# Report : Groupy: a group membership service

Zhanbo Cui

October 6, 2021

## 1 Introduction

In this assignment, I implemented a small group service which involves the concept of group communication, leader election and failure detection in distributed systems. In this report, I will explain the principle and key implementation of this system. Finally, the issues mentioned in optional task also be discussed.

## 2 Main problems and solutions

### 2.1 Groupy architecture

There are many nodes in the group, each of them contain application layer, group layer and network layer. The *application layer processes* are responsible for sending or receiving state change message from *group layer processes*, meanwhile it also maintains a GUI. As for the *group layer processes*, they are responsible for forwarding message to other *group layer processes*. There are two different roles in the system: *leader node* and *slave node*, *slave node* can only communicate with *leader node*, and *leader node* can multicast message to all nodes in the group. Finally, the Erlang virtual machine can help us complete the work of the network layer fortunately.

### 2.2 Handling failure

As mentioned in the tutorial, all nodes use the Erlang built in support to detect and report that processes have crashed.

```
init(Id, Group, Master, Rnd)->
    ...
 %% this node begain to monitor the Leader
```
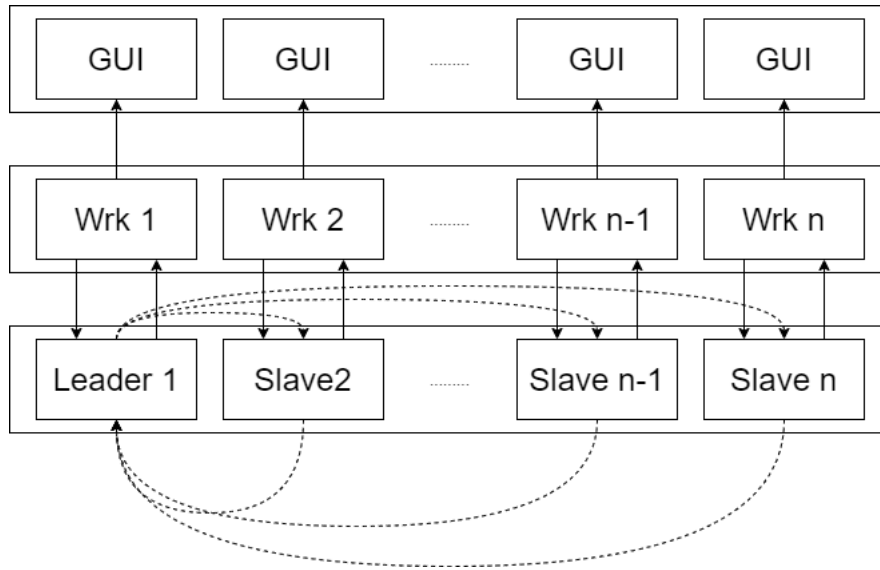
Figure 1: Groupy architecture

```
    erlang:monitor(process, Leader),
%% enter receive mode
    slave(Id, Master, Leader, N+1, Last, NewSlaves, NewGroup);
    ...
```

When a leader crashed(receive the DOWN message), the first node in the ordered slave list is selected as the new leader and multicast the new view to all slave nodes, then the *group layer processes* will deliver the new view to *application layer processes*.

## 2.3   Reliable multicast

If a leader crashes while doing a multicast, only some of the nodes will receive the message they all should have received. For more reliable multicast, the new leader need to re-multicast the message to all nodes in the group. Old leader send the message by following the order of node in the list, so the new leader as the second one in the list defiantly has received the message. New leader only need to multicast the last message to all the nodes when it was appointed to leader. Meanwhile, each message is attached with a serial number, so nodes don't care the message which had been received.

```
election(Id, Master, N, Last, Slaves, [_|NewGroup])->
    ...
```

```
%%if the first candidate is itself
   [Self|Rest]->
     %%resend the last message to all nodes in the group
       bcast(Id, Last, Rest),
             ...
     %%update the sequence number
       leader(Id, Master, N+1, Rest, NewGroup);
   ...
```

## 2.4   Optional task

We design system under the condition that message is not lost when forwarding between the nodes. It depends on the Erlang system to guarantee all of them. However, in a real world environment, messages can be lost in delivery.

### 2.4.1   Strong synchronization

We can ask an acknowledgement when the message delivered between nodes or even between *application layer processes* and *group layer processes*. If the *ack* message is not received over timeout, the message will be resent.



Figure 2: Strong synchronization

We can use this method during the connection establishment phase, because a lost message during connection establishment will cause the connection to fail.

### 2.4.2 Weak synchronization

If use *Strong synchronization* during the normal message exchange phase ensures that each node receives a consistent history of messages. But I don't think it's necessary, because we won't always receive no message, we just need to make sure that we receive new information as much as possible. In previous pattern matching in system, node receive the the message which sequence number is equal to the sequence number kept by node and than update the sequence number. However, if some message lost, the node will receive a message with a larger sequence number, at which point node should receive this new message and update its sequence number based on the sequence number of this message.

```
%%old message
    {msg, I, _} when I < N->
        slave(Id, Master, Leader, N, Last, Slaves, Group);
%%new message
    {msg, I, Msg} ->
        ...
    slave(Id, Master, Leader, I + 1, NewLast, Slaves, Group);
```

## 3 Evaluation

We start with four nodes and output their connection establishment and message passing process as far as possible.



Figure 3: Nodes boost

Figure 4: How a node join a existing group

Close the node and output the election process. The previous leader is node 1, now node 3 is elected as a leader. Node 3 Begin to "bcast".



Figure 5: Election



Figure 6: Another election

# 4 Conclusion

Through this assignment we can see that with a feasible election method we can ensure that a system based on group communication will work in the event of some errors. But group communication is a kind of indirect communication, it is very dependent on the message, so it is necessary to devise some mechanisms to ensure that message can be delivered correctly.