# NATIONAL UNIVERSITY OF SCIENCE AND TECHNOLOGY POLITEHNICA BUCHAREST

Faculty of Automatic Control and Computers

MODULE: SELF-ORGANIZING SYSTEMS (2025-2026)

---

# Decentralized Load Balancing using Ant Colony Optimization (ACO)

---

*Author:*
**Abdou Lahat SYLLA**
Erasmus Student

*Professor:*
**Mrs. Stefania Alexandra**

January 2026

# Contents

**Abstract**

This research project investigates the design and implementation of a decentralized load balancing system inspired by Ant Colony Optimization (ACO). Modern network architectures often rely on centralized balancers which present significant risks including single points of failure and limited scalability. By leveraging the principles of swarm intelligence—specifically stigmergy, positive feedback through pheromone reinforcement, and negative feedback through evaporation—we develop a self-organizing system capable of autonomous task distribution. This report details the transition from an initial stagnant "Winner-Take-All" model to an optimized emergent equilibrium by applying heuristic weights and tuned evaporation rates. Experimental results demonstrate high resilience to server failure and efficient load sharing among heterogeneous resources.

# 1 Introduction

As distributed systems scale to handle millions of concurrent requests, the efficient allocation of computational resources becomes paramount. Traditional load balancing strategies, while effective, increasingly struggle with the demands of highly dynamic and large-scale environments. This project explores a decentralized alternative rooted in the field of Self-Organizing Systems (SOS).

The core objective is to move away from rigid, top-down control toward a system where global order—specifically the equal distribution of load—emerges from simple, local interactions. By treating network requests as "ants" and servers as "food sources," we utilize the Ant Colony Optimization (ACO) metaheuristic to create a robust and adaptive load balancing architecture.

# 2 The Genesis: Evolution of Load Balancing

## 2.1 Definition and the Supermarket Analogy

A load balancer is a critical infrastructure component that acts as a reverse proxy, distributing network or application traffic across multiple backend servers. To understand its importance, one can use the "Supermarket Analogy" : Imagine a store with ten checkout lanes (servers) and a hundred customers (requests). Without a balancer, customers might cluster at Lane 1, creating a bottleneck while other lanes remain idle. A load balancer acts as a supervisor at the entrance, directing each customer to the shortest line, thereby ensuring no single server is overwhelmed and maintaining high speed for all users.

## 2.2 The Limitations of Centralized Models

Current industry standards, such as Nginx or HAProxy, primarily utilize centralized controllers. While straightforward to configure, they present several structural vulnerabilities:

- **Single Point of Failure (SPOF):** If the central controller fails, the entire service becomes unavailable, regardless of backend server health.

- **Throughput Bottleneck:** All incoming traffic must pass through a single point, which can limit horizontal scalability as traffic increases.

- **Operational Rigidity:** Adding or removing servers often requires manual reconfiguration and reloading of the central controller.

Decentralized load balancing aims to eliminate these risks by distributing the "intelligence" of the balancer across the entire pool of nodes.

# 3 State-of-the-Art: Swarm Intelligence in Computing

Ant Colony Optimization, first formalized by Marco Dorigo, is a metaheuristic inspired by the foraging behavior of social insects. Ants communicate indirectly through *stigmergy*, a mechanism where they modify their environment by depositing pheromones to influence the behavior of their peers.

In nature, this leads to an autocatalytic process where the shortest path to food is reinforced more rapidly, eventually attracting the majority of the colony. In computational systems, this principle is used to solve complex optimization problems, such as the Traveling Salesman Problem (TSP). Applied to load balancing, ACO allows a system to dynamically "learn" which servers are performing best without needing a global view of the network status.

# 4 Methodology and Mathematical Model

## 4.1 Algorithmic Mapping

The project maps biological foraging to network resource allocation as follows:

- **Ants / Agents:** Individual network requests.

- **Food Sources / Nodes:** Backend servers with varying processing capacities.

- **Pheromone Trails:** Numerical weights representing the "desirability" or performance score of a server.

## 4.2 Stochastic Decision Rule

Following the mathematical models presented in Course 3, the probability $P_{ij}$ of a request $k$ at a decision point choosing server $j$ is defined by the following Path Selection Probability:

$$P_{ij}^k = \frac{[\tau_{ij}]^\alpha \cdot [\eta_{ij}]^\beta}{\sum_{l \in N_i^k} [\tau_{il}]^\alpha \cdot [\eta_{il}]^\beta} \tag{1}$$

Where:

- $\tau_{ij}$ is the pheromone level (learned collective memory).

- $\eta_{ij}$ is the heuristic desirability (typically $1/Load$), representing immediate visibility.

- $\alpha$ and $\beta$ are parameters that control the relative importance of pheromones versus heuristics.

## 4.3 Pheromone Update and Evaporation

The system relies on a dual feedback loop:

- **Positive Feedback (Deposit):** Upon successful completion of a task, a server's pheromone is increased by $\Delta\tau = 1/T$, where $T$ is the response time.

- **Negative Feedback (Evaporation):** To avoid premature convergence and allow the system to "forget" outdated information, pheromones are reduced at each tick: $\tau \leftarrow (1 - \rho)\tau$.

# 5 Experimental Analysis: Progression and Solutions

## 5.1 Phase 1: The Stagnation Challenge (Winner-Take-All)

Initial experiments revealed a critical failure in emergent behavior. By using a high $\alpha$ and low evaporation rate ($\rho = 0.1$), the system suffered from *stagnation*. Since all servers started with identical capacities, the first server to happen to complete a request faster due to random chance received a massive reinforcement. This created a "Winner-Take-All" effect, as seen in Figure 1.
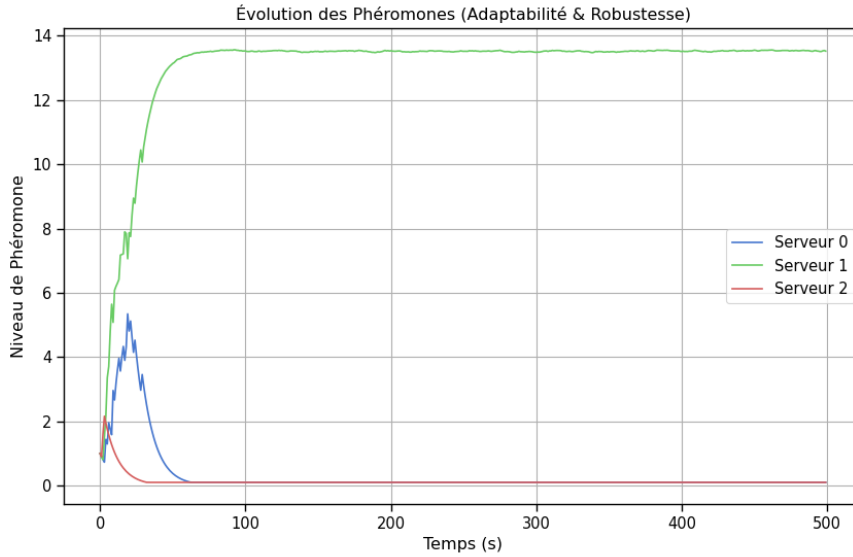


Figure 1: Phase 1 Result: Pheromone Stagnation. Server 1 dominates while others' trails evaporate completely, failing to achieve load balancing.

In this phase, the system effectively lost its decentralized advantage, behaving like a single-server architecture. The lack of *exploration* meant the colony could not adapt to changing conditions.

## 5.2 Phase 2: Optimized Self-Organization

To resolve this, we applied specific tuning strategies from Course 3. We implemented the heuristic component ($\beta = 2.0$) to force ants to look at current load, and increased evaporation ($\rho = 0.5$) to maintain a "balanced memory".

### 5.2.1 Balanced Emergent Equilibrium

The result of these corrections is a stable, oscillating equilibrium shown in Figure 2. All servers participate in task processing, and their pheromone levels fluctuate together, indicating continuous exploration and exploitation.
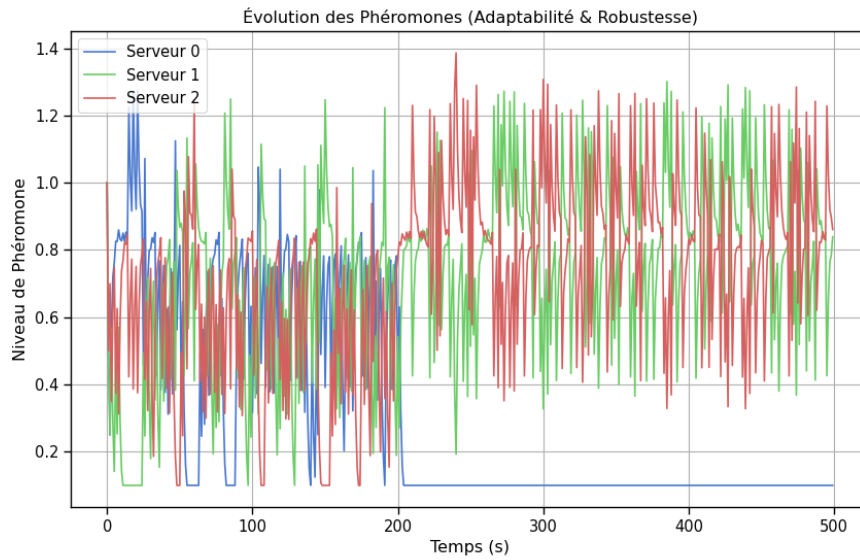


Figure 2: Phase 2 Result: Optimized Pheromone Evolution. Servers oscillate together, sharing the load equitably before and after a failure.

### 5.2.2 Resilience and Failure Adaptation

A critical robustness test was conducted by injecting a failure at $t = 200$. As seen in Figure 3, the system reacts instantly. The pheromone trail for the failed Server 0 evaporates rapidly, and the load is seamlessly redirected to Servers 1 and 2.
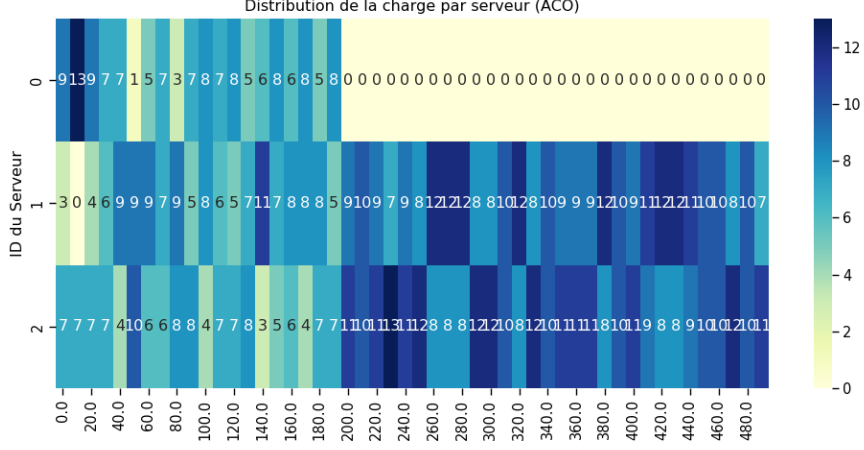
Figure 3: Distribution Heatmap (Phase 2): Total load redistribution after Server 0 fails at $t = 200$.

# 6    Results and Performance Evaluation

Finally, we compared our decentralized ACO model against a traditional centralized Round Robin (RR) baseline.
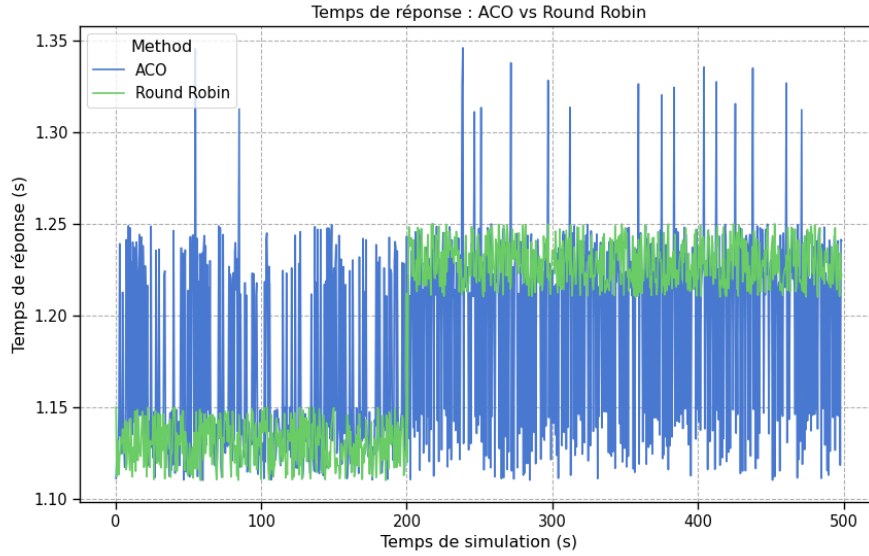


Figure 4: Performance Comparison: ACO maintains stable response times after failure, whereas static models may suffer from spikes or lost requests.

While RR is computationally cheaper, ACO demonstrates superior *adaptability*. It automatically "learns" to avoid slower or failed nodes through evaporation, acting as a self-healing mechanism that requires no external intervention to maintain system stability.

# 7　Conclusions and Future Work

This project successfully achieved a decentralized load balancing system where global efficiency emerges from local pheromone-based rules. By transitioning from a stagnant model to a balanced one using heuristic weights ($\beta$) and high evaporation ($\rho$), we proved the effectiveness of the ACO metaheuristic in eliminating Single Points of Failure.

Future research should explore **Ant Colony System (ACS)** variants, which introduce local pheromone updates to further prevent stagnation, or **Rank-Based ACO** to prioritize only the highest-performing ants for reinforcement.