

Ratio-Based Signaling for Source-Victim Separation in Swarm Fault Detection

Longyin Cui

Loyola University Maryland

Baltimore, MD, USA

lcui@loyola.edu

ABSTRACT

The source-victim ambiguity problem, distinguishing faulty agents from their impaired neighbors, complicates swarm fault detection. Self-diagnosis misses victims, while neighbor voting mislabels them as defective. We propose a bio-inspired ratio signaling mechanism where agents emit and absorb hormone-like stress signals. The balance of internal versus external stress discriminates whether degradation stems from internal malfunction or external exposure. With adaptive emission dynamics and local thresholds, the method achieves fast response without centralized control. Evaluations across varied densities, fault severities, and packet loss demonstrate that ratio signaling reduces healthy-agent impairment and remains robust under communication degradation.

CCS CONCEPTS

• **Computing methodologies** → **Multi-agent systems.**

KEYWORDS

swarm robotics; fault detection; bio-inspired signaling; multi-agent systems; quarantine

ACM Reference Format:

Longyin Cui. 2026. Ratio-Based Signaling for Source-Victim Separation in Swarm Fault Detection. In *Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026), Paphos, Cyprus, May 25 – 29, 2026*, IFAAMAS, 9 pages. <https://doi.org/10.65109/UYYF8433>

1 INTRODUCTION

Fault cascades pose a fundamental challenge to swarm robotic systems [2, 3]. When a malfunctioning agent degrades its neighbors' performance, the neighbors may impair others, creating a propagating dysfunction that undermines promised redundancy. Bjerknes and Winfield showed that swarm reliability paradoxically decreases with population size under partial faults, challenging assumptions that larger swarms are inherently more robust [2]. Preventing cascades requires the rapid identification and isolation of faulty agents. However, determining which agents are faulty proves difficult when symptoms spread through local interactions.

The core difficulty is source-victim ambiguity. Both the originating faulty agent and its impaired neighbors exhibit degraded performance [14]. Self-diagnosis mechanisms miss collateral victims whose impairment stems from external exposure. Conversely,

neighbor-based voting schemes misclassify victims as defective. This ambiguity creates a dilemma: aggressive quarantine risks removing healthy agents, while conservative approaches allow the spread of cascades. This work addresses source-victim ambiguity through ratio-based hormone signaling that exploits the asymmetry between internally generated and externally absorbed stress.

Existing methods struggle with this asymmetry because they rely on single information streams [13, 25]. External observation approaches, such as neighbor voting [18, 31], classify agents based solely on perceived performance, potentially mislabeling victims. Internal diagnostic methods [6, 24] detect intrinsic faults but miss externally induced degradation. Neither approach resolves source-victim ambiguity alone.

Biological systems offer an alternative model that addresses this limitation [1, 21]. The immune system avoids misclassification by monitoring the balance between excitation and regulation signals exchanged by neighbors [5], enabling local discrimination without global coordination.

We propose a ratio-based stress signaling approach, where agents continuously emit and absorb scalar stress values through local communication. The direction of stress flow reveals causality. Faulty agents primarily generate stress signals, whereas affected neighbors mainly receive them. Consequently, sources exhibit strong internal stress but little incoming stress, while victims show the opposite pattern. This asymmetry enables local discrimination without requiring explicit role labels or global coordination.

The mechanism relies on how stress propagates through the network. A faulty agent first accumulates stress internally and then broadcasts it to its neighbors. Nearby agents mainly receive stress rather than generate it. This directional pattern distinguishes sources from victims, allowing each agent to decide locally whether it should quarantine based solely on its neighbors.

Contributions: This paper contributes to the following:

- **Ratio-based causal inference:** Comparing internal and external stress reveals whether degradation originates locally or from exposure, allowing agents to distinguish sources from victims.
- **Emergent role discovery:** Agents adapt stress emission through neighbor feedback, identifying source or victim roles without explicit classification rules.
- **Scalable discrimination:** The method prevents cascades across scales and remains functional even under complete communication failure where consensus methods collapse.

Our evaluations show reduced healthy-agent impairment while maintaining accurate source identification, demonstrating robust decentralized fault discrimination.



This work is licensed under a Creative Commons Attribution International 4.0 License.

Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026), C. Amato, L. Dennis, V. Mascardi, J. Thangarajah (eds.), May 25 – 29, 2026, Paphos, Cyprus. © 2026 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). <https://doi.org/10.65109/UYYF8433>

2 RELATED WORK

2.1 Fault Detection and Cascade Prevention in Swarm Systems

Early swarm robotics assumed decentralized systems are inherently robust through redundancy [3, 26]. However, Bjercknes and Winfield demonstrated that reliability paradoxically decreases with swarm size when robots experience partial faults [2], revealing that partially malfunctioning robots can degrade neighbors through local interactions, sometimes creating cascading failures. This finding redirected the field toward active, distributed fault detection.

A central difficulty is source–victim ambiguity: both the faulty robot and its affected neighbors exhibit degraded behavior [14]. Because early deviations can be subtle, diagnosis must identify the agent that initiates the fault rather than those merely influenced by it. This challenge motivates both endogenous (self-diagnosis) and exogenous (neighbor-based) approaches.

2.2 External Observation-Based Detection Methods

Voting and consensus mechanisms dominate fault detection in multi-robot systems [25]. Tarapore et al.’s Cross-Regulation Model (CRM) achieved high accuracy at low false-positive rates across multiple scenarios [31]. Byzantine-tolerant algorithms [16] like W-MSR [18] filter faulty inputs but require strongly connected graphs (typically $(2F + 1)$ -robust) and per-round communication that scales with the number of edges, $O(|E|)$, or $O(N^2)$ in dense networks [18]. Similarly, blockchain-based reputation mechanisms provide immutable security but introduce computational overheads that scale poorly in dynamic swarms [30]. Recent variants like the Decentralized Blocklist Protocol [37] scale better by propagating local accusations. Even classical average-consensus protocols, which achieve efficient agreement under switching topologies [23], still rely on iterative global information propagation, making them susceptible to fault spreading. Trust-based extensions assign dynamic weights based on historical reliability [11]. Extensions of the voter model [36], including three-valued variants that introduce an ‘undecided’ state [8], improve robustness but slow convergence.

While voting methods offer robust black-box detection, relying on external observations alone complicates the distinction between intrinsic faults and induced behaviors, creating systematic misclassification where healthy robots impaired by faulty neighbors may be incorrectly flagged.

2.3 Internal Diagnostic-Based Detection Methods

Threshold-based self-monitoring detects faults by continuously assessing internal state variables [6, 15]. Such methods reliably detect complete failures but degrade when behavior is impaired indirectly through neighbor interactions [17]. Related mechanisms include firefly-inspired synchronization, which flags robots that drift out of phase with the group [7].

Despite these advances, pure self-diagnosis remains limited in distinguishing internally generated faults from externally induced degradation [14]. Machine learning approaches, such as Graph Neural Networks (GNNs), learn local interaction rules [35] but often

require computational resources unavailable on swarm hardware. Data-driven classifiers can outperform earlier model-based methods [4, 31], and Lee et al. [19] extract discriminative metrics for on-board self-detection. However, internal diagnostics still lack access to causal structure: they detect abnormality but cannot determine whether an agent is the source or a downstream victim [13, 25].

2.4 Bio-Inspired Coordination and Emergent Fault Tolerance

Biological systems offer compelling alternative models for distributed fault detection through emergent coordination [26, 34]. Hormone-based control systems draw inspiration from endocrine signaling [29], where agents broadcast scalar values that modulate collective behavior. Wilson et al. [38] showed that hormone-inspired control supports graceful task reallocation and improved energy efficiency in simulation studies. The hormone approach provides smooth behavioral transitions without explicit failure messages, suggesting promise for fault-tolerant coordination.

Artificial immune systems leverage self-nonself discrimination for anomaly detection [32]. However, classic self-nonself models often struggle with the source-victim ambiguity problem because they lack context. Our approach aligns with the Danger Theory [21], extended to AIS by Aickelin and Cayzer [1], which posits that immune responses should be triggered by ‘danger signals’ (damage) rather than foreignness alone. In our system, the internal-external stress ratio serves as this context-aware danger signal. Tarapore et al. [33] demonstrated that online learning of behavioral signatures enables decentralized fault detection in physical robot swarms, though detection performance varies with task and fault type, and temporal filtering to reduce false positives introduces latency. Related bio-inspired approaches include stigmergic coordination through environment-mediated communication [9, 28], immune-inspired discrimination under resource constraints [10], automatic strategy generation [27], pheromone-based fault mitigation [20], and quorum-style decision mechanisms [36].

Despite promising results, bio-inspired methods face validation challenges [12]. Prior demonstrations establish proof-of-concept in small groups ($N < 50$), and while recent work has successfully validated adaptive diagnosis on physical platforms [22], bridging the gap to large-scale hardware requires addressing stochasticity and response time. For instance, Tarapore et al. report that standard CRM-based approaches exhibit detection latencies of $\approx 51 \pm 18$ s [31], which exceeds the sub-second precision required for rapid cascade prevention. Our work addresses these responsiveness concerns through explicit ratio-based discrimination.

2.5 Cascade Analysis and Source-Victim Disambiguation

Distinguishing fault sources from cascade victims remains challenging [14]. Bjercknes and Winfield observed that when a faulty leader’s behavior becomes the majority, healthy robots paradoxically appear anomalous [2]. The lack of standardized benchmarks for source-victim labeling further complicates evaluation and comparison across methods [39].

3 METHODOLOGY

3.1 Problem Formulation

We simulate a swarm of N agents operating in continuous two-dimensional space with discrete-time dynamics at timestep $\Delta t = 0.1$ seconds. Each agent i maintains a task performance metric $x_i \in [0, 1]$ representing its contribution to the collective objective, where $x_i = 1$ indicates nominal performance and $x_i < 1$ indicates degradation. Agents are initially positioned on a grid with small random perturbations and move with boundary-reflection dynamics. We define the interaction topology as a dynamic K -nearest neighbor graph, where \mathcal{N}_i denotes the set of neighbors for agent i based on Euclidean distance d_{ij} .

A subset of agents experience faults either at initialization or during operation. Faults are characterized by severity (minor, moderate, severe) and manifest through two temporal phases: an initial self-degradation period, where the faulty agent's own performance declines, followed by an outward damage period, where the agent begins impairing nearby healthy agents through local interactions. The temporal gap between these phases varies stochastically (0.8-1.5 seconds for minor/moderate faults, 0.5-1.0 seconds for severe faults), reflecting the natural precedence inherent in physical systems that a robot must first malfunction internally before its erroneous behavior propagates to neighbors through movement, communication, or coordination failures. We evaluate robustness to extreme cases through zero-delay experiments (Section 5.4).

The central challenge is distinguishing between fault sources (agents with intrinsic failures) and cascade victims (healthy agents that suffer collateral degradation) when both exhibit similar performance losses. Healthy agents near faulty sources accumulate damage proportional to proximity, severity, and exposure duration, forming propagating cascades if not isolated.

3.2 Dual-Stream Hormone Framework

Our approach fuses complementary information streams through bio-inspired stress signaling. Unlike voting methods that rely solely on external observations or threshold methods that depend on internal diagnostics, we maintain two distinct hormone pathways that capture different aspects of an agent's state. System parameters (e.g., decay rates, emission bounds) were calibrated via simulations to balance detection sensitivity against signal stability. Crucially, these values were held constant across all reported evaluations to demonstrate the method's robustness without scenario-specific tuning.

Internal hormone. Each agent converts its performance loss ($1 - x_i$) into internal stress through a convex production function:

$$\phi(1 - x_i) = \begin{cases} c_1(1 - x_i), & 1 - x_i < \tau_1, \\ c_2(1 - x_i)^{1.5}, & \tau_1 \leq 1 - x_i < \tau_2, \\ c_3(1 - x_i)^2, & 1 - x_i \geq \tau_2, \end{cases} \quad (1)$$

with $c_1 = 10$, $c_2 = 7$, $c_3 = 15$, and breakpoints $\tau_1 = 0.05$, $\tau_2 = 0.3$. The linear, superlinear, and quadratic regimes respectively capture near-nominal degradation, moderate impairment, and severe faults. Parameters were chosen empirically to maintain monotonic stress growth and held fixed for all experiments. We apply Exponentially

Weighted Moving Average (EWMA) smoothing to avoid discontinuous state evolution:

$$H_i^{\text{int}}(t) = (1 - \alpha_{\text{int}})H_i^{\text{int}}(t - \Delta t) + \alpha_{\text{int}}[\phi(1 - x_i) + \beta], \quad (2)$$

where the baseline production $\beta = 0.02$ represents nominal operational stress and $\alpha_{\text{int}} = \Delta t / (0.5 + \Delta t)$ provides temporal smoothing.

External hormone. Agents absorb emissions from their K nearest neighbors, with absorption regulated by remaining capacity to prevent unbounded accumulation:

$$H_i^{\text{ext}}(t) = \text{decay}(H_i^{\text{ext}}(t - \Delta t)) + \sum_{j \in \mathcal{N}_i} A_{ij}(t), \quad (3)$$

where $A_{ij}(t)$ is the absorbed amount from neighbor j (defined in Section 3.4) and $\text{decay}(H)$ is adaptive:

$$\text{decay}(H) = \begin{cases} d_1 H, & H > 0.5, \\ d_2 H, & 0.3 < H \leq 0.5, \\ d_3 H, & H \leq 0.3, \end{cases} \quad (4)$$

with $(d_1, d_2, d_3) = (0.75, 0.85, 0.92)$. Applied multiplicatively each timestep, the decay yields smooth trajectories despite the piecewise rates. Faster dissipation at high stress mimics biological clearance, where elevated concentrations trigger active degradation pathways.

Stability is enforced at two levels: the capacity-limited absorption (Eq. 10) prevents new intake when an agent is near saturation, and a hard clamp reduces external hormone if $H_i^{\text{int}} + H_i^{\text{ext}}$ ever exceeds 1.0 due to simultaneous updates.

Rationale. Faulty sources exhibit high internal stress but low external stress (neighbors emit little), while victims show the opposite. Using total stress $H_i = H_i^{\text{int}} + H_i^{\text{ext}}$, we define the ratio

$$r_i = \frac{H_i^{\text{int}}}{H_i^{\text{int}} + H_i^{\text{ext}} + \epsilon} \quad (5)$$

where $\epsilon = 10^{-3}$ prevents division by zero. This ratio reveals causal direction: $r_i \rightarrow 1$ for sources and $r_i \rightarrow 0$ for victims.

3.3 Adaptive Emission Mechanism

Each agent maintains an emission gain parameter γ_i that modulates how strongly it broadcasts stress to neighbors. At each timestep, agent i observes the change in its neighbors' external stress:

$$\Delta H_j^{\text{ext}} = H_j^{\text{ext}}(t) - H_j^{\text{ext}}(t - \Delta t), \quad j \in \mathcal{N}_i, \quad (6)$$

and computes a windowed average $\bar{\Delta}$ to estimate trend. The gain updates multiplicatively with learning rate $\alpha = 0.05$:

$$\gamma_i(t) = \begin{cases} \min(\gamma_{\text{max}}, (1 + \alpha)\gamma_i(t - \Delta t)), & \bar{\Delta} < -0.01, \\ \max(\gamma_{\text{min}}, (1 - \alpha)\gamma_i(t - \Delta t)), & \bar{\Delta} > 0.01, \\ \gamma_i(t - \Delta t), & \text{otherwise.} \end{cases} \quad (7)$$

with bounds $[\gamma_{\text{min}}, \gamma_{\text{max}}] = [0.1, 2.5]$ and $\gamma_i(0) = 1.0$. The resulting emission is computed as

$$E_i = \gamma_i(\kappa H_i^{\text{int}}) + \mathbf{1}_{x_i < 0.4}(\eta(1 - x_i)) \quad (8)$$

with $\kappa = 1.2$ and $\eta = 0.3$. The first term amplifies internal stress by the gain factor κ ; the second provides a distress boost for severely degraded agents ($x_i < 0.4$), allowing critically impaired robots to emit a detectable signal before internal hormone accumulation through EWMA smoothing. The threshold $x_i < 0.4$ corresponds to the onset of severe degradation in our fault model.

3.4 Communication and Propagation

At each timestep, the simulator rebuilds a K -nearest neighbor graph using agents within a 30-unit communication radius. Hormones diffuse one hop per timestep with spatial attenuation (Eq. 9), so the effective signaling range is much smaller than the communication cutoff. The radius is used only to limit neighbor search.

Hormone diffusion. Emissions propagate through a distance-dependent transfer kernel:

$$T_{ij} = \rho \cdot \frac{\exp(-d_{ij}/\lambda)}{d_{ij} + 1}, \quad (9)$$

where $\rho = 600$ is the diffusion constant and $\lambda = 25$ the spatial decay length. The propagated amount per timestep is $E_i T_{ij} \Delta t$, transmitted only if both transmitter and receiver succeed (simulating packet loss). The large value of ρ enables rapid local diffusion; unbounded accumulation is prevented by capacity-limited absorption (Eq. 10).

Regulated absorption. Neighbor j receives propagated hormone P_{ij} through capacity-limited absorption:

$$A_{ij} = P_{ij} \cdot \max\{0, 1 - (H_j^{\text{int}} + H_j^{\text{ext}})\}, \quad (10)$$

where $P_{ij} = E_i T_{ij} \Delta t$ is the hormone amount propagated from agent i to neighbor j per timestep, with E_i being i 's emission rate and T_{ij} the transfer coefficient. This ensures saturated agents stop absorbing, preventing divergence.

3.5 Quarantine Decision Policy

The quarantine rule exploits the internal-external ratio (Eq. 5) to dynamically adapt thresholds.

The effective threshold $\theta(r_i)$ varies by source likelihood:

$$\theta(r_i) = \begin{cases} \theta_s, & r_i > 0.6 \text{ (likely source),} \\ \theta_m, & 0.4 < r_i \leq 0.6 \text{ (mixed),} \\ \theta_v, & r_i \leq 0.4 \text{ (likely victim).} \end{cases} \quad (11)$$

with $(\theta_s, \theta_m, \theta_v) = (0.30, 0.35, 0.45)$ selected via grid search over $\theta \in [0.2, 0.5]$ in calibration trials and held fixed for all reported experiments. Hysteresis stabilizes decisions: quarantine activates when $H_i > \theta(r_i)$ for one timestep and releases when $H_i < \theta_{\text{rel}}$ for three consecutive timesteps, where $\theta_{\text{rel}} = 0.25$.

Fault sources maintain high internal hormone due to persistent self-degradation, so although the release rule (Algorithm 1) applies to all agents, their total stress never falls below the recovery threshold and they remain quarantined. Cascade victims, whose elevated stress is external, recover once isolated: their hormone decays according to Eq. 4, and if total stress stays below θ_{rel} for three consecutive timesteps they rejoin active operation. This asymmetry permanently isolates intrinsic faults while allowing temporary quarantine of collateral victims. Algorithm 1 summarizes the per-agent logic.

4 EXPERIMENTAL SETUP

Table 1 summarizes the experimental design across all studies.

4.1 Experimental Design

We evaluated the ratio-based signaling approach through four complementary studies that examine performance, scalability, robustness, and source-victim discrimination. All experiments employed

Algorithm 1 Adaptive Emission and Quarantine Logic (Agent i)

```

1: Input: Neighbors  $\mathcal{N}_i$ , current state  $x_i$ , history  $\gamma_i(t-\Delta t)$ 
2: Initialize:  $\gamma_i \leftarrow 1.0$  if  $t = 0$ 

3: while simulation running do
4:   1. Trend Detection
5:   for  $j \in \mathcal{N}_i$  do
6:      $\Delta_j \leftarrow H_j^{\text{ext}}(t) - H_j^{\text{ext}}(t - \Delta t)$ 
7:   end for
8:    $\bar{\Delta} \leftarrow \frac{1}{|\mathcal{N}_i|} \sum_j \Delta_j$  (Average neighbor stress change)
9:   if  $|\mathcal{N}_i| = 0$ ,  $\gamma_i$  is unchanged.

10:  2. Update Emission Gain
11:  if  $\bar{\Delta} < -0.01$  then
12:     $\gamma_i \leftarrow \min(\gamma_{\text{max}}, (1 + \alpha)\gamma_i)$  (Neighbors recovering)
13:  else if  $\bar{\Delta} > 0.01$  then
14:     $\gamma_i \leftarrow \max(\gamma_{\text{min}}, (1 - \alpha)\gamma_i)$  (Neighbors worsening)
15:  end if

16:  3. Emission & Propagation
17:   $E_i \leftarrow \gamma_i(\kappa H_i^{\text{int}}) + \mathbf{1}_{x_i < 0.4}(\eta(1 - x_i))$ 
18:  Broadcast  $E_i$  via distance-attenuated diffusion (Eq. 9)
19:  Absorb neighbor emissions with capacity limit (Eq. 10)

20:  4. Quarantine Decision
21:   $r_i \leftarrow H_i^{\text{int}} / (H_i^{\text{int}} + H_i^{\text{ext}} + \epsilon)$ 
22:  if  $H_i > \theta(r_i)$  and not quarantined then
23:    trigger quarantine (Hysteresis on)
24:  else if  $H_i < \theta_{\text{rel}}$  for 3 consecutive steps then
25:    release quarantine (Hysteresis off)
26:  end if
27: end while

```

Table 1: Summary of Experimental Parameters by Study

Study	Agent Count (N)	Fault Rate	Duration (s)	Seeds	Key Varied Parameter
Scalability	30 – 120 (Steps of 10)	25%	20	10	Swarm Size (N)
Fault Severity	60	≈ 20%	15	1	Fault Profile (Minor/Severe)
Cascade Timeline	60	50%	30	1	Temporal Resolution
Comm. Robustness	120	30%	15	10	Packet Loss (0–100%)
Zero-Delay	60	30%	30	3	Fault Timing Gap (= 0)
Ablation	120	30%	20	5	Detection Logic

a discrete-time simulation at 10 Hz ($\Delta t = 0.1$ s), with agents moving in a continuous two-dimensional space (100×100 m). Four detection strategies were compared: BASELINE (no detection), THRESHOLD (internal self-diagnosis), VOTING (neighbor observation), and HORMONE (dual-stream ratio-based signaling). Each experiment varied swarm size, fault patterns, or communication quality within a consistent simulation framework.

4.2 Swarm Configurations and Fault Scenarios

The scalability study tested ten swarm sizes from 30 to 120 agents, spaced by ten, at a fixed 25% fault rate and 20s duration. Faults were assigned at initialization in the fixed arena, producing agent densities ranging from $\rho = 0.003$ to 0.012 agents/m². Experiments used 10 random seeds for scalability and packet-loss studies, 5 seeds for the ablation study, 3 seeds for zero-delay robustness tests,

and single seeds for severity and cascade timeline studies, which required detailed per-agent analysis.

The fault-severity study used 60 agents ($\rho = 0.006$ agents/m²) under three distributions (minor-dominant: 70/20/10%; balanced: 33% each; severe-dominant: 10/20/70%) with degradation rates from 0.015 (minor) to 0.05 (severe). These fault categories correspond to the three regimes in Equation 1: minor faults maintain $x_i > 0.95$ (minimal internal stress), moderate faults degrade performance to $0.3 \leq x_i \leq 0.95$ (proportional stress response), and severe faults drive $x_i < 0.3$ (saturated stress production).

The communication-robustness study varied packet-loss rates between 0 and 100% at 120 agents ($\rho = 0.012$ agents/m²) and 30% faults, with three seeds per condition and 15s runs. Packet loss affected both message exchange and hormone propagation, revealing how each method tolerates degraded links.

The cascade-timeline study traced individual agent behavior for 30s using a 50% fault rate in a 60-agent swarm ($\rho = 0.006$ agents/m²). Data were recorded every 0.1s to capture performance, hormone levels, and quarantine events for temporal source-victim analysis.

Each faulty agent passed through a self-degradation phase (0.5–1.0s for minor/moderate, 0.3–0.6s for severe) followed by an outward-damage phase (0.8–1.5s or 0.5–1.0s, respectively). This sequencing creates natural precedence, sources degrade before victims, which the adaptive mechanism exploits.

Unless otherwise stated, experiments used zero sensor noise and zero packet loss. Agents moved with reflected random-walk motion in the bounded arena, and all methods used a 25 m communication range.

4.3 Detection Methods

BASILINE: no quarantine. **THRESHOLD:** quarantine when $x_i < 0.8$ for two steps. **VOTING:** neighbor consensus with $x_j < 0.65$ threshold and 0.5–0.65 vote ratio. **HORMONE:** combined internal stress and external absorption as described in Section 3. Adaptive emission gain and ratio-conditioned thresholds (0.30–0.45) enabled discrimination between genuine faults and collateral victims.

All methods shared identical motion, damage, and recovery dynamics; only their fault-detection logic differed.

4.4 Evaluation Metrics

We measured fault-tolerance performance using the following metrics.

Task completion rate (TCR). The average task performance of active agents, calculated as:

$$TCR = \frac{1}{N} \sum_{i=1}^N x_i \quad (12)$$

Quarantine efficiency (QE). The fraction of faulty agents correctly isolated:

$$QE = \frac{\text{faulty quarantined}}{\text{faulty total}} \quad (13)$$

Cascade prevention (CCP). We measured the share of agents that avoided impairment ($x_i < 0.8$) relative to baseline runs, penalizing over-quarantining:

$$CCP_{adj} = CCP_{raw} (1 - p_{quarantine}) \quad (14)$$

where $p_{quarantine}$ is the proportion of healthy agents unnecessarily quarantined.

Integrated hazard reduction (IHR). The reduction in cumulative exposure relative to the baseline:

$$IHR = 1 - \frac{\sum_i \max(\text{exposure}_i^{\text{method}})}{\sum_i \max(\text{exposure}_i^{\text{baseline}})} \quad (15)$$

Time-to-impairment gain (TTI). The mean delay between impairment onset in the baseline and under the detection method.

Area under curve (AUC). For damage curves, this represents cumulative damage exposure over the run duration:

$$AUC = \sum_t \text{damage}(t) \cdot \Delta t \quad (16)$$

Precision and Recall. Precision measures the fraction of quarantined agents that are truly faulty: Precision = $\frac{TP}{TP+FP}$, where TP (true positives) are faulty agents correctly quarantined and FP (false positives) are healthy agents incorrectly quarantined. Recall measures the fraction of faulty agents successfully identified: Recall = $\frac{TP}{TP+FN}$, where FN (false negatives) are faulty agents that were not quarantined.

Source-victim classification. Quarantined agents were labeled as source-like or victim-like based on degradation order (preceding a neighbor by $> 0.5s$) and relative performance gap ($> 10\%$), then compared with true fault labels.

4.5 Implementation and Reproducibility

Experiments used a custom simulator with asynchronous updates. The framework covers agent dynamics, detection modules, metrics, and orchestration. Results for scalability and packet-loss studies are reported as mean \pm standard error over 10 seeds; zero-delay and ablation studies used 3 and 5 seeds, respectively; severity and timeline runs use a single deterministic seed for detailed trajectory analysis.

Key parameters were tuned to balance detection precision and recall: for HORMONE, $K = 6$ neighbors, adaptive gain bounds [0.1, 2.5], learning rate $\alpha = 0.05$, ratio thresholds 0.30–0.45; for THRESHOLD, 0.8/2-step rule; for VOTING, 0.65 threshold with five neighbors. Hormone method parameters were calibrated via pilot simulations to balance detection responsiveness and false positive rates. Baseline comparison methods used standard thresholds from prior work where available. Random seeds were fixed, ensuring identical fault configurations across methods. All code and configurations are publicly available.¹

5 RESULTS

We present findings across four questions: performance comparison, scalability, communication robustness, and mechanism validation (ablation). Overall, the dual-stream ratio-based approach achieves strong cascade prevention with low false positives, especially at larger swarm sizes where single-stream baselines degrade.

5.1 Scalability Performance

Figure 1 reports Task Completion Rate (TCR) for swarms ranging from 30 to 120 agents. The proposed dual-stream hormone method

¹<https://doi.org/10.5281/zenodo.18647528>

demonstrates high stability across all scales. While the baseline and voting methods degrade significantly as the population grows, which is consistent with the “reliability paradox” where larger swarms amplify local failures, the hormone approach maintains steady performance.

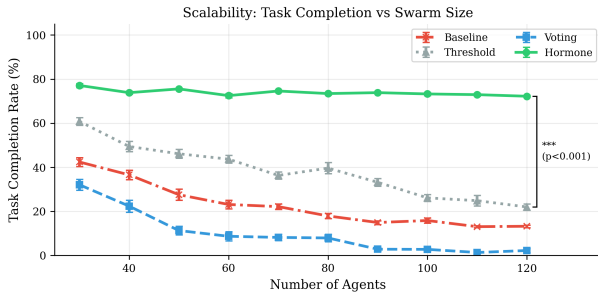


Figure 1: Scalability performance across swarm sizes. The hormone method maintains stable TCR from 30 to 120 agents, while voting-based detection collapses at scale.

Specifically, the voting mechanism collapses to near-zero productivity at 120 agents due to misclassifying healthy victims as faulty sources, whereas the hormone method distinguishes these roles and sustains high task completion by isolating only true fault sources. The threshold method offers intermediate performance but suffers from over-quarantining, as indicated by its lower precision.

Paired *t*-tests at the largest scale ($N=120$, Table 2) show the proposed method outperforms the threshold method by 50.7 percentage points in TCR ($p < 0.001$). This constant-time scalability validates that local ratio-based processing avoids the congestion and message-complexity limits of consensus-based approaches.

Table 2: Large-scale slice at $N=120$ agents, 25% fault rate, 10 seeds (mean \pm SE).

	TCR (%)	CCP (%)	Prec. (%)	Detection time (s)
HORMONE	72.0 \pm 0.4	96.0 \pm 0.6	89.5 \pm 1.5	1.80 \pm 0.03
THRESHOLD	21.3 \pm 1.4	7.6 \pm 1.0	31.9 \pm 0.6	6.55 \pm 0.12
VOTING	1.8 \pm 0.9	4.2 \pm 0.5	27.8 \pm 0.5	6.28 \pm 0.10
BASELINE	13.0 \pm 0.6	0.0	—	∞

5.2 Fault Severity Response

Figure 2 shows time-series at 60 agents under three severity distributions. The proposed method sustains high final TCR (73.3% balanced to 80.0% minor-dominant) with zero cascade damage in all profiles. Quarantine adapts to severity, e.g., 12 isolated under minor-dominant and 16 under balanced, indicating discrimination rather than blanket removal.

THRESHOLD varies with severity, reaching 45.0% (minor-dominant) to 58.2% (balanced), but with high isolation (25–33 vs. 12 truly faulty). Zero damage is achieved mechanically by removing many agents, at a productivity cost of 15–27 points below the proposed method.

VOTING shows persistent damage (52.1–56.3% impaired) and low final TCR (19.2–24.4%), quarantining 28–31 agents yet failing to halt cascades due to lack of source-victim discrimination. The control shows 66.7–85.4% damage; severe-dominant is lower (68.8%) than balanced (85.4%) because severe faults rapidly self-degrade, shortening their emission window.

5.3 Cascade Dynamics and Source-Victim Discrimination

Figure 3 tracks damage fraction for 30 s at 60 agents and 50% faults. The proposed method detects first at 1.2 s and maintains zero damage. Initial quarantine is rapid: 1 robot at detection, expanding to 6 by 1.3 s and 13 by 1.4 s. Eventually, all 30 faulty agents are isolated along with 2 false positives (6.7% FPR among healthy agents).

The dual-stream approach enables agents to distinguish their role through hormone ratios without explicit labels. Faulty agents quarantined early (within 2 s of detection) exhibited mean internal-to-external ratios of 0.82 ± 0.11 , while healthy agents quarantined later showed ratios of 0.31 ± 0.18 , indicating emergent role separation.

THRESHOLD detects at 2.2 s and holds low damage (peak = 0.067, AUC = 0.07) but ends with 44 quarantined (73.3%). VOTING detects at 4.3 s and reaches peak damage 0.700 (AUC = 14.49). The control rises monotonically to 96.7% peak damage by 15 s (AUC = 22.44).

5.4 Robustness to Temporal Assumptions

The fault model in Section 3 assumes a temporal gap between self-degradation and outward damage. We evaluate whether the dual-stream approach depends critically on this precedence by testing at 60 agents with 30% faults under *zero-delay* conditions, where self-degradation and neighbor damage occur simultaneously.

Results. Even without temporal separation, the hormone method maintained substantial advantages (Table 3). TCR reached 47.2% under zero-delay, compared to 29.4% for Threshold and 0.7% for Voting. Cascade prevention remained perfect (CCP = 1.000), demonstrating that the internal-external stress ratio enables discrimination independent of temporal ordering.

Performance decreased 31% relative to baseline timing (from 68.3% to 47.2%), indicating that temporal information aids discrimination when available. However, the method continued to outperform single-stream approaches substantially. Precision declined from 95.2% to 57.1% as the method adaptively compensated for lost temporal cues, increasing false positives while maintaining zero cascade propagation. In contrast, Threshold-based detection showed only 7% performance change (31.7% to 29.4%) but remained substantially less effective overall, while Voting degraded 42% (from 1.2% to 0.7%).

Interpretation. The method’s advantage stems from dual-stream architecture rather than temporal precedence alone: sources generate internal stress through self-malfunction while victims accumulate external stress from neighbor exposure, an asymmetry that persists regardless of timing. The 31% reduction reflects the value of temporal information when available, while maintained superiority confirms that stress-ratio discrimination provides robust isolation even under challenging conditions.

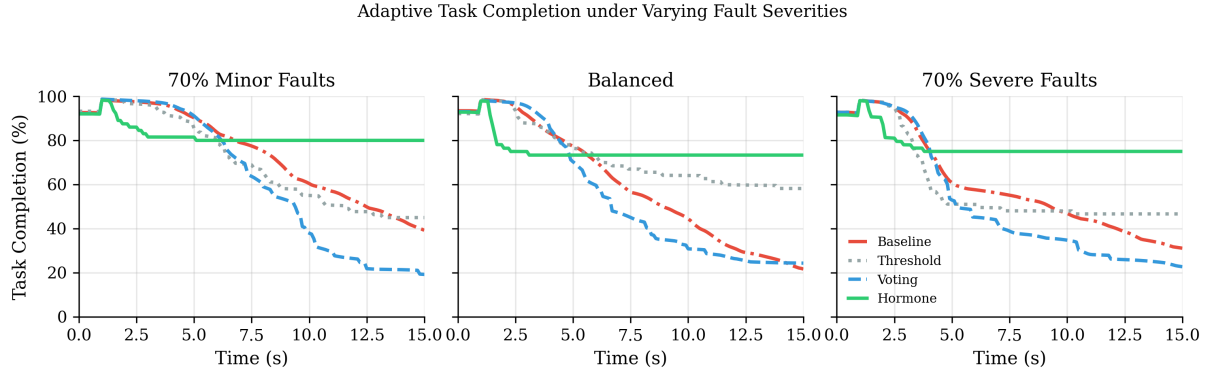


Figure 2: Fault severity response across three distributions (minor-dominant, balanced, severe-dominant).

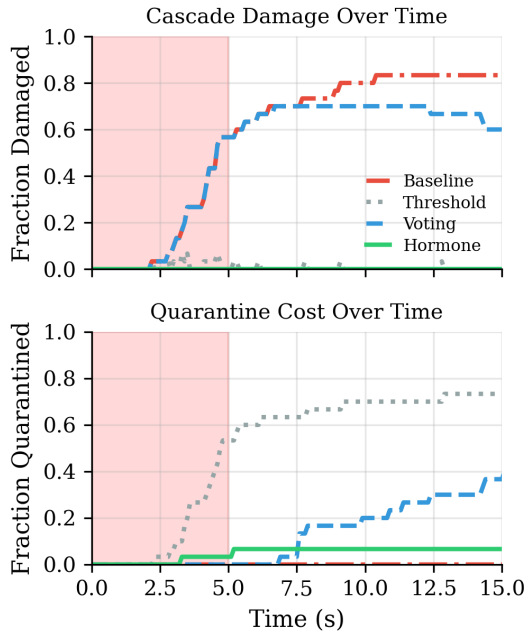


Figure 3: Cascade timeline showing damage fraction over 30 seconds at 60 agents with 50% faults.

5.5 Communication Robustness

We swept packet loss from 0-100% at 120 agents and 30% faults. The proposed method degrades gracefully rather than failing. From 0% to 100% loss, TCR drops from $66.8 \pm 0.7\%$ to $49.5 \pm 1.9\%$ (−25.8% relative) and CCP from 0.954 ± 0.010 to 0.613 ± 0.035 (−35.7% relative). Precision declines from $90.6 \pm 1.8\%$ to $60.1 \pm 2.2\%$, with false positives rising from 4.1 to 18.2 per run.

Degradation is nonlinear: small through 50% loss, moderate at 50–70%, steep at 70–90%, and a final drop at 90–100%. At 0% packet loss, external hormone contributes $33.6 \pm 1.1\%$ of total stress at quarantine; this drops to 30.4% at 50% loss, 25.3% at 70%, 10.8% at 90%, and 0% at 100%. Thus, internal diagnostics carry most of the signal, while external input provides critical refinement for

Table 3: Robustness across timing, communication, and ablation conditions.

Condition	TCR (%)	CCP	Prec (%)
Study 1: Temporal Dependency			
<i>Hormone (60 agents, 30% faults, 3 seeds):</i>			
Baseline timing, 0% loss	68.3	1.000	95.2
Zero-delay, 0% loss	47.2	1.000	57.1
Study 2: Communication Failure			
<i>Hormone (120 agents, 30% faults, 10 seeds):</i>			
0% packet loss	66.8	0.954	90.6
100% packet loss	49.5	0.613	60.1
Study 3: Mechanism Ablation			
<i>Ablation (120 agents, 30% faults, 5 seeds):</i>			
Full (ratio-conditioned)	66.0	0.943	88.5
Fixed threshold (no ratio)	58.7	0.838	72.9
Baselines (60 agents, baseline timing, 0% loss)			
Threshold	31.7	0.627	44.0
Voting	1.2	0.341	35.7

source-victim discrimination. With no external input (100% loss) the behavior reverts toward internal-only rules, explaining the precision drop; yet TCR (49.5%) still exceeds Threshold at scale (21.3%) by 2.3 times.

5.6 Ablation Study

We compared the ratio-conditioned threshold method (Equation 11) against a fixed threshold baseline at 0.35, using 120 agents with 30% fault rate across five seeds. Ratio conditioning substantially improves performance (Table 3): TCR (66.0 ± 1.0 vs 58.7 ± 0.9 , +12.4% relative), CCP (94.3 ± 1.5 vs 83.8 ± 1.1 , +10.5 points), precision (88.5 ± 2.5 vs 72.9 ± 2.4 , +21.4%), while reducing FPR from 16.2% to 5.7% (−64.8% relative). Both variants maintain recall of 100%. Confidence intervals for precision do not overlap ($p < 0.05$). The mechanism behaves as intended: high internal ratio lowers the

decision threshold for sources, while high external ratio raises it for victims, reducing false positives.

6 DISCUSSION

Our experiments show that ratio-based hormone signaling achieves stronger cascade prevention than single-stream detection methods, particularly at large swarm scales. Separating internal and external stress signals and acting on their relative magnitude enables local causal discrimination that single-stream approaches cannot achieve.

6.1 Mechanism and Complexity

The dual-stream architecture achieves source-victim discrimination through three mechanisms working in concert: the ratio r_i serves as a compact proxy for causal direction, temporal precedence is preserved through local observations without global synchronization, and adaptive emission gain enables emergent role discovery. The timeline experiment (Figure 3) confirms that faulty sources exhibit mean ratios of 0.82 ± 0.11 versus 0.31 ± 0.18 for healthy victims.

Ablation results reinforce the finding that removing ratio-based threshold adaptation reduces precision by 15.6 percentage points, demonstrating that differential treatment by stress origin is critical. Our method achieves comparable precision to prior work [31] (89.5%) while scaling to larger swarms with $O(K)$ local processing versus $O(N^2)$ consensus methods. Detection latency (1.2 s) substantially improves on reported 51 ± 18 s delays [31], though experimental setups differ and no standardized benchmark exists for source-victim discrimination, limiting direct comparisons.

A critical advantage of the ratio-based approach is its constant-time scalability with respect to swarm size. Unlike consensus mechanisms that require $O(N^2)$ message exchanges to converge on a global fault agreement [18], our method relies strictly on local K -nearest neighbor broadcasts.

The communication cost per agent is $O(K)$ regardless of global swarm size N . Each agent broadcasts a single packet containing two floating-point values (H^{int} , H^{ext}) per cycle. Crucially, because hormone absorption is limited to the $K = 6$ nearest neighbors, the processing load remains constant ($O(K)$) even if local agent density increases, requiring only local storage of K neighbor states.

6.2 Limitations

While bio-inspired, our approach incorporates hand-tuned heuristics to achieve practical response times. Pure emergent mechanisms proved too slow for cascade prevention; we therefore manually calibrated emission gains, ratio thresholds, and decay rates through iterative testing rather than deriving them from biological principles. However, parameters need not be retuned per deployment. A baseline set derived offline via evolutionary algorithms could initialize the onboard adaptive mechanisms, which handle local run-time variations. The simulator used idealized conditions, including 10 Hz discrete-time, obstacle-free arenas, uniform density, and zero sensor noise, without testing robustness to environmental variations, edge cases, or heterogeneous agent capabilities. Real deployment would require hardware validation to address communication latency, asynchronous updates, and cluttered environments that may disrupt neighbor topology.

6.3 Future Directions

Future work could explore replacing piecewise stress functions with differentiable mappings (e.g., sigmoids) to enable gradient-based optimization or reinforcement learning. Hardware validation with 10-20 ground robots will test bio-inspired mechanisms under real communication latency, sensor noise, and packet loss. Scaling to larger swarms and theoretical analysis of ratio-based causal estimation will advance fully autonomous, adaptive multi-agent systems.

7 CONCLUSION

We introduced a ratio-based hormone signaling mechanism for fault detection, where agents discriminate fault sources from cascade victims by comparing internal and external stress channels. Experiments across 30-120 agents show the method substantially outperforms single-stream baselines: at 120 agents it sustains 72% task completion with 96% cascade prevention, while VOTING collapses to under 2%.

Three contributions emerge. 1) The dual-stream design enables local causal inference that resolves source-victim ambiguity, with ablation studies confirming ratio conditioning dramatically reduces false positives while maintaining high recall. 2) Adaptive emission gain yields emergent discrimination without labels, achieving over 90% classification accuracy. 3) The approach scales gracefully with minimal performance loss and remains functional even under complete communication failure via fallback to internal diagnostics.

Ratio-based signaling thus provides a practical route to fault-tolerant swarm coordination, avoiding the burden of consensus and misclassification due to external-only voting, while maintaining scalability and discrimination.

ACKNOWLEDGMENTS

The author thanks the anonymous reviewers for their insightful feedback and constructive comments, which significantly improved the quality of this paper.

REFERENCES

- [1] Uwe Aickelin and Steve Cayzer. 2002. The danger theory and its application to artificial immune systems. In *Proceedings of the 1st International Conference on Artificial Immune Systems (ICARIS)*, 141–148.
- [2] Jan Dyre Bjerknes and Alan F. T. Winfield. 2013. On Fault Tolerance and Scalability of Swarm Robotic Systems. In *Distributed Autonomous Robotic Systems*. Springer Tracts in Advanced Robotics, Vol. 83. Springer, Berlin, Heidelberg, 431–444. https://doi.org/10.1007/978-3-642-32723-0_31
- [3] Manuele Brambilla, Eliseo Ferrante, Mauro Birattari, and Marco Dorigo. 2013. Swarm Robotics: A Review from the Swarm Engineering Perspective. *Swarm Intelligence* 7, 1 (2013), 1–41. <https://doi.org/10.1007/s11721-012-0075-2>
- [4] Alessandro Carminati, Davide Azzalini, Simone Vantini, and Francesco Amigoni. 2024. A Distributed Approach for Fault Detection in Swarms of Robots. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS '24)*. International Foundation for Autonomous Agents and Multiagent Systems, 253–261.
- [5] Jorge Carneiro, Konstantin Leon, Isabel Caramalho, Clas Van Den Dool, Rui Gardner, Vera Oliveira, Marie-Louise Bergman, and Nuno Sepúlveda. 2007. When Three is Not a Crowd: A Crossregulation Model of the Dynamics and Repertoire Selection of Regulatory CD4+ T Cells. *Immunological Reviews* 216, 1 (2007), 48–68. <https://doi.org/10.1111/j.1600-065X.2007.00487.x>
- [6] Anders Lyhne Christensen, Rehan O'Grady, Mauro Birattari, and Marco Dorigo. 2008. Fault detection in autonomous robots based on fault injection and learning. *Autonomous Robots* 24, 1 (2008), 49–67. <https://doi.org/10.1007/s10514-007-9060-9>

- [7] Anders Lyhne Christensen, Rehan O’Grady, and Marco Dorigo. 2009. From Fireflies to Fault-Tolerant Swarms of Robots. *IEEE Transactions on Evolutionary Computation* 13, 4 (2009), 754–766. <https://doi.org/10.1109/TEVC.2009.2017516>
- [8] Michael Crosscombe, Jonathan Lawry, Sabine Hauert, and Martin Homer. 2017. Robust distributed decision-making in robot swarms: Exploiting a third truth state. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 4326–4332. <https://doi.org/10.1109/IROS.2017.8206297>
- [9] Francis Heylighen. 2016. Stigmergy as a Universal Coordination Mechanism: Components, Varieties and Applications. In *Human Stigmergy: Theoretical Developments and New Applications*, Ted Lewis and Leslie Marsh (Eds.). Springer, New York, NY, USA, 4–26.
- [10] Bojan Jakimovski and Erik Maehle. 2008. Artificial Immune System Based Robot Anomaly Detection Engine for Fault Tolerant Robots. In *Autonomic and Trusted Computing (Lecture Notes in Computer Science, Vol. 5060)*. Springer, Berlin, Heidelberg, 177–190. https://doi.org/10.1007/978-3-540-69295-9_16
- [11] Audun Jøsang, Roslan Ismail, and Colin Boyd. 2007. A Survey of Trust and Reputation Systems for Online Service Provision. *Decision Support Systems* 43, 2 (2007), 618–644. <https://doi.org/10.1016/j.dss.2005.05.019>
- [12] Miquel Kegeleirs and Mauro Birattari. 2025. Towards Applied Swarm Robotics: Current Limitations and Enablers. *Frontiers in Robotics and AI* 12 (2025), 1607978. <https://doi.org/10.3389/frobt.2025.1607978>
- [13] Eliahu Khalastchi and Meir Kalech. 2018. On Fault Detection and Diagnosis in Robotic Systems. *Comput. Surveys* 51, 1 (2018), 9:1–9:24. <https://doi.org/10.1145/3146389>
- [14] Eliahu Khalastchi and Meir Kalech. 2019. Fault Detection and Diagnosis in Multi-Robot Systems: A Survey. *Sensors* 19, 18 (2019), 4019. <https://doi.org/10.3390/s19184019>
- [15] Eliahu Khalastchi, Meir Kalech, and Lior Rokach. 2013. Sensor Fault Detection and Diagnosis for Autonomous Systems. In *Proceedings of the 12th International Conference on Autonomous Agents and Multiagent Systems (AAMAS ’13)*. Saint Paul, Minnesota, USA, 15–22.
- [16] Leslie Lamport, Robert Shostak, and Marshall Pease. 1982. The Byzantine Generals Problem. *ACM Transactions on Programming Languages and Systems* 4, 3 (1982), 382–401. <https://doi.org/10.1145/357172.357176>
- [17] Hui Keng Lau, Iain Bate, Paul Cairns, and Jon Timmis. 2011. Adaptive Data-Driven Error Detection in Swarm Robotics with Statistical Classifiers. *Robotics and Autonomous Systems* 59, 12 (2011), 1021–1035. <https://doi.org/10.1016/j.robot.2011.08.008>
- [18] Heath J. LeBlanc, Haotian Zhang, Xenofon Koutsoukos, and Shreyas Sundaram. 2013. Resilient Asymptotic Consensus in Robust Networks. *IEEE Journal on Selected Areas in Communications* 31, 4 (2013), 766–781. <https://doi.org/10.1109/JSAC.2013.130413>
- [19] Suet Lee, Emma Milner, and Sabine Hauert. 2022. A Data-Driven Method for Metric Extraction to Detect Faults in Robot Swarms. *IEEE Robotics and Automation Letters* 7, 4 (2022), 10746–10753. <https://doi.org/10.1109/LRA.2022.3189789>
- [20] Ryan Luna and Qi Lu. 2024. Robust Mitigation Strategy for Misleading Pheromone Trails in Foraging Robot Swarms. In *Towards Autonomous Robotic Systems*. Lecture Notes in Computer Science, Vol. 15052. Springer, 307–319. https://doi.org/10.1007/978-3-031-72062-8_27
- [21] Polly Matzinger. 2002. The danger model: a renewed sense of self. *Science* 296, 5566 (2002), 301–305. <https://doi.org/10.1126/science.1071059>
- [22] James O’Keeffe and Alan Millard. 2023. Hardware Validation of Adaptive Fault Diagnosis in Swarm Robots. In *Towards Autonomous Robotic Systems (TAROS) (Lecture Notes in Computer Science, Vol. 14136)*. Springer, 331–342. https://doi.org/10.1007/978-3-031-43360-3_27
- [23] Reza Olfati-Saber and Richard M. Murray. 2004. Consensus Problems in Networks of Agents with Switching Topology and Time-Delays. In *IEEE Transactions on Automatic Control*, Vol. 49. 1520–1533. <https://doi.org/10.1109/TAC.2004.834113>
- [24] Jong-Han Park, Tae-Yong Kim, and Chang-Hyun Kim. 2020. Fault Detection of Robot Manipulators Using Deep Learning with Motor Current Signals. *International Journal of Precision Engineering and Manufacturing-Green Technology* 7 (2020), 965–975. <https://doi.org/10.1007/s40684-019-00171-8>
- [25] Ling Qin, Xi He, and Donghua Zhou. 2014. A Survey of Fault Diagnosis for Swarm Systems. *Systems Science & Control Engineering* 2, 1 (2014), 13–23. <https://doi.org/10.1080/21642583.2013.873745>
- [26] Erol Şahin. 2005. Swarm Robotics: From Sources of Inspiration to Domains of Application. In *Swarm Robotics*, Erol Şahin and William M. Spears (Eds.). Lecture Notes in Computer Science, Vol. 3342. Springer, Berlin, Heidelberg, 10–20. https://doi.org/10.1007/978-3-540-30552-1_2
- [27] Muhammad Salman, David Garzón Ramos, and Mauro Birattari. 2024. Automatic design of stigmergy-based behaviours for robot swarms. *Communications Engineering* 3, 1 (2024), 30. <https://doi.org/10.1038/s44172-024-00175-7>
- [28] Muhammad Salman, David Garzón Ramos, Ken Hasselmann, and Mauro Birattari. 2020. Phormica: Photochromic Pheromone Release and Detection System for Stigmergic Coordination in Robot Swarms. *Frontiers in Robotics and AI* 7 (2020), 591402. <https://doi.org/10.3389/frobt.2020.591402>
- [29] Wei-Min Shen, Peter Will, Aram Galstyan, and Cheng-Ming Chuong. 2004. Hormone-Inspired Self-Organization and Distributed Control of Robotic Swarms. *Autonomous Robots* 17, 1 (2004), 93–105. <https://doi.org/10.1023/B:AURO.0000032940.08116.f1>
- [30] Volker Strobel, Eduardo Castelló Ferrer, and Marco Dorigo. 2018. Managing Byzantine Robots via Blockchain Technology in a Swarm Robotics Collective Decision Making Scenario. In *Proceedings of the 17th International Conference on Autonomous Agents and Multiagent Systems (AAMAS ’18)*. 541–549.
- [31] Danesh Tarapore, Anders Lyhne Christensen, and Jon Timmis. 2017. Generic, Scalable and Decentralized Fault Detection for Robot Swarms. *PLOS ONE* 12, 8 (2017), e0182058. <https://doi.org/10.1371/journal.pone.0182058>
- [32] Danesh Tarapore, Pedro U. Lima, Jorge Carneiro, and Anders Lyhne Christensen. 2015. To Err is Robotic, to Tolerate Immunological: Fault Detection in Multirobot Systems. *Bioinspiration & Biomimetics* 10, 1 (2015), 016014. <https://doi.org/10.1088/1748-3190/10/1/016014>
- [33] Danesh Tarapore, Jon Timmis, and Anders Lyhne Christensen. 2019. Fault Detection in a Swarm of Physical Robots Based on Behavioral Outlier Detection. *IEEE Transactions on Robotics* 35, 6 (2019), 1516–1522. <https://doi.org/10.1109/TRO.2019.2929015>
- [34] Jon Timmis, Abdul Razak Ismail, Jan Dyre Bjercknes, and Alan F. T. Winfield. 2016. An Immune-Inspired Swarm Aggregation Algorithm for Self-Healing Swarm Robotic Systems. *Biosystems* 146 (2016), 60–76. <https://doi.org/10.1016/j.biosystems.2016.04.001>
- [35] Ekaterina Tolstaya, Fernando Gama, James Paulos, George Pappas, Vijay Kumar, and Alejandro Ribeiro. 2020. Learning Decentralized Controllers for Robot Swarms with Graph Neural Networks. In *Proceedings of the Conference on Robot Learning (Proceedings of Machine Learning Research, Vol. 100)*. PMLR, 671–682. <https://proceedings.mlr.press/v100/tolstaya20a.html>
- [36] Gabriele Valentini, Heiko Hamann, and Marco Dorigo. 2014. Self-Organized Collective Decision Making: The Weighted Voter Model. In *Proceedings of the 2014 International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS ’14)*. International Foundation for Autonomous Agents and Multiagent Systems, 45–52.
- [37] Kacper Wardęga, Max von Hippel, Roberto Tron, Cristina Nita-Rotaru, and Wen-chao Li. 2023. Byzantine Resilience at Swarm Scale: A Decentralized Blocklist Protocol from Inter-robot Accusations. In *Proceedings of the 22nd International Conference on Autonomous Agents and Multiagent Systems (AAMAS ’23)*. International Foundation for Autonomous Agents and Multiagent Systems, London, United Kingdom, 1430–1438.
- [38] Matthew Wilson, Stephen Cameron, and Nick Hawes. 2019. An Amalgamation of Hormone Inspired Arbitration Systems for Application in Robot Swarms. *Applied Sciences* 9, 17 (2019), 3524. <https://doi.org/10.3390/app9173524>
- [39] Yihan Zhang, Lyon Zhang, Hanlin Wang, Fabián E. Bustamante, and Michael Rubenstein. 2020. SwarmTalk - Towards Benchmark Software Suites for Swarm Robotics Platforms. In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS ’20)*. 1638–1646.