

Distributed Systems and Cloud Computing

B3 Project: Election Distributed Algorithms

Mattia Di Battista (0304938)
Università degli studi di Roma "Tor Vergata"
Rome, Italy
mattia.dibattista@alumni.uniroma2.eu

Abstract—Election distributed algorithms are a specific implementation of consensus distributed algorithms in which the aim is to find a leader. In this work, Chang and Roberts, and Bully algorithms are implemented with several services. Furthermore, we deploy the whole decentralized network on Docker containers and execute it on an AWS EC2 instance.

Index Terms—Ring-based algorithm, Bully algorithm, TCP, Docker, AWS EC2, Ansible

I. INTRODUCTION

This work aims to implement two election distributed algorithms, make a decentralized network of nodes using Docker containers and execute the whole application on an AWS EC2 instance. In the following, we describe the services used by the algorithms, the implementation of the algorithms (i.e., **Chang and Roberts** and **Bully** algorithms), the whole architecture where to perform them, and three tests used to demonstrate the operation.

II. SERVICES

A. Register Service

The **register service** provides knowledge of the entire network to all nodes belonging to it. Moreover, it generates a unique random ID for each member of the topology (see VII).

The register node listens by default on the **TCP** port 1234 (it can be changed from *SDCC/sdcc/config.json*). The listening window is kept open for a *SOCKET_TIMEOUT* period (defined in *SDCC/sdcc/register/src/constants.py*) after which a node receives the member's network list.

Initially, every node has an ID equal to *DEFAULT_ID* (by default set to -1). Instead, the register is identified by value 0. After the register phase, a unique ID to each node is associated.

A node generates multiple sockets for whole its activity. During the register phase, two sockets are created: the first used to communicate with the register node and the second one used later to listen to eventual packets from other nodes. IP address, port, and ID of this latest must be sent before its use to make known to other nodes that information.

B. Heartbeat

The **heartbeat service** allows detaching crashes by the coordinator nodes.

The service is kept active using a thread that sends heartbeat messages to the leader through a dedicated **TCP** socket (different from the two defined in II-A). After the heartbeat

message the thread waits for a while (i.e., *TOTAL_DELAY* period defined in *SDCC/sdcc/node/src/constants.py*). If the timeout occurs a crash is found out and consequently a new election is started. Only the ACK message indicates a running coordinator.

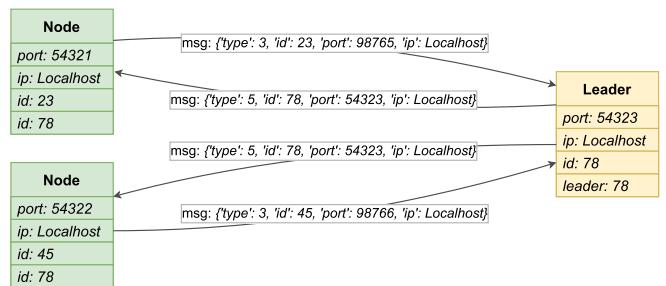


Fig. 1. Heartbeat service invoked by two nodes.

Heartbeat messages are sent periodically based on the *HEARTBEAT_TIME* constant. If the received message has a different type from HEARTBEAT or the sender is not the current coordinator the packet is ignored and the thread keeps listening for the remaining time.

C. Verbose

The **verbose flag** shows all messages exchanged (i.e., received, sent) throughout the node lifetime¹ (see VII). To activate the verbose flag is necessary to pass the -v parameter from the command line (see VI).

The main info shown are:

- Timestamp
- Node info (i.e., IP address, port number, ID)
- Receiver/Sender info (i.e., IP address, port number, ID)
- Message content

The **Logging** library is used to define the message syntax.

```
def set_logging() -> logging:  
    logging.basicConfig(  
        level=logging.DEBUG,  
        format="[%(levelname)s]  
              %(asctime)s\n%(message)s",  
        datefmt='%b-%d-%y %I:%M:%S'  
    )  
    return logging
```

¹Also register node implements this service.

Method to instantiate a logging object.

D. Delay

To stress more the system, the **delay** method is introduced in *SDCC/sdcc/node/src/helper.py*. It generates a period waited by the sender to forward the current packet. This may cause the receiver timeout to expire.

```
def delay(flag: bool, ub: int):
    if flag:
        delay = randint(0, floor(ub*1.5))
        time.sleep(delay)
```

It is activated by default when tests are performed (see. IV), but is also configured by the command line through the *-d* flag.

III. ALGORITHM IMPLEMENTATION

What follows does not describe how the algorithms work, but only how particular aspects are implemented (reference at [1]).

The implementation consists of the abstract class² and election distributed algorithms classes that extend the first one.

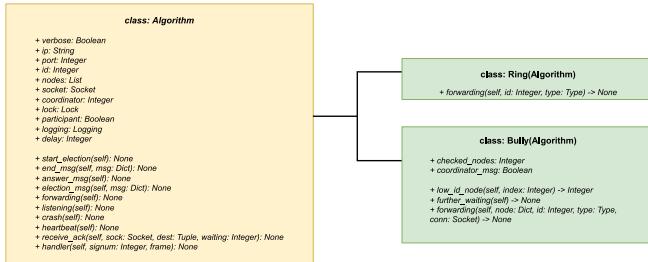


Fig. 2. Logic implementation of the classes.

Six types of messages can be exchanged between nodes:

```
class Type(Enum):
    ELECTION = 0
    END = 1
    ANSWER = 2
    HEARTBEAT = 3
    REGISTER = 4
    ACK = 5
```

ANSWER type is used only by the **Bully algorithm**, whereas REGISTER is sent during register phase.

Both algorithms begin with run the listening thread after which they start an election. Only after completing these two phases, the heartbeat can start.

```
def __init__(self, ...):
    ...
    self.lock = Lock()

    thread = Thread(target=self.listening)
    thread.daemon = True
    thread.start()
```

²Defined in *SDCC/sdcc/node/src/Algorithm.py*

```
self.start_election()
Algorithm.heartbeat(self)
```

Many data are accessed simultaneously from multiple threads, thus a **lock** is defined to manage shared resources.

```
self.lock.acquire()
if self.participant or (self.coordinator in
    [self.id, DEFAULT_ID]):
    self.lock.release()
    continue
```

Example of **lock** management in heartbeat method.

A. Chang and Roberts Algorithm

The algorithm is suitable for a collection of processes arranged in a logical ring. Each process has a communication channel to the next one in the ring. All messages are sent clockwise around the ring. The ID's node is used to define the **ring network**: the next node is the one with the greatest ID than current and the lowest among others.

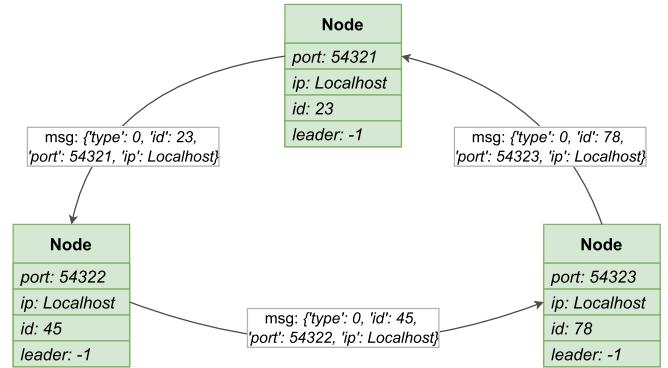


Fig. 3. Election started by node with *id*=23 in Ring topology.

When a leader crash occurs or the node's timer associated with a heartbeat message elapses, the leader is removed from the nodes list, therefore remaining nodes can not interact with it even if it is still active.

B. Bully Algorithm

With different respect to the **Ring-based algorithm**, the **Bully algorithm** assumes that each process knows which processes have higher identifiers and that it can communicate with all such processes. That information is sent by the register node (as described in II-A).

1) **Scenarios:** Unlike in the III-A, for the **Bully algorithm**, the current leader is not removed from the nodes list during leader crash or timeout scenarios.

The considered cases when the leader is executing are:

1) If the leader delays sending the ACK packet, other nodes start a new election that will produce the next coordinator also if the previous one is running³.

³In the next election it will be elected again.

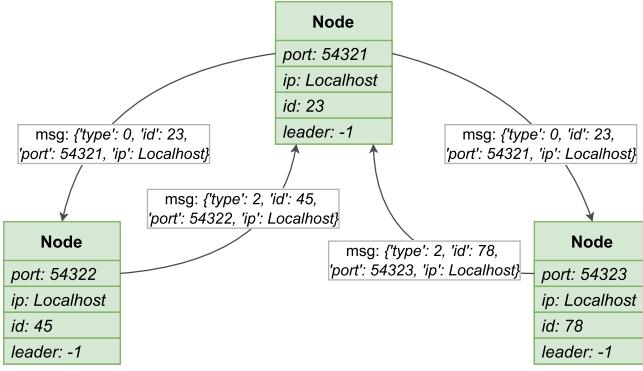


Fig. 4. Election started by node with $id=23$ using **Bully algorithm**.

- 2) If the leader is stopped⁴ a new election will start. In the meanwhile all messages sent to the sleeping node are queued, so when it wakes up⁵ will receive those messages and begin the leader again.

IV. TESTS

The following tests are performed:

- 1) *Test A*: a generic node fails
- 2) *Test B*: the coordinator fails
- 3) *Test C*: both generic and coordinator nodes fail

To interrupt a specific node **psutil** library is used. It kills a process that is listening on a certain **TCP** port⁶. That mechanism exploits the sorted list of nodes send by the register. E.g., the fact that the coordinator node occupies the last position in the list is used in *Test B*.

```

def kill_node(self, port: int):
    for proc in process_iter():
        for conn in
            proc.connections(kind='inet'):
                if conn.laddr.port == port:
                    proc.send_signal(signal.SIGINT)

```

The network is keep running for a while before the test and after the interruption to show its activity.

Test execution is interactive (see VII) means that the user chooses which test executes sees logging information on the terminal and sets which algorithm has to be performed.

As described in II-D to stress more the application, the delay option is active.

V. DEPLOYMENT

The network is deployed on an **AWS EC2** instance where every node runs on a **Docker** container. **Docker Compose** is used to automate the container's creation⁷.

By default, **Docker Compose** creates a network where containers can communicate with each other, so that is used.

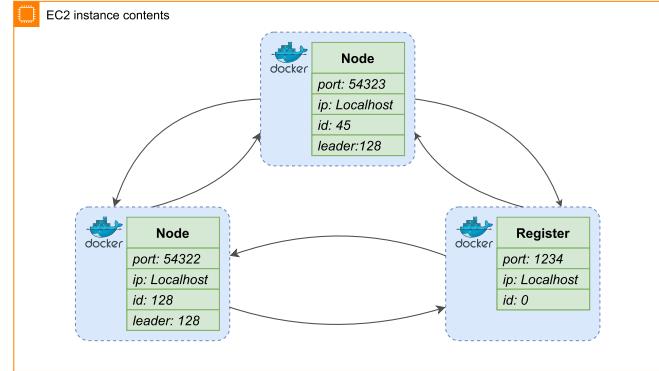


Fig. 5. Deployment using **AWS EC2** instance and **Docker** containers.

To automate the deployment procedure **Ansible** is used to install **Docker** and to forward the code application on an **EC2** instance. See VII.

VI. HOW TO USE

The application can be run in two different ways:

- 1) Local execution without **Docker** containers
- 2) Remote execution on **AWS EC2** instance using **Docker** containers⁸(as shown in V)

The complete list of commands is available here.

VII. RUNNING EXAMPLES

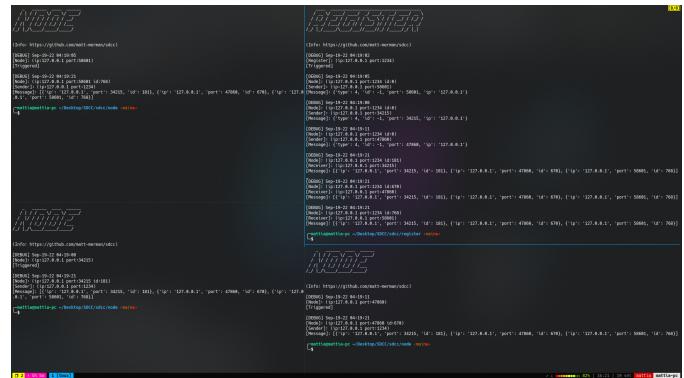


Fig. 6. Register phase from three generic nodes and register node.

REFERENCES

- [1] George Coulouris, Jean Dollimore, and Tim Kindberg. *Distributed Systems: Concepts and Design (International Computer Science)*. 4th rev. ed. Addison-Wesley Longman, Amsterdam, 2005. ISBN: 0321263545.

⁴Using **Ctrl + z** combination.

⁵Using **fg** command.

⁶It requires root privileges (see VI).

⁷See SDCC/sdcc/docker-compose.yml, SDCC/sdcc/node/Dockerfile and SDCC/sdcc/register/Dockerfile.

⁸That execution requires an **AWS** account.

```

[14/14]

(Info: https://github.com/matt-merman/sdcc)

[DEBUG] Sep-27-22 03:30:20
[Register]: (ip:register port:1234)
[Triggered]

[DEBUG] Sep-27-22 03:30:24
[Node]: (ip:register port:1234 id:8)
[Sender]: (ip:127.0.0.1 port:41477)
[Message]: {'type': 4, 'id': -1, 'port': '0.0.0.0'}

[DEBUG] Sep-27-22 03:30:26
[Node]: (ip:register port:1234 id:8)
[Sender]: (ip:127.0.0.1 port:37865)
[Message]: {'type': 4, 'id': -1, 'port': '0.0.0.0'}

[DEBUG] Sep-27-22 03:30:27
[Node]: (ip:register port:1234 id:8)
[Sender]: (ip:127.0.0.1 port:36679)
[Message]: {'type': 4, 'id': -1, 'port': '0.0.0.0'}

[DEBUG] Sep-27-22 03:30:37
[Node]: (ip:register port:1234 id:16)
[Sender]: (ip:127.0.0.8 port:45587)
[Message]: [{"ip": "0.0.0.0", "port": 45587, "id": 16}, {"ip": "0.0.0.0", "port": 59713, "id": 521}, {"ip": "0.0.0.0", "port": 59475, "id": 987}]

0 1 0m 24m 1 [emux]

```

Fig. 7. Example of the message showed by register node.

```

(Info: https://github.com/matt-merman/sdcc)

1. Node Failure
2. Coordinator Failure
3. Two nodes Fall (Include the Coordinator)
4. Change Algorithm (Bully by default)
5. Exit
(NOTICE: 4 nodes except the register is spawning)

0 1 0d 3h 2m 1 sudo

```

Fig. 8. User interface to tests execution

```

ubuntu@ip-172-31-23-55:~$ ll
total 10
drwxr-x--- 9 ubuntu ubuntu 4096 Sep 27 07:32 .
drwxr-xr-x 3 root  root  4096 Sep 27 07:18 ..
drwxr-xr-x 10 ubuntu ubuntu 4096 Sep 27 07:32 ansible/
-rw-r--r-- 1 ubuntu ubuntu 228 Sep 27 07:32 .bash_logout
-rw-r--r-- 1 ubuntu ubuntu 3771 Sep 27 07:32 .bashrc
drwxr-xr-x 2 ubuntu ubuntu 4096 Sep 27 07:26 .cache/
drwxr-xr-x 2 ubuntu ubuntu 4096 Sep 27 07:26 .config/
drwxr-xr-x 2 ubuntu ubuntu 4096 Sep 27 07:18 .ssh/
-rw-r--r-- 1 ubuntu ubuntu 1  Sep 27 07:26 sudo_as_admin.successful
drwxr-xr-x 2 ubuntu ubuntu 4096 Sep 27 07:32 test/
drwxrwxr-x 2 ubuntu ubuntu 4096 Sep 27 07:31 ansible/
-rw-r--r-- 1 ubuntu ubuntu 112 Sep 27 07:27 config.json
-rw-r--r-- 1 ubuntu ubuntu 184 Sep 27 07:27 docker-compose.yml
drwxrwxr-x 3 ubuntu ubuntu 4096 Sep 27 07:28 node/
drwxrwxr-x 3 ubuntu ubuntu 4096 Sep 27 07:30 register/
-rw-r--r-- 1 ubuntu ubuntu 1  Sep 27 07:30 requirements.txt
-rw-r--r-- 1 ubuntu ubuntu 184 Sep 27 07:27 tests.py
-rw-rw-r-- 1 ubuntu ubuntu 3596 Sep 27 07:28 tests.py
ubuntu@ip-172-31-23-55:~$ docker --version
Docker version 20.10.12, build 39ca2be
ubuntu@ip-172-31-23-55:~$ docker-compose --version
docker-compose version 1.29.2, build unknown
ubuntu@ip-172-31-23-55:~$ _

0 1 0m 20m 1 ssh

```

Fig. 9. Docker, Docker Compose and application code info on an AWS EC2 instance.