

ASSIGNMENT № 2

Akshat Khare, IIT Delhi

14/04/2020

Introduction to GHS

There exists centralized algorithms to compute minimum spanning trees, example: Kruskal, Prim's algorithm, etc. In 1983, Gallager, Humblet and Spira proposed an algorithm which computes minimum spanning tree from a distributed system perspective. MST is required in many applications involving networking overlay and other distributed applications. GHS calculates MST in at max $5N\log_2 N + 2E$ messages without every node just knowing weight of edges it is connected to. It is proved that if the edge weights are unique MST is unique [1]. If the edges are not unique we can establish a trichotomy by enforcing an order based on the unique ids of the nodes across the edge.

GHS works by merging fragments till there is a single connected fragment. It is also noteworthy that once edge is designated as a *Branch* edge it stays that way till end. Hence the algorithm proceeds progressively marking required edges as Branch edges till MST is found. More information about the algorithm can be found in original paper[1]. I also referred Guy Flysher and Amir Rubinshtein's explanation[2] of GHS because original paper failed to explain many concepts in required brevity.

Details of Implementation

I used *python3* for my implementation. I made use of *multiprocessing* module for ipc procedures and communications. Since communication channels need to be multiple producer and single consumer (a node will receive message from it's neighbouring nodes), I made use of *queues* instead of *pipes* as pipes are single producer single consumer.

I made use of **blocking receive calls** as node essentially is triggered by incoming messages and there is no background task to be done which will get affected by process going to sleep. Since number of processes (number of nodes) is at max 100. I used *multiprocessing.Process* class as I tested that 100 processes can be handled by Linux OS. If the number of processes would be higher I would have made use of *multiprocessing.Pool* as then OS would not be able to handle that high number of processes.

At last the main process which parsed the input file and spawned N process sends a message to each process asking about their branch edges. Main process compiles the results and prints on the console in the format required. This additional process which was essential to assignment involved $2N$ messages and didn't affect the message complexity of $5N\log_2 N + 2E$.

Packages **numpy**, **multiprocessing** and **sys** need to be installed.

The skeleton of Node class is:

Listing 1: Node Class.

```

1 class Node:
2     """docstring for Node"""
3     def __init__(self, infoStart):
4         self.uid = infoStart.uid #unique identifier
5         self.edges = infoStart.edges #list of tuples of the form (nodeId, ↵
            weight)
6         self.edgeToWeight = {} #a dictionary for the fast retrieval of ↵
            edgeWeight from NodeId
7         for i,j in self.edges:
8             self.edgeToWeight[i] = j
9         self.queues = infoStart.queues #a dictionary of the messageQueue ↵
            of the form (nodeId, queue)
10        self.queue = infoStart.queue #queue of the self
11        self.masterQueue = infoStart.masterQueue #queue of the master ↵
            process which spawned all processes
12        self.SN = "Sleeping" #State of the Node
13        self.SE = {} #Dictionary of the form (nodeId, status of edge) used↵
            to get the status of edges
14        self.test_edge = None
15        for i,j in self.edges:
16            self.SE[i] = "Basic"

```

About Listing 1... Each node maintains its own set of variables, consisting of its state (denoted by SN and assuming possible values Sleeping, Find, and Found) and the state of the adjacent edges. The state of edge j is denoted by $SE(j)$ and can assume the possible values Basic, Branch, and Rejected. It is possible for the edge states at the two nodes adjacent to the edge to be temporarily inconsistent[1]. Initially for each node, $SN = \text{Sleeping}$ and $SE(j) = \text{Basic}$ for each adjacent edge j . Each node also maintains a fragment identity FN, a level LN, and variables best-edge, best-wt, test-edge, in-branch, and find-count, all of whose initial values are immaterial.[1]

Listing 2: Node's knowledge of the graph

```

1 class InfoStart:
2     """docstring for InfoStart"""
3     def __init__(self, uid, edges, queues, queue, masterQueue):
4         self.uid = uid #Unique id
5         self.edges = edges #neighbouring edge and their weights
6         self.queues = queues #neighbouring edges queues
7         self.queue = queue #nodes own queue
8         self.masterQueue = masterQueue #master process queue

```

About Listing 2... A node is spawned with only above information.

Listing 3: Master spawning nodes and waking up

```
1 if __name__ == '__main__':
2     ....
3     nodesQueues = [multiprocessing.Queue() for i in range(numNodes)] #↔
4     masterQueue = multiprocessing.Queue() #master's own queue
5     processes = []
6     for i in range(numNodes):
7         queuedic = {} #node's relevent queue data
8         for j,k in adjacencyList[i]:
9             queuedic[j] = nodesQueues[j] #only providing neighbouring ↔
10            edges queues
11            infoStart = InfoStart(i, adjacencyList[i], queuedic, nodesQueues[i↔
12            ], masterQueue)
13            p = multiprocessing.Process(target=nodecode, args=(infoStart,))
14            p.start()
15            processes.append(p)
16            nodesQueues[0].put(Message("wakeup", [], -1))#waking up nodeId 0
```

About Listing 3... Main function (process) spawns node processes with relevent information embedded in *Infostart* class. It wakes up 0th node.

Listing 4: Function executed by every Node

```
1 def nodecode(infoStart):
2     node = Node(infoStart)
3     node.receiveAndProcess()
```

About Listing 4... Node initiates its relevent information table from infostart which is a instance of *Infostart* class. Internals of receiveAndProcess are described below. It does a blocking receive message call and processes the messages endlessly until it halts itself or receive a halting reply from master.

Listing 5: Receive Message and Process

```
1 class Node:
2     ....
3     def receiveAndProcess(self):
4         while(True):
5             message = self.queue.get() #blocking receive
```

```

6         self.processMessage(message)
7     def processMessage(self, message):
8         typemessage = message.typemessage
9         senderid = message.senderid
10        metadata = message.metadata
11        if DEBUG: print("Process ", self.uid, "received ", typemessage, "↔
            from " , senderid, "with metadata ", metadata)
12        if typemessage=="wakeup":
13            self.wakeup()
14        elif typemessage=="connect":
15            self.connect(metadata[0], senderid, message)
16        elif typemessage=="initiate":
17            self.initiate(metadata[0], metadata[1], metadata[2], senderid)
18        elif typemessage=="test":
19            self.testResponse(metadata[0], metadata[1], senderid, message)
20        elif typemessage=="accept":
21            self.accept(senderid)
22        elif typemessage=="reject":
23            self.reject(senderid)
24        elif typemessage=="report":
25            self.reportResponse(metadata[0], senderid, message)
26        elif typemessage=="changeRoot":
27            self.changeRootResponse()
28        elif typemessage=="queryStatus":
29            self.queryStatusResponse()
30        else:
31            if DEBUG: print("Unrecognised message")

```

About Listing 5... Handling of the messages is as described by paper[1]. Message has *metadata* which consists of extra data bundled up with message.

Listing 6: Message Class

```

1 class Message():
2     """docstring for Message"""
3     def __init__(self, typemessage, metadata, senderid):
4         self.typemessage = typemessage #Type of message like wakeup, ↔
            connect, etc
5         self.metadata = metadata #Any extra information bundled up with ↔
            data
6         self.senderid = senderid #Id of whoever sent the message

```

About Listing 6... This message is passed via *multiprocessing.Queue*.

Listing 7: Core nodes inform main process that MST is computed in last elif

```

1 class Node:

```

```

2     ....
3     def reportResponse(self, weightparam, senderEdge, message):
4         ....
5         elif weightparam == self.bestWeight and self.bestWeight == float('↵
            inf'):
6             self.masterQueue.put(Message("done", [self.SE, self.inBranch],↵
                self.uid)) #GHS is computed, inform master

```

About Listing 7... When the bestWeight received by core nodes is received and it is infinite, core nodes inform main process that MST is computed.

Listing 8: Main process gathers MST information

```

1 if __name__ == '__main__':::
2     ....
3     recvmesssage = masterQueue.get()
4     if recvmesssage.typemessage=="done":
5         for i in range(numNodes):
6             nodesQueues[i].put(Message("queryStatus",[], -1)) #get Branch ↵
                edges information
7 class Node:
8     ....
9     def queryStatusResponse(self):
10         self.masterQueue.put(Message("queryAnswer", [self.SE, self.↵
            inBranch], self.uid))
11         sys.exit() #MST computed and sent information to master, time to ↵
            exit

```

About Listing 8... Main process gathers information about branch edges and prints the output according to specification.

Running time analysis

Upper limit of messages is $5N\log_2 N + 2E$.

How to run the code

Listing 9: Execute this on console

```

1 #! /bin/bash
2 pip install multiprocessing numpy
3 python main.py [name_of_input_file]

```

messages and upperlimit vs edges (nodes=100)

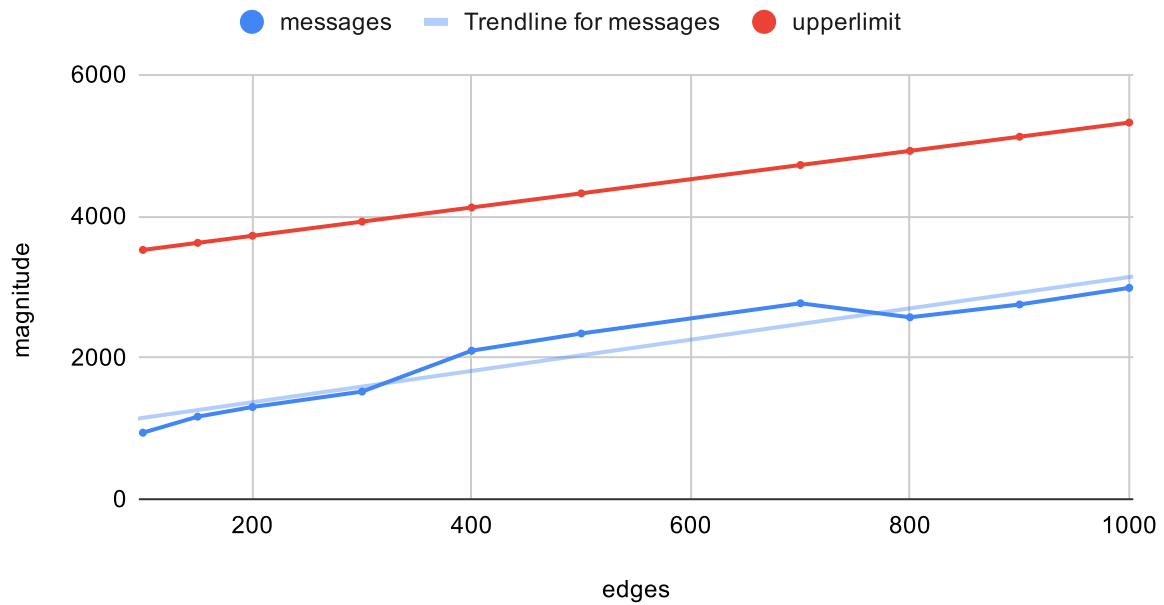


Figure 1: Keeping nodes constant we see that number of messages indeed increase linearly.

user+sys vs edges (nodes = 100)

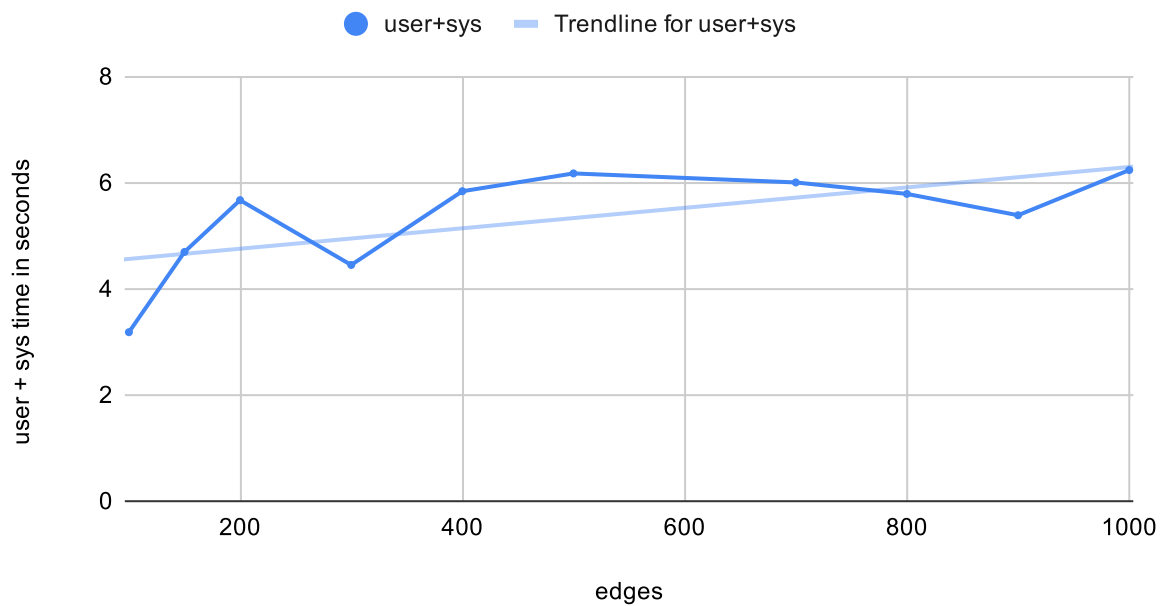


Figure 2: Time does not exactly follow a linear trend. This might be because of my code being executed on virtual machine (virtualbox)

messages and upperlimit vs nodes(log(nodes)) (edges=100)

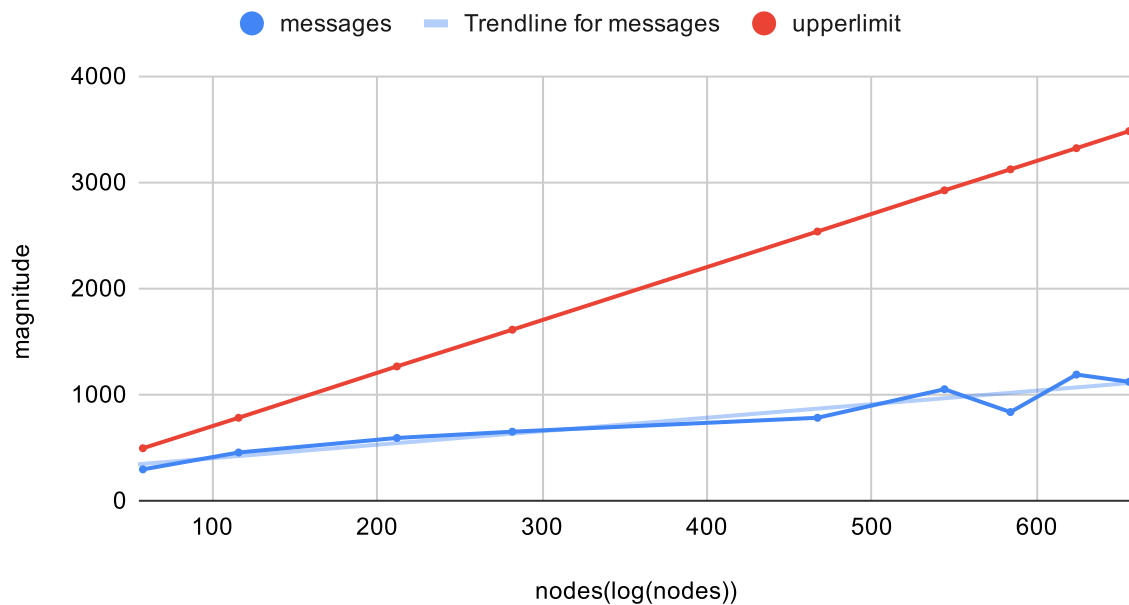


Figure 3: Number of messages is indeed linear with $N \log_2 N$

user+sys vs. n

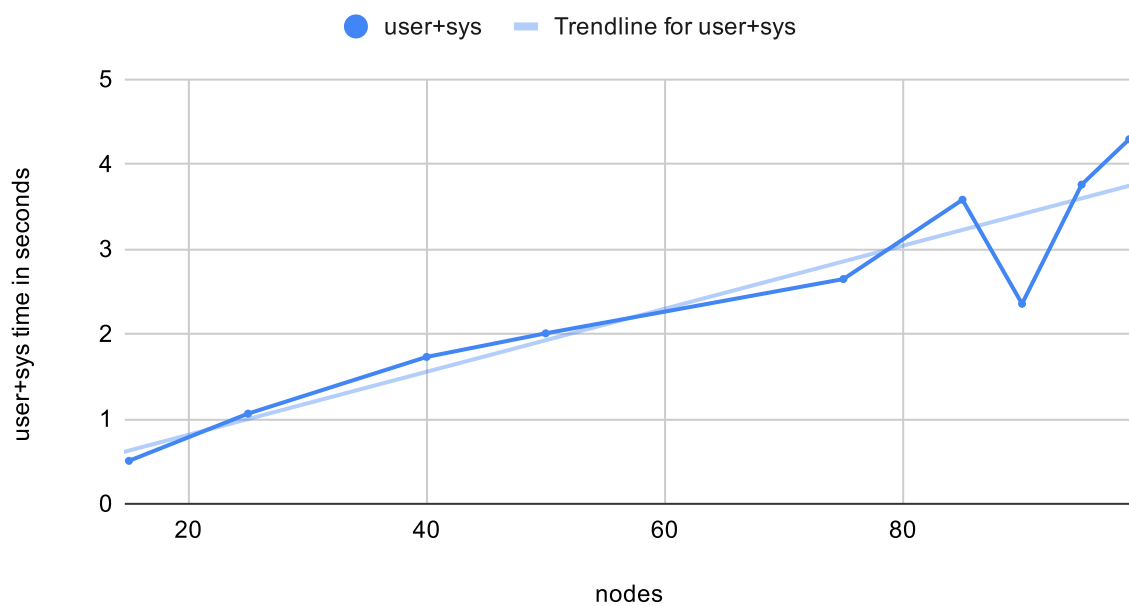


Figure 4: Time does not exactly follow a linear trend. This might be because of my code being executed on virtual machine

user+sys time vs messages

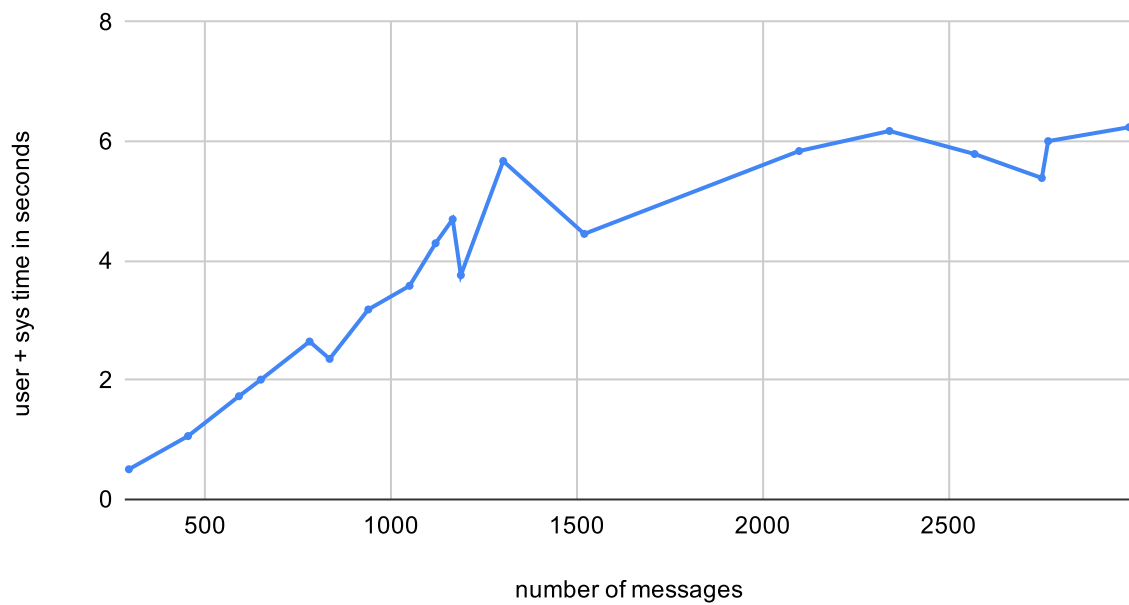


Figure 5: Time does not exactly follow a linear trend. This might be because of my code being executed on virtual machine

About Listing 9... My code requires *multiprocessing* and *numpy* packages to work. Please install them and then simply run *main.py* with inputfile. The output in desired format will be printed in console.

References

- [1] R. G. Gallager, P. A. Humblet and P. M. Spira. *A distributed algorithm for minimum weight spanning trees*. Massachusetts Institute of Technology, 1983.
- [2] Guy Flysher and Amir Rubinshtein adaptation. *A distributed algorithm for minimum weight spanning trees*.