

CSE 586

DISTRIBUTED SYSTEMS

PROJECT PHASE 3

Authors: Sarveshwar Singhal(50418642) | UBIT - sarveshw
Abhishek Thosar (50394602) | UBIT - athosar

Date of submission: April 8, 2022

Introduction:

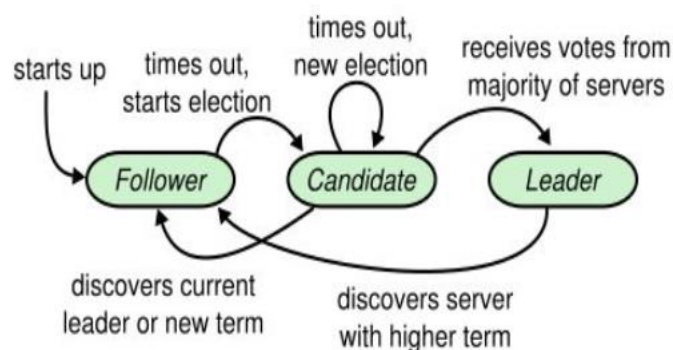
The phase 3 of the project deals with understanding and learning RAFT protocol, socket programming, server-side programming. We learned socket programming, multi-threading, asynchronous call in python to perform the Leader election using the RAFT protocol. Multiple Nodes/servers were created in Docker along with controller to create a distributed system. All the nodes are identical with respect to code.

Initially when the system starts every node is a candidate and when the timeout happens for any one of the node election starts. Post that a leader is elected. Once that leader is elected it periodically sends heartbeats to all other nodes to maintain it's leadership. In the network if any nodes fail or status of a node changes (through controller) appropriate action is triggered automatically. In case if a leader node is shutdown elections are called or appropriate action is taken based on the controller command.

The main motivation came from looking at the services of big tech giants, their systems never went down, so they must be following distributed systems approach to have data redundancy in case if a server fails, they have other servers to cater the user requests with persistent data.

Design Overview:

The below image is the design overview of our Leader election system implementing RAFT algorithm.



Implementation:

As part of the project, we created 5 nodes and 1 controller node in the docker container. To perform leader election within the nodes we implemented RAFT algorithm. All the 5 nodes have identical code. All the operations are done on server side. Initially when the system starts (i.e. we start the docker container) all the nodes are candidates as there is no leader. As soon as one of the nodes timed out it request for elections and elections starts.

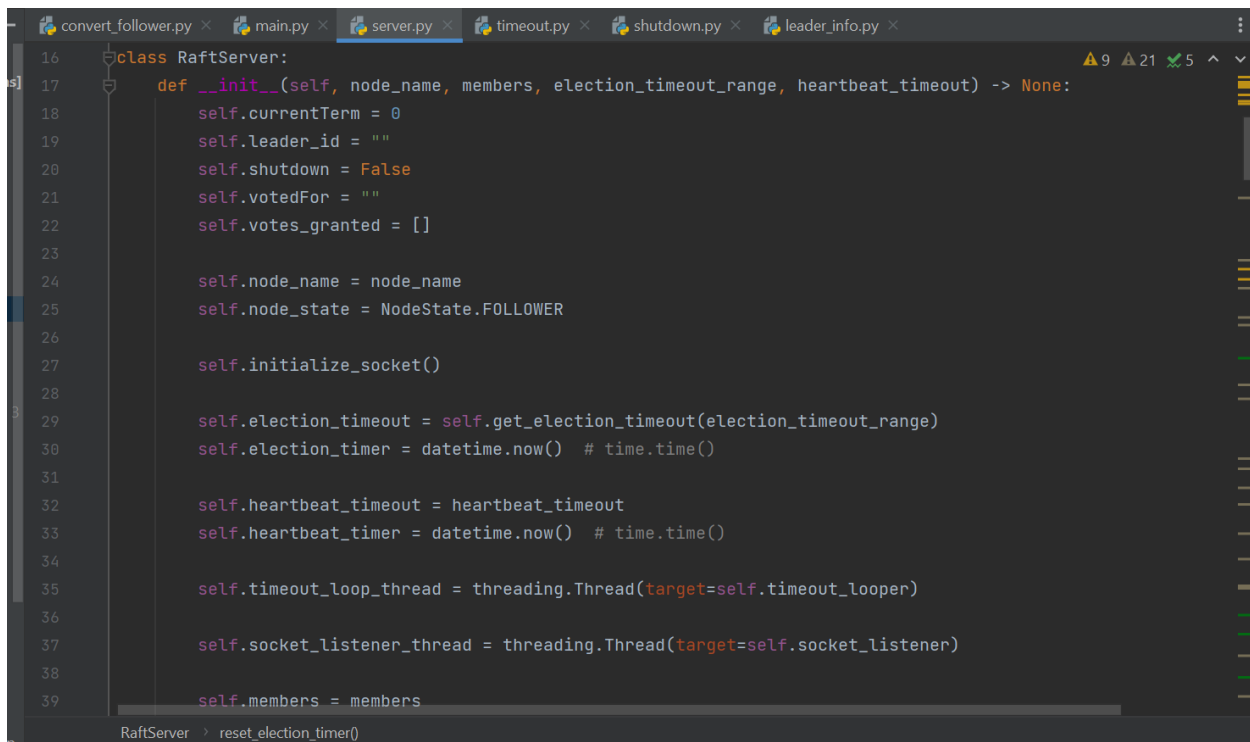
We have used python to code the server-side logic. Communication within the nodes is done with the help of web sockets. Server-side code used multi-threading to communicate within nodes asynchronously.

Classes/Methods:

1. RaftServer: This class returns a RAFT server instance.
 - a. *get_election_timeout*: This method returns the election timeout.
 - b. *socket_listener*: Thread to actively listen to socket communication.
 - c. *become_candidate*: Method to change node state to candidate and start election process by timeout.
 - d. *Become_leader*: After majority vote this method changes server state to leader.
 - e. *send_heartbeat*: This method is used to send heartbeat to all other nodes in the distributed network.
 - f. *send_append_entry*: This method is used to send append entry with data to all other nodes in our distributed system.
 - g. *handle_vote_request*: This method is used to handle the voting procedure.
 - h. *Handle_vote_request_response*: This method parses vote request response *granted or not granted*
 - i. *Handle_append_rpc*: This method handles *heartbeat* signals from leader node

- j. *convert_to_follower*: This method is used to convert a node to follower state.
 - k. *Check_majority_votes*: Verifies node majority nodes
 - l. *Set_election_term*: Updates node's current term
 - m. *reset_election_timer*: Reset election timer after starting election timeout
 - n. *Reset_heartbeat_timer*: Resets heartbeat timer clock.
2. *Controller*: This method is used to test the implementation of Leader Election.

Validations/Screenshots:



```
16 class RaftServer:
17     def __init__(self, node_name, members, election_timeout_range, heartbeat_timeout) -> None:
18         self.currentTerm = 0
19         self.leader_id = ""
20         self.shutdown = False
21         self.votedFor = ""
22         self.votes_granted = []
23
24         self.node_name = node_name
25         self.node_state = NodeState.FOLLOWER
26
27         self.initialize_socket()
28
29         self.election_timeout = self.get_election_timeout(election_timeout_range)
30         self.election_timer = datetime.now() # time.time()
31
32         self.heartbeat_timeout = heartbeat_timeout
33         self.heartbeat_timer = datetime.now() # time.time()
34
35         self.timeout_loop_thread = threading.Thread(target=self.timeout_looper)
36
37         self.socket_listener_thread = threading.Thread(target=self.socket_listener)
38
39         self.members = members
```

The screenshot shows a code editor with several tabs open: convert_follower.py, main.py, server.py (active), timeout.py, shutdown.py, and leader_info.py. The code in the active tab defines the RaftServer class. The __init__ method initializes various attributes: currentTerm, leader_id, shutdown, votedFor, votes_granted, node_name, node_state, election_timeout, election_timer, heartbeat_timeout, heartbeat_timer, timeout_loop_thread, socket_listener_thread, and members. The node_state is set to NodeState.FOLLOWER. The election_timer and heartbeat_timer are set to the current datetime. The timeout_loop_thread and socket_listener_thread are created as new threads. The RaftServer class is instantiated with the reset_election_timer() method.

The above image shows the initialization of RAFT server class, where all the variable like heartbeat_time, node_state, election_timeout etc. are set.

```

65     def socket_listener(self):
66         while True:
67             try:
68                 msg, addr = self.socket.recvfrom(1024)
69             except:
70                 print(f"NO DATA TO READ .....")
71
72             if not msg:
73                 continue
74
75             # Decoding the Message received from nodes
76             decoded_msg = json.loads(msg.decode('utf-8'))
77
78             # Check for Controller request
79             contReqType = decoded_msg['request']
80             if contReqType == 'CONVERT_FOLLOWER':
81                 self.convert_to_follower()
82                 continue
83             if contReqType == 'TIMEOUT':
84                 self.become_candidate()
85                 continue
86             if contReqType == 'SHUTDOWN':
87                 self.shutdown = True

```

The above image shows the socket_listener method which is used to listen all the incoming request on socket.

```

103     def become_candidate(self):
104         self.node_state = NodeState.CANDIDATE
105         self.reset_election_timer()
106         # increase the term
107         self.currentTerm += 1
108         self.votedFor = self.node_name
109         self.votes_granted.append(self.node_name)
110
111         # Get the votes for each member
112
113         for member in self.members:
114             # request vote
115             messageObj = RequestVoteMessage(self.node_name, RequestType.VOTE_REQUEST, self.currentTerm, self)
116             msg_bytes = json.dumps(messageObj.__dict__).encode()
117             print(f"VOTE REQUEST :: SENDING TO :: {member}")
118             try:
119                 self.socket.sendto(msg_bytes, (member, 5555))
120             except:
121                 if member in self.members:
122                     self.members.remove(member)
123                     print(f"VOTE REQUEST ERROR:: SENDING TO :: {member}")
124                     continue
125
126         self.check_majority_votes()

```

The above image shows become_candidate method which is used to send vote request to other candidates (i.e. nodes)

```

127
128 def become_leader(self):
129     self.node_state = NodeState.LEADER
130     self.reset_heartbeat_timer()
131     # send each node update about becoming leader
132     # send appendentries
133     self.send_append_entry()
134
135 def send_heartbeat(self):
136     self.reset_heartbeat_timer()
137     # send append entries request
138     self.send_append_entry()
139
140 def send_append_entry(self):
141     for member in self.members:
142         messageObj = AppendRPCMessage(self.node_name, RequestType.APPEND_RPC, self.currentTerm, self.node_
143         msg_bytes = json.dumps(messageObj.__dict__).encode()
144         try:
145             self.socket.sendto(msg_bytes, (member, 5555))
146         except:
147             if member in self.members:
148                 self.members.remove(member)
149             print(f"send_append_entry ERROR:: SENDING TO :: {member}")
150             continue

```

RaftServer > reset_election_timer()

The image shows various methods like `become_leader`, `send_heartbeat`, `send_append_entry`. These methods are used change a node to state to leader, send heartbeat to other nodes and send data(append) entry to all other nodes.

```

151
152 def initialize_socket(self):
153     self.socket = socket.socket(family=socket.AF_INET, type=socket.SOCK_DGRAM)
154     # Bind the node to sender ip and port
155     self.socket.bind((self.node_name, 5555))
156
157 def handle_vote_request(self, msg):
158     # obj = json.loads(msg)
159     requestObj = RequestVoteMessage(**msg)
160     # parse the message
161     self.set_election_term(requestObj.term)
162     # Check if node term is less than or greater than received term
163     granted = self.currentTerm <= requestObj.term
164
165     print(f"self.currentTerm < requestObj.term :: {self.currentTerm} < {requestObj.term}")
166
167     if granted: # request granted
168         self.votedFor = requestObj.candidate_id
169
170     messageObj = RequestVoteResponse(self.node_name, RequestType.VOTE_ACK, self.currentTerm, granted)
171     msg_bytes = json.dumps(messageObj.__dict__).encode()
172     self.socket.sendto(msg_bytes, (requestObj.sender_name, 5555))
173
174     # SEND:-> voting request response as granted

```

RaftServer > reset_election_timer()

The above image shows `handle_vote_request` method which is used to check the voting condition for a node. If a node will vote in favor or in against.

```

13 # Initialize
14 FORMAT = 'utf-8'
15 sender = "controller"
16 request = str(input())
17 input_target = str(input())
18 port = 5555
19 # Request
20 msg['sender_name'] = sender
21 # msg['request'] = "SHUTDOWN"#"CONVERT_FOLLOWER"
22 msg['request'] = request
23 print(f"Request Created : {msg}")
24
25 # Socket Creation and Binding
26 skt = socket.socket(family=socket.AF_INET, type=socket.SOCK_DGRAM)
27 skt.bind((sender, port))
28
29 # Send Message
30 try:
31     # Encoding and sending the message
32     skt.sendto(json.dumps(msg).encode(FORMAT), (input_target, port))
33     if request == "LEADER_INFO":
34         msg, addr = skt.recvfrom(1024)
35         res = msg.decode('utf-8')
36         print(res)

```

The above image shows the implementation of controller.

References:

<https://raft.github.io/>

<https://raft.github.io/raft.pdf>

<http://thesecretlivesofdata.com/raft/> (Highly recommended)

<https://www.youtube.com/watch?v=YbZ3zDzDnrw>

<https://docs.python.org/3/library/socket.html>

<https://docs.python.org/3/library/threading.html>