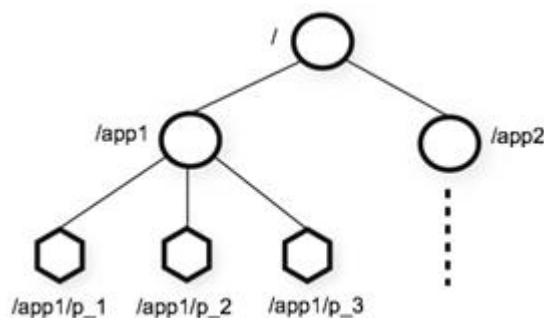


1. zookeeper的数据模型

- ZooKeeper 的数据模型，在结构上和标准文件系统的非常相似，拥有一个层次的命名空间，都是采用**树形层次结构**。



- ZooKeeper 树中的每个节点被称为一个Znode。和文件系统的目录树一样，ZooKeeper 树中的每个节点可以拥有子节点。

但也有不同之处：

1. Znode **兼具文件和目录两种特点**。既像文件一样维护着数据、元信息、ACL、时间戳等数据结构，又像目录一样可以作为路径标识的一部分，并可以具有子 Znode。用户对 Znode 具有增、删、改、查等操作（权限允许的情况下）。
2. Znode 存**储数据大小有限制**。ZooKeeper 虽然可以关联一些数据，但并没有被设计为常规的数据库或者大数据存储，相反的是，它用来管理调度数据，比如分布式应用中的配置文件信息、状态信息、汇集位置等等。这些数据的共同特性就是它们都是很小的数据，通常以 KB 为大小单位。ZooKeeper 的服务器和客户端都被设计为严格检查并限制每个 Znode 的数据大小至多 1M，常规使用中应该远小于此值。
3. Znode **通过路径引用**，如同 Unix 中的文件路径。**路径必须是绝对的**，因此他们必须由斜杠字符来开头。除此以外，他们必须是唯一的，也就是说每一个路径只有一个表示，因此这些路径不能改变。在 ZooKeeper 中，路径由 Unicode 字符串组成，并且有一些限制。字符串"/zookeeper"用以保存管理信息，比如关键配额信息。
4. 每个 Znode 由 3 部分组成：
 - **stat**：此为状态信息,描述该 Znode 的版本,权限等信息
 - **data**：与该 Znode 关联的数据
 - **children**：该 Znode 下的子节点

2. Znode节点类型

2.1 Znode 有两种，分别为**临时节点**和**永久节点**。节点的类型在创建时即被确定，并且不能改变。

- **临时节点**：该节点的生命周期依赖于创建它们的会话。一旦会话结束，临时节点将被自动删除，当然可以也可以手动删除。临时节点不允许拥有子节点。
- **永久节点**：该节点的生命周期不依赖于会话，并且只有在客户端显示执行删除操作的时候，他们才能被删除。

2.2 Znode 还有一个序列化的特性，如果创建的时候指定的话，该 Znode 的名字后面会自动追加一个不断增加的序列号。序列号对于此节点的父节点来说是唯一的，这样便会记录每个子节点创建的先后顺序。它的格式为“%10d”(10 位数字，没有数值的数位用 0 补充，例如“0000000001”)。

2.3 这样便会存在四种类型的 Znode 节点，分别对应：

- PERSISTENT：永久节点
- EPHEMERAL：临时节点
- PERSISTENT_SEQUENTIAL：永久节点、序列化
- EPHEMERAL_SEQUENTIAL：临时节点、序列化

3.Zookeeper的Shell 客户端操作

3.1 登录Zookeeper客户端

```
1 bin/zkCli.sh -server node01:2181
```

3.2 Zookeeper客户端操作命令

| 命令 | 说明 | 参数 |
|---|----------------------------|-------------------------------|
| <code>create [-s] [-e] path data acl</code> | 创建Znode | -s 指定是顺序节点 -e 指定是临时节点 |
| <code>ls path [watch]</code> | 列出Path下所有子Znode | |
| <code>get path [watch]</code> | 获取Path对应的Znode的数据和属性 | |
| <code>ls2 path [watch]</code> | 查看Path下所有子Znode以及子Znode的属性 | |
| <code>set path data [version]</code> | 更新节点 | version 数据版本 |
| <code>delete path [version]</code> | 删除节点, 如果要删除的节点有子Znode则无法删除 | version 数据版本 |
| <code>rmr path</code> | 删除节点, 如果有子Znode则递归删除 | |
| <code>setquota -n -b val path</code> | 修改Znode配额 | -n 设置子节点最大个数 -b 设置节点数据最大长度 |
| <code>history</code> | 列出历史记录 | |

3.3 操作实例

列出Path下的所有Znode

```
ls /
```

创建永久节点

```
create /hello world
```

创建临时节点:

```
create -e /abc 123
```

创建永久序列化节点:

```
create -s /zhangsan boy
```

创建临时序列化节点:

```
create -e -s /lisi boy
```

修改节点数据

```
set /hello zookeeper
```

删除节点, 如果要删除的节点有子Znode则无法删除

```
delete /hello
```

删除节点, 如果有子Znode则递归删除

```
rmr /abc
```

列出历史记录

```
histroy
```

3.4 节点属性

每个 znode 都包含了一系列的属性，通过命令 get，可以获得节点的属性。

```
[zk: node-22(CONNECTED) 2] get /aaa0000000001
hello22
cZxid = 0x2000000003
ctime = Fri Sep 22 16:47:35 CST 2017
mZxid = 0x2000000007
mtime = Fri Sep 22 17:26:15 CST 2017
pZxid = 0x2000000003
cversion = 0
dataVersion = 2
aclVersion = 0
ephemeralOwner = 0x0
dataLength = 7
numChildren = 0
```

dataVersion：数据版本号，每次对节点进行 set 操作，dataVersion 的值都会增加 1（即使设置的是相同的数据），可有效避免了数据更新时出现的先后顺序问题。

cversion：子节点版本号。当 znode 的子节点有变化时，cversion 的值就会增加 1。

aclVersion：ACL 的版本号。

cZxid：Znode 创建的事务 id。

mZxid：Znode 被修改的事务 id，即每次对 znode 的修改都会更新 mZxid。

- 对于 zk 来说，每次的变化都会产生一个唯一的事务 id，zxid (ZooKeeper Transaction Id)。通过 zxid，可以确定更新操作的先后顺序。例如，如果 zxid1 小于 zxid2，说明 zxid1 操作先于 zxid2 发生，zxid 对于整个 zk 都是唯一的，

ctime：节点创建时的时间戳。

mtime：节点最新一次更新发生时的时间戳。

ephemeralOwner:如果该节点为临时节点, ephemeralOwner 值表示与该节点绑定的 session id. 如果不是,ephemeralOwner 值为 0.

3.5 Zookeeper的watch机制

- 通知类似于数据库中的触发器, 对某个Znode设置 **Watcher** , 当Znode发生变化的时候, **WatchManager** 会调用对应的 **Watcher**
- 当Znode发生删除, 修改, 创建, 子节点修改的时候, 对应的 **Watcher** 会得到通知
- **Watcher** 的特点
 - **一次性触发** 一个 **Watcher** 只会被触发一次, 如果还需要继续监听, 则需要再次添加 **Watcher**
 - **事件封装**: **Watcher** 得到的事件是被封装过的, 包括三个内容 **keeperState**, **eventType**, **path**

| KeeperState | EventType | 触发条件 | 说明 |
|---------------|------------------|-------------------|------------------------------|
| | None | 连接成功 | |
| SyncConnected | NodeCreated | Znode被创建 | 此时处于连接状态 |
| SyncConnected | NodeDeleted | Znode被删除 | 此时处于连接状态 |
| SyncConnected | NodeDataChanged | Znode数据被改变 | 此时处于连接状态 |
| SyncConnected | NodeChildChanged | Znode的子Znode数据被改变 | 此时处于连接状态 |
| Disconnected | None | 客户端和服务端断开连接 | 此时客户端和服务端处于断开连接状态 |
| Expired | None | 会话超时 | 会收到一个SessionExpiredException |
| AuthFailed | None | 权限验证失败 | 会收到一个AuthFailedException |

4: zookeeper的JavaAPI操作

这里操作Zookeeper的JavaAPI使用的是一套zookeeper客户端框架 Curator，解决了很多Zookeeper客户端非常底层的细节开发工作。

Curator包含了几个包：

- **curator-framework**：对zookeeper的底层api的一些封装
- **curator-recipes**：封装了一些高级特性，如：Cache事件监听、选举、分布式锁、分布式计数器等

Maven依赖(使用curator的版本：2.12.0，对应Zookeeper的版本为：3.4.x，如果跨版本会有兼容性问题，很有可能导致节点操作失败)：

4.1、创建java工程，导入jar包

创建maven java工程，导入jar包

```
1 <!-- <repositories>
```



```
2
3     <repository>
4
5         <id>cloudera</id>
6
7         <url>https://repository.cloudera.com/artifactory/cloudera-
repos/</url>
8
9     </repository>
10
11 </repositories> -->
12
13 <dependencies>
14
15     <dependency>
16
17         <groupId>org.apache.curator</groupId>
18
19         <artifactId>curator-framework</artifactId>
20
21         <version>2.12.0</version>
22
23     </dependency>
24
25     <dependency>
26
27         <groupId>org.apache.curator</groupId>
28
29         <artifactId>curator-recipes</artifactId>
30
31         <version>2.12.0</version>
32
33     </dependency>
34
35     <dependency>
36
37         <groupId>com.google.collections</groupId>
38
39         <artifactId>google-collections</artifactId>
40
41         <version>1.0</version>
42     </dependency>
43     <dependency>
44         <groupId>junit</groupId>
```



```
45         <artifactId>junit</artifactId>
46         <version>RELEASE</version>
47     </dependency>
48     <dependency>
49         <groupId>org.slf4j</groupId>
50         <artifactId>slf4j-simple</artifactId>
51         <version>1.7.25</version>
52     </dependency>
53 </dependencies>
54
55 <build>
56
57     <plugins>
58
59         <!-- java编译插件 -->
60
61         <plugin>
62
63             <groupId>org.apache.maven.plugins</groupId>
64
65             <artifactId>maven-compiler-plugin</artifactId>
66
67             <version>3.2</version>
68
69             <configuration>
70
71                 <source>1.8</source>
72
73                 <target>1.8</target>
74
75                 <encoding>UTF-8</encoding>
76
77             </configuration>
78
79         </plugin>
80
81     </plugins>
82
83 </build>
```

4.2 节点的操作



创建永久节点

```
1  @Test
2  public void createNode() throws Exception {
3
4      RetryPolicy retryPolicy = new ExponentialBackoffRetry(1000, 1);
5      //获取客户端对象
6      CuratorFramework client =
7      CuratorFrameworkFactory.newClient("192.168.174.100:2181,192.168.174.110:2181,192.168.174.120:2181", 1000, 1000, retryPolicy);
8
9      //调用start开启客户端操作
10     client.start();
11
12     //通过create来进行创建节点，并且需要指定节点类型
13     client.create().creatingParentsIfNeeded().withMode(CreateMode.PERSISTENT).forPath("/hello3/world");
14
15     client.close();
16 }
```

创建临时节点

```
1  public void createNode2() throws Exception {
2
3      RetryPolicy retryPolicy = new ExponentialBackoffRetry(3000, 1);
4
5      CuratorFramework client =
6      CuratorFrameworkFactory.newClient("node01:2181,node02:2181,node03:2181", 3000, 3000, retryPolicy);
7
8      client.start();
9
10     client.create().creatingParentsIfNeeded().withMode(CreateMode.EPHEMERAL).forPath("/hello5/world");
11
12     Thread.sleep(5000);
13
14     client.close();
15 }
```

修改节点数据

```
1  /**
2
3      * 节点下面添加数据与修改是类似的，一个节点下面会有一个数据，新的数据会覆盖旧的数据
4
5      * @throws Exception
6
7      */
8
9      @Test
10     public void nodeData() throws Exception {
11
12         RetryPolicy retryPolicy = new ExponentialBackoffRetry(3000, 1);
13
14         CuratorFramework client =
15         CuratorFrameworkFactory.newClient("node01:2181,node02:2181,node03:2181",
16         3000, 3000, retryPolicy);
17
18         client.start();
19
20         client.setData().forPath("/hello5", "hello7".getBytes());
21
22         client.close();
23     }
```

节点数据查询

```
1
2
3  /**
4
5      * 数据查询
6
7      */
8
9      @Test
10
11     public void updateNode() throws Exception {
```



```
12
13     RetryPolicy retryPolicy = new ExponentialBackoffRetry(3000, 1);
14
15     CuratorFramework client =
16     CuratorFrameworkFactory.newClient("node01:2181,node02:2181,node03:2181",
17     3000, 3000, retryPolicy);
18
19     client.start();
20
21     byte[] forPath = client.getData().forPath("/hello5");
22
23     System.out.println(new String(forPath));
24
25     client.close();
26 }
```

节点watch机制

```
1     /**
2
3     * zookeeper的watch机制
4
5     * @throws Exception
6
7     */
8     @Test
9     public void watchNode() throws Exception {
10
11         RetryPolicy policy = new ExponentialBackoffRetry(3000, 3);
12
13         CuratorFramework client =
14         CuratorFrameworkFactory.newClient("node01:2181,node02:2181,node03:2181",
15         policy);
16
17         client.start();
18
19         // ExecutorService pool = Executors.newCachedThreadPool();
20
21         //设置节点的cache
```



```
21         TreeCache treeCache = new TreeCache(client, "/hello5");
22
23         //设置监听器和处理过程
24
25         treeCache.getListenable().addListener(new TreeCacheListener()
26     {
27
28         @Override
29
30         public void childEvent(CuratorFramework client,
31     TreeCacheEvent event) throws Exception {
32
33             ChildData data = event.getData();
34
35             if(data !=null){
36
37                 switch (event.getType()) {
38
39                     case NODE_ADDED:
40
41                         System.out.println("NODE_ADDED : "+
42     data.getPath() +" 数据:"+ new String(data.getData()));
43
44                         break;
45
46                     case NODE_REMOVED:
47
48                         System.out.println("NODE_REMOVED : "+
49     data.getPath() +" 数据:"+ new String(data.getData()));
50
51                         break;
52
53                     case NODE_UPDATED:
54
55                         System.out.println("NODE_UPDATED : "+
56     data.getPath() +" 数据:"+ new String(data.getData()));
57
58                         break;
59
60                     default:
61
62                         break;
63                 }
64             }
65         }
66     }
```



```
60
61         }
62
63         }else{
64
65             System.out.println( "data is null : "+
event.getType());
66
67         }
68
69     }
70
71     });
72
73     //开始监听
74
75     treeCache.start();
76
77     Thread.sleep(50000000);
78
79 }
```