

开源大数据技术专场（上午）:Spark、HBase、JStorm应用与实践

转载

[aibiba0894](#) 于 2016-10-24 15:29:00 发布 34 收藏

文章标签: [大数据](#) [运维](#) [java](#)

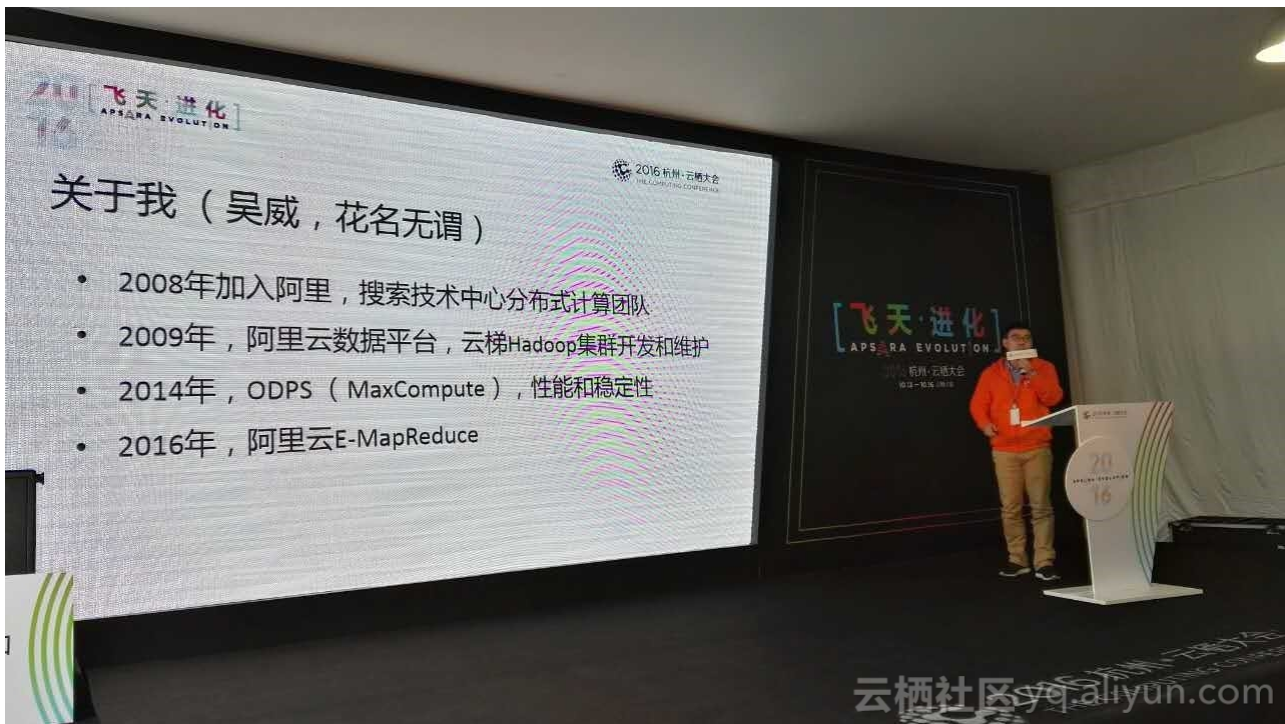
原文链接: <http://www.cnblogs.com/onetwo/p/5993158.html>

版权

16日上午9点, 2016云栖大会“开源大数据技术专场”(全天)在阿里云技术专家封神的主持下开启。通过封神了解到, 在上午的专场中, 阿里云高级技术专家无谓、阿里云技术专家封神、阿里巴巴中间件技术部高级技术专家天梧、阿里巴巴中间件技术部资深技术专家纪君祥将给大家带来Hadoop、Spark、HBase、JStorm Turbo等内容。



无谓: Hadoop过去现在未来, 从阿里云梯到E-MapReduce

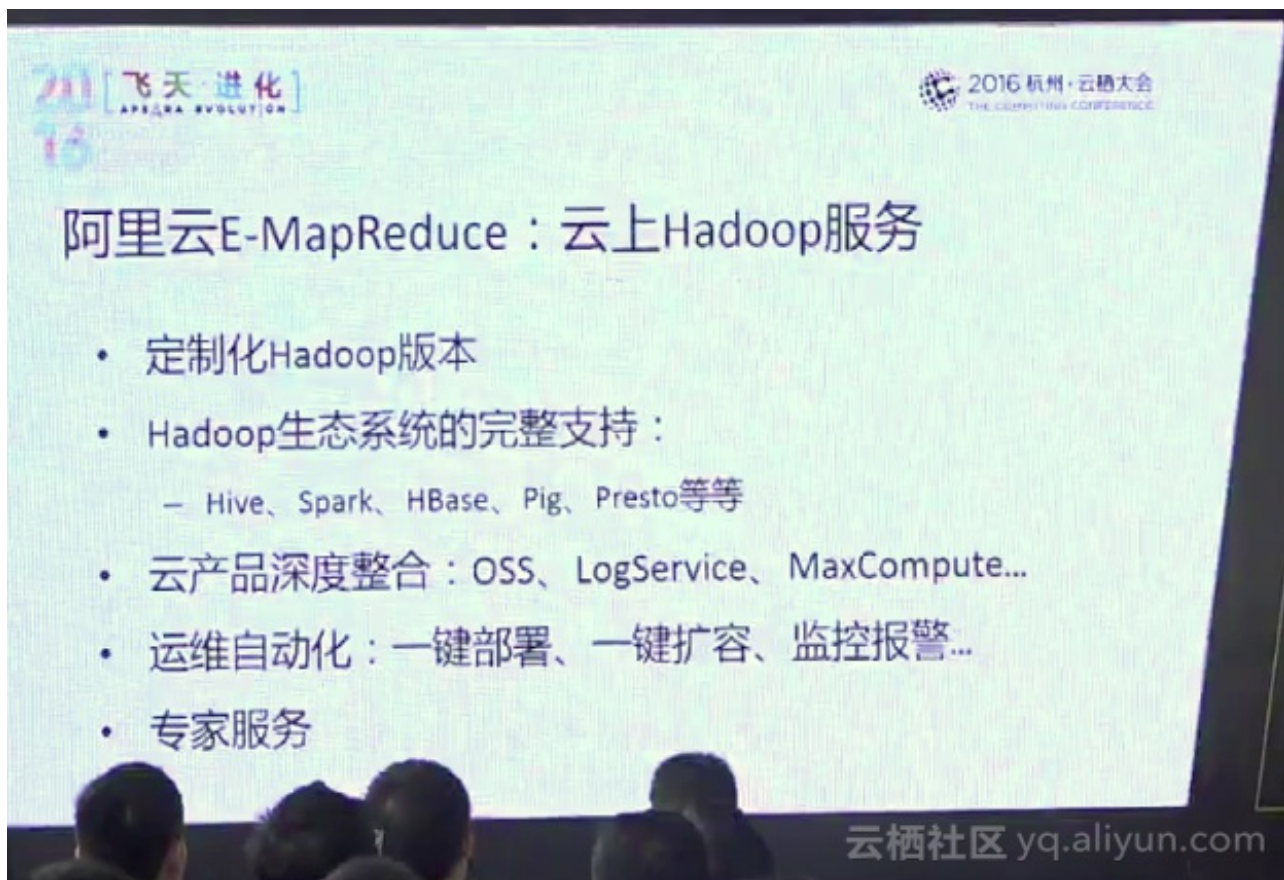


阿里云高级技术专家 无谓

从开辟大数据先河至现在，风雨十年，Hadoop已成为企业的通用大数据框架。而作为上午的第一个演讲，无谓首先给我们总结了Hadoop这十年，也是从离线到在线的十年，其中意义重大的事情有：YARN成为大数据操作系统；Hadoop成为企业级解决方案，涵盖数据可视化工具、存储、计算、数据管理等；机器学习和人工智能的支持；Mahout->oryx，批处理到实时处理的学习工具。



而在这段时间，阿里从2008年就已经参与到Hadoop中，其主要阶段可概括为：2008-2009期间，建立了多部门独立的Hadoop集群；2009-2015，主要做云梯集群和服务，包括：集群统一运维，专业的开发团队；数据统一管理，集团层面的全局视图；资源错峰分配，整体成本最优；2015-至今，阿里云E-MapReduce，阿里云对外的Hadoop基础服务。



2016 [飞天·进化] AWS/AZURE EVOLUTION

2016 杭州·云栖大会 THE CLOUDPIPING CONFERENCE

阿里云E-MapReduce：云上Hadoop服务

- 定制化Hadoop版本
- Hadoop生态系统的完整支持：
 - Hive、Spark、HBase、Pig、Presto等等
- 云产品深度整合：OSS、LogService、MaxCompute...
- 运维自动化：一键部署、一键扩容、监控报警...
- 专家服务

云栖社区 yq.aliyun.com

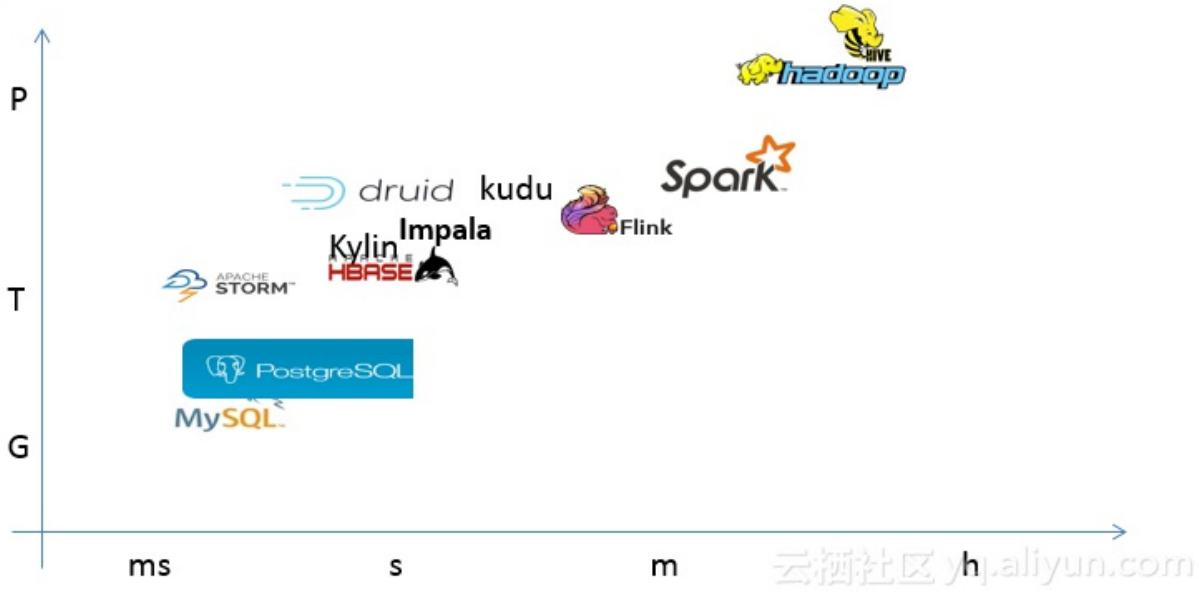
随后，无谓还重点分享了阿里内部的Hadoop服务云梯：全局资源调度：支持业务优先级（基于Fair Scheduler）；安全性，HDFS上的扩展ACL，Hive安全认证和授权；稳定性，消除异常作业对全局的影响 Master HA；扩展性：Master节点的单点性能压力，跨机房的部署架构；云梯医生：集群诊断系统，最后，通过无谓，我们还体会了阿里云分享的技术红利E-MapReduce。

封神：Spark实践与探索



阿里云技术专家 封神

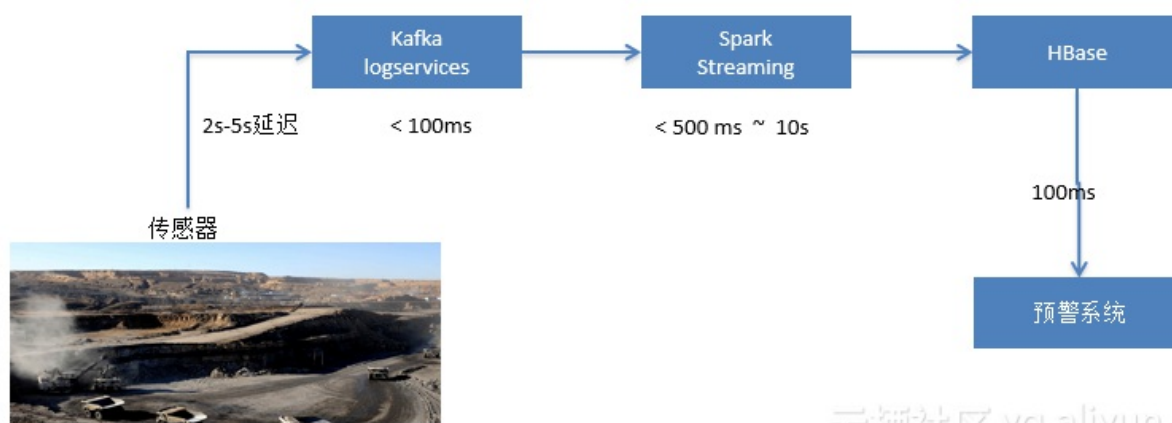
封神专注于大数据领域，拥有7年的分布式引擎开发经验，先后参与了上万台Hadoop、ODPS集群的开发。在本次演讲中，他主要介绍了数据处理技术、About Spark、阿里的Spark历程、Spark与云，及Spark未来多个方面。



在时下流行大数据技术对比中，封神首先从数据处理时间与数据量两个方面维度进行了切入，在这个过程中，我们会发现，没有哪个软件能解决所有的问题，能解决问题也是在一个范围内，即使是Spark、Flink等。目前存在有意思的事情是：Greenplum类似的MPP引擎想处理大数据的需求，Hadoop等被定位为大数据的引擎也想解决小数据的问题（列式存储、或者也加入一些索引）。图中右上角的想往左边靠，减少延迟，图中左下角的想往上面靠，增大能处理的数据量。此外在DB/MPP与Hadoop的对比上，Hadoop生态圈为何如何火爆也能有所体现：首先，在硬件需求上，DB/MPP可能需要小型机和高端存储，同时也需要RAID，而Hadoop只需要普通的PC机；容错性上，DB/MPP重跑即可，而Hadoop则需要容错；在调度模型上，DB/MPP使用了基于线程的调度，而Hadoop则需要做CPU/Memory的调度；最后，在衡量指标上DB/MPP一般以QPS为标准，而Hadoop相关系统一般更看重吞吐。

✓ Batch	✓ Batch ✓ Interactive	✓ Batch ✓ Interactive ✓ Memory ✓ Near-Real Time Streaming ✓ Full Stack	✓ Hybrid (Batch+Streaming) ✓ Interactive ✓ Memory ✓ Near-Real Time ✓ Streaming ✓ Full Stack	✓ Hybrid(Batch+Streaming) ✓ Interactive ✓ Real-Time Streaming ✓ Native Iterative Processing ✓ Full Stack
MapReduce	DAG: Direct Acyclic Graphs	RDD: Resilient Distributed Datasets	RDD: Resilient Distributed Datasets	Cyclic Dataflows
1G	2G	3G	3.8G	4G

随后，封神详细的介绍了Spark的各个部分，更重点介绍了Spark链路、Spark Core、Spark弹性伸缩，并结合业务介绍了Spark的相关实践。而在阿里的使用上，封神表示，在10年至12年，阿里就对Spark进行了初步尝试，其中主要使用了Standalone方式，主要做了Spark Mllib机器学习的探索；在12-14年，Spark on YARN成为主要探索部分，至此集群的规模已达数百台，其中14还做了基于内存计算的研究；而从15年至今，团队已对公共云提供服务，主要是作业平台、运维平台。



云栖社区 yq.aliyun.com

演讲期间，封神还就机器学习、即席查询、流式计算3个时下热门的大数据实用场景进行了解析，在机器学习场景中，他对比了Mahout on mr、MLlib On Spark，以及MPICH2三个方式，其中在易用性和稳定性上，Spark都有一个非常好的表现，但是在性能上MPICH2则更胜一筹，但是后者不管稳定和易用都与Spark相去甚远。

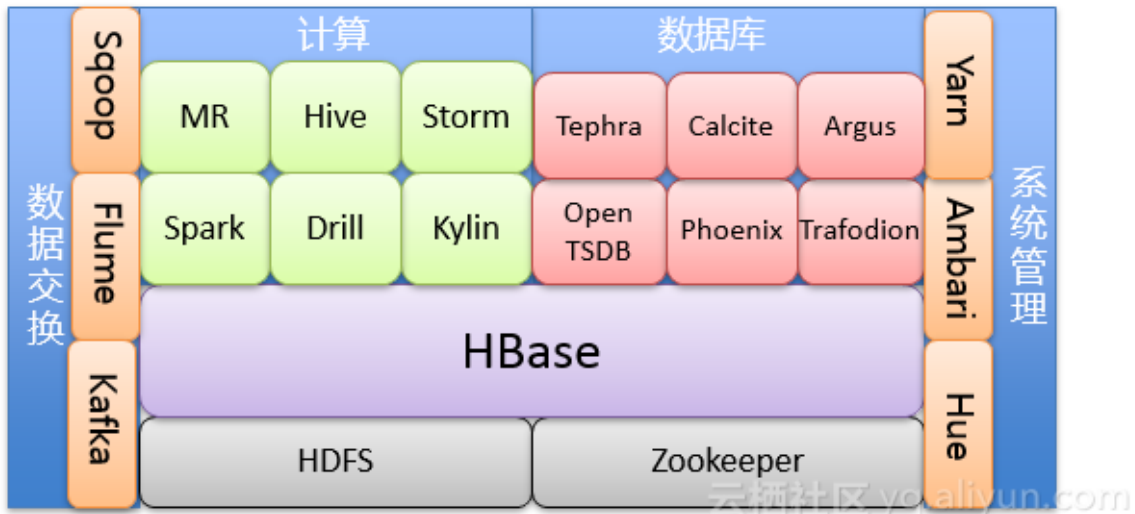
最后，在Spark未来的方向上，封神表示，支持ANSI SQL，性能接近MPP数据仓库，一切基于优化（Catalyst），新硬件的支持（比如：大内存、GPU），更加友好的支持云显然更为重要。

天梧：HBase的一些实践与探索



阿里巴巴中间件技术部高级技术专家 天梧

时至今日，数据已经在各行各业产生价值，在天梧的演讲中，他首先为大家总结了大数据的应用形式，主要包括：万物万面，精准定像；数据赋能，运筹帷幄；智能生活。然而，机构如果想真正地享受数据带来的红利，机构仍然需要拼尽全力去克服大数据场景的数据特点，比如：基础量大、增长快、计算与存储的实时性要求迫切、时效性短、易发散、易产生脏数据等。随后，天梧结合各个应用场景开始了HBase实践的讲解。

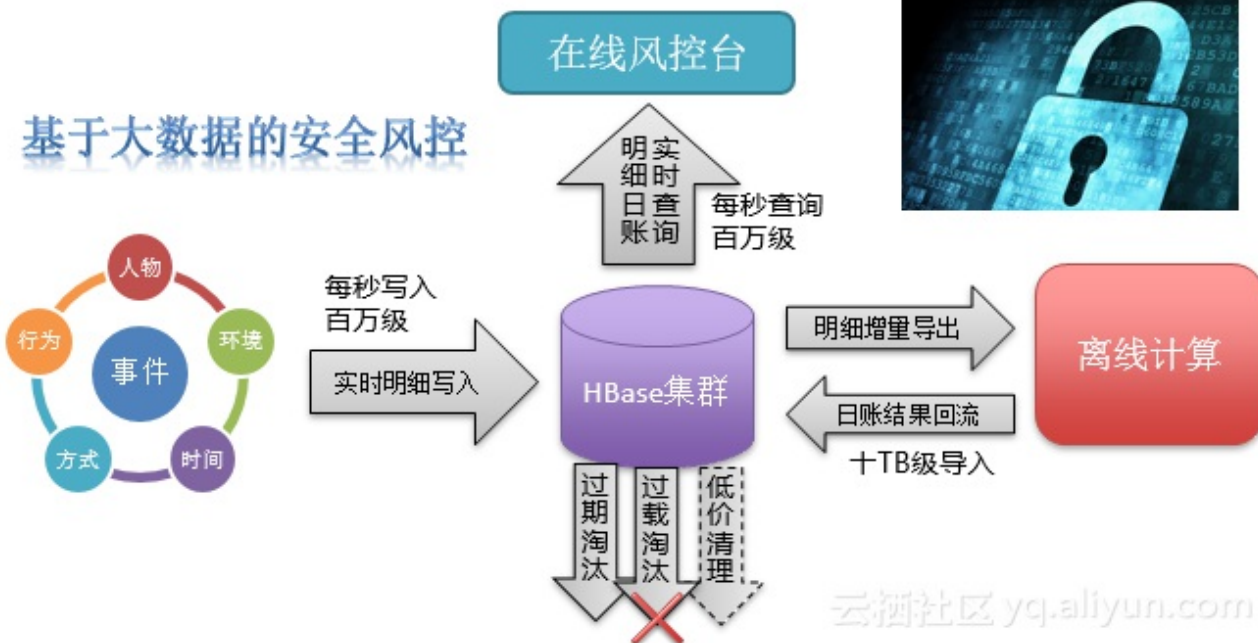


Hadoop Database，是一个基于Google BigTable论文设计的高可靠性、高性能、可伸缩的分布式存储系统，它的具体特性有：松散表，实时更新、增量导入、多维删除，随机查询、范围查询，高伸缩、高可用、高可靠、高性能、高适应，在线分布式NOSQL数据库。

与Hadoop的天然集成让HBase天生具备了很多优势，在阿里之外，同样得到了 Intel、Facebook、Cludera、Hortonworks、小米等公司的支持。而在此之外，HBase的其他基因同样深受大数据玩家的喜爱，包括：自动分区，分区自动分裂，分区在线Merge，可应对数据爆发式增长和访问爆发式增长；LSM，写吞吐高，不受SSD随机写入放大干扰，不受空间放大干扰；存储计算分离，负载均衡更高效，资源扩容更节省，存储优化更便捷（非对称副本冗余：异构介质、Erasure Code等）。



基于大数据的安全风控



可以说，HBase为大数据而生。然而就如任何开源软件，HBase的使用同样需要大量的研发投入。在这里，阿里也基于阿里巴巴/蚂蚁的环境和业务需求，对社区HBase进行深度定制与改进，从内核引擎、解决方案、稳定护航、发展支撑等全方位提供一站式大数据基础存储服务，就拿灾备体系来说，包括集群数据复制的诉求、多集群数据复制、流量切换、跨集群一致性保证、深度优化的宕机恢复能力等方面。集群数据复制的诉求，数据一致，延迟低，吞吐大，多源多目标，链路粒度细，异构系统，可视可追踪等；多集群数据复制，异步模式，同步模式，支持多地多单元、表级复制、循环流动，支持延迟/拓扑/复制详情可视，支持数据的链路追踪，支持实时复制到异构系统，并发、吞吐、实时的有效权衡异步模式；流量切换，虚拟地址映射，支持一键切换、自动切换；跨集群一致性保证，基于读写保护的强一致；深度优化的宕机恢复能力。

天梧表示，在此之外，在HBase上阿里还做了调整、报警、健康等各个方面的工作。而在未来，更大硬件支持、容器化部署也将是一大研究的方向。

纪君祥：阿里巴巴实时计算平台 JStorm Turbo



阿里巴巴中间件技术部资深技术专家 纪君祥

通过纪君祥了解到，从2013年4月3日起，JStorm已经发布了25个版本，部署方式包括Standalone、JStorm-on-yarn、JStorm-on-docker等方式，部署超过4000台主机，支撑了1500以上的应用，拥有超过2000+的 topologies。

- Fraud Detection: Nut/Velocity
- Auditing: Alimama Ads auditing, Alimama P4P auditing, AMG repository auditing, ...
- Statistics: Eagleeye, AE pv – uv statistics, BI statistics...
- Monitoring: Tlog, Rds-monitor, Oceanbase-monitor, Cainiao Raidar, QuaMonitor, YunOSMonitor, ...
- Data Transport: RDS-log-sync, Unify-Log
- Realtime Recommendation: Alipay 1315, Alipay Hyperloop, Tpp
- Scheduling Apps: AE mail realtime analyzer, 1688 Wangxiaobao

云栖社区 yq.aliyun.com

在JStorm与Storm区别上，纪君祥提到JStorm更是一个流处理生态系统，而不是简单的一个流计算框架。同时，对于企业来说JStorm还是一个成熟的Java版Storm，它不仅运营更快、更稳定，也具备了更多的功能。

转载于:<https://www.cnblogs.com/onetwo/p/5993158.html>



[创作打卡挑战赛](#) >
[赢取流量/现金/CSDN周边激励大奖](#)