

Apache Hadoop 2.8.0正式发布

时隔两年，Apache Hadoop终于又有大改版，Apache基金会近日发布了Hadoop 2.8版，一次新增了2,919项更新功能或新特色。不过，Hadoop官网建议，2.8.0仍有少数功能在测试，要等到释出2.8.1或是2.8.2版才适合用于正式环境。

在2.8.0版众多更新，主要分布于4大套件分别是：

- 共用套件 (Common)
- 底层分散式档案系统HDFS套件(HDFS)
- MapReduce运算套件(MapReduce)
- YARN分析框架(YARN)



如果想及时了

解Spark、Hadoop或者Hbase相关的文章，欢迎关注微信公共帐号：iteblog_hadoop

例如，共用套件中，可无上限存取S3档案的Hadoop内建S3A机制，可直接外挂上任何AWS验证的资料来源，也可用hadoop验证API直接取得S3A验证，取代透过XML配置档来设定的方式，还能支持Amazon STS临时验证。另外，新版也开始支持直接存取 Azure Data Lake 文件系统。这两项功能都有助于简化Hadoop从大型云端储存服务中取得大量分析资料。

除此之外，DFS分布式文件系统，支持异步调用重试 (Async Call Retry) 和容错转移 (Failover) ，来降低DFS文件系统重连的门槛。安全性方面，新版可透过Servlet过滤器，来防护XFS攻击 (跨框架脚本攻击， Cross Frame Scripting) 。构建方面，2.8则用Yetus取代了wrapper版本发布方式，还新增了Docker软体包的发布方式，更容易打包建置环境来测试或交换。

HDFS套件则是强化了WebHDFS对伪造跨站请求 (Cross-site Request Forgery) 的检查来提高安

全性，也支持OAuth2验证方式。HDFS还新增了多层式巢状加密区机制，不再只能指定单一个目录加密，也能对目录底下的不同目录，分别建立加密区来强化控管。另外，还采用了新的DataNode协定，可以避免NameNode因为心跳的延迟而导致不正确地处理DataNodes的状态。

虽然大多数Hadoop使用者不会在Windows环境执行程序，但YARN套件还是新增对Windows环境的CPU资源监控。而MapReduce套件的新特色之一则是，发布MapReduce运算任务时，可以顺便加上标签以便于管理。

完整的更新如下：

Common

- Support async call retry and failover which can be used in async DFS implementation with retry effort.
- Cross Frame Scripting (XFS) prevention for UIs can be provided through a common servlet filter.
- S3A improvements: add ability to plug in any AWSCredentialsProvider, support read s3a credentials from hadoop credential provider API in addition to XML configuraiton files, support Amazon STS temporary credentials
- WASB improvements: adding append API support
- Build enhancements: replace dev-support with wrappers to Yetus, provide a docker based solution to setup a build environment, remove CHANGES.txt and rework the change log and release notes.
- Add posixGroups support for LDAP groups mapping service.
- Support integration with Azure Data Lake (ADL) as an alternative Hadoop-compatible file system.

HDFS

- WebHDFS enhancements: integrate CSRF prevention filter in WebHDFS, support OAuth2 in WebHDFS, disallow/allow snapshots via WebHDFS
- Allow long-running Balancer to login with keytab
- Add ReverseXML processor which reconstructs an fsimage from an XML file. This will make it easy to create fsimages for testing, and manually edit fsimages when there is corruption
- Support nested encryption zones
- DataNode Lifeline Protocol: an alternative protocol for reporting DataNode liveness. This can prevent the NameNode from incorrectly marking DataNodes as stale or dead in highly overloaded clusters where heartbeat processing is suffering delays.
- Logging HDFS operation's caller context into audit logs
- A new Datanode command for evicting writers which is useful when data node decommissioning is blocked by slow writers.

YARN

- NodeManager CPU resource monitoring in Windows.
- NM shutdown more graceful: NM will unregister to RM immediately rather than waiting for timeout to be LOST (if NM work preserving is not enabled).
- Add ability to fail a specific AM attempt in scenario of AM attempt get stuck.
- CallerContext support in YARN audit log.
- ATS versioning support: a new configuration to indicate timeline service version.

MAPREDUCE

- Allow node labels get specified in submitting MR jobs
- Add a new tool to combine aggregated logs into HAR files

本博客文章除特别声明，全部都是原创！
转载本文请加上：转载自过往记忆（<https://www.iteblog.com/>）
本文链接：【】（）