

Trust-Aware Behavior Reflection for Robot Swarm Self-Healing*

Rui Liu
Robotics Institute
Carnegie Mellon University
rui.liu2@cs.cmu.edu

Fan Jia
Robotics Institute
Carnegie Mellon University
fanjia1996@gmail.com

Wenhao Luo
Robotics Institute
Carnegie Mellon University
whluo12@gmail.com

Meghan Chandarana
Dept. of Mech. Eng.
Carnegie Mellon University
mchandar@cmu.edu

Changjoo Nam
Ctr. for Robotics Research
Korean Inst. of Tech.
cnam@kist.re.kr

Michael Lewis
School of Comp. and Info.
University of Pittsburgh
ml@sis.pitt.edu

Katia Sycara
Robotics Institute
Carnegie Mellon University
katia@cs.cmu.edu

ABSTRACT

The deployment of robot swarms is influenced by real-world factors, such as motor issues, sensor failure, and wind disturbances. These factors cause the appearance of faulty robots. In a decentralized swarm, sharing incorrect information from faulty robots will lead to undesired swarm behaviors, such as swarm disconnection and incorrect heading directions. We envision a system where a human operator is exerting supervisory control over a remote swarm by indicating changes in trust to the swarm via a "trust-signal". By correcting faulty behaviors, trust between the human and the swarm is maintained to facilitate human-swarm cooperation. In this research, a trust-aware behavior reflection method – **Trust-R** – is designed based on a weighted mean subsequence reduced algorithm (WMSR). By using Trust-R, detected faulty behaviors are automatically corrected by the swarm in a decentralized fashion by referring to the motion status of their trusted neighbors and isolating failed robots from the others. Based on real-world scenarios, three types of robot faults – degraded performance caused by motor wear, abnormal motion caused by system uncertainty and motion deviation caused by an external disturbance such as wind – were simulated to test the effectiveness of Trust-R. Results show that Trust-R is effective in correcting swarm behaviors for swarm self-healing.

KEYWORDS

Trust-R; WMSR; Trust; Behavior Reflection; Swarm Self-Healing

ACM Reference Format:

Rui Liu, Fan Jia, Wenhao Luo, Meghan Chandarana, Changjoo Nam, Michael Lewis, and Katia Sycara. 2019. Trust-Aware Behavior Reflection for Robot Swarm Self-Healing. In *Proc. of the 18th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2019), Montreal, Canada, May 13–17, 2019*, IFAAMAS, 9 pages.

1 INTRODUCTION

Robot swarms use simple, local control laws to achieve a desired global emergent behavior over time. In using only local information, these systems are flexible to changes in the environment conditions and swarm size. The scalability of robot swarms leads to their use

*This work has been funded by AFOSR award FA9550-15-1-0442 and AFOSR/AFRL award FA9550-18-1-0251.

Proc. of the 18th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2019), N. Agmon, M. E. Taylor, E. Elkind, M. Veloso (eds.), May 13–17, 2019, Montreal, Canada. © 2019 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

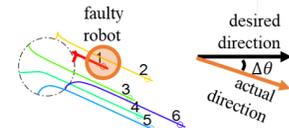


Figure 1: Undesired swarm flocking caused by faulty robots sharing incorrect information with their neighbors.

in a variety of applications, such as search and rescue [1], disaster relief [17], and environmental monitoring [8].

In many of these applications, human operators use supervisory control interfaces to remotely adjust mission goals and requirements. Human trust in the swarm is critical for effective human-swarm cooperation [21][6]. When a swarm is untrusted due to its unsatisfied performance, unnecessary interventions, such as diverting swarms' paths through the mission space and assigning additional intermediary spots to pass by, from human operators will increase. Unnecessary interventions are often time-consuming and require the attention of a group of robots to receive new inputs from the human operator, leading to delayed goal attainments or a decreased efficiency in human-swarm cooperation [18]. While, when a swarm is trusted, human operators are more willing to rely on automation to perform tasks, thereby reducing unnecessary interventions [12].

However, real-world faults, such as motor degradation, sensor failure or wind disturbance, make maintaining trust between humans and swarms challenging [13][5]. These factors can cause undesirable and uncontrollable robot behavior, such as a robot or group of robots getting disconnected as shown by robot 1 in Figure 1. In addition, faulty robots may share incorrect information with other members of the swarm leading to incorrect behaviors of the swarm as a whole (Figure 1). The unpredictable nature of these real-world faults can reduce human trust in the reliability of the swarm. Unlike centralized systems where these faults can be directly corrected by centralized control commands, decentralized systems have a difficult time identifying these issues and are more susceptible to the effect faults have on the overall system behavior [7][2].

We envision a system where the human operator assumes a supervisory control role over the remote swarm. In such a scenario, the tolerance levels for faulty and failed robot behaviors are prescribed by the human's interpretation of the application requirements. Using these prescribed tolerance levels, the human operator

monitors the performance of the swarm. It has been found in the literature that if the swarm performance decreases e.g. due to faulty robots, the operator's trust also decreases [18][12]. The operator indicates changes in trust (in particular trust decrease) to the swarm via, a so-called "trust-signal". The trust signal also contains information about expected swarm behavior and any current swarm deviations from that behavior. Therefore, it is critical to correct these faulty behaviors in a timely manner to ensure a high level of trust is maintained between the human and swarm. In this paper, a decentralized trust-aware behavior reflection (Trust-R) method is proposed to correct swarms' faulty behaviors. With Trust-R correction, swarms with faulty robots can repair themselves to attain human-assigned goals and still receive high-level human trust with minimum unnecessary interventions. The contributions are two-fold. First, a novel trust-aware reflection (Trust-R) algorithm is presented to help robots with a semi-automated, self-behavior diagnosis. Instead of merely judging whether it is normal or abnormal, each robot identifies its level of faultiness from a human trust perspective. Second, a reflective correction method is designed. Robots leverage the communicated levels of faultiness from their neighbors to update their motion status using only the information received from their trusted neighbors. As a result, information exchange is encouraged with trusted robots and discouraged with untrusted robots resulting in behavior correction of the whole swarm.

In this paper, similar to previous self-healing work [16], we assume the faulty robots are a minority. Thus, it is possible to correct the swarm behavior by following the trusted robots.

2 RELATED WORK

Recent swarm self-healing research has focused on simulating faulty robot behaviors. In [22], faulty robots were defined and simulated as the robots not located in the desired position defined by a swarm's network topology. In [16], faulty robots were defined and simulated as robots with incorrect heading directions. By comparing observed behavior with ideally designed robot behaviors, the faulty behaviors were detected and corrected. These two methods are effective in swarm healing. However, neither considers the presence of a range of real-world factors, such as sensor failures or wind disturbance. These factors can greatly influence swarm behaviors in the real world and cannot be ignored. In our trust-R method, frequently observed faults, such as degraded motors, system uncertainty, and wind disturbance, were considered, showing our method to be general and suitable for real world environments. Moreover, Trust-R has the potential to support adaptive swarm deployments.

Additional emphasis has been placed on passive healing strategies that increase the swarm resilience. [16][23][14][9][10] increased the tolerance of faulty robots in the swarm by encouraging larger network robustness. As a result the negative influence of faulty robots on the swarm can be limited. Although this method is able to dilute the negative influence, the passive strategies usually require relatively high swarm connectivity and required specification of tolerable speed values which may be difficult to specify in advance. Therefore, it is necessary to actively correct these faulty behaviors when they appear. When faulty robots mislead the normal robots in a swarm, the proposed Trust-R method corrects the faulty robots by

referring to those who are trusted. The failed robots are isolated by lowering the communication quality between the failed robots and the others. Trust-R can correct faulty behaviors that cannot be prevented by other techniques that aim to increase swarm resilience. Combining the passive prevention methods and the Trust-R active correction method, would make a swarm more robust.

In [3][20], faulty behaviors, such as fixed heading directions, were identified by tracking the temporal motion trajectory of a robot. However, the severity (i.e., effect on the whole swarm's behavior) of the faulty robots in different scenarios was not assessed. Without assessment of the severity, appropriate control strategies are difficult to design. In the proposed Trust-R method, degree of severity is diagnosed. Different control strategies, such as "take trusted robots as reliable information sources", "correct the robots with mild fault", "isolate the failed robots from the swarm", are then designed for exchanging information and adjusting connectivity among robots.

3 ILLUSTRATIVE SCENARIO FOR SWARM SELF-HEALING: UNTRUSTED FLOCKING

The task scenario for the swarm in an obstacle free environment is selected as distributed, biased swarm flocking. In these scenarios, controllers are tasked with ensuring the required coordinated motion necessary to reach a desired motion consensus. During each update step, robots estimate the global variable by exchanging and averaging the motion statuses of their neighbors. Using local robot interactions and updates, agreed global variables, such as motion direction and velocity, will be achieved to guide a consensus motion for the whole swarm.

Consider a robot swarm of n holonomic robots with positions $X_i \in \mathbb{R}^3$, where $X_i = (\mathbf{x}_{i,h}, \mathbf{x}_{i,v}, \theta_i)$. Each robot is assigned a unique identifier (UID) $i \in \{1, 2, \dots, n\}$. The communication graph is given by $G = (\mathcal{V}, \mathcal{E})$. Every node $v \in \mathcal{V}$ represents a robot. Every robot i only communicates with its direct neighbors $j \in N_i$, where N_i is the set of all neighbors of i within the communication radius, R . If robot j is a neighbor of i , then edge $(v_i, v_j) \in \mathcal{E}$. The connectivity graph is connected and undirected (i.e., $(v_i, v_j) \in \mathcal{E} \Rightarrow (v_j, v_i) \in \mathcal{E}$). The dynamic model [11] for each robot is defined as follows. A robot i is controlled by the linear velocity \mathbf{u}_i^v and angular velocity \mathbf{u}_i^w generated by motors. \mathbf{x}_i, θ_i denotes horizontal and vertical positions, and orientation state respectively.

$$\dot{\mathbf{x}}_{i,h} = \mathbf{u}_i^v \cos(\theta^i) \quad (1)$$

$$\dot{\mathbf{x}}_{i,v} = \mathbf{u}_i^v \sin(\theta^i) \quad (2)$$

$$\dot{\theta}_i = \mathbf{u}_i^w \quad (3)$$

Bearing vector $\mathbf{b}^i \in \mathbb{R}^2 : \|\mathbf{b}^i\|_2 = 1$ denotes the heading direction of robot i . \mathbf{b}^{ij} denotes the bearing vector between robot i and j . $\phi(\mathbf{c}_1, \mathbf{c}_2)$ provides a general equation for the smallest angle rotating from a motion direction \mathbf{c}_1 to direction \mathbf{c}_2 , where $\hat{\mathbf{a}}_3$ is the unit vector along the positive z-axis.

$$\mathbf{b}_{ij} = \frac{\mathbf{x}_j - \mathbf{x}_i}{\|\mathbf{x}_j - \mathbf{x}_i\|_2} \quad (4)$$

$$\phi(\mathbf{c}_1, \mathbf{c}_2) = \text{sgn}((\mathbf{c}_1 \times \mathbf{c}_2)^T \hat{\mathbf{a}}_3) \cos^{-1}\left(\frac{\mathbf{c}_1^T \mathbf{c}_2}{\|\mathbf{c}_1\|_2 \|\mathbf{c}_2\|_2}\right) \quad (5)$$

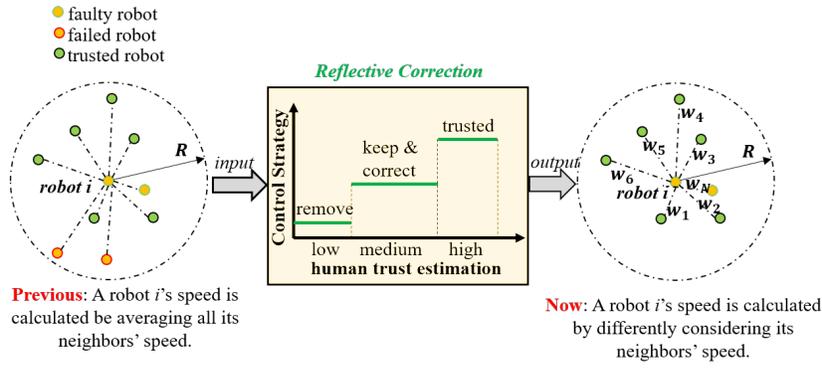


Figure 2: Illustration of change in speed calculation in the Trust-Aware Behavior Reflection for Swarm Self-Healing.

The distributed control for biased flocking is shown below. The heading direction for the swarm is specified by a given direction $\mathbf{q}_0 = \sum_{v \in \mathcal{V}} \theta_i$. N_i^r denotes neighbors of robot i within the repulsion radius r . N_i^R denotes neighbors outside of r but within the communication radius R . \mathbf{q}_{N_i} denotes the average direction of a robot and its neighbors within R . \mathbf{v}_i and \mathbf{y}_i are speed vectors.

$$\mathbf{u}_i^v = K_v(\mathbf{v}_i + \mathbf{q}_{N_i})^T \mathbf{b}_i \quad (6)$$

$$\mathbf{u}_i^w = K_w(\mathbf{y}_i + \phi(\mathbf{b}_i, \mathbf{q}_{N_i})) \quad (7)$$

$$\mathbf{v}_i[t+1] \leftarrow \frac{1}{N_i+1} \left(\mathbf{v}_i + \sum_{j \in N_i^r} \frac{-\mathbf{b}_{ij}}{\|\mathbf{x}_j - \mathbf{x}_i\|_2^2} + \sum_{j \in N_i^R} \mathbf{v}_j \right) \quad (8)$$

$$\mathbf{y}_i[t+1] \leftarrow \frac{1}{N_i+1} \left(\mathbf{y}_i + \sum_{j \in N_i^r} \phi(\mathbf{b}_i, -\mathbf{b}_{ij}) + \sum_{j \in N_i^R} \phi(\mathbf{b}_i, \mathbf{b}_{ij}) \right) \quad (9)$$

The speed \mathbf{u}_i of the robot i is updated using Equation 10.

Equations 6 – 9 can be simplified to equation 10. At each time step t , a robot i update its motion status by averaging its neighbors' motion status.

$$\mathbf{u}_i[t+1] = \frac{1}{N_i+1} (\mathbf{u}_i[t] + \sum_{j \in N_i} \mathbf{u}_j[t]) \quad (10)$$

As seen from the distributed update method above, faulty robots will be able to relay unreliable motion information to their neighbors which in turn will mislead their neighbors' motions.

Definition I (Faulty Robots and Failed Robots): “*Faulty robot*” refers to a robot with undesired behaviors, due to propagation of faulty data from a failed robot, environmental disturbance, etc., i.e. the faulty behavior is correctable. “*Failed robot*” refers to a robot with -undesired behaviors, which are not correctable.

Definition II (Untrusted Swarm): During the swarm deployments – influenced by faulty and failed robots – a swarm shows abnormal behaviors, such as partial disconnection or heading deviation. This decreases human trust in the swarm’s performance. This type of swarms is defined as an “untrusted swarm”.

Definition III (Influential Factors and Robot Faults): The real-world factors, such as degraded motors on a robot, uncertainty in

sensors and mechanical systems, or wind/rain disturbances from environments can cause abnormal robot behaviors and impair robot performance. These factors are defined as “influential factors”. Abnormal robot behaviors, such as degraded performance or abnormal motions, are defined as “robot faults”.

4 TRUST-AWARE BEHAVIOR REFLECTION (TRUST-R) FOR SWARM SELF-HEALING

The overall architecture of our self-healing method is shown in Figure 2. Based on a human’s trust signal that also indicates human’s diagnosis and level of faults, (e.g. low, medium or high, of the fault), each robot determines its strategy of communication between itself and its neighbors. When faulty robots appear in a swarm, it becomes unreliable to update a robot’s status by considering in the faulty robots’ motion status (calculated by Equation 10) [4]. Instead, it is more reliable to constrain information sharing between a faulty robot and its neighbors. In particular, if the trust level is high (faultiness is low) then the strategy “accept high-trust information” is employed. On the other hand, if trust level is medium (fault level is medium) then “reduce middle-trust information” is employed; and if trust level is low (faultiness is high) then “refuse low-trust information”. We propose a *novel information updating method* based on the weighted mean subsequence reduced algorithm (WMSR) [15]. Instead of merely averaging values as in the previous update method, our Trust-R method updates information differentially based on the communication quality (Equation 11). Weights w_i are calculated in Sections 3.2.2 and 3.2.3.

$$\mathbf{u}_i[t+1] = w_i[t] \mathbf{u}_i[t] + \sum_{j \in N_i} w_j[t] \mathbf{u}_j[t] \quad (11)$$

4.1 Human Trust in Faulty Behaviors

Abnormal robot behaviors inside a swarm decrease human trust in the swarm [12]. A human operator knows the swarm behavior requirements of the mission she is pursuing, such as requirements of connectivity and heading direction, and therefore can estimate a relation between current performance of the swarm and expected performance. $\delta(u_{actual}, u_{expect})$ is defined as the difference between expected speed/heading direction and actual speed/heading direction. u_{expect} is calculated using Equation 11 by referring to a robot’s neighbors. u_{actual} is read directly from a robot’s motion

sensors. δ is used to calculate the trust score. If $\delta(u_{actual}, u_{expect})$ is smaller than a small percentage β_1 , robot behaviors are normal with “high trust”. If $\delta(u_{actual}, u_{expect})$ is larger than a small percentage β_1 and smaller than a big percentage β_2 , robot behaviors are faulty with “medium trust”. If $\delta(u_{actual}, u_{expect})$ is larger than the big percentage β_2 , robot behaviors have completely failed with “low trust”. The β values are found by examining the differences between the speed and heading-direction according to UAV control requirements in different scenarios.

4.2 Trust-Aware Communication Quality Assessment

The overall communication graph for robot i is $\mathcal{E} = \{(i, j) \mid j \in N_i\}$. Based on the estimated trust levels of the two robots $\{i, j\}$, communication quality, $f_{ij} \in [0, 1]$, is used to measure the reliability of exchanged information. The trust-aware communication quality is dynamically updated to reflect the changing communication graph using Equation 12. The best communication distance between two robots i and j is ρ . Communication within ρ is considered as the communication with the best quality. The communication radius is R . The parameter, η , is used as a weighting factor to discourage the impact of faulty robots on their neighbors.

$$f_{ij} = \begin{cases} 0 & \|\mathbf{x}_i - \mathbf{x}_j\| \geq R \\ \frac{1}{2}(g_i + g_j)\eta & \|\mathbf{x}_i - \mathbf{x}_j\| \leq \rho \\ \frac{(g_i + g_j)\eta}{2} \exp \frac{-\gamma(\|\mathbf{x}_i - \mathbf{x}_j\| - \rho)}{R - \rho} & \text{otherwise} \end{cases} \quad (12)$$

where g_i is the trust level of robot i . The above communication quality evaluation method implies that within the communication range, the communication reliability is the average of the two robots’ trust values. If both robots are trusted, their communication is the most reliable; if one robot is faulty, the most reliable communication under that connection is the communication from the trusted robot.

The quality assessment for robot communications is visualized in Figure 3. The rationale of designing the trust-aware communication quality is to encourage information sharing with trusted robots by using higher upper limits on their communication quality, while discouraging information sharing with untrusted robots by using lower upper limits on the communication quality. Meanwhile, to encourage a compact swarm with closer distances among robots, the communication quality is decreased if the robot distance increases. Figure 3 shows that the communication quality among trusted robots is close to 1, while the quality among failed robots is 0.

For the curves shown in Figure 3, the g values are (1, 0.5, 0) for trusted robots, faulty robots and failed robots, respectively. η values are (1, 1, 0.4, 0.3, 0.2, 0.2) and γ values are (0.1, 0.5, 1, 3, 5, 7) for communications between trusted-trusted robots (trust-trust), trusted-faulty robots (trust-faulty), trusted-failed robots (trust-failed), faulty-faulty robots (faulty-faulty), faulty-failed robots (faulty-failed), failed-failed robots (failed-failed). g and η are used to set upper limits on the communication quality. γ defines the sensitivity of quality to mutual distance. For the remainder of the paper we set the communication radius to be $R = 12m$ and the best communication distance to be $\rho = 4m$.

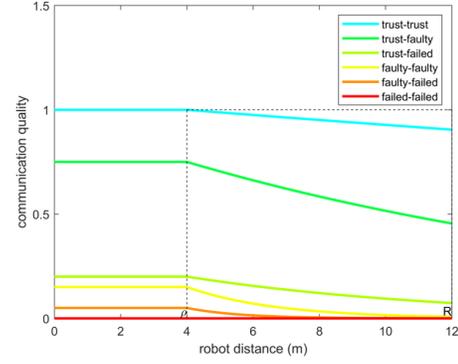


Figure 3: Illustration of the Trust-aware communication quality assessment. Information shared by trusted robots is encouraged with higher upper limits, while untrusted information is discouraged with lower upper limits.

The adjacency matrix, A , that describes the communication graph is given by:

$$[A]_{ij} = \begin{cases} 0 & i \neq j \\ f_{ij} & i = j. \end{cases} \quad (13)$$

The degree matrix, D , is:

$$[D]_{ij} = \begin{cases} 0 & i \neq j \\ \sum_j f_{ij} & i = j. \end{cases} \quad (14)$$

The novel trust-weighted Laplacian matrix, $[L]_{ij}$, calculated as $[L]_{ij} = [D]_{ij} - [A]_{ij}$ can then be defined as:

$$[L]_{ij} = \begin{cases} -f_{ij} & i \neq j \\ \sum_j f_{ij} & i = j. \end{cases} \quad (15)$$

The eigenvalues $\{\lambda_i \mid i = 1, 2, \dots, n\}$ of L are real and they satisfy $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. The connectivity measure λ_2 is estimated by the equation $L e_2 = \lambda_2 e_2$ and the eigenvector e_2 .

4.3 Trust-Aware Swarm Behavior Correction

A swarm proactively corrects its faulty behaviors using a two step process. First, it corrects the faulty robots by restraining the negative influence from faulty robots and referring to trusted robots for behavior correction. The failed robots are isolated from other trusted robots, preventing the sharing of unreliable motion information. The connectivity control in Section 5 is then used to reduce the distance between robots and their “normal” neighbors. In doing so, a robot adjusts its behavior – heading direction and speed – using a larger amount of trusted motion information.

$$w_k[t] = \frac{\hat{f}_k[t]}{\hat{f}_i[t] + \sum_{j \in N_i} \hat{f}_j[t]}, k \in [i, N_i] \quad (16)$$

Weights for updating each robot’s status are calculated by Equations 11 and 16. The result of the weighted update mechanism is shown on the right side of Figure 2. For updating a robot i , weights w_k are calculated by normalizing all the communication quality values in a communication range, shown in equation 16. When $k = i$, $\hat{f}_k = g_i$ (i.e., the trust level of itself). If $k = j \in N_i$ then

$\hat{f}_k = f_{ij}$ (i.e., the communication quality between robots i and j).
 $\hat{f}_i = g_i$ for all values of k .

With the trust-weighted update, the control input \mathbf{u}_v^i and \mathbf{u}_w^i for robot motors are changed to $\mathbf{u}_{v,trust}^i$ and $\mathbf{u}_{w,trust}^i$. The gains K_v and K_w are parameters for adjusting the motor output.

$$\mathbf{u}_{i,trust}^v = (K_v + K_{v,trust})(\mathbf{v}_i + \mathbf{q}_{N_i})^T \mathbf{b}_i \quad (17)$$

$$\mathbf{u}_{i,trust}^w = (K_w + K_{w,trust})(\mathbf{y}_i + \phi(\mathbf{b}_i, \mathbf{q}_{N_i})) \quad (18)$$

Let $\mathbf{u}_i[t+1]$ denote the actual speed of a robot with abnormal behaviors at the moment $t+1$, then the expected speed calculated by referring to its neighbors is denoted by $\mathbf{u}_{i,trust}[t+1]$. The extra trust-gain $K_{v,trust}$ and $K_{w,trust}$ can then be solved to adjust the control output of robot motors. The gains are updated based on the difference between the actual and the human-trusted robot speeds.

$$K_{v,trust}[t+1] = \frac{\mathbf{u}_{i,trust}^v[t] - \mathbf{u}_i^v[t]}{\mathbf{u}_i^v[t]} \quad (19)$$

$$K_{w,trust}[t+1] = \frac{\mathbf{u}_{i,trust}^w[t] - \mathbf{u}_i^w[t]}{\mathbf{u}_i^w[t]} \quad (20)$$

To avoid collision, the safe distance (repulsion radius) for separating robots is set to r . For a pair of robots i and j , their positions at the moment t are \mathbf{x}_i and \mathbf{x}_j . The overall swarm safety is maintained during the correction period $[0, T]$ by maintaining safety distance $h_{i,j}^{safe}$ for any robot pair i and j . $\mathcal{H}_{i,j}^{safe}$ is the set of all safe distances.

$$h_{i,j, safe}(t) = \|\mathbf{x}_i - \mathbf{x}_j\|^2 - \rho^2, \forall i, j \quad (21)$$

$$\mathcal{H}_{i,j, safe}[t] = \{\mathbf{x} \in \mathbb{R}^2, t \in [0, T] : h_{i,j}^{safe}(t) \geq 0\} \quad (22)$$

5 TRUST-AWARE CONNECTIVITY MAINTENANCE FOR MOTION CONSENSUS

To further correct faulty swarm behaviors, connectivities between a faulty robot and the other trusted robots are strengthened by Equation 23. The graph's Laplacian matrix is L and the algebraic connectivity is λ_2 . The connectivity (described by the second smallest eigenvalue λ_2 of L , \mathbf{e}_2 is the corresponding eigenvector) between a faulty – correctable – robot and a neighboring – trusted – robot is improved by reducing the distance between the faulty robot and the trusted neighbor. The more reliable information provides a mechanism for correcting the faulty robot's behavior. $\mathbf{x}_{i,\psi}$ is the position component on direction ψ (horizontal or vertical direction) for the robot i . $\frac{\alpha L(\mathbf{x})}{\alpha \mathbf{x}_{i,\psi}}$ is computed by calculating the difference between the reliability values, f_{ij} , at adjacent time steps as shown in Equation 24.

$$\mathbf{u}_i = \nabla_{i,\psi} \lambda_2 \quad (23)$$

$$= \frac{\alpha \lambda_2(L)}{\alpha \mathbf{x}_{i,\psi}} = \frac{\alpha \lambda_2(L)}{\alpha L(\mathbf{x})} \frac{\alpha L(\mathbf{x})}{\alpha \mathbf{x}_{i,\psi}} = \text{Trace} \left\{ \left[\frac{e_2 e_2^T}{e_2^T e_2} \right]^T \left[\frac{\alpha L(\mathbf{x})}{\alpha \mathbf{x}_{i,\psi}} \right] \right\} \quad (24)$$

THEOREM 5.1. *The novel method, Trust-R, reduces untrusted information and encourages the trusted information among robots.*

PROOF. For two robots $a, b \in G$, a is an abnormal robot, while b is a normal robot. To update the motion status of a target robot i , weights of its neighbors a and b ($a, b \in N_i$) for information exchange are w_a and w_b respectively.

$$w_a[t] = \frac{\hat{f}_a[t]}{\hat{f}_i[t] + \sum_{j \in N_i} \hat{f}_j[t]}, w_b[t] = \frac{\hat{f}_b[t]}{\hat{f}_i[t] + \sum_{j \in N_i} \hat{f}_j[t]}$$

An abnormal robot's trust level is lower than that of a normal robot. As a result, $\hat{f}_a[t] \leq \hat{f}_b[t] \Rightarrow w_a[t] < w_b[t]$. Therefore, with trust awareness, Trust-R reduces the untrusted information given by abnormal robots, while encouraging the sharing of trustworthy information from normal robots. \square

THEOREM 5.2. *The novel Trust-R method encourages a relatively closer distance between a robot and other trusted robots, and encourages a relatively farther distance between a robot and other untrusted robots. The adjustment will be reduced to zero once the flocking consensus is reached.*

PROOF. When using the trust-aware communication quality to adjust the distance of robot i to other robots, the adjustment along a direction ψ is

$$\mathbf{u}_i = \text{Trace} \left\{ \left[\frac{e_2 e_2^T}{e_2^T e_2} \right]^T \left[\frac{\alpha L(\mathbf{x})}{\alpha \mathbf{x}_{i,\psi}} \right] \right\} = \text{Trace} \left\{ \left[\frac{e_2 e_2^T}{e_2^T e_2} \right]^T \left[\frac{\alpha [L]_{ij}}{\alpha \mathbf{x}_{i,\psi}} \right] \right\}$$

For the off-diagonal elements in the Laplacian matrix, L , $\frac{\alpha [L]_{ij}}{\alpha \mathbf{x}_{i,\psi}}$ is solved by

$$\sum_K - \frac{\alpha f_{ij}}{\alpha \mathbf{x}_{k,\psi}} \mathbf{u}_{k,\psi} = \frac{\alpha f_{ij}}{\alpha \mathbf{x}_{i,\psi}} (\mathbf{u}_{j,\psi} - \mathbf{u}_{i,\psi})$$

For the diagonal elements in L , $\frac{\alpha [L]_{ij}}{\alpha \mathbf{x}_{i,\psi}}$ is solved by

$$\sum_k \left(\sum_j \frac{\alpha f_{ij}}{\alpha \mathbf{x}_{k,\psi}} \right) \mathbf{u}_{k,\psi} = \sum_j \frac{\alpha f_{ij}}{\alpha \mathbf{x}_{i,\psi}} (\mathbf{u}_{i,\psi} - \mathbf{u}_{j,\psi})$$

Since

$$\frac{\alpha f_{ij}}{\alpha \mathbf{x}_{i,\psi}} = - \frac{\gamma \eta (g_i + g_j) (\mathbf{x}_{i,\psi} - \mathbf{x}_{j,\psi})}{2(R - \rho) \|\mathbf{x}_i - \mathbf{x}_j\|} \exp \frac{-\gamma (\|\mathbf{x}_i - \mathbf{x}_j\| - \rho)}{R - \rho}$$

$\frac{\alpha f_{ij}}{\alpha \mathbf{x}_{i,\psi}}$ is bounded by the robot distance which is smaller than communication radius R . For a desired flocking direction \mathbf{q}_0 , the adjustment degree \mathbf{u}_i , between two robots i and j , is positively correlated with their average trust score $\frac{(g_i + g_j) \eta}{2}$. A larger trust score leads to a larger adjustment. Therefore, Trust-R encourages a relatively closer distance between a robot and other trusted robots, and encourages relatively farther distances to abnormal robots. The abnormal faulty robots are gradually abandoned by the swarm.

As robots reach the motion consensus along the heading direction, \mathbf{u}_i will be equal to \mathbf{u}_j within a limited time. Therefore, $\frac{\alpha f_{ij}}{\alpha \mathbf{x}_{i,\psi}}$ will be 0, stopping the adjustment when the consensus is reached. \square

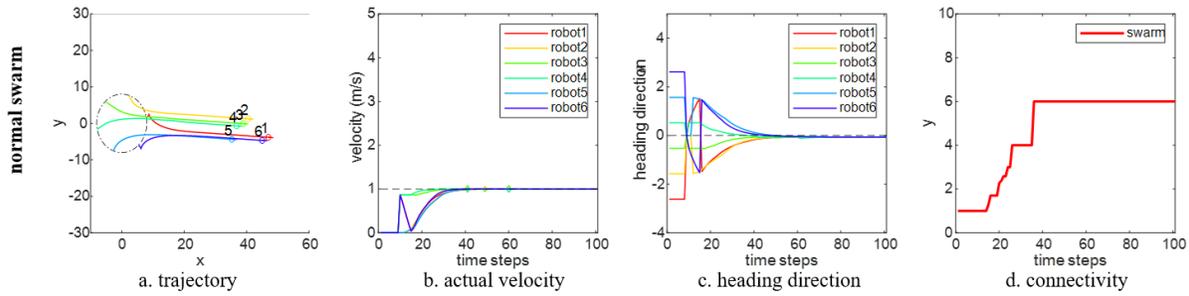


Figure 4: System response given by normal flocking. With a distributed control method, the swarm without a faulty robot flocks to the human assigned direction: “East”. with consensus on both heading and motion.

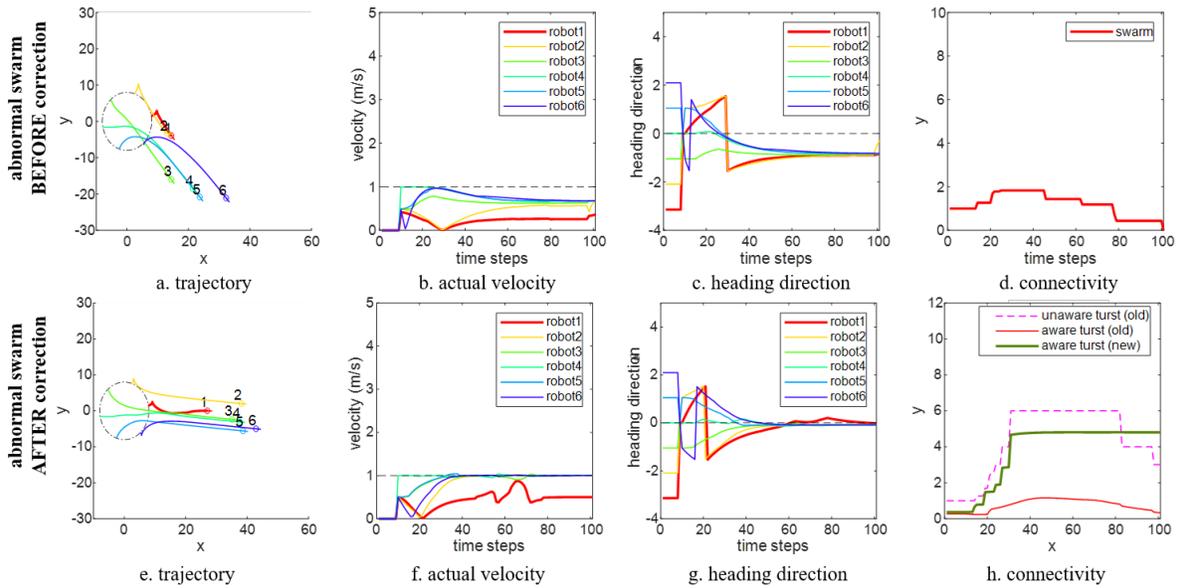


Figure 5: System response given by untrusted flocking caused by a robot with faulty motor. Under the influence of faulty, robot 1, the motion consensus cannot be achieved.

6 EVALUATION

To validate the effectiveness of Trust-R in helping the swarm self-heal, three real-world faults were simulated using MATLAB: a degraded motor, system uncertainty (sensor dysfunction) and wind disturbance. These faults commonly happen in complex environments, such as densely distributed forests/buildings and extreme weather conditions, which can affect robot communication, spatial distributions and system reliability [5][19]. Our goal in using Trust-R is to repair the untrusted swarm misled by faulty robots by improving the swarms’ environmental adaptation. The task for the swarm in all experiments is distributed biased flocking. All results reported are the performance of the system for a single run. The non-stochastic nature of the algorithm results in the same behavior each time a specific parameter configuration is run.

To focus on “correction” of faulty swarms with different faults and to reduce the difficulties in analyzing the behaviors of individual robots, the number of robots was set to a small number – 6 – and the biased heading-direction was fixed to “East”. The initial number

of faulty/failed robots for each scenario was chosen to be either 1 or 2 robots. Under the influence of these abnormal robots, several of the neighboring robots can also become faulty/failed. The map size for the flocking was 60m×60m. The velocity for each robot was set as 1.0m/s. To observe the misleading effect of one faulty robot on its neighbors, robot locations were initialized in a circle with radius of 8m. The heading direction of all the robots pointed to the circle center. To avoid collision, the repulsion radius securing robot safety was set as 2m. For all conducted experiments $\beta_1 = 10\%$ and $\beta_2 = 50\%$ were used for the faulty behavior detection.

6.1 Limited Speed – Degraded Motor(s)

Due to a degraded motor, the speed of a robot was constrained such that the velocity and angular speeds could not reach the designed speeds. In this case, the speed of the faulty robots was lower than the normal robots. Because of the exchange of motion statuses with the faulty robots, the speed of some robots was influenced. The upper limits of linear velocity and angular velocity were set

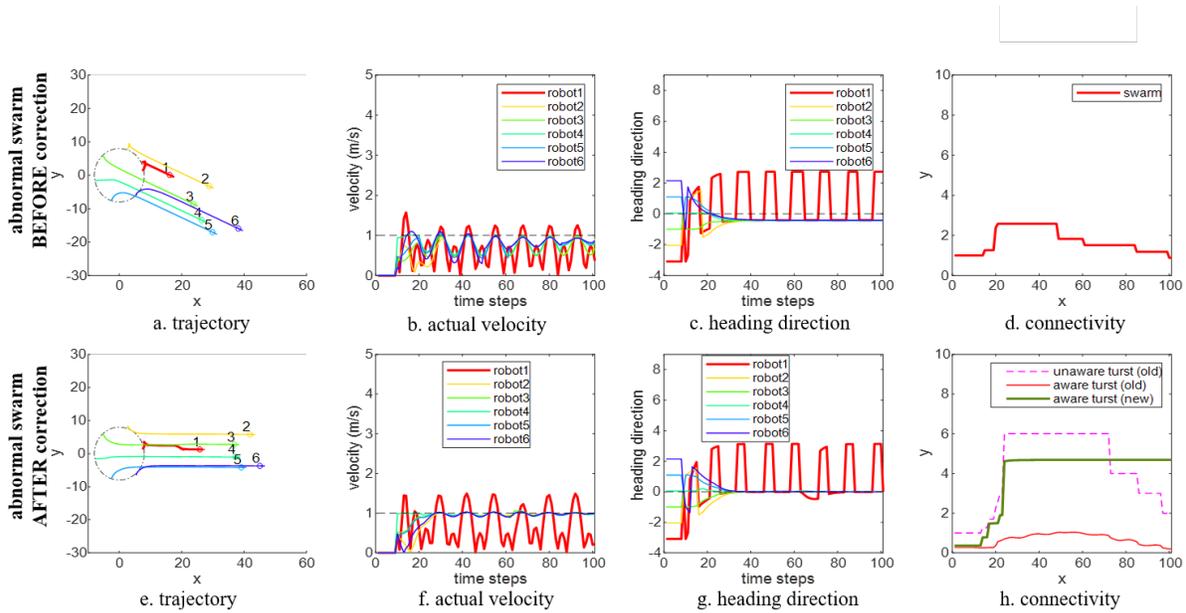


Figure 6: System response given by untrusted behaviors caused by system uncertainty (faulty sensor).

as $0.5m/s$ and $0.5rad/s$. For comparison, normal biased flocking without faulty robots was simulated as a baseline.

For a normal swarm (Figure 4), after about 38 time steps (19s), the velocity of all 6 robots achieved the desired consensus of $1.0m/s$; After about 50 time steps (25s), the heading of all 6 robots achieved consensus on the “East” direction. The connectivity λ_2 was 6, which means all robots achieved the best communication in this scenario.

Figure 5 shows a scenario in which Robot 1 had a degraded motor. As shown, with the faulty robot in the swarm, the velocity consensus was not achieved within 100 time steps (50s), and the faulty Robot 1 was disconnected from the swarm (Figure 5(a,b)). The heading direction of the swarm shifted to $1rad$ after 64 time steps (32s) (Figure 5(c)). Connectivity with the faulty robot was decreased to 0 after 100 time steps (50s)(Figure 5(d)).

Trust-R provided robots in the swarm additional awareness by assessing the motion statuses of a robot and its neighbors. In this case, Robot 1, whose speed is 70% lower than the expected speed, was considered an untrusted and failed robot (Figure 5e), thereby decreasing swarm performance and human trust. The communication quality between Robot 1 and other normal robots decreased as calculated by the “trust-failed” curve in Figure 3. With the Trust-R correction, the information exchanged with Robot 1 was tightly constrained. After 70 time steps (35s), Robot 1 was disconnected from the normal robots. The swarm with only trusted robots achieved velocity consensus after 32 time steps (16s), and achieved consensus on heading after 50 time steps (25s) with only a -0.1 rad deviation (Figure 5(f,g)). This demonstrates that Trust-R was effective in correcting the faulty behaviors of the swarm. Shown as Figure 5(h), the connectivity, λ_2 , of the old swarm without Trust-R maintained a low-level of connectivity and decreased to 0 after 80 time steps (40s). In contrast, the swarm which constrained the information exchanged with the faulty Robot 1 had connectivity that increased

to a high level of 4.8, showing the effectiveness of Trust-R in encouraging connectivity among trusted robots.

6.2 Abnormal Motion – System Uncertainty

Due to the system uncertainties such as sensor failures, lost GPS signals and internal disturbances from the mechanical systems, robots may show abnormal behaviors such as sinusoidal motion, random motion, or fixed-direction motion. For this case study, a sinusoidal motion was investigated. Robot 1 had abnormal sinusoidal motion and angular velocity (shown as Figure 6(a,b)), with amplitude of $1.5m/s$. Without correction, the motion consensus was not achieved (Figure 6(b,c)). The connectivity decreased to 0.8 after about 100 time steps (50 seconds), shown in Figure 6(d). With the Trust-R correction, misleading information from the untrusted robot 1 was quickly constrained. The new swarm without faulty robots achieved velocity consensus after 30 time steps (15s) and achieved heading direction consensus after 30 time steps (15s) with $0rad$ deviation, shown in Figure 6(e,f,g). As shown in Figure 6(h), connectivity of the swarm without Trust-R correction remained low as a farther distance between the faulty Robot 1 and other trusted robots was encouraged. On the other hand, the connectivity of the swarm isolating the faulty Robot 1 increased to 4.9 by using Trust-R, showing again the effectiveness of Trust-R in correcting abnormal swarm behaviors.

6.3 Motion Deviation – Wind Disturbance

When some robots in a swarm cross into a wind zone, the wind will give the robots extra linear and angular velocity. For this experiment, a wind region with size of 15×15 was located in the convex hull formed by the following set of vertices ((15,4), (30,4), (30,19),(15,19)). Before reaching the region, the swarm had already

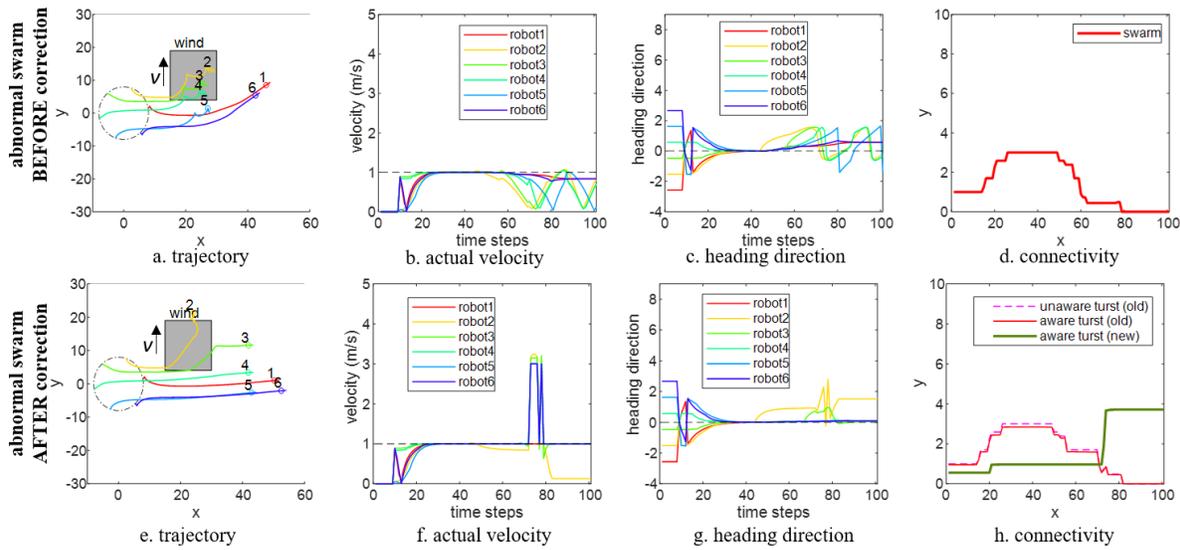


Figure 7: System response given by untrusted flocking caused by a wind disturbance.

achieved motion consensus. Some robots will cross the wind region and gain an extra $0.25m/s$ linear velocity along the “North” direction and an angular deviation of $0.1rad/s$.

As Figure 7(a) shows, Robots 2 and 3 crossed the wind region first. They then attracted Robots 4 and 5 into the wind zone. Without correction, a motion consensus was not achieved (Figure 7(a-c)). The connectivity decreased to 0 after about 78 time steps (39s)(Figure 7d). With the Trust-R correction, misleading information from the untrusted Robots 2 and 3 was quickly constrained and their influence on the other robots was largely reduced. The new swarm without faulty robots achieved velocity consensus after about 30 time steps (15s) and achieved heading direction consensus after 40 time steps (about 20s), shown in Figure 7(e,f,g). As shown in Figure 7(h), connectivity of the old swarm without the Trust-R correction was decreased to 0 after 80 time steps (40s), due to the disconnection of the faulty Robot 2. In contrast, the connectivity of the new swarm (swarm after removing the faulty Robot 2) increased to 4, showing the effectiveness of Trust-R in correcting abnormal swarm behaviors caused by disturbances such as wind.

Once the wind disturbance has passed, the difference between the actual velocity and the expected velocity will decrease. If the robot is already disconnected, given the characteristics of the distributed control, the robot is no longer reachable and will be ignored by the swarm (behavior correction of robot 2 in Case Study III). If the robot is still within the communication range of the other robots after the wind disturbance has passed, the previously faulty robot will then be considered a normal robot with a high communication quality. The new high-level communication can then be used to correct the previously faulty robot’s behavior (robots 3 and 4).

7 CONCLUSION & FUTURE WORK

We presented a trust aware behavior reflection method – Trust-R – to help a swarm self-heal. This allows the swarm to repair its

overall behavior when faulty behaviors of its members occur. Three types of robot faults – limited performance caused by a degraded motor, abnormal motion caused by system uncertainty such as sensor failure and motion deviation caused by wind disturbance – were simulated. With Trust-R, the motion status of the robots are estimated and the corresponding communication quality is determined. The robots are encouraged to communicate with trusted robots and discouraged from communicating with untrusted robots. In doing so the negative influence caused by misleading information is largely reduced and swarm behaviors are corrected. The simulation results for the three faulty scenarios demonstrate the effectiveness of Trust-R in correcting a range of faulty swarm behaviors.

In this paper, the relation between human trust and swarm performance was presumed linear. Future research will focus on user studies to assess the effects on trust of self-healing swarm strategies which preserve swarm performance despite observably faulty behavior of individual robots. Additional studies will be conducted to characterize the type and severity of faulty behaviors necessary to decrease a human’s trust to the point of intervention despite self-healing behavior. In this paper faults from different categories were remediated using the same mechanism in Trust-R which leads to expulsion of faulty robots from the swarm. In cases such as the wind disturbance where the robot’s communications are faulty but not the robot itself or where neighboring robots might assist we hope to develop compensatory strategies to salvage those faulty robots we can. Lastly, we plan to develop signatures from faulty behavior detection data to analyze the temporal status of robot motion and interactions. Based on these models, faulty robot behaviors can be automatically detected during the swarm deployments and the time needed to mitigate or expel faulty robots can be reduced. More complex scenarios (e.g., obstacles) will also be considered. The system’s ability to mitigate the effect of faulty and failed robots on the performance of the swarm in these complex scenarios will be evaluated.

REFERENCES

- [1] Hai Chen, Xin-min Wang, and Yan Li. 2009. A survey of autonomous control for UAV. In *Artificial Intelligence and Computational Intelligence, 2009. AICT'09. International Conference on*, Vol. 2. IEEE, 267–271.
- [2] Rafer Cooley, Shaya Wolf, and Mike Borowczak. 2018. Secure and Decentralized Swarm Behavior with Autonomous Agents for Smart Cities. *arXiv preprint arXiv:1806.02496* (2018).
- [3] Yuan-Shun Dai, Michael Hinchey, Manish Madhusoodan, James L Rash, and Xukai Zou. 2006. A prototype model for self-healing and self-reproduction in swarm robotics system. In *Dependable, Autonomic and Secure Computing, 2nd IEEE International Symposium on*. IEEE, 3–10.
- [4] Ali Jadbabaie, Jie Lin, and A Stephen Morse. 2003. Coordination of groups of mobile autonomous agents using nearest neighbor rules. *IEEE Transactions on automatic control* 48, 6 (2003), 988–1001.
- [5] Belkacem Khaldi, Fouzi Harrou, Foudil Cherif, and Ying Sun. 2017. Monitoring a robot swarm using a data-driven fault detection approach. *Robotics and Autonomous Systems* 97 (2017), 193–203.
- [6] Andreas Kolling, Phillip Walker, Nilanjan Chakraborty, Katia Sycara, and Michael Lewis. 2016. Human interaction with robot swarms: A survey. *IEEE Transactions on Human-Machine Systems* 46, 1 (2016), 9–26.
- [7] Eric J Leaman, Brian Q Geuther, and Bahareh Behkam. 2018. Hybrid Centralized/Decentralized Control of Bacteria-Based Bio-Hybrid Microrobots. In *2018 International Conference on Manipulation, Automation and Robotics at Small Scales (MARSS)*. IEEE, 1–6.
- [8] Ali Marjovi and Lino Marques. 2013. Optimal spatial formation of swarm robotic gas sensors in odor plume finding. *Autonomous Robots* 35, 2-3 (2013), 93–109.
- [9] Nithin Mathews, Anders Lyhne Christensen, Rehan Oa'Grady, Francesco Mondada, and Marco Dorigo. 2017. Mergeable nervous systems for robots. *Nature communications* 8, 1 (2017), 439.
- [10] Nithin Mathews, Anders Lyhne Christensen, Alessandro Stranieri, Alexander Scheidler, and Marco Dorigo. 2019. Supervised morphogenesis: Exploiting morphological flexibility of self-assembling multirobot systems through cooperation with aerial robots. *Robotics and Autonomous Systems* 112 (2019), 154–167.
- [11] Sasanka Nagavalli, Nilanjan Chakraborty, and Katia Sycara. 2017. Automated sequencing of swarm behaviors for supervisory control of robotic swarms. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2674–2681.
- [12] Changjoo Nam, Phillip Walker, Michael Lewis, and Katia Sycara. 2017. Predicting trust in human control of swarms via inverse reinforcement learning. In *Robot and Human Interactive Communication (RO-MAN), 2017 26th IEEE International Symposium on*. IEEE, 528–533.
- [13] James O'Keefe, Danesh Tarapore, Alan G Millard, and Jon Timmis. 2017. Towards fault diagnosis in robot swarms: An online behaviour characterisation approach. In *Conference Towards Autonomous Robotic Systems*. Springer, 393–407.
- [14] Andrzej Pelc and David Peleg. 2005. Broadcasting with locally bounded byzantine faults. *Inform. Process. Lett.* 93, 3 (2005), 109–115.
- [15] David Saldana, Amanda Prorok, Shreyas Sundaram, Mario FM Campos, and Vijay Kumar. 2017. Resilient consensus for time-varying networks of dynamic agents. In *American Control Conference (ACC), 2017*. IEEE, 252–258.
- [16] Kelsey Saulnier, David Saldana, Amanda Prorok, George J Pappas, and Vijay Kumar. 2017. Resilient flocking for mobile robot teams. *IEEE Robotics and Automation Letters* 2, 2 (2017), 1039–1046.
- [17] Nathan Schurr, Janusz Marecki, Milind Tambe, Paul Scerri, Nikhil Kasinadhuni, and John P Lewis. 2005. The Future of Disaster Response: Humans Working with Multiagent Teams using DEFACTO.. In *AAAI spring symposium: AI technologies for homeland security*. 9–16.
- [18] Tina Setter, Andrea Gasparri, and Magnus Egerstedt. 2016. Trust-based interactions in teams of mobile agents. In *American Control Conference*. 6158–6163.
- [19] Andreas Steyven, Emma Hart, and Ben Paechter. 2017. An investigation of environmental influence on the benefits of adaptation mechanisms in evolutionary swarm robotics. In *Proceedings of the Genetic and Evolutionary Computation Conference*. ACM, 155–162.
- [20] Danesh Tarapore, Anders Lyhne Christensen, and Jon Timmis. 2017. Generic, scalable and decentralized fault detection for robot swarms. *PLoS one* 12, 8 (2017), e0182058.
- [21] Xiaotian Wang and Yue Wang. 2017. Co-design of Control and Scheduling for Human-Swarm Collaboration Systems Based on Mutual Trust. In *Trends in Control and Decision-Making for Human-Robot Collaboration Systems*. Springer, 387–413.
- [22] Fei Zhang and Weidong Chen. 2007. Self-healing for mobile robot networks with motion synchronization. In *Intelligent Robots and Systems, 2007. IROS 2007. IEEE/RSJ International Conference on*. IEEE, 3107–3112.
- [23] Haotian Zhang and Shreyas Sundaram. 2012. Robustness of information diffusion algorithms to locally bounded adversaries. In *American Control Conference (ACC), 2012*. IEEE, 5855–5861.