

相比于Hadoop1.0，Hadoop 2.0中的HDFS增加了两个重大特性，HA和Federaion。HA即为High Availability，用于解决NameNode单点故障问题，该特性通过热备的方式为主NameNode提供一个备用者，一旦主NameNode出现故障，可以迅速切换

Secondary NameNode至备NameNode，从而实现不间断对外提供服务。Federation即为“联邦”，该特性允许一个HDFS集群中存在多个NameNode同时对外提供服务，这些NameNode分管一部分目录（水平切分），彼此之间相互隔离，但共享底层的数据节点存储资源。

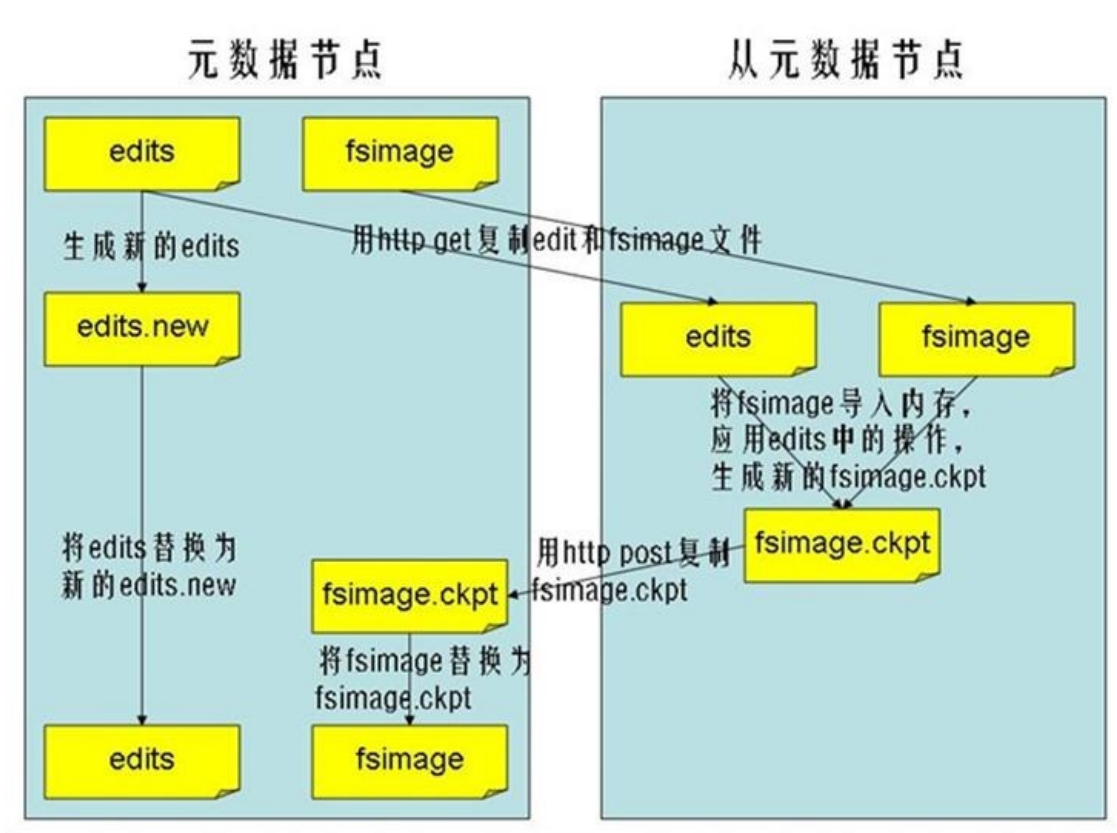
Secondary NameNode的检查点进程启动，是由两个配置参数控制的：

fs.checkpoint.period，指定连续两次检查点的最大时间间隔，默认值是1小时。

fs.checkpoint.size定义了edits日志文件的最大值，一旦超过这个值会导致强制执行检查点（即使没到检查点的最大时间间隔）。默认值是64MB。

日志与镜像的定期合并总共分五步：

- 1、SecondaryNameNode通知NameNode准备提交edits文件，此时主节点产生edits.new
- 2、SecondaryNameNode通过http get方式获取NameNode的fsimage与edits文件（在SecondaryNameNode的current同级目录下可见到 temp.check-point或者previous-checkpoint目录，这些目录中存储着从namenode拷贝来的镜像文件）
- 3、SecondaryNameNode开始合并获取的上述两个文件，产生一个新的fsimage文件fsimage.ckpt
- 4、SecondaryNameNode用http post方式发送fsimage.ckpt至NameNode
- 5、NameNode将fsimage.ckpt与edits.new文件分别重命名为fsimage与edits，然后更新fstime，整个checkpoint过程到此结束。



四、HDFS中文件读写操作流程

在HDFS中，文件的读写过程就是client和NameNode以及DataNode一起交互的过程。我们已经知道NameNode管理着文件系统的元数据，DataNode存储的是实际的数据，那么client就会联系NameNode以获取文件的元数据，而真正的文件读取操作是直接和DataNode进行交互的。

写文件的过程：

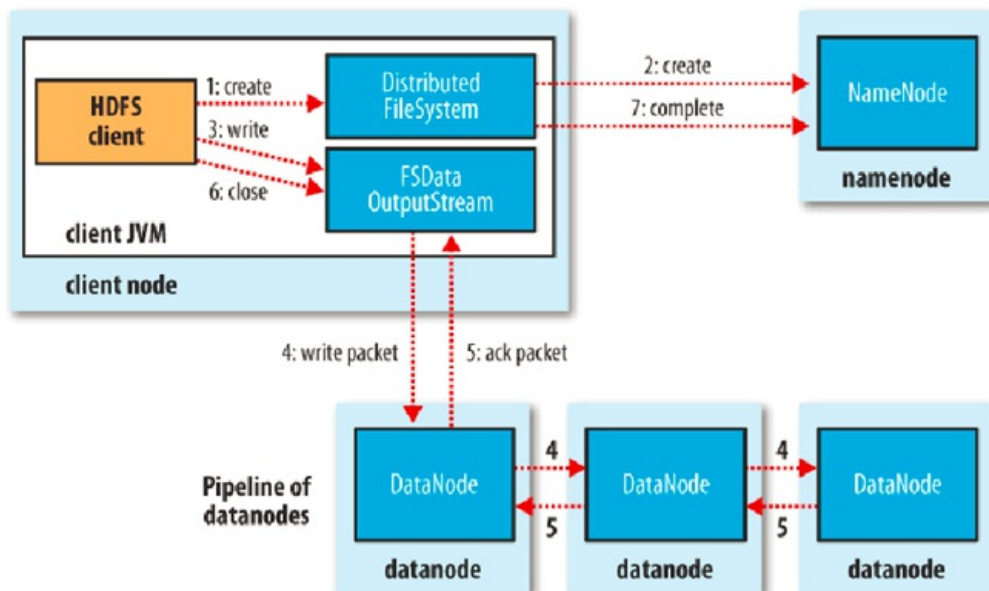
- 1.初始化FileSystem，客户端调用create()来创建文件
- 2.FileSystem用RPC调用元数据节点，在文件系统的命名空间中创建一个新的文件，元数据节点首先确定文件原来不存在，并且客户端有创建文件的权限，然后创建新文件。
- 3.FileSystem返回DFSOutputStream，客户端用于写数据，客户端开始写入数据。
- 4.DFSOutputStream将数据分成块，写入data queue。data queue由Data Streamer读取，并通知元数据节点分配数据节点，用来

存储数据块(每块默认复制3块)。分配的数据节点放在一个pipeline里。Data Streamer将数据块写入pipeline中的第一个数据节点。第一个数据节点将数据块发送给第二个数据节点。第二个数据节点将数据发送给第三个数据节点。

5.DFSOutputStream为发出去的数据块保存了ack queue，等待pipeline中的数据节点告知数据已经写入成功。

6.当客户端结束写入数据，则调用stream的close函数。此操作将所有数据块写入pipeline中的数据节点，并等待ack queue返回成功。最后通知元数据节点写入完毕。

7.如果数据节点在写入的过程中失败，关闭pipeline，将ack queue中的数据块放入data queue的开始，当前的数据块在已经写入的数据节点中被元数据节点赋予新的标示，则错误节点重启后能够察觉其数据块是过时的，会被删除。失败的数据节点从pipeline中移除，另外的数据块则写入pipeline中的另外两个数据节点。元数据节点则被通知此数据块是复制块数不足，将来会再创建第三份备份。



读文件的过程：

1.初始化FileSystem，然后客户端(client)用FileSystem的open()函数打开文件

2.FileSystem用RPC调用元数据节点，得到文件的数据块信息，对于每一个数据块，元数据节点返回保存数据块的数据节点的地址。

3.FileSystem返回FSDataInputStream给客户端，用来读取数据，客户端调用stream的read()函数开始读取数据。

4.DFSInputStream连接保存此文件第一个数据块的最近的数据节点，data从数据节点读到客户端(client)

5.当此数据块读取完毕时，DFSInputStream关闭并此数据节点的连接，然后连接此文件下一个数据块的最近的数据节点。

6.当客户端读取完毕数据的时候，调用FSDataInputStream的close函数。

7.在读取数据的过程中，如果客户端在与数据节点通信出现错误，则尝试连接包含此数据块的下一个数据节点。

8.失败的数据节点将被记录，以后不再连接。

