

47号Gamer

博客园 首页 新随笔 联系 订阅 管理

随笔 - 249 文章 - 0 评论 - 1 阅读 - 74738

昵称：47号Gamer
园龄：2年10个月
粉丝：4
关注：3
+加关注

2021年12月						
日	一	二	三	四	五	六
28	29	30	1	2	3	4
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30	31	1
2	3	4	5	6	7	8

搜索

常用链接

我的随笔
我的评论
我的参与
最新评论
我的标签

随笔分类

Java8(11)
Java基础(31)
JVM(7)
MyCat(6)
MySQL(21)
Netty(5)
Spring(24)
SpringBoot(29)
zookeeper(17)
分布式(12)
其他(15)
设计模式(22)
数据结构(6)
算法(5)
线程(27)
消息中间件(10)
转载(1)

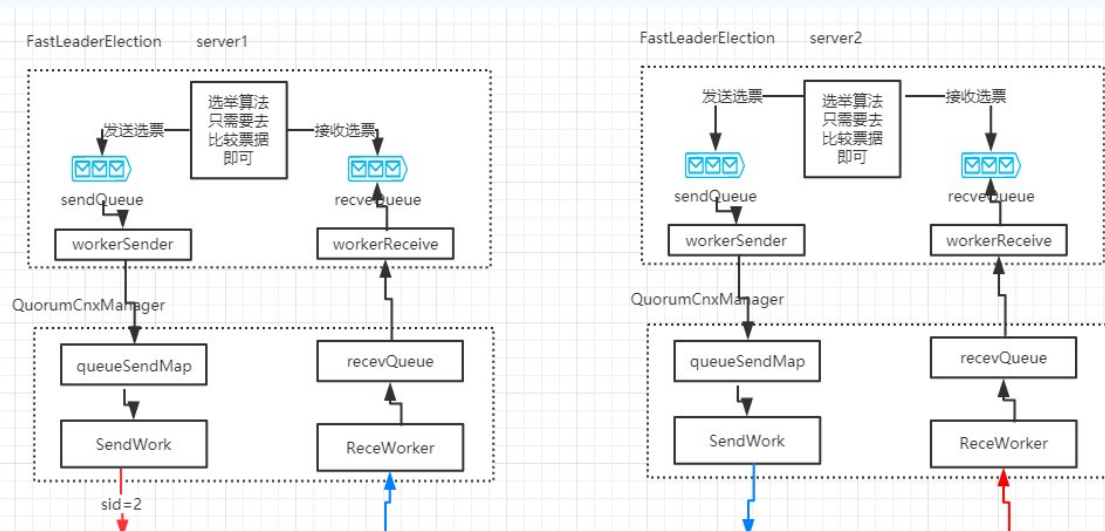
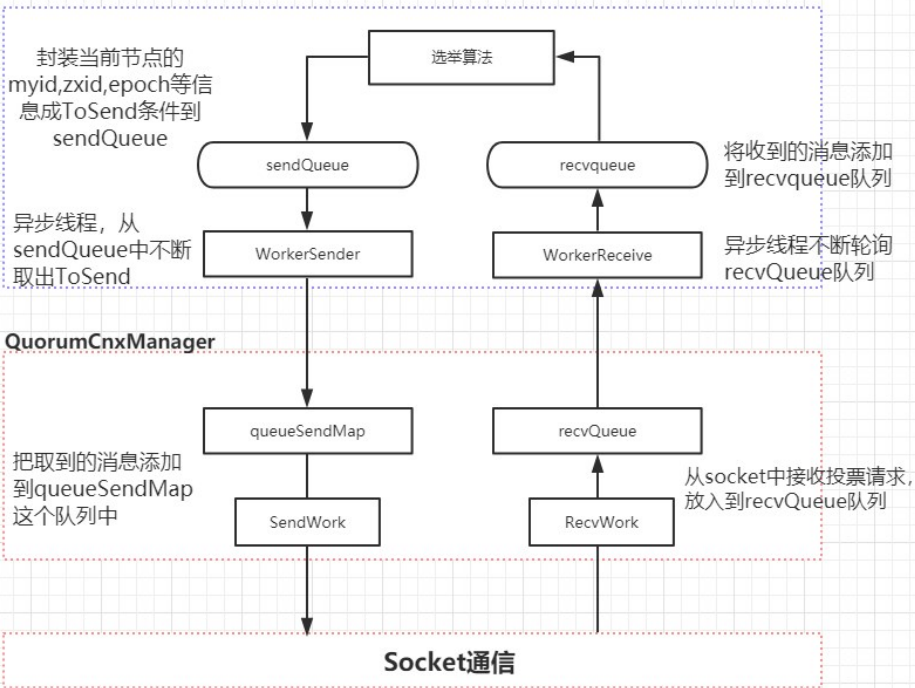
随笔档案

2021年8月(2)
2021年7月(7)
2021年4月(1)
2021年3月(3)
2021年2月(6)
2021年1月(5)
2020年12月(11)
2020年11月(17)
2020年10月(59)
2020年9月(63)
2020年8月(27)
2020年7月(13)
2020年6月(35)

zookeeper原理之投票的网络通信流程

通信流程图：

FastLeaderElection



阅读排行榜

1. java8 Stream对List<Map>的分组合并操作(5614)
2. Hash表之ASL和不成功ASL的计算 (平均查找长度)(2954)
3. Java中的CompletableFuture超时使用(2904)
4. RabbitMQ消息最终一致性解决方案 (TCC方式) (2295)
5. Java8 中map中删除元素的简单方法(2069)

评论排行榜

1. Spring IOC基于注解容器的初始化(1)

推荐排行榜

1. Kettle大量数据快速导出的解决方案 (利用SQL导出百万级数据, 挺快的) (1)
2. 内部类与静态内部类详解(1)
3. Java8的CompletableFuture在方法内使用不当, 导致局部变量出现线程安全问题(1)
4. Java8 中map中删除元素的简单方法(1)
5. ArrayBlockingQueue原理分析讲解(1)

最新评论

1. Re:Spring IOC基于注解容器的初始化
前面基于注解的初始化一直说着好好的, 小结变成了基于Xml的初始化。
--吴小破

Socket

接收数据 Notification 和发送 ToSend

ToSender	Notification
leader; 被推荐的服务器 sid zxid; 被推荐的服务器当前最新的事务 id peerEpoch; 被推荐的服务器当前所处的 epoch electionepoch; 当前服务器所处的 epoch stat 当前服务器状态 sid 接收消息的服务器 sid (myid)	leader; //被推荐的服务器 sid zxid; 被推荐的服务器最新事务 id peerEpoch; 被推荐的服务器当前所处的 epoch electionEpoch 选举服务器所处的 epoch stat; 选举服务器当前的状态 sid; 选举服务器的 sid

通信过程源码分析

每个 zk 服务启动后创建 socket 监听

```
1  protected Election createElectionAlgorithm(int electionAlgorithm){
2      //...
3      case 3:
4          qcm = createCnxnManager();
5          QuorumCnxManager.Listener listener =
6              qcm.listener;
7          if(listener != null){
8              listener.start();
9          }
10         // 启动监听listener 实现了线程, 所以在 run 方法中可以看到构建ServerSocket 的请求, 这里专门用来接收其他zkServer
11         // 的投票请求
12         // 这块后续再分析
13         @Override
14         public void run() {
15             int numRetries = 0;
16             InetAddress addr;
17             while(!isshutdown) && (numRetries < 3)){
18                 try {
19                     ss = new ServerSocket();
20                 }
21             }
22         }
23     }
```

FastLeaderElection.lookForLeader

这个方法在前面分析过, 里面会调用 sendNotifications 来发送投票请求

```
1  public Vote lookForLeader() throws InterruptedException {
2      //省略部分代码
3      sendNotifications(); //这个方法, 会把当前zk 服务器的信息添加到 sendqueue
4      /*
5       * Loop in which we exchange
6       * notifications until we find a leader
7       */
8      while ((self.getPeerState() == ServerState.LOOKING) &&
9          //省略部分代码
10     }
```

FastLeaderElection.sendqueue

sendQueue 这个队列的数据, 是通过 WorkerSender 来进行获取并发送的。而这个 WorkerSender 线程, 在构建 fastLeaderElection 时, 会启动

```
1  class WorkerSender extends ZooKeeperThread {
2      public void run() {
3          while (!stop) {
4              try { //从队列中获取 ToSend 对象
5                  ToSend m = sendqueue.poll(3000, TimeUnit.MILLISECONDS);
6                  if(m == null) continue;
7                  process(m);
```

```

8      //省略部分代码
9      void process(ToSend m) {
10         ByteBuffer requestBuffer = buildMsg(m.state.ordinal(),
11         m.leader, m.zxid,
12         m.electionEpoch,
13         m.peerEpoch);
14         managerToSend(m.sid, requestBuffer); // 这里就是调用 QuorumCnxManager
15         // 进行消息发送
16     }
17 }
18 }
19 }
20 }

```

QuorumCnxManagerToSend

```

1  public void toSend(Long sid, ByteBuffer b) {
2
3      if (this.mySid == sid) { // 如果接受者是自己, 直接放置到接收队列
4          b.position(0);
5          addToRecvQueue(new Message(b.duplicate(), sid));
6      } else { // 否则发送到对应的发送队列上
7          ArrayBlockingQueue<ByteBuffer> bq = new ArrayBlockingQueue<ByteBuffer>(SEND_CAPACITY);
8          // 判断当前的 sid 是否已经存在于发送队列, 如果是, 则直接把已经存在的数据发送出去
9          ArrayBlockingQueue<ByteBuffer> bqExisting = queueSendMap.putIfAbsent(sid, bq);
10         if (bqExisting != null) {
11             addToSendQueue(bqExisting, b);
12         } else {
13             addToSendQueue(bq, b);
14         }
15         connectOne(sid); // 连接申请调用链 connectOne-->initiateConnection-
16         // ->startConnection . startConnection 就是发送方启动入口
17     }
18 }

```

startConnection

```

1  private boolean startConnection(Socket sock, Long sid) {
2      // 省略部分代码
3      if (sid > this.mySid) {
4          // 为了防止重复建立连接, 只允许 sid 大的主动连接 sid 小的
5          closeSocket(sock);
6      } else {
7          // 构建一个发送线程和接收线程, 负责针对当前连接的数据传递, 后续的逻辑比较简单, 就不做分析
8          SendWorker sw = new SendWorker(sock, sid);
9          RecvWorker rw = new RecvWorker(sock, din, sid, sw);
10         sw.setRecv(rw);
11     }
12 }

```

SendWorker 会监听对应 sid 的阻塞队列, 启动的时候回如果队列为空时会重新发送一次最前最后的消息, 以防上一次处理是服务器异常退出, 造成上一条消息未处理成功; 然后就是不停监听队里, 发现有消息时调用send 方法RecvWorker : RecvWorker 不停监听 socket 的 inputStream, 读取消息放到消息接收队列中,消息放入队列中, qcm 的流程就完毕了。

QuorumCnxManager.Listener

listener 监听到客户端请求之后, 开始处理消息

```

1  public void run() {
2      // 省略部分代码
3      while (!shutdown) {
4          Socket client = ss.accept();
5          setSockOpts(client);
6          LOG.info("Received connection request" + client.getRemoteSocketAddress());
7          if (quorumSaslAuthEnabled) {
8              receiveConnectionAsync(client);
9          } else {
10             receiveConnection(client); // 接收客户端请求
11         }
12     }
13 }

```

```
12     }
13 }

QuorumCnxManager.receiveConnection

1 public void receiveConnection(final Socket sock) {
2     DataInputStream din = null;
3     try {
4         // 获取客户端的数据包
5         din = new DataInputStream(new BufferedInputStream(sock.getInputStream()));
6         handleConnection(sock, din);// 调用 handle 进行处理
7     } catch (IOException e) {
8         LOG.error("Exception handling connection, addr: {}, closing server connection",
9             sock.getRemoteSocketAddress());
10        closeSocket(sock);
11    }
12 }
```

handleConnection

```
1 private void handleConnection(Socket sock, DataInputStream din)throws IOException {
2     Long sid = null;
3     try {
4         //获取客户端的 sid,也就是 myid
5         sid = din.readLong();
6         if (sid < 0) {
7             sid = din.readLong();
8             if (sid < this.mySid) {
9                 //为了防止重复建立连接,只允许 sid 大的主动连接 sid 小的
10                SendWorker sw = senderWorkerMap.get(sid);
11                if (sw != null) {
12                    sw.finish();//关闭连接
13                }
14                LOG.debug("Create new connection to server: " + sid);
15                closeSocket(sock);//关闭连接
16                connectOne(sid);//向 sid 发起连接
17            } else { //同样, 构建一个 SendWorker 和RecvWorker 进行发送和接收数据
18                SendWorker sw = new
19                    SendWorker(sock, sid);
20                RecvWorker rw = new
21                    RecvWorker(sock, din, sid, sw);
22                sw.setRecv(rw);
23            }
24        }
25    }
26 }
```

分类: [zookeeper](#)

好文要顶 关注我 收藏该文



47号Gamer

关注 - 3

粉丝 - 4

+加关注

« 上一篇: [zookeeper原理之选举流程分析](#)
» 下一篇: [zookeeper原理之Leader选举完成之后的处理逻辑](#)

0

推荐

0

反对

编辑推荐：

- 带团队后的日常思考（六）
- 聊聊云原生和微服务架构
- 异或运算的巧用：不用额外的变量，如何交换两个变量的值？
- OAuth 2.1 带来了哪些变化
- 理解微前端技术原理



最新新闻：

- 夺走 20 多亿用户「唱反调」的声音，全球最大的视频网站凭什么？（2021-12-02 17:33）
 - vivo T1 体验：高能实力派，还有颜，妥妥的C位（2021-12-02 17:22）
 - 2021 互联网文学盘点：年轻人集体发疯，中年人爱凡尔赛（2021-12-02 17:11）
 - Meta 开源全新移动端AI生成神器PyTorch Live，造个AI应用5分钟，安卓iOS都支持（2021-12-02 17:00）
 - 和「咸鱼」做邻居！从虹宇宙看国内元宇宙应用的真相（2021-12-02 16:53）
- » 更多新闻...