

百度运维大数据存储平台 设计与实践

李玉明 (liyuming01@baidu.com)

百度智能运维数据平台负责人

QCon

全球软件开发大会

10月17-19日 上海·宝华万豪酒店



扫码锁定席位

九折即将结束

团购还享更多优惠，折扣有效期至9月17日

扫描右方二维码即可查看大会信息及购票



如果在使用过程中遇到任何问题，可联系大会主办方，欢迎咨询！

微信：qcon-0410

电话：010-84782011

ArchSummit

全球架构师峰会 2017



扫码锁定席位

12月8-9日 北京·国际会议中心

七折即将截止立省2040元

使用限时优惠码AS200，

以目前最优惠价格报名ArchSummit

仅限前20名用户，优惠码有效期至9月19日，

扫描右方二维码即可使用



如果在使用过程中遇到任何问题，可联系大会主办方，欢迎咨询！

微信：aschina666

电话：15201647919

极客搜索

全站干货，一键触达，只为技术

s.geekbang.org



扫描二维码立即体验

有没有一种搜索方式，能整合 InfoQ 中文站、极客邦科技旗下12大微信公众号矩阵的全部资源？

极客搜索，这款针对极客邦科技全站内容资源的轻量级搜索引擎，做到了！

扫描上方二维码，极客搜索！

这里只有 技术领导者

EGO会员第二季招募季正式开启



E小欧

报名时间：9月1日-9月15日

扫描添加E小欧，
邀您进入EGO会员预报名群

立即报名



TABLE OF CONTENTS

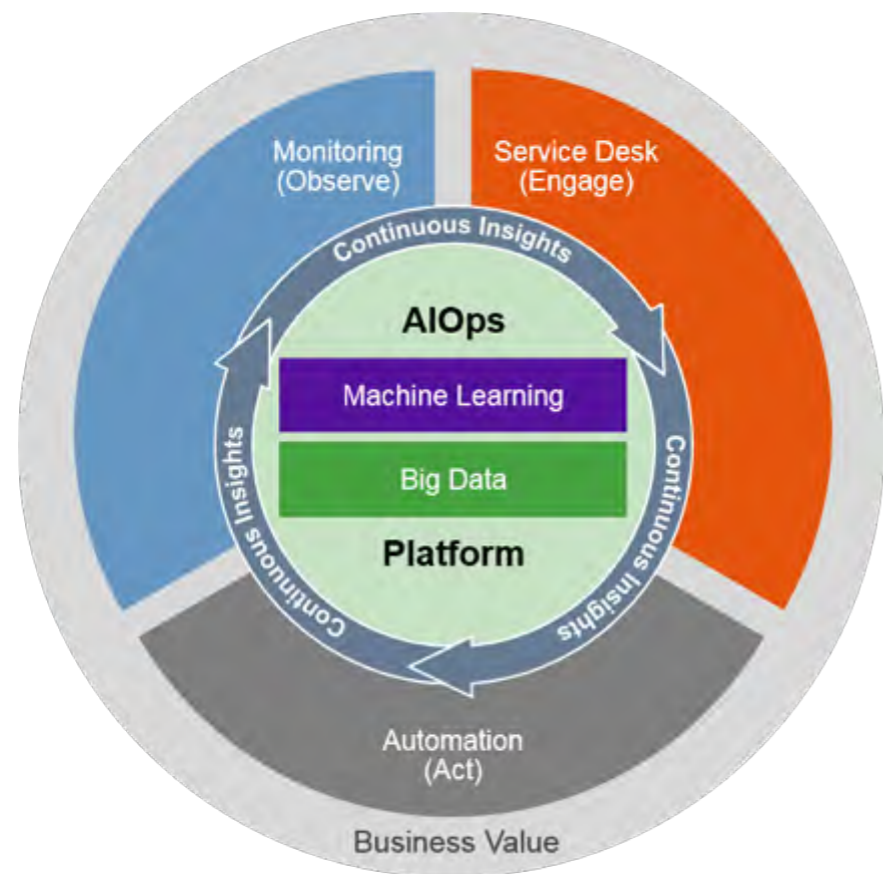
智能运维与大数据

大规模时序数据存储

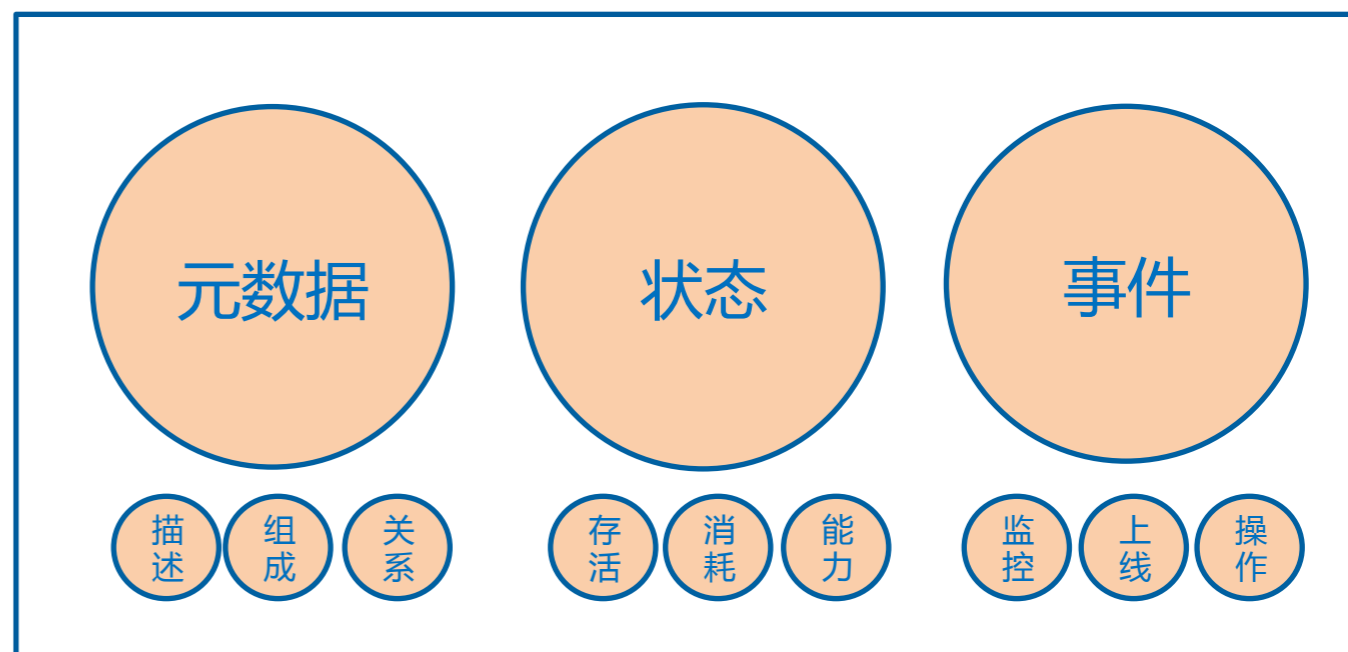
海量运维事件数据存储

运维知识库

智能运维与大数据



- 运维大数据
 - 元数据：服务/实例/机房/网络等及其关联
 - 状态数据：资源和业务等监控指标（时序）
 - 事件数据：异常事件/上线事件等



Source : Gartner Report
 IT Operations Analytics Must Be Placed Within
 an AIOps Context.
 Will Cappelli (Research VP) | 26 August 2016

TABLE OF CONTENTS

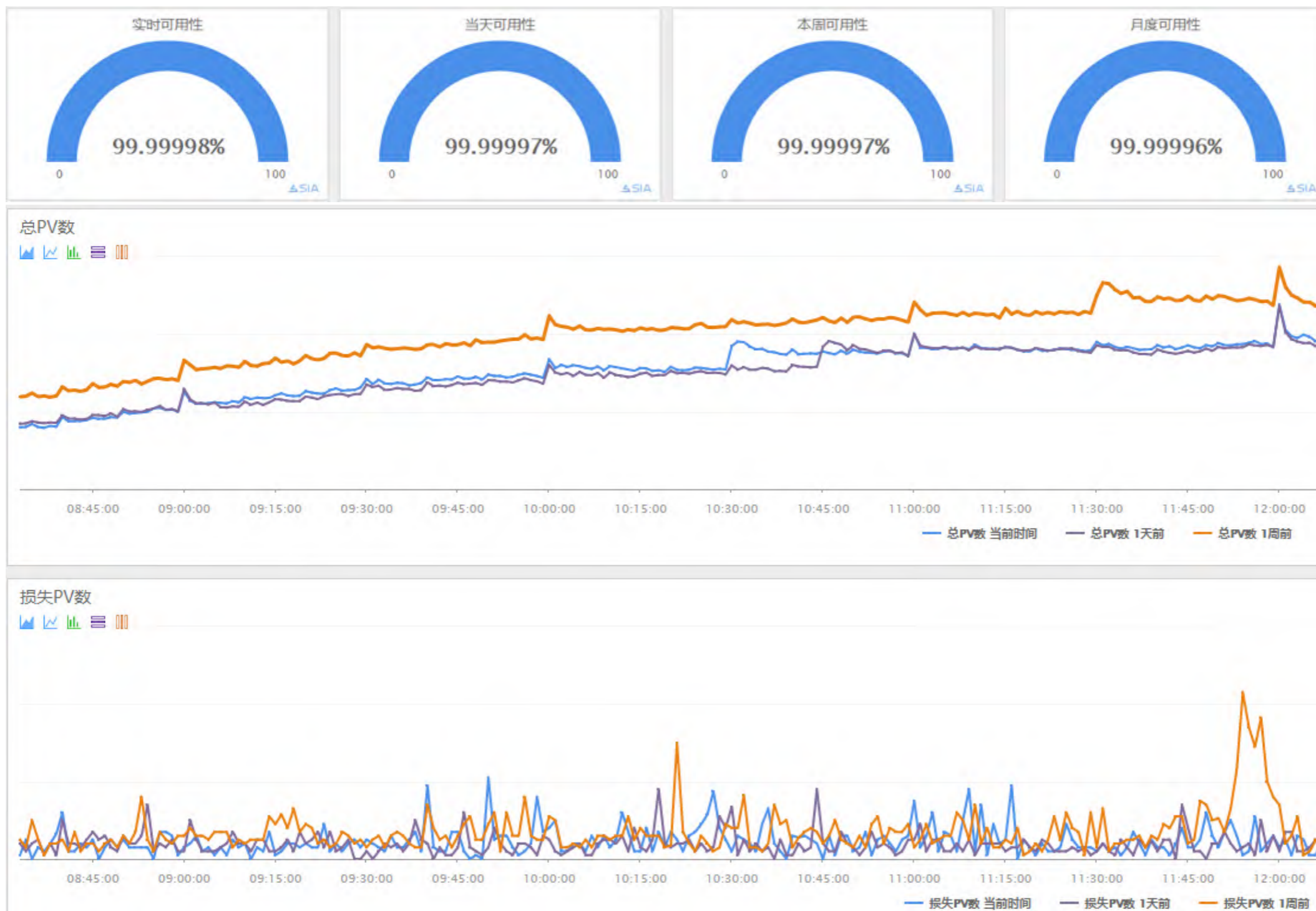
智能运维与大数据

大规模时序数据存储

海量运维事件数据存储

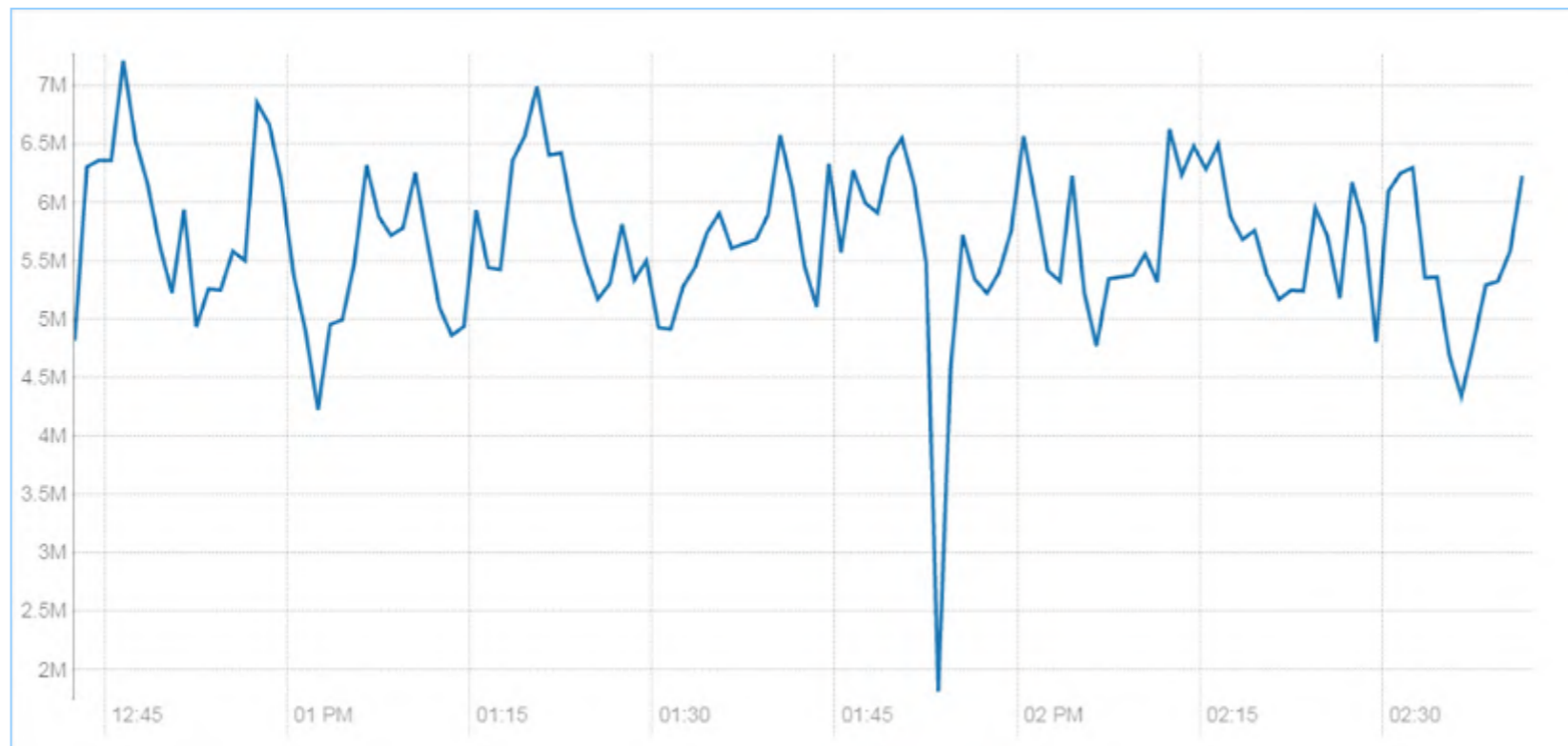
运维知识库

服务状态的监控 (from 赛亚)



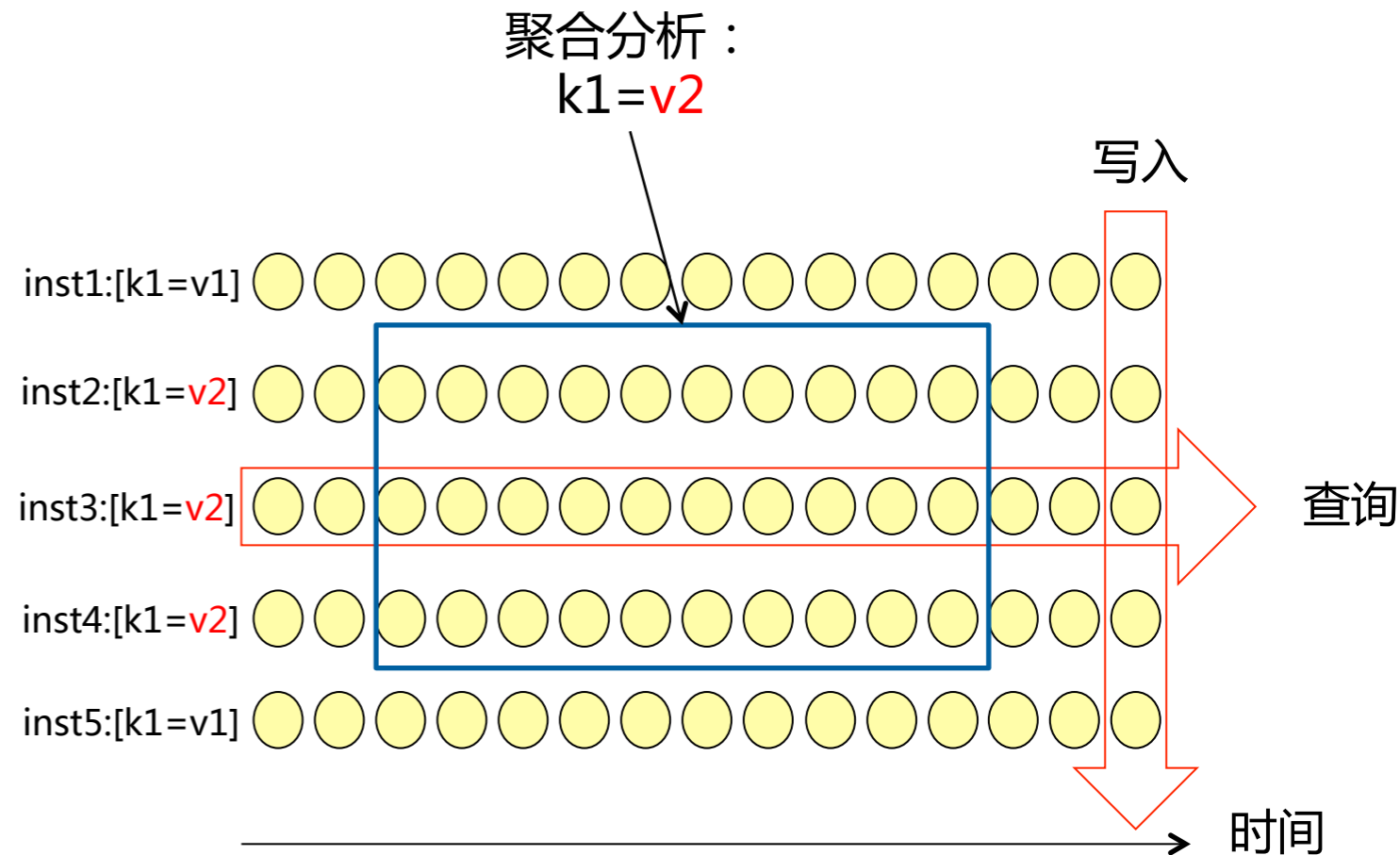
时间序列数据

- 时序 (Time Series) : object + tags + metric + datapoints
 - 监控对象(object): 1.nginx.www.tc 【机器/实例/服务/机房/网络等】
 - 监控项(metric) : pv 【资源消耗/业务指标/服务状态等】
 - 标签(tag , 或称维度) : province=shanghai, isp=unicom
 - 数据点(datapoint): (timestamp, value)=(1504594100, 5000)



时序数据读写特征

- 写多读少：95% ~ 99%写入
- 读写正交：从时间维度看
- 按标签（tag）对数据聚合分析
- 新近数据更重要，在线频繁访问
- 历史数据长期保存做离线计算
- 事务特性（ACID）需求不强



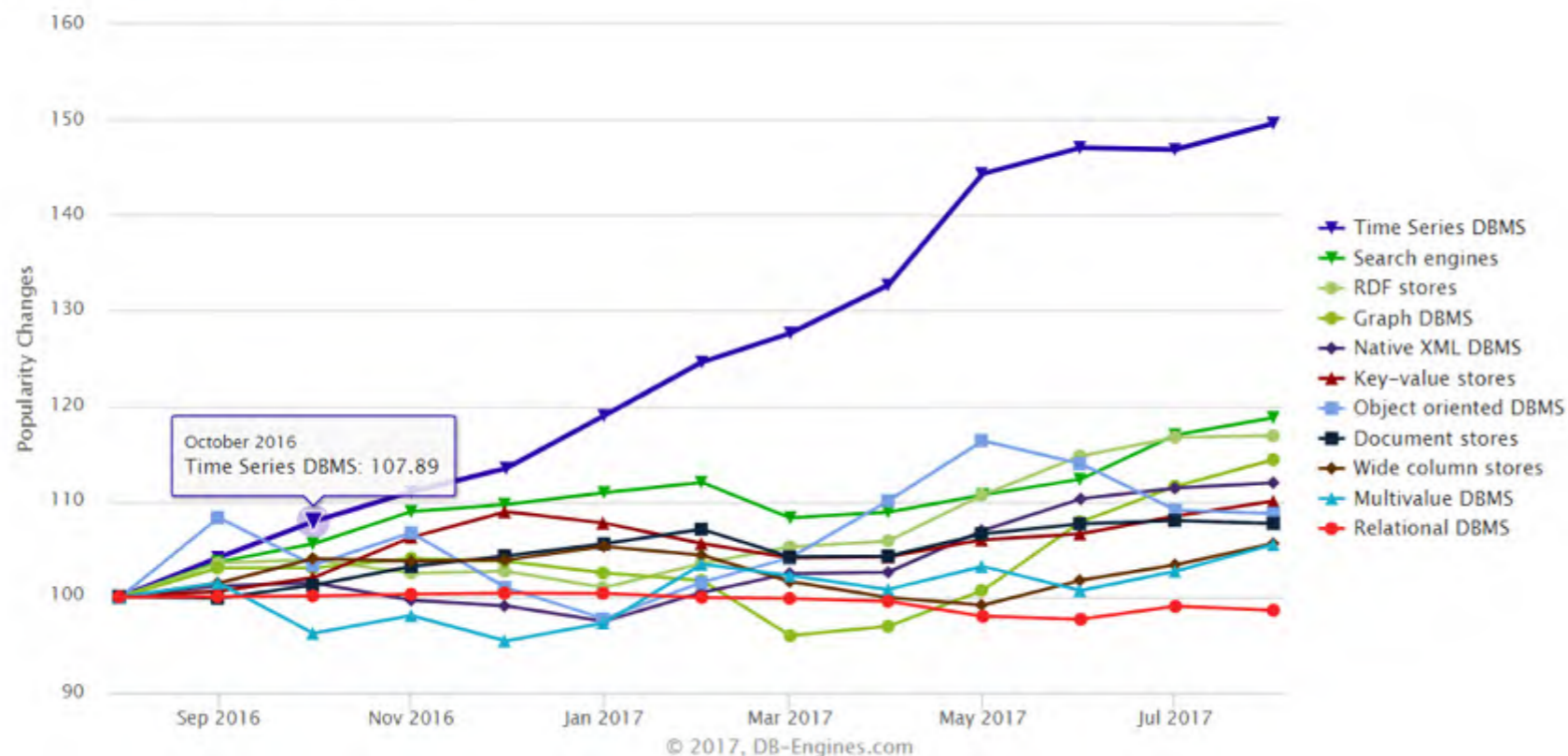
时序数据存储的技术挑战（百度）

- 监控规模
 - 监控对象：千万级
 - 监控项和标签组合：平均百级别
 - 指标规模（曲线数）：**10亿级**
- 系统吞吐量
 - 采集周期：10s，甚至5s
 - 监控数据点写入：**几千万/每秒**，达**几万亿/每天**
 - 查询量：**几万次/每秒**，达**几十亿/每天**
- 可用性/性能
 - 持续高负载：**7 * 24** 小时
 - 可用性：**99.99%**
 - 查询性能：500ms p99th

业界与开源

- 开源
 - InfluxDB : 底层是自研TSM存储引擎
 - OpenTSDB : 底层使用HBase
 - KairosDB : 底层使用Cassandra
 - Prometheus : 基于LevelDB存储引擎
- 业界
 - Facebook : Gorilla(自研)+ HBase
 - Twitter : Cache+Manhattan(自研)
- 共同点
 - 底层存储基于**LSM**(Log-structured Merge-Tree)

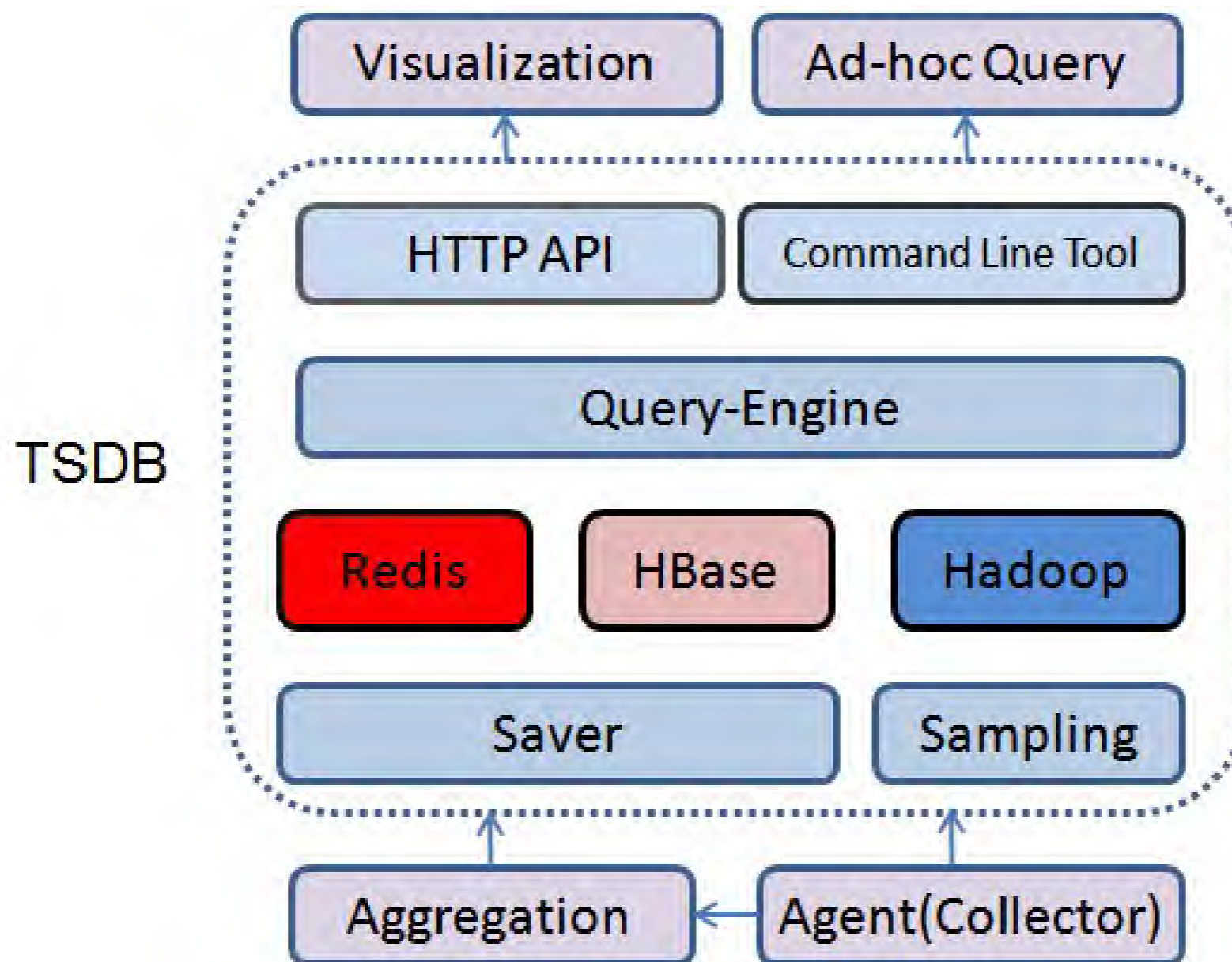
Trend of the last 12 months



Source : DB-Engines

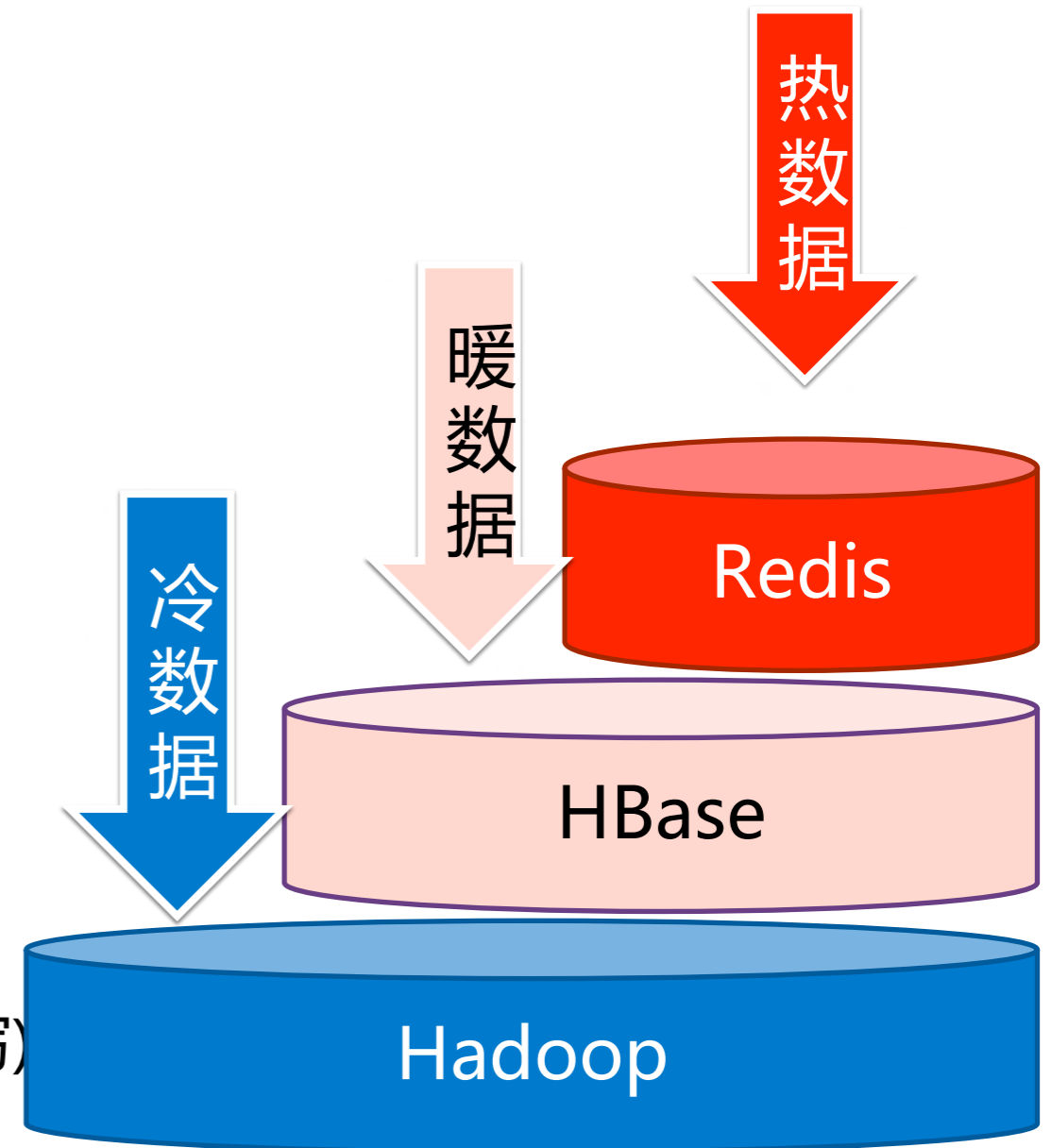
流行趋势 : 最近1年

时序数据 (TSDB) 的存储架构



存储层次

- 层次化存储
 - 关键/热数据缓存到Redis
 - 暖数据存储到 HBase
 - 冷数据存储到 Hadoop
- Redis
 - 成熟的缓存存储
 - 百度自研的分布式Redis
- HBase
 - 分布式表格存储，可扩展性
 - **写优化**(基于LSM将随机写转化为批量顺序写)
 - 业界有成功案例
- Hadoop



存储结构之HBase

- 行键：`<entity_id> <metric_id> <timestamp_hour>`
 - 监控对象+标签：`object:[tagk1=tagv1,...,tagkN=tagvN]` → `entity_id`
 - 监控项：`metric` → `metric_id`
 - 时间戳：`timestamp` → `timestamp_hour + timestamp_offset`
- 数据格式：借鉴OpenTSDB，按小时的存成一行，如下表格
 - 列为`timestamp_offset`，单元格存储`value`
- 按维度的二级索引
 - 建立维度到`entity_id`的二级索引

Row Key	+0	+10	+20	...	+3590
0x01037130068293841292148000	82	78	84	...	89

数据压缩

- 压缩算法 : Facebook Gorilla 时序数据压缩
 - Delta-of-delta 压缩时间戳
 - XOR 压缩数值
- 压缩效果
 - 在原有基础上**80%**的存储节省 (内存)
 - <8%的CPU成本

N-2 timestamp	02:00:00	-
N-1 timestamp	02:01:02	Delta: 62
timestamp	02:02:02	Delta: 60
		Delta of deltas: -2

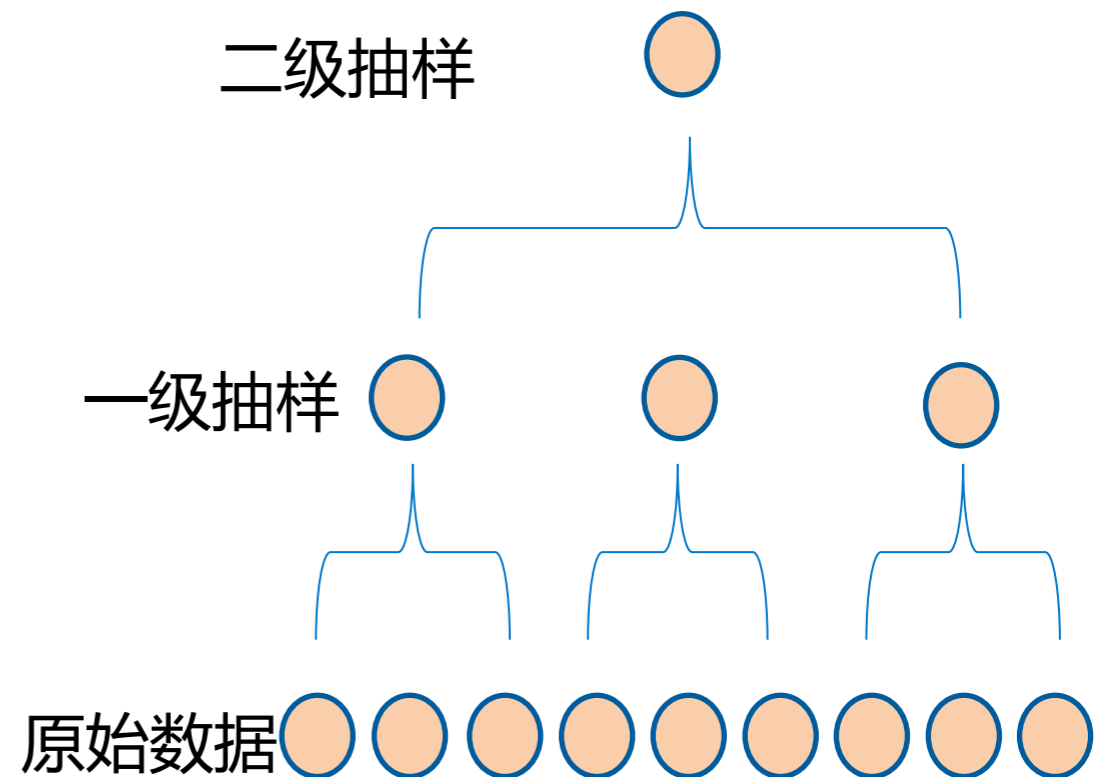
Previous Value	12.0	0x4028000000000000
Value	24.0	0x4038000000000000
XOR	-	0x0010000000000000

11 leading zeros, # of meaningful bits is 1

Source : Gorilla
2015 VLDB

分级统计抽样

- 统计抽样
 - 统计产生 {max, min, sum, count}
 - 查询上对用户透明
- 分级数据保存策略【HBase中】
 - 原始数据在线保存数十天
 - 一级抽样保存数月
 - 二级抽样存储数年



分集群/分库/分表

- 分集群
 - 按不同的监控系统分集群，比如网络监控独立集群
- 分库
 - 按产品线的来分库
 - 定制数据的存储特征，比如存储时长
- 分表
 - 按时间分表
 - 数据的批量过期清除
 - 减小HBase的compaction压力

高可用/可扩展

- 可用性问题
 - HBase的可用性存在不足
 - 单机房故障容灾
- 高可用方案
 - 跨地域双活集群
 - 故障自动切换
- 可扩展
 - 虚拟化实例部署
 - 弹性扩缩容

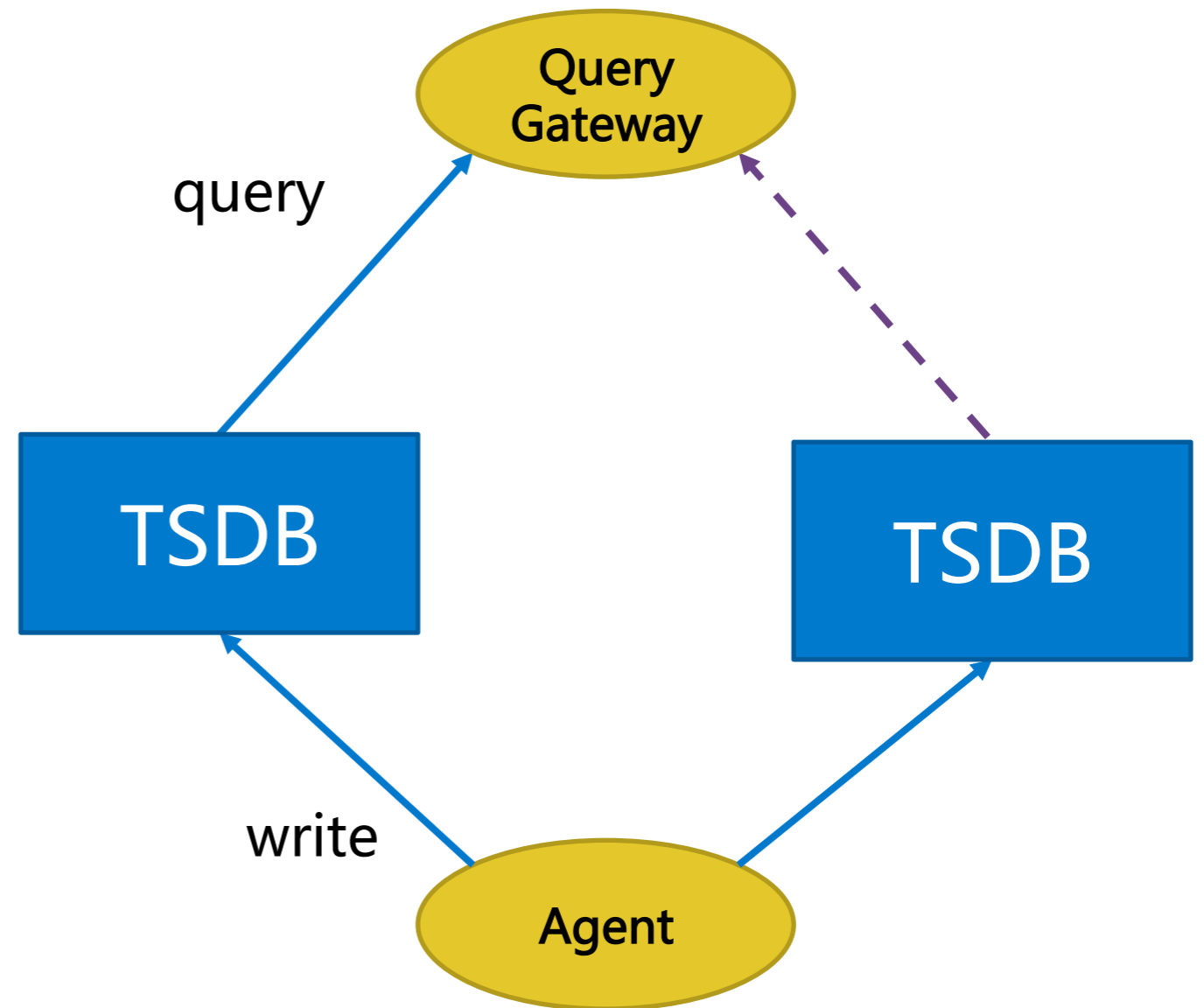


TABLE OF CONTENTS

智能运维与大数据

大规模时序数据存储

海量运维事件数据存储

运维知识库

事件数据库：EventDB

- 运维事件
 - 服务变更、运营活动和服务异常等跟运维相关的事件
- 事件数据模型
 - JSON格式，带有时间域，可扩展
 - 按事件类型分别建立schema：变更、运营和异常等
- 系统规模
 - 日增长量：千万数量级
 - 总事件数：百亿数量级
- 技术方案：ElasticSearch + OpenResty
 - 底层存储：ElasticSearch，分布式的全文搜索引擎，按时间、按类型分表
 - 前端Proxy：OpenResty，双写主备ES集群，自动故障切换

TABLE OF CONTENTS

智能运维与大数据

大规模时序数据存储

海量运维事件数据存储

运维知识库

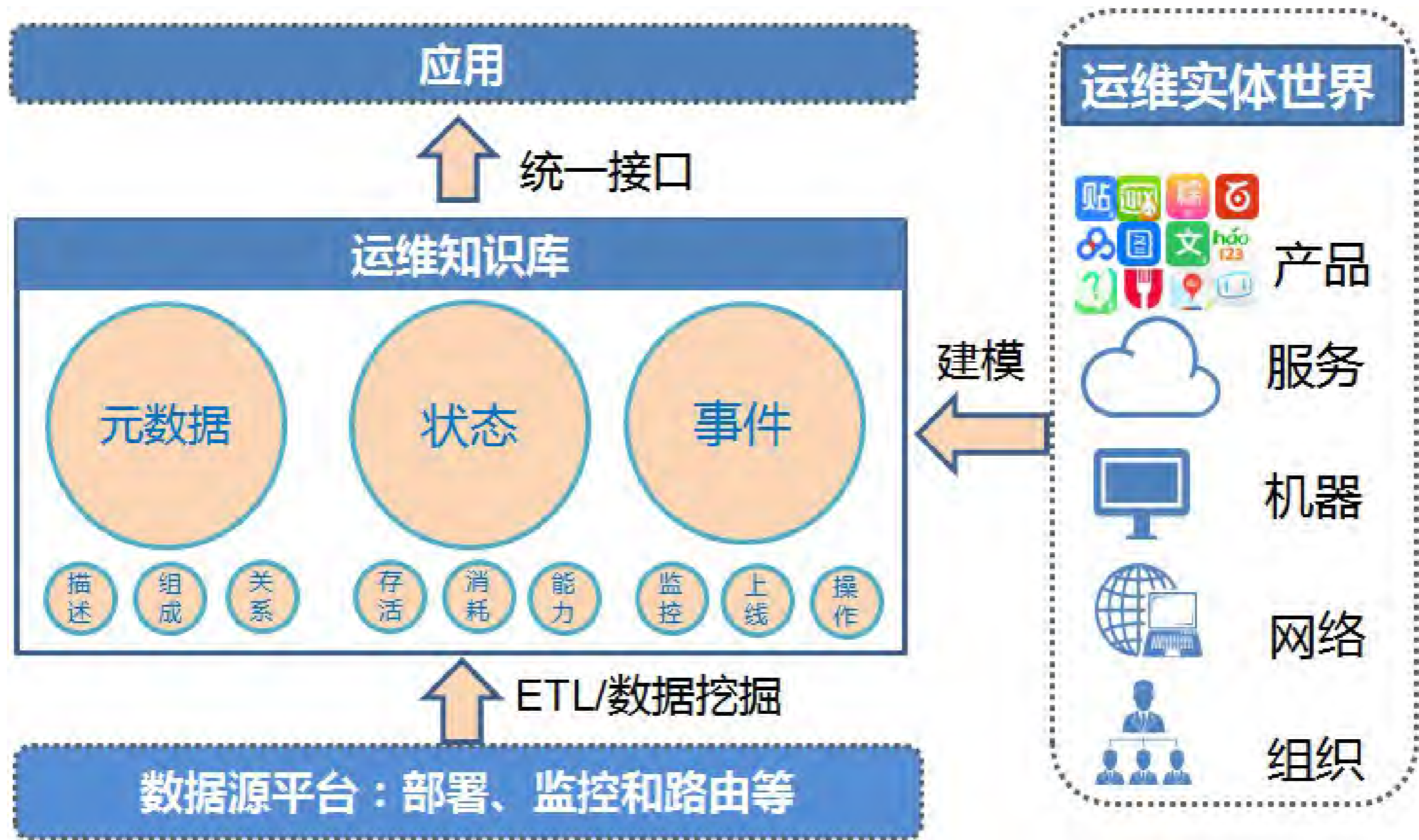
运维知识库

- 背景
 - 同场景：平台多，数据格式不同，数据分散
 - 跨场景：监控、部署和调度等数据未关联
 - 部分数据未结构化（以文档/网页形式存在）
- 目标
 - 建立对“运维世界”的统一、数字化描述
 - 打通数据的壁垒，形成智能运维闭环
 - 积累数据，挖掘运维知识

书同文：一致运维“语言”



运维知识库的系统架构



总结

- 运维大数据存储平台
 - 统一数据模型，可扩展
 - 拥抱开源，对技术持开放心态
 - 量材适用，合理设计
- 大规模时序数据存储
 - 针对写入优化（LSM树）
 - 按数据冷热做分层存储
 - 按时序特点做极致压缩：Delta/XOR
 - 分治：分集群/分库/分表
 - 可用性：主备，自动故障切换



AIOps智能运维

THANKS!

智能时代的新运维

CNUTCon 2017