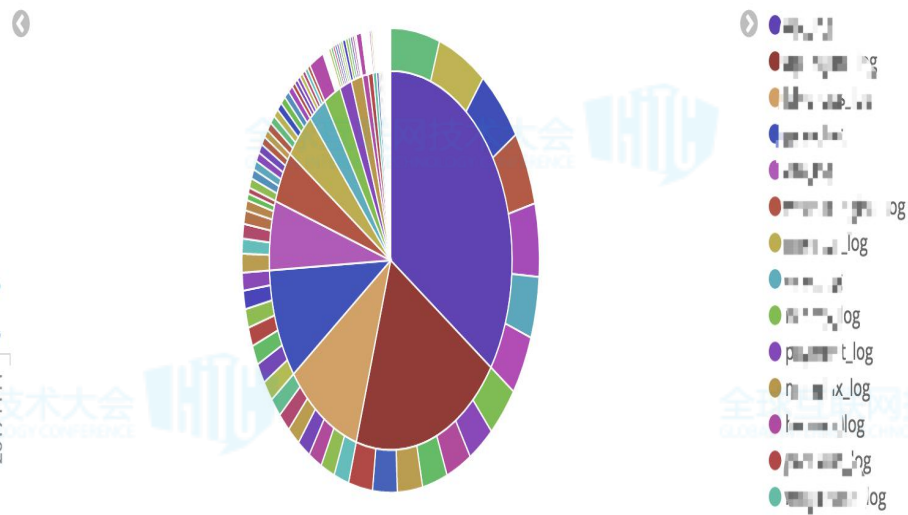
全链路监控



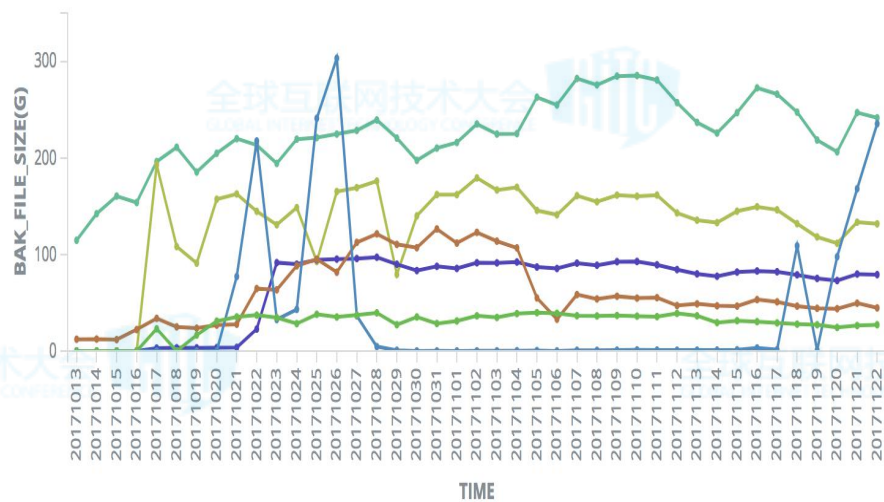
BAK_FILE_TOTAL_NUM_SIZE



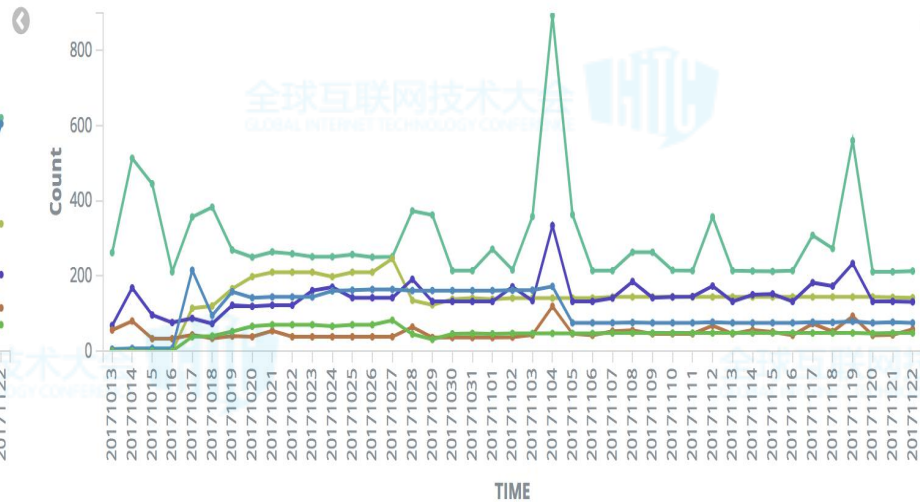
BAK_FILE_TYPE_TOTAL_SIZE



BAK_FIKE_TYPE_SIZE



BAK_FIKE_TYPE_NUM



平台建设

- Yum 源
- Puppet
- Ansible
- Zabbix
- Ganglia

平台建设

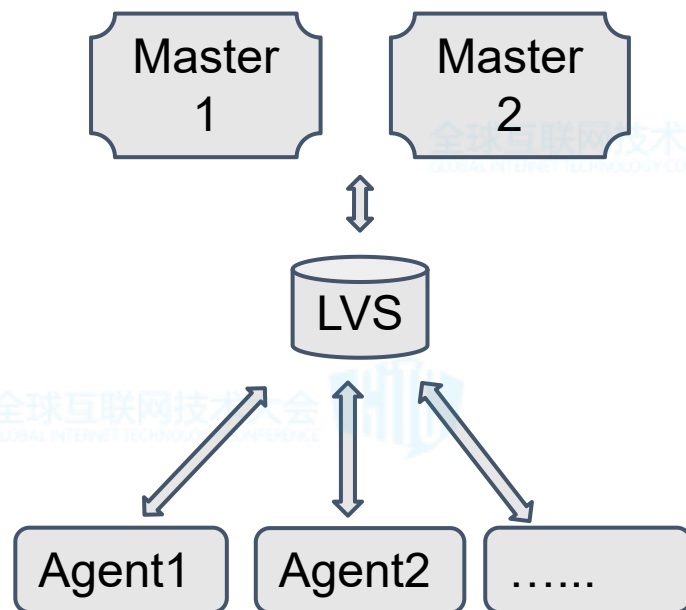
- Puppet

配置统一管理

- Ansible

自动化部署

Puppet (Hadoop)



平台建设

- HDFS

- 约束目录规范，严禁私自创建目录
- Quota限额，文件数、存储

- YARN

- 按作业类型划分固定的Queue（online/offline/bi/dw/default）

- Jobname 规范格式 zhangsan—xxx

- 集群账号

- 按部门、按业务分配用户
- 个人帐号仅供测试，组帐号上生产调度，Job提交需指定queue

- 数据权限

- HDFS Acl

- Job管理

- workflow定义方式，通过配置文件，GitLab CI 自动提交

未来展望

- 平台安全
- 自动化建设
- 数据质量

全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



平台安全

- 身份认证

- Kerberos

- 身份管理

- LDAP

- 授权访问

- 数据授权访问:SENTRY

- 安全审计

- 数据加密

自动化建设

- 更自动化
- 更智能化

全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



全球互联网技术大会



数据质量

- 数据生命周期管理
- 数据间的血缘关系
- 数据格式标准化



Thanks

加入我们一起做有意义的事
caoyongpeng@mobike.com