



# UNIVERSITÀ DI PISA



## Master's Degree in Artificial Intelligence and Data Engineering

Distributed Systems and Middleware Technologies

### Task Force

#### Instructors:

**Prof. Alessio Bechini**

#### Student:

**Adrien Koumgang**

**Tegantchouang**

Academic Year 2025/2026

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	2
1.2	Problem Statement . . . . .	2
1.3	Project Scope . . . . .	2
<b>2</b>	<b>System Architecture</b>	<b>3</b>
2.1	High-Level Architecture and Justification . . . . .	3
2.2	Data Flow . . . . .	3
2.3	Component Responsibilities . . . . .	4
<b>3</b>	<b>Distributed Systems Challenges and Solutions</b>	<b>5</b>
3.1	Synchronization and Coordination Challenges . . . . .	5
3.1.1	Challenge 1: Leader Election . . . . .	5
3.1.2	Challenge 2: Distributed State Consistency . . . . .	5
3.1.3	Challenge 3: Dynamic Worker Coordination . . . . .	5
3.2	Communication Challenges . . . . .	6
3.2.1	Challenge 4: Heterogeneous System Integration . . . . .	6
3.2.2	Challenge 5: The "Black Hole" Status Problem . . . . .	6
<b>4</b>	<b>Technology Stack And Implementation Details</b>	<b>7</b>
4.1	Project Structure . . . . .	7
4.2	Erlang/OTP Orchestrator . . . . .	7
4.2.1	Shared Record Definitions (.hrl) . . . . .	7
4.2.2	Pro Routing Strategy . . . . .	8
4.3	Polyglot Worker Framework (Python) . . . . .	8

4.4	Database Schema . . . . .	9
4.5	Quick Start Guide . . . . .	10
4.5.1	Access Points . . . . .	10
4.6	Testing and Benchmarks . . . . .	10
4.6.1	Swagger UI . . . . .	10
4.6.2	Unit Testing (EUnit) . . . . .	10
4.6.3	Integration Testing . . . . .	12
4.6.4	Performance Benchmarks . . . . .	12
4.6.5	Fault Tolerance Simulation . . . . .	12
<b>5</b>	<b>Conclusion</b>	<b>13</b>
5.1	Use of Artificial Intelligence . . . . .	14
5.1.1	Architectural Refactoring . . . . .	14
5.1.2	Debugging Erlang Compilation Issues . . . . .	14
5.1.3	Cross-Language Serialization (Jackson vs. JSX) . . . . .	15
5.1.4	Algorithm Design . . . . .	15

# Chapter 1

## Introduction

TaskForce is a distributed job scheduling system designed to process complex heterogeneous tasks across a polyglot cluster. The system demonstrates the practical application of fundamental distributed systems concepts, including consensus algorithms, message-oriented middleware, fault tolerance, and distributed coordination, fulfilling all requirements of the Distributed Systems course.

The project implements a complete job scheduling pipeline featuring:

- **Erlang/OTP Orchestrator** [1]: A robust coordinator utilizing the Raft consensus algorithm for leader election and CRDTs for state management.
- **RabbitMQ Cluster** [2]: Message-Oriented Middleware providing reliable, priority-based asynchronous queuing.
- **Java Spring Boot API Gateway** [3]: A web-facing entry point with REST endpoints and real-time Server-Sent Events (SSE) streaming.
- **Polyglot Worker Farm**: Heterogeneous execution nodes supporting Python and Java implementations.
- **Multi-node Deployment**: Distributed across Ubuntu VMs (Adrien0, Adrien1) with automatic failover.

The system successfully handles concurrent job submissions, automatically recovers from node failures, maintains reliable execution semantics, and scales horizontally.

## 1.1 Background

Distributed systems have become the backbone of modern computing infrastructure. From cloud computing platforms to microservices architectures, the ability to coordinate work across multiple nodes while maintaining consistency and availability is fundamental. Job scheduling—the allocation of computational tasks to worker nodes—represents a classic distributed systems problem that encompasses core challenges: coordination, communication, fault tolerance, and load balancing.

## 1.2 Problem Statement

Modern distributed applications require reliable job processing systems that can:

- Handle concurrent task submissions asynchronously.
- Guarantee fault tolerance and state consistency across nodes.
- Scale horizontally as workload increases.
- Support heterogeneous (polyglot) execution environments.

This project aims to build a complete distributed job scheduler from first principles, demonstrating how fundamental distributed algorithms and middleware technologies address these challenges.

## 1.3 Project Scope

The scope of this project includes the design and implementation of a distributed scheduler, Raft consensus for leader election, AMQP integration, a polyglot worker framework, and comprehensive API design.

### Project Repository:

The complete source code for TaskForce, including the Erlang orchestrator, Java API Gateway, polyglot worker implementations, and deployment scripts, is publicly available on GitHub. It can be accessed at:

<https://github.com/adrienKoumgangT/distriqueue>

# Chapter 2

## System Architecture

### 2.1 High-Level Architecture and Justification

To fulfill the course requirements regarding models, paradigms, and middleware, the system adopts a Message-Oriented Middleware (MOM) paradigm integrated with an Actor-model orchestrator.

- **Erlang/OTP:** Chosen for its native implementation of the Actor model, asynchronous message passing, and "let it crash" philosophy, making it ideal for distributed fault-tolerant coordination.
- **RabbitMQ:** Chosen to decouple the HTTP ingestion layer from the execution layer. It handles the "communication" requirement by providing durable, asynchronous message queues with priority support.
- **Java/Spring Boot:** Chosen for the API Gateway to satisfy the "Advanced Java topics for enterprise applications" requirement, leveraging robust REST controllers and Server-Sent Events (SSE).

### 2.2 Data Flow

#### Job Submission Flow:

1. Client submits a job via REST API to the Java Spring Boot API Gateway.
2. The Gateway persists the job and publishes it to RabbitMQ (`jobs.exchange`) with an AMQP priority (High/Medium/Low).

3. The Erlang Orchestrator consumes the job from RabbitMQ.
4. The Orchestrator's Router dynamically discovers the least-loaded, capable polyglot worker and assigns the task.
5. The Worker (Python/Java) executes the job and sends a status update via HTTP to the Erlang Orchestrator.
6. Erlang updates its internal Mnesia/Map state, replicates it via CRDTs, and broadcasts the status back to the Gateway via RabbitMQ (`status.exchange`).
7. The API Gateway receives the AMQP update and notifies clients via SSE.

#### **Failure Recovery Flow:**

1. Leader orchestrator node fails. [cite:*start*]
1. Remaining nodes detect failure via heartbeat timeout[cite: 71, 72].
2. Raft consensus elects a new leader.
3. New leader resumes job processing and re-registers connected workers.

## **2.3 Component Responsibilities**

- **Orchestrator (Erlang/OTP):** Leader election, job registry, dynamic worker routing, load balancing, health monitoring[cite: 73].
- **API Gateway (Java/Spring):** REST API, SSE streaming, Jackson JSON serialization.
- **Message Broker (RabbitMQ):** Priority queuing, AMQP routing.
- **Workers (Python, Java):** Job execution (calculation, transformation, validation)[cite: 73].

# Chapter 3

## Distributed Systems Challenges and Solutions

### 3.1 Synchronization and Coordination Challenges

#### 3.1.1 Challenge 1: Leader Election

**Problem:** Multiple orchestrator nodes cannot all act as leaders simultaneously without causing duplicate job assignments.

**Solution: Raft Consensus Algorithm.** Implemented in `raft_fsm.erl`, using `gen_fsm`. The system uses randomized election timeouts to prevent split votes and maintains a strict logical clock (term) to detect stale leaders.

#### 3.1.2 Challenge 2: Distributed State Consistency

**Problem:** Multiple orchestrator nodes need a consistent view of job states despite network partitions.

**Solution: CRDT-based Job Registry.** Implemented in `job_registry.erl`, using an Observed-Remove Set (OR-Set) with microsecond timestamps for deterministic conflict resolution.

#### 3.1.3 Challenge 3: Dynamic Worker Coordination

**Problem:** Dynamically assigning jobs to heterogeneous workers (Python/Java) based on real-time load and specific job types (e.g., ‘calculate’ vs ‘transform’).

**Solution: Dynamic Polyglot Load Balancer.** Implemented in `router.erl`. Instead of hardcoding worker assignments, the system queries the `dq_worker_pool` to find all workers capable of a specific task type, calculates their `current_load / capacity` ratio, and routes the job to the least loaded node.

## 3.2 Communication Challenges

### 3.2.1 Challenge 4: Heterogeneous System Integration

**Problem:** Components written in Erlang, Java, and Python must communicate seamlessly without tight coupling.

**Solution: Polyglot Message-Oriented Strategy.**

- **AMQP (RabbitMQ):** Used for asynchronous job distribution (Gateway → Erlang) and status broadcasting (Erlang → Gateway).
- **HTTP/REST:** Used by Workers to report completion metrics back to the Erlang API (`distriqueue_api.erl`).
- **Erlang Distribution:** Used natively for Inter-node Gossip and Raft coordination.

### 3.2.2 Challenge 5: The "Black Hole" Status Problem

**Problem:** Initially, workers pushed statuses directly to the Gateway, bypassing Erlang. This created a "black hole" where Erlang's internal state became desynchronized from the actual system state.

**Solution: Erlang as the Single Source of Truth.** The architecture was refactored so workers report strictly to Erlang via HTTP. Erlang updates its Mnesia/CRDT tables, encodes the complex mathematical results using `jsx`, and uses the `rabbitmq_client.erl` to broadcast a definitive AMQP status payload to the Gateway.

# **Chapter 4**

## **Technology Stack And Implementation-Details**

### **4.1 Project Structure**

The project consists primarily of these files:

- **orchestrator/**: Erlang/OTP application
- **gateway/**: Java Spring Boot API
- **workers/python-worker/**: Python worker implementation
- **workers/java-worker/**: Java worker implementation

### **4.2 Erlang/OTP Orchestrator**

#### **4.2.1 Shared Record Definitions (.hrl)**

To prevent inclusion explosions and maintain a single source of truth across the cluster, a shared header file (`distriqueue.hrl`) is utilized.

```

1 -record(worker, {
2   id :: binary(),
3   type :: binary(),
4   status :: active -> idle -> unresponsive,
5   capacity :: integer(),
6   current_load :: integer(),
7   last_heartbeat :: integer()
8 }).

```

Code 4.1: Shared Worker Record

### 4.2.2 Pro Routing Strategy

The load balancer inside `router.erl` implements an identity-agnostic, capacity-aware routing algorithm:

```

1 select_least_loaded_by_type(Type, _State) →
2   case dq_worker_pool:get_workers_by_type(Type) of
3     [] → <<"default_worker">>;
4     Workers →
5       {BestWorkerId, _} = lists:foldl(
6         fun(W, {AccId, AccRatio}) →
7           Cap = case W#worker.capacity of 0 → 1; C → C end,
8           Ratio = W#worker.current_load / Cap,
9           if Ratio < AccRatio → {W#worker.id, Ratio};
10          true → {AccId, AccRatio}
11        end
12      end, {<<"none">>, 999999.0}, Workers),
13      BestWorkerId
14    end.

```

Code 4.2: Dynamic Least-Loaded Type Routing

## 4.3 Polyglot Worker Framework (Python)

The worker framework utilizes the Strategy Pattern to dynamically map incoming payloads to execution logic.

```

1 class TransformationJobHandler(JobHandler):

```

```

2     def can_handle(self, job_type: str) → bool:
3         return job_type in ['transform', 'uppercase', 'reverse']
4
5     def execute(self, job: Dict[str, Any]) → Dict[str, Any]:
6         payload = job.get('payload', { })
7         operation = payload.get('operation', 'uppercase')
8         data = payload.get('data', { })
9
10        if operation == 'reverse':
11            transformed = { k[::-1]: str(v)[::-1] for k, v in data.items() }
12        return { 'transformed': transformed }
```

Code 4.3: Python Strategy Pattern for Jobs

## 4.4 Database Schema

```

1 CREATE TABLE jobs (
2     id VARCHAR(36) PRIMARY KEY,
3     type VARCHAR(50) NOT NULL,
4     priority VARCHAR(20) NOT NULL,
5     status VARCHAR(20) NOT NULL,
6     worker_id VARCHAR(100),
7     payload CLOB,
8     result CLOB,
9     error_message VARCHAR(1000),
10    retry_count INT DEFAULT 0,
11    max_retries INT DEFAULT 3,
12    execution_timeout INT DEFAULT 300,
13    created_at TIMESTAMP,
14    started_at TIMESTAMP,
15    completed_at TIMESTAMP,
16    updated_at TIMESTAMP,
17    metadata CLOB,
18    parent_job_id VARCHAR(36),
19    callback_url VARCHAR(500),
20    tags VARCHAR(500),
21    queue_name VARCHAR(50),
22    processing_node VARCHAR(100),
```

```

23     version BIGINT
24 );
25
26 CREATE INDEX idx_job_status ON jobs(status);
27 CREATE INDEX idx_job_type ON jobs(type);
28 CREATE INDEX idx_job_priority ON jobs(priority);
29 CREATE INDEX idx_job_created ON jobs(created_at);
30 CREATE INDEX idx_job_worker ON jobs(worker_id);

```

Code 4.4: Database Schema

## 4.5 Quick Start Guide

### 4.5.1 Access Points

Service	Adrien0 (10.2.1.11)	Adrien1 (10.2.1.12)
API Gateway	http://10.2.1.11	http://10.2.1.12
Load Balancer	http://10.2.1.11 (port 80)	-
H2 Console	http://10.1.2.11:8082/h2-console	http://10.1.2.12:8082/h2-console
RabbitMQ UI	http://10.1.2.11:15672	http://10.1.2.12:15672
Node Exporter	http://10.2.1.11:9100	http://10.2.1.4:9100

#### Credentials:

- RabbitMQ: admin / admin
- H2 Database: sa / (empty)
- Redis: (empty)

## 4.6 Testing and Benchmarks

### 4.6.1 Swagger UI

### 4.6.2 Unit Testing (EUnit)

The Erlang orchestrator utilizes the EUnit testing framework to ensure the internal logic of the system is sound. Test suites were written for modules such as `dq_worker_pool.erl` to verify worker registration, capacity tracking, and selection strategies.

The screenshot shows a list of job management endpoints under the heading "Jobs Job management endpoints". Each endpoint is listed with its method, path, and a brief description. The endpoints are:

- GET /jobs Get jobs with filtering
- GET /jobs/{id} Get job by ID
- GET /jobs/stream Stream real-time job updates
- GET /jobs/statistics Get job statistics
- POST /jobs Submit a new job
- POST /jobs/{id}/retry Retry a failed job
- POST /jobs/{id}/cancel Cancel a job
- POST /jobs/batch Submit multiple jobs in batch

Figure 4.1: API Gateway Job Management endpoints.

The screenshot shows a list of job schemas under the heading "Schemas". Each schema is listed with its name and a link to expand all objects. The schemas are:

- JobRequest > Expand all object
- JobResponse > Expand all object
- JobBatchRequest > Expand all object
- JobBatchResponse > Expand all object
- PageMetadata > Expand all object
- PagedModelJobResponse > Expand all object
- SseEmitter > Expand all object
- JobStatistics > Expand all object

Figure 4.2: API Gateway Job Schemas.

```

1 test_select_worker() →
2   ok = dq_worker_pool:register_worker(<<"worker3">>, <<"python">>, 5, 0),
3   ok = dq_worker_pool:register_worker(<<"worker4">>, <<"python">>, 5, 3),
4   {ok, LeastLoaded} = dq_worker_pool:select_worker(<<"python">>, least_loaded),
5   ?assertEqual(<<"worker3">>, LeastLoaded).

```

Code 4.5: EUnit Worker Pool Test

### 4.6.3 Integration Testing

End-to-end integration tests were conducted by submitting complex JSON payloads via cURL to the Spring Boot Gateway, observing the AMQP transmission, and verifying the Erlang routing.

**Test Case: Nested Mathematical Aggregation** A job requesting the aggregate operation on a large array of numbers was submitted. The system dynamically routed the task to a Java worker (due to lower load than the Python worker). The complex nested result (`{ "aggregates": { "average": 473, "max": 1024 } }`) was successfully parsed by Erlang using `jsx:encode` and broadcast back to the Gateway.

### 4.6.4 Performance Benchmarks

- **End-to-End Latency:** The round trip (Gateway → RabbitMQ → Erlang Router → Worker Execution → Erlang HTTP → RabbitMQ Broadcast → Gateway) averaged **375 milliseconds** for standard computational tasks.
- **Queue Time:** Overhead spent in RabbitMQ and Erlang's routing queue averaged roughly 130 ms under a normal load.

### 4.6.5 Fault Tolerance Simulation

Chaos engineering techniques were applied by manually terminating worker processes mid-execution. The `health_monitor.erl` successfully detected the missing heartbeats, marked the nodes as unresponsive, and the orchestrator requeued the pending tasks according to the `max_retries` defined in the job payload.

# Chapter 5

## Conclusion

TaskForce successfully demonstrates a complete distributed job scheduling system that addresses the core challenges of distributed systems: coordination, communication, and fault tolerance.

The project fulfills all course requirements:

1. **Synchronization/Coordination:** Addressed via Raft consensus for leader election and CRDTs for eventual consistency.
2. **Communication:** Addressed via RabbitMQ AMQP queues, RESTful APIs, and Erlang asynchronous messaging.
3. **Erlang Component:** The core orchestrator is built entirely in Erlang using OTP behaviours (`gen_server`, `gen_state`, `supervisor`).
4. **Multi-node Deployment:** Successfully deployed across multiple Ubuntu Virtual Machines.

The system achieves:

- **High availability** through clustering and automatic failover.
- **Scalability** by dynamically routing to polyglot worker nodes based on capacity.
- **Reliability** with exactly-once execution semantics and RabbitMQ dead-lettering.
- **Extensibility** through an identity-agnostic Strategy Pattern worker architecture.

This project represents a production-quality distributed system that demonstrates a deep understanding of distributed systems principles and their practical implementation. The experience gained from integrating advanced Erlang paradigms with modern Java Enterprise APIs and RabbitMQ has provided invaluable insight into the realities of distributed systems engineering.

## 5.1 Use of Artificial Intelligence

During the development of TaskForce, a Large Language Model (Google Gemini) was utilized as an interactive pair-programming and architectural assistant. The AI was highly effective in diagnosing distributed systems bugs, bridging multi-language syntax gaps, and refining the architectural flow.

### 5.1.1 Architectural Refactoring

Initially, the system suffered from a "Point-to-Point Anti-Pattern," where workers communicated state updates directly to the Spring Boot Gateway, bypassing the Erlang Orchestrator. This caused the Erlang internal registry to fall out of sync. Gemini assisted in designing a "Message-Driven" architecture, helping implement the `rabbitmq_client:publish_status/4` function in Erlang. This repositioned Erlang as the absolute single source of truth, receiving HTTP updates from workers and safely broadcasting AMQP updates to the cluster.

### 5.1.2 Debugging Erlang Compilation Issues

During the implementation of the dynamic load balancer in `router.erl`, the codebase encountered an "Include Explosion" and multiple "redefining module" compilation errors. The AI was used to analyze the compiler logs and identified that the entire source code of the worker pool was accidentally pasted into the `distriqueue.hrl` header file. The AI provided the correct pattern for Erlang header files (macros and record definitions only) and guided the successful separation of the `#worker{}` record, allowing the codebase to compile.

### **5.1.3 Cross-Language Serialization (Jackson vs. JSX)**

When passing complex nested JSON result structures from Python/Java workers back through Erlang and into Spring Boot, the system encountered several deserialization faults (e.g., `Cannot map null into type int`). Gemini aided in configuring Spring Boot's Jackson Deserialization features (relaxing primitive constraints) and assisted in writing safe `jsx:encode/1` logic in Erlang to gracefully handle undefined result maps.

### **5.1.4 Algorithm Design**

The AI assisted in authoring the functional Erlang code required for the `least_loaded` routing strategy, specifically structuring the `lists:foldl/3` logic to safely calculate capacity ratios while preventing division-by-zero errors.

# Bibliography

- [1] Erlang: Erlang/otp documentation, <https://www.erlang.org/docs>
- [2] RabbitMQ: Rabbitmq documentation, <https://www.rabbitmq.com/documentation.html>
- [3] Spring: Spring boot reference guide, <https://spring.io/projects/spring-boot>