



UNIVERSITÀ DI PISA

Distributed Systems and Middleware Technologies

## **Distributed File System with Pastry**

**Iacopo Canetta**

**Francesco Berti**

---

ANNO ACCADEMICO 2024/2025

# Contents

# 1 Introduction

This project consists in the implementation of a peer-to-peer distributed file storage using a distributed hash table for locating the files among the nodes. The distributed hash table is based on Pastry protocol.

## 2 Pastry

Pastry implements a logical ring of nodes, and each node has an ID. The files are hashed to get a key, and the key is used to identify which node is responsible for that file, checking which one has the longest common prefix. When a file is accessed or stored, a request is sent to the current node. When a node receives a request, it uses a routing table, structured as a tree, for finding the next hop for the request. Every time a row of the routing table is accessed,  $b$  bits of the prefix of the key are used to identify the correct column. If the  $b$  bits of the prefix are owned by the current node, the next row is visited, otherwise the address of the next hop is retrieved.

Whenever an address is not found or the whole table is crossed, an extra table is checked, called leaf set and containing the nodes with id close to the current node; if no node in the leaf set has a longer shared prefix with the key than the current node, then the current node is responsible for the file. This routing table, using  $N$  nodes, guarantees a maximum number of hops equal to  $\log n$ .

## 3 Implementation

The implementation of the system is divided in two parts: one for the distributed file system, created in Erlang, and a webserver realized in Java, also thanks to the Frameworks and implementations of the Java Servlets.

### 3.1 Erlang

The Erlang code is divided in multiple modules, in order to isolate functions and keep them short enough.

#### 3.1.1 Util.erl, File\_handler.erl, Network.erl

These files are created for having simple utility functions. *File\_handler.erl* contains the utility functions for dealing with files, *Network.erl* contains simple functions for starting the Erlang nodes and setting the cookies, in order to receive messages from different machines, and *Util.erl* contains generic utility functions.

#### 3.1.2 Key\_gen.erl

*Key\_gen.erl* is the file that implements the key system for Pastry. The keys are generated hashing strings and represented as *bitstrings*; this type is chosen because it can be used for representing streams of any number of bits, differently from the *binary* type, which is exclusively for multiples of 8 bits.

*Bitstrings*, however, have little documentation and almost no function dedicated, so *Key\_gen.erl* also implements the functions for comparing these numbers and handling them as hexadecimal, since Pastry requires groups of 4 bits for the routing.

#### 3.1.3 Routing.erl

*Routing.erl* represents the core of Pastry protocol, implementing the Routing Table. The routing table is represented as a Tuple containing the TableKey, so the key of the node owning the table, and the table. The table is a list of 32 lists, representing the rows; the single entries of the rows are represented by a tuple containing the hexadecimal identifying the column and the node information.

During the routing or any table modification, instead of doing the whole routing, this approach requires of detecting the number of hexadecimal in common in the prefix, which represents the index of the row, and the first hexadecimal not in common, which represent the Column key.

This module contains also the functions for modifying the table and some utility functions relative to it.

#### 3.1.4 Leaf\_set.erl

*Leaf\_set.erl* implements the Leaf Set; the Leaf Set is represented as a Tuple of two lists, *Left* and *Right*, representing the two halves of it. The two list contain at most *L2* elements, represented as tuples key-node information.

The module also offers functions for modifying the leaf set and some utility relative to it.

#### 3.1.5 Node\_actions.erl

*Node\_actions.erl* implements some utility functions used by the upper modules to interact with files, sending messages, or broadcasting messages to lists, routing tables, or leaf sets.

The most important function is *full\_route()*, which implements the full navigation of the routing table, and, in case of no match or it reaches the bottom, the leaf set is analyzed; the function returns the information of the routed node, or *route\_end* if the current node is responsible for the key in input.

#### 3.1.6 Pastry\_actions.erl

*Pastry\_actions.erl* implements all the functions for handling the pastry requests and responses to the other nodes in the network, plus the handling of the self messages for peridoc events. The main functions are:

- *join()*: function used for responding to the *join* message with the routing table row, and for forwarding the message to the next node of the routing; function used by the node to join the network through another known node;
- *join\_res\_handle()*: function used for updating the routing table and leaf set after a *join\_response*;
- *exit()*: function for handling the exit of another node from the network;
- *keepalive()*: sends keepalive to all known nodes;
- *update\_keepalive()*: updates the keepalive of given nodes;
- *keepalive\_res()*: responses to the keepalive to make sure that, even if a node knows another but not viceversa, the one knowing still receives keepalives from the other;
- *check\_expired\_nodes()*: checks which nodes are disconnected and updates the tables;
- *share\_info()*: shares the info about the nodes with the known nodes;

#### 3.1.7 Backup\_actions.erl

*Backup\_actions.erl* implements all the functions for sending and handling the backups. For design choice, the backups are kept only by the nodes in the leaf<sub>s</sub>et.

- *join()*: function used for responding to the *join* message with the routing table row, and for forwarding the message to the next node of the routing; function used by the node to join the network through another known node;
- *join\_res\_handle()*: function used for updating the routing table and leaf set after a *join\_response*;
- *exit()*: function for handling the exit of another node from the network;
- *keepalive()*: sends keepalive to all known nodes;
- *update\_keepalive()*: updates the keepalive of given nodes;
- *keepalive\_res()*: responses to the keepalive to make sure that, even if a node knows another but not viceversa, the one knowing still receives keepalives from the other;

- *check\_expired\_nodes()*: checks which nodes are disconnected and updates the tables;
- *share\_info()*: shares the info about the nodes with the known nodes;

## 4 Testing and Results

Once Mininet and ONOS are initialized, the network is tested to guarantee that it works as intended by the VPLS. The network is tested both with the *ping* and both with a configured broadcast ethernet frame. The broadcast is executed using the python script *frameBroadcaster.py*.

For checking that the VPLS are working as intended, we used *Wireshark*, checking which interfaces and devices receive the ethernet frames broadcasted. Different topologies were initialized as *net.json* and run. The topologies and some of the screenshots of wireshark are shown below, to prove how only the interfaces in the same VPLS are reached by the broadcast.

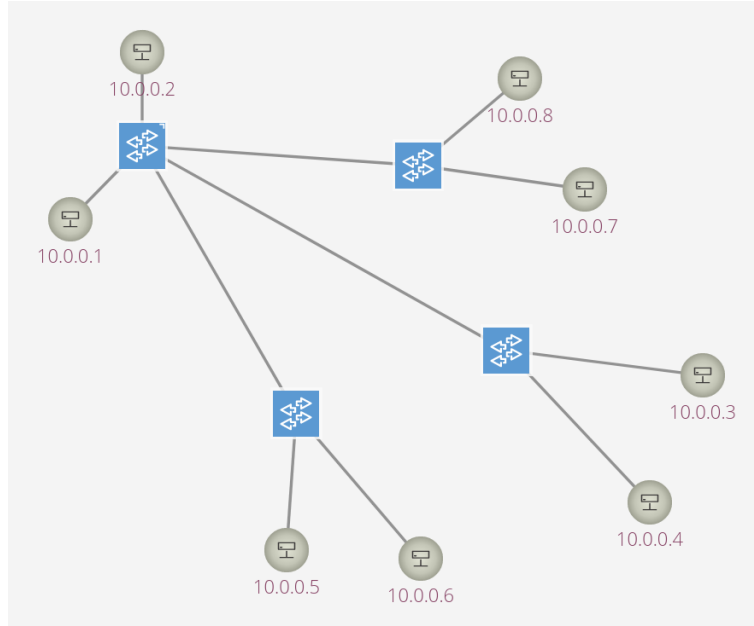


Figure 1: Topology 1

No.	Time	Source	Destination	Protocol	Length	Info
22	3.071330667	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
25	3.071334214	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
28	3.071338301	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
54	8.096294238	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
55	8.096301181	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
56	8.096297244	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
89	13.120884527	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
90	13.120896748	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
91	13.120899178	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
110	18.152811841	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
111	18.152816666	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
112	18.152822182	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
146	23.179633081	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
147	23.179638231	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
148	23.179644553	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
182	28.211282709	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
183	28.211300242	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
184	28.211290784	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
202	33.248498672	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
203	33.248507639	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64
204	33.248502870	10.0.0.1	10.255.255.255	ICMP	47	Echo (ping) request id=0x0000, seq=0/0, ttl=64

Figure 2: Broadcast example 1: servers reached

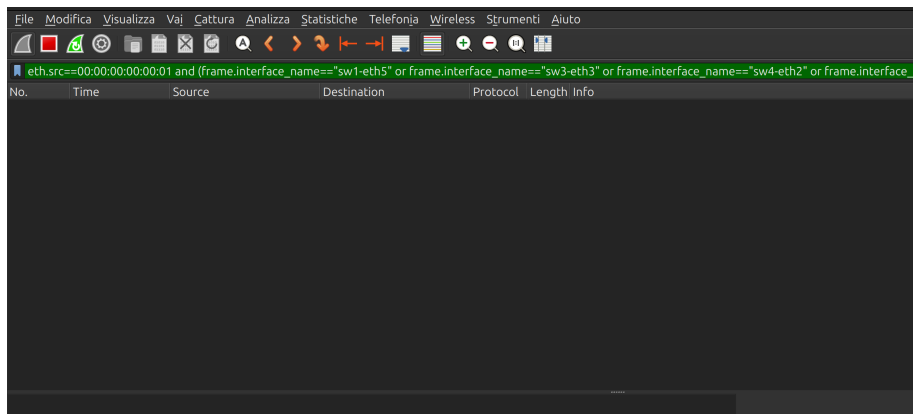


Figure 3: Broadcast example 1: all other interfaces not reached

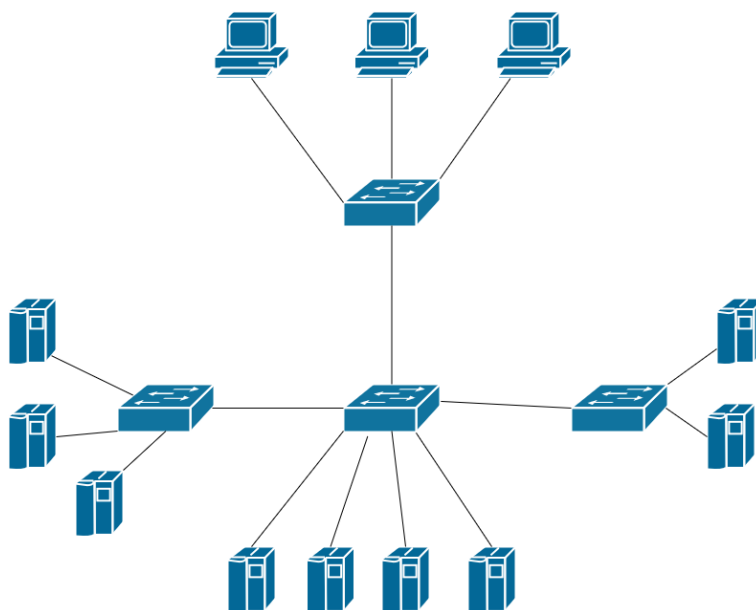


Figure 4: Topology 2



No.	Time	Source	Destination	Protocol	Length	Info
64	4.522976433	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
67	4.522974751	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
70	4.522980352	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
103	9.557274631	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
104	9.557278019	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
105	9.557282678	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
156	14.579213250	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
157	14.579219762	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
158	14.579227537	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
184	19.606206576	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
185	19.606211977	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
186	19.606218429	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
236	24.639032979	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
237	24.639038459	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
238	24.639045061	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64
288	29.676294408	10.0.0.3	10.255.255.255	ICMP	51	Echo (ping) request id=0x0000, seq=0/0, ttl=64

Figure 8: Broadcast example 3: servers reached

No.	Time	Source	Destination	Protocol	Length	Info
-----	------	--------	-------------	----------	--------	------

Figure 9: Broadcast example 3: all other interfaces not reached