

CSE5020 Assignment 2

1 问题描述

Chandy-Lamport 算法是一个用于记录分布式系统全局一致状态的快照算法。在异步分布式系统中，由于缺乏全局时钟，直接获取系统的全局状态是不可能的。该算法通过 marker 消息的传播，实现了在不停止系统运行的情况下捕获一个一致的全局快照。

本次实验基于普林斯顿大学 COS418 课程的作业模板，使用 Go 语言实现了完整的 Chandy-Lamport 分布式快照协议。

2 算法原理

Chandy-Lamport 算法的核心思想如下：

1. **初始化快照**: 当一个进程开始快照时，它记录自己的本地状态，并向所有出站通道发送 marker 消息。
2. **首次接收 marker**: 当进程第一次收到 marker 时，它立即记录自己的本地状态，将接收到 marker 的通道状态记录为空，并向所有出站通道发送 marker 消息。
3. **后续接收 marker**: 当进程已经记录过状态后再次收到 marker，它记录该通道从记录状态到收到 marker 之间的所有消息。
4. **完成条件**: 当进程从所有入站通道都接收到 marker 后，该进程的快照完成。

算法保证了捕获的全局状态是一致的，即不存在“孤儿消息”（在快照中出现但发送者未记录发送的消息）。

3 实现细节

3.1 数据结构设计

为了实现 Chandy-Lamport 算法，在 Server 结构体中添加了以下字段：

```

1 type Server struct {
2     Id          string
3     Tokens      int
4     sim         *Simulator
5     outboundLinks map[string]*Link
6     inboundLinks map[string]*Link
7     // 新增字段
8     snapshots    *SyncMap // snapshotId -> *SnapshotState
9     markersReceived *SyncMap // snapshotId -> map[src]bool
10 }

```

在 Simulator 结构体中添加了：

```

1 type Simulator struct {
2     time        int
3     nextSnapshotId int
4     servers     map[string]*Server
5     logger      *Logger
6     // 新增字段
7     snapshotComplete *SyncMap // snapshotId -> map[serverId]bool
8     snapshotDone    *SyncMap // snapshotId -> chan bool
9 }

```

使用 SyncMap 而非普通 map 是为了保证并发安全，避免数据竞争。

3.2 核心函数实现

3.2.1 StartSnapshot 函数

该函数在服务器上启动快照过程：

```

1 func (server *Server) StartSnapshot(snapshotId int) {
2     // 记录本地状态
3     snapshot := &SnapshotState{
4         id:      snapshotId,
5         tokens:  make(map[string]int),
6         messages: make([]*SnapshotMessage, 0),
7     }
8     snapshot.tokens[server.Id] = server.Tokens
9     server.snapshots.Store(snapshotId, snapshot)
10
11    // 初始化marker追踪

```

```

12     markers := make(map[string]bool)
13     server.markersReceived.Store(snapshotId, markers)
14
15     // 向所有出站链路发送marker
16     server.SendToNeighbors(MarkerMessage{snapshotId})
17
18     // 检查是否立即完成（无入站链路）
19     if len(server.inboundLinks) == 0 {
20         server.sim.NotifySnapshotComplete(server.Id, snapshotId)
21     }
22 }
```

3.2.2 HandlePacket 函数

该函数处理接收到的消息：

```

1 func (server *Server) HandlePacket(src string, message interface{}) {
2     switch msg := message.(type) {
3         case TokenMessage:
4             server.Tokens += msg.numTokens
5             // 记录尚未收到marker的通道上的消息
6             server.markersReceived.Range(func(key, value interface{}) {
7                 bool {
8                     snapshotId := key.(int)
9                     markers := value.(map[string]bool)
10                    if !markers[src] {
11                        snapVal, _ := server.snapshots.Load(snapshotId)
12                        snapshot := snapVal.(*SnapshotState)
13                        snapshot.messages = append(snapshot.messages,
14                            &SnapshotMessage{src, server.Id, msg})
15                    }
16                    return true
17                })
18        case MarkerMessage:
19            snapshotId := msg.snapshotId
20            if _, exists := server.markersReceived.Load(snapshotId); !
21                exists {
22                    server.StartSnapshot(snapshotId) // 首次收到marker
23                }
24            markersVal, _ := server.markersReceived.Load(snapshotId)
```

```

23     markers := markersVal.(map[string]bool)
24     markers[src] = true
25
26     // 检查是否所有入站链路都收到marker
27     if len(markers) == len(server.inboundLinks) {
28         server.sim.NotifySnapshotComplete(server.Id, snapshotId)
29     }
30 }
31 }
```

3.3 并发控制

由于分布式系统中存在多个并发的快照过程和消息传递，必须确保线程安全。本实现采用了 SyncMap 来保护共享数据结构，避免了 data race 问题。

4 测试结果

实现的算法通过了所有测试用例，包括：

- Test2NodesSimple: 两节点简单拓扑
- Test2NodesSingleMessage: 两节点单消息传递
- Test3NodesMultipleMessages: 三节点多消息传递
- Test3NodesMultipleBidirectionalMessages: 三节点双向消息传递
- Test8NodesSequentialSnapshots: 八节点顺序快照
- Test8NodesConcurrentSnapshots: 八节点并发快照
- Test10NodesDirectedEdges: 十节点有向边拓扑

5 总结

本次实验成功实现了 Chandy-Lamport 分布式快照算法，深入理解了分布式系统中全局状态捕获的原理和实现细节。通过处理并发控制、消息传递、状态记录等问题，掌握了分布式算法实现的关键技术。该算法在分布式系统的检查点、故障恢复、调试等场景中有重要应用价值。

```
PASS
ok      src/chandy-lamport      0.300s
● PS C:\Users\Administrator\Desktop\cse5020\assignments_template\assignment2\src\chandy-lamport> go test -v
== RUN Test2NodesSimple
Running test '2nodes.top', '2nodes-simple.events'
--- PASS: Test2NodesSimple (0.00s)
== RUN Test2NodesSingleMessage
Running test '2nodes.top', '2nodes-message.events'
--- PASS: Test2NodesSingleMessage (0.00s)
== RUN Test3NodesMultipleMessages
Running test '3nodes.top', '3nodes-simple.events'
--- PASS: Test3NodesMultipleMessages (0.00s)
== RUN Test3NodesMultipleBidirectionalMessages
Running test '3nodes.top', '3nodes-bidirectional-messages.events'
--- PASS: Test3NodesMultipleBidirectionalMessages (0.00s)
== RUN Test8NodesSequentialSnapshots
Running test '8nodes.top', '8nodes-sequential-snapshots.events'
--- PASS: Test8NodesSequentialSnapshots (0.00s)
== RUN Test8NodesConcurrentSnapshots
Running test '8nodes.top', '8nodes-concurrent-snapshots.events'
--- PASS: Test8NodesConcurrentSnapshots (0.00s)
== RUN Test10NodesDirectedEdges
Running test '10nodes.top', '10nodes.events'
--- PASS: Test10NodesDirectedEdges (0.00s)
PASS
ok      src/chandy-lamport      0.305s
```

图 1: 所有测试用例通过的结果截图