# Swarm Oracle:

# Trustless Blockchain Agreements through Robot Swarms

Alexandre Pacheco[1,4,*] Hanqing Zhao[2], Volker Strobel[1], Tarik Roukny[6], Gregory Dudek[3], Andreagiovanni Reina[1,4,5], Marco Dorigo[1]

[1]IRIDIA, Université Libre de Bruxelles, Brussels, Belgium

[2]Department of Electrical and Computer Engineering, Université Laval, Québec, Canada

[3]School of Computer Science, McGill University, Montréal, Canada

[4]Centre for the Advanced Study of Collective Behaviour, Universität Konstanz, Germany

[5]Department of Collective Behaviour, Max Planck Institute of Animal Behavior, Konstanz, Germany

[6]Faculty of Economics and Business, KU Leuven, Belgium

[*]To whom correspondence should be addressed; `alexandre.melo.pacheco@gmail.com`

## Abstract

Blockchain consensus, rooted in the principle "don't trust, verify", limits access to real-world data, which may be ambiguous or inaccessible to some participants. Oracles address this limitation by supplying data to blockchains, but existing solutions may reduce autonomy, transparency, or reintroduce the need for trust. We propose Swarm Oracle: a decentralized network of autonomous robots—that is, a robot swarm—that use onboard sensors and peer-to-peer communication to collectively verify real-world data and provide it to smart contracts on public blockchains. Swarm Oracle leverages the built-in decentralization, fault tolerance and mobility of robot swarms, which can flexibly adapt to meet information requests on-demand, even in remote locations. Unlike typical cooperative robot swarms, Swarm Oracle integrates robots from multiple stakeholders, protecting the system from single-party biases but also introducing potential adversarial behavior. To

1

ensure the secure, trustless and global consensus required by blockchains, we employ a Byzantine fault-tolerant protocol that enables robots from different stakeholders to operate together, reaching social agreements of higher quality than the estimates of individual robots. Through extensive experiments using both real and simulated robots, we showcase how consensus on uncertain environmental information can be achieved, despite several types of attacks orchestrated by large proportions of the robots, and how a reputation system based on blockchain tokens lets Swarm Oracle autonomously recover from faults and attacks, a requirement for long-term operation.

# 1 Introduction

Blockchain technology facilitates the creation and governance of public digital resources through peer-to-peer collaboration over the internet. Bitcoin [1] exemplifies this as the first successful public monetary ledger, enabling censorship-resistant transactions worldwide. The key ideas behind Bitcoin's success were later adopted and generalized by Ethereum [2], a blockchain envisioned as a "world computer" that hosts smart contracts—programs written in a Turing-complete language and executed collectively by the Ethereum network. To operate and scale globally, blockchain systems are designed to be decentralized, transparent, and trustless, relying on large networks of nodes to independently verify that new blocks follow consensus rules and that user transactions are valid. However, this design imposes a fundamental limitation: blockchains cannot securely incorporate information unless it is algorithmically verifiable by every node. This limitation extends to smart contracts which rely on accurate real-world data to support decision-making, such as exchange rates for decentralized finance, seismic activity measurements for disaster insurance, or environmental indicators for managing smart cities and natural resources [3].

External entities that provide data to blockchains are called *oracles*, a term borrowed from Alan Turing's "oracle-machines"—hypothetical devices that can solve problems beyond the limits of algorithmic computation [4]. However, relying on oracles to obtain accurate real-world data reintroduces elements of trust and centralization, potentially undermining blockchain's core properties such as automation, censorship resistance, and transparency.

This issue is known as blockchains' "oracle problem" [5]. As an illustrative example, consider a smart contract that distributes rewards to local landowners for maintaining a sustainable forest (e.g., using funds originating from the sale of carbon credits). In principle, such a contract could implement transparent, immutable rules that guarantee landowners direct and equitable access to funds, independent of shifting politics, opaque bureaucracies and financial intermediaries. However, to determine rewards, the contract requires reliable information about the state of the forest: Are biodiversity levels increasing? At what rate are trees being felled? Has a fire occurred? Prospective blockchain projects could, for instance, obtain this data from a trusted environmental agency or autonomous software that processes satellite images—both of which acting as oracles.

To reduce the need for trust, decentralized oracle networks such as Chainlink [6] aggregate data from multiple software oracles that retrieve and process online information—such as the satellite images in our example, but also asset prices and news events. Beniiche [7] proposes a taxonomy of existing oracle solutions based on their underlying trust models, ranging from authoritative to decentralized consensus-based approaches, and on the nature of their participants (e.g., software, hardware, and even human). In contrast with software oracles, which rely on online sources, hardware oracles can collect data directly from the physical environment using sensors. A prominent example is Helium [8], whose IoT network provides wireless coverage while gathering geolocation and connectivity data for on-chain applications. However, despite growing demand for transparent and trustworthy data, existing decentralized oracles, hardware ones in particular, lack the autonomy, flexibility, and scalability required for most applications.

In this paper, we propose Swarm Oracle, a decentralized oracle network composed of autonomous robots that act as hardware oracles, directly sensing and processing data from the physical world. By leveraging swarm robotics' peer-to-peer interactions, decentralization, and self-organization [9], Swarm Oracle can potentially scale or adapt robots' tasks to flexibly meet the demands of real-world applications. Even when composed of simple robots equipped with inexpensive, noisy or biased sensors, Swarm Oracle harnesses the wisdom of the crowd [10, 11]

3

through a consensus protocol, aggregating data from multiple robots to reach agreements that are more accurate and reliable than individual robots' sensor readings. Swarm Oracle can, for example, monitor fish stocks, wildfires, traffic conditions, or trash deposits, providing important information required by smart contracts to enact transparent regulations to manage shared resources [12]. In this context, Swarm Oracle functions as a Decentralized Physical Infrastructure Network (DePIN) [13], retrieving and serving real-world data to blockchains on demand, as part of the growing decentralized economy.

We envision Swarm Oracle as a collective of robots from different stakeholders, a prerequisite for a *trustless* oracle network that prevents biases during the design, manufacturing, or programming of the robots. Achieving consensus in a mobile multi-robot system with unreliable and sparse communication is a major challenge, further complicated by the possibility that stakeholders hold conflicting interests or that malicious actors attempt to manipulate on-chain data for personal gain. Even under these adversarial conditions, Swarm Oracle consensus must be both *fault-tolerant* and *global*—otherwise, discrepancies among the robots could compromise the integrity and reproducibility of the blockchain's deterministic state.

Unfortunately, these requirements are not met by typical swarm robotics consensus methods, where robots gradually adjust their individual beliefs (e.g., through local averaging of sensor readings [14]) to converge toward an *approximate* collective value, possibly applying local rejection of extreme values to enhance robustness (e.g., the W-MSR rule [15]). While some strategies can withstand adversarial robots [16, 17, 18], security hinges on robots' remaining sufficiently connected to honest robots. In sparsely connected robot swarms, these approaches may fail when networks partition—especially if malicious robots target specific neighborhoods. Most critically, they allow for persistent (small) discrepancies in opinions across the swarm, falling short of Swarm Oracle's requirements.

To reach trustless, fault-tolerant and global consensus without requiring external infrastructure, we implement Swarm Oracle through a blockchain hosted by the robots themselves. This second-layer blockchain (also called "sidechain") is adapted to the limitations of the robots, employing a lightweight consensus protocol to avoid the high energy demands of

protocols based on proof-of-work [1]. The Swarm Oracle protocol, which includes data aggregation methods, is deployed as a smart contract on this blockchain, as illustrated in Figure 1. The Swarm Oracle smart contract interacts with external infrastructure through events, receiving queries from public blockchains and publishing consensus results once agreements are reached. A general taxonomy on information exchange methods between oracle networks and smart contracts is available in [19]. While our framework centers on blockchain as both the application and the deployment vehicle, the trustworthy information produced by Swarm Oracle can also support off-chain decision processes, particularly when transparency and impartiality are required, such as applications for non-profit or government institutions. The Swarm Oracle protocol could equally be deployed as a smart contract on a public blockchain, though this could introduce transaction costs, especially when swarm sizes are large or when robots collect high volumes of data [20].

Previous research [21, 22, 23, 24] showed how robot swarms can benefit from blockchains to detect inconsistent or extreme values reported by some robots, and how to use blockchain tokens as participation credentials to protect the system from "double-spending" and "Sybil" attacks.[1] However, the "oracle agreement" in these studies, in which robots vote and agree on the fraction of white floor tiles, offered no formal guarantees against attacks, and only a few malicious behaviors were studied. A small group of coordinated robots could potentially synchronize their votes to exploit the system, thereby gaining cryptotokens and voting power. Zhao et al. [25] presented a blockchain-based Byzantine fault-tolerant algorithm, however, this preliminary work lacked validation on physical robots, leaving open questions about the practicality of the approach and its ability to withstand coordinated attacks under realistic constraints, without incurring excessive delays or communication overhead.

In this work, we deploy a Byzantine fault-tolerant consensus protocol on a swarm of 12 physical robots, to form an experimental Swarm Oracle. We stress-test the system at its theoretical limit of fault tolerance: when one third of the robots coordinate attacks to compromise

---

[1] Attacks where participants send conflicting information to different peers ( "double-spending"), and where they generate false identities ("Sybil").
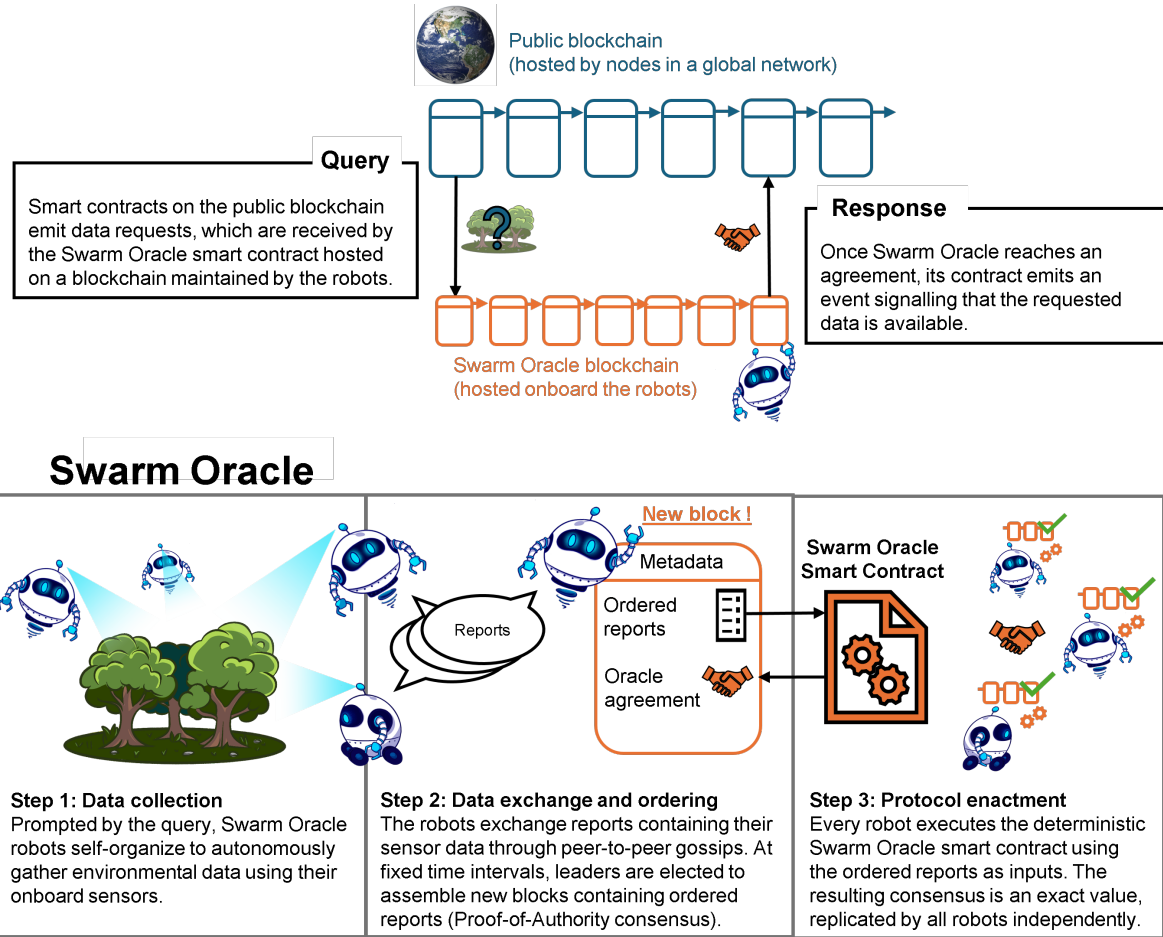
5

Figure 1: Deployment of the Swarm Oracle as a service that delivers reliable data to a public blockchain. In this deployment scenario, the protocol is implemented as a smart contract on a dedicated second-layer blockchain maintained by the robots. By hosting the contract on a dedicated blockchain tailored to their capabilities, the robots can reach a deterministic global (i.e., swarm-wide) agreement from the aggregation of individual robots' data, with lower latency and transaction costs than a public blockchain deployment.

the Swarm Oracle's *safety* and *liveness* properties [26].[2] We also integrate a reputation system that enables Swarm Oracle to self-heal from faults and attacks, as this is a crucial requirement for continuous and autonomous operations. Unlike previous reputation systems [23, 24, 25] that penalized only robots providing unreliable information, Swarm Oracle also safeguards liveness by penalizing robots that become inactive, either due to accumulating faults or attacks. In our experimental scenario, robots use onboard cameras to identify the colors of landmarks in the environment. This scenario reflects real-world conditions where robots have varying levels of accuracy and individual biases in their sensing, as color perception is highly influenced by environmental and hardware variability (e.g., camera calibration and lighting conditions). We perform an extensive experimental study using real robots (approximately 22 hours of runtime) and simulated robots (approximately 278 hours of runtime), showcasing how the Swarm Oracle provides correct agreements and recovers from attacks and faults.

## 2 Results

### 2.1 Experimental Scenario

Our experimental scenario, shown in Figure 2, contains three colored landmarks (red, green, and blue), each marked with AprilTags indicating whether they are "valuable" or not. In this study, the red landmark is designated as valuable, while green and blue are not. The oracle query under consideration is: "What are the RGB color values of the valuable landmark?" This setup illustrates a potential real-world deployment of a Swarm Oracle, composed of a swarm of 12 mobile robots that explore the environment and monitor its physical properties using onboard sensors and peer-to-peer communications. Each robot is equipped with an inexpensive monocular camera used to detect the colors of the landmarks and to read the AprilTags. However, color readings are highly noisy due to occlusions, sensor quality, and light conditions. Additionally, individual differences in each robot's color perception introduce systematic biases. As a Swarm Oracle, robots combine their readings to compute an accurate

---

[2]The safety property states that *bad things won't happen*—in our case, that the robots will not agree on incorrect information. The liveness property states that *good things will happen*—in our case, that agreements are reached within feasible time and resource constraints.
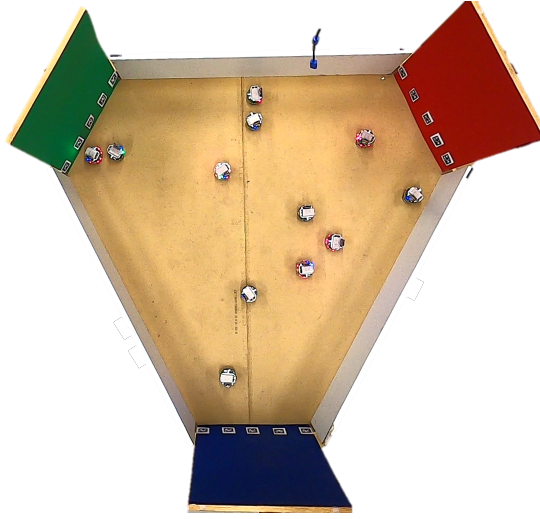
Figure 2: The arena contains three colored landmarks that can be detected by the robot cameras from a distance (potentially from any position in the arena if there is a clear line of sight) and used for navigation. Each landmark has AprilTags, readable from a short distance (approximately 10 cm), which encode "not valuable" on green and blue landmarks, and "valuable" on red.

consensus value for the RGB color of the landmark (a three-dimensional array of continuous values). Robots can navigate toward visible color blobs, but can only determine whether a landmark is valuable or not by reading the AprilTags at close range (approximately 10 cm). Adversary robots can use the green and blue landmarks as distractors to deliver some of their attacks.

## 2.2   The Swarm Oracle

The design and properties of Swarm Oracle are presented in detail in Methods. Here, we provide a brief, intuitive overview in the context of our experimental scenario.

Robots move toward large color blobs they detect and, within 10 cm of an AprilTag, they stop to sample RGB values. These values (or *observations*) are then reported to nearby robots (i.e., robots not farther away than approximately 15 cm). Each report additionally contains the sending robot's identity, a deposit of reputation tokens, and a vote indicating the value of the resource. A blockchain, maintained by the robots in their peer-to-peer network, hosts Swarm Oracle's consensus protocol as a smart contract, orders reports, and tracks the

reputation tokens held by each robot—which add up to a total circulating supply of $T$ tokens.

The smart contract aggregates reports into clusters via k-means clustering, and computes proposals which are the centroids of the clusters. A new cluster is created when the distance between an observation and existing proposals exceeds a *clustering threshold R*.

To send reports, robots must deposit a fixed fraction of the reputation tokens it owns, which we call *deposit quota* $K \in ]0, 1]$. Swarm Oracle's consensus protocol requires that proposals accumulate a pool of deposited tokens larger than $\frac{2}{3}KT$ before a consensus decision is made, therefore, the value $K$ regulates the maximum number of proposals that Swarm Oracle can run in parallel. Once $\lfloor K^{-1} \rfloor$ proposals are awaiting a decision, no more clusters are generated and reports that do not match existing ones are dropped.

Once the $\frac{2}{3}KT$ pool of tokens is reached, the decision is made in favor of the majority votes, individually weighted by deposits. If the threshold exceeds $\frac{2}{3}$, an adversary with less than $\frac{1}{3}$ of the tokens could cause a deadlock simply by not sending any reports. If the threshold is below $\frac{2}{3}$, the adversary could force arbitrary agreements with under $\frac{1}{3}$ of the tokens (i.e., half of $\frac{2}{3}$) by acquiring a majority stake in a proposal. As such, requiring a threshold of $\frac{2}{3}$ of $KT$ tokens guarantees both liveness and safety, as long as adversaries control no more than one third of the tokens [27, 28].

Autonomous Swarm Oracle deployments require recovery mechanisms to prevent faults accumulating over time. To this end, we implement a *reputation system* which rewards reputation tokens to robots in the majority group and penalizes the others, according to Equation 1 in Methods. Each time a consensus is reached, the deposits in reports contradicting the majority are redistributed to the majority (improving safety), and a fixed amount of tokens—established by the issuance constant $I_c$—is awarded to the majority (improving liveness). By using reputation tokens as both participation credentials and weights for future contributions, Swarm Oracle heals itself from the damage caused by Byzantine faults and attacks.

## 2.3    Byzantine faults and attacks

We use the term *Byzantine* to refer to both faults and attacks, emphasizing that Swarm Oracle makes no distinction between the two in its design.

To evaluate Swarm Oracle's fault tolerance, we systematically introduce *Byzantine attackers* into the swarm. The attackers are robots with deliberately modified behaviors designed to degrade Swarm Oracle's safety and/or liveness. In each experiment, all attackers employ the same modified behavior, since a coordinated strategy (collusion) maximizes their chances of success. We investigate four distinct attack types, each targeting a critical aspect of the system's operation:

- **Safety attack**: The attackers report green and blue landmarks as "valuable" and red as "not valuable".

- **Liveness attack**: The attackers abstain from voting.

- **Combined attack**: The attackers combine the two previous strategies: they report for green and blue landmarks, but abstain from red.

- **Physical attack**: The attackers attempt to disrupt other robots' camera readings by physical blocking their line of sight and activating blue LEDs.

Robots also exhibit *Byzantine faults* that occur unintentionally, i.e., without explicit fault injection. Faults can affect safety, for example when obstructions or unusual lighting conditions lead to incorrect proposals. They also affect liveness, for example when robots become stuck due to low-quality obstacle sensors or wheel actuators. Unlike attacks, Byzantine faults are considered transient, since we expect that the robots eventually recover and resume normal operations. If persistent faults occur (e.g., broken motors, damaged sensors or disconnected cables), we halt the experiment and repair the affected robots to prevent faults from accumulating and biasing results over time.

|  |  | Clustering Threshold $R$ (with f = 0) | | | | # of Byzantine Attackers $f$ (with R = 60) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  |  | 20 | 40 | 60 | 80 | 0 | 1 | 2 | 3 | 4 |
|  | Safety | - | - | - | - | - | 10 | 10 | 10 | 10 |
| **Type** | Liveness | - | - | - | - | - | 10 | 10 | 10 | 10 |
| **of** | Combined | - | - | - | - | - | 10 | 10 | 10 | 10 |
| **Attack** | Physical | - | - | - | - | - | 10 | 10 | 10 | 10 |
|  | No attack | 10 | 10 | 10* | 10 | 10* | - | - | - | - |

Table 1: Short-run experiments (real robots), showing the number of trials per configuration. The star symbol (*) indicates when parameter configurations are repeated. In total, we performed 200 trials over a total duration of $22\,\mathrm{h}\,32\,\mathrm{m}$.

## 2.4 Experiments

The experiments were conducted in two batches: *short-run* and *long-run*. In the short-run experiments, performed on a physical swarm of 12 Pi-puck robots [29], we study the effect that the clustering threshold $R$ and the number of Byzantine attackers $f$ have on the quality and costs of agreements (Table 1). The experiments terminate once the first agreement on the valuable landmark is reached. We analyze the costs, defined as the time and number of reports to reach consensus, and the consensus error, measured as a Euclidean distance in RGB space.

In the long-run experiments, simulated in ARGoS [30], we let multiple agreements occur, evaluating how our reputation system performs over time. We vary two key design parameters: the deposit quota $K$ and the issuance constant $I_c$, as shown in Table 2. The simulations accurately replicate key physical parameters of the environment and the robots, including physical collisions, field-of-view occlusions, motion, and communication range. To simulate noisy color perception accurately, each time a robot detects a landmark, it samples at random an RGB observation of that landmark made by a corresponding real robot. A depiction of all RGB observations and Swarm Oracle agreements generated during the short-run experiments is shown in Figure 3.

| | | Deposit Quota K (with $I_c = KT_0$) | | Issuance Constant $I_c$ (with K = $^1/_3$) | |
|---|---|---|---|---|---|
| | | 1 | $^1/_3$ | 0 | $KT_0$ |
| **Type of Attack** | Safety | 40 | 30* | 30 | 30* |
| | Liveness | 40 | 30* | 30 | 30* |
| | Combined | 40 | 30* | 30 | 30* |
| | Physical | 40 | 30* | 30 | 30* |

Table 2: Long-run experiments (in simulation), showing the number of trials per configuration. The star symbol (*) indicates when parameter configurations are repeated. In all long-run configurations we use $f = 4$ Byzantine attackers and $R = 60$ clustering threshold. In total, we performed 400 trials over a total duration of $278\,\mathrm{h}\,52\,\mathrm{m}$.
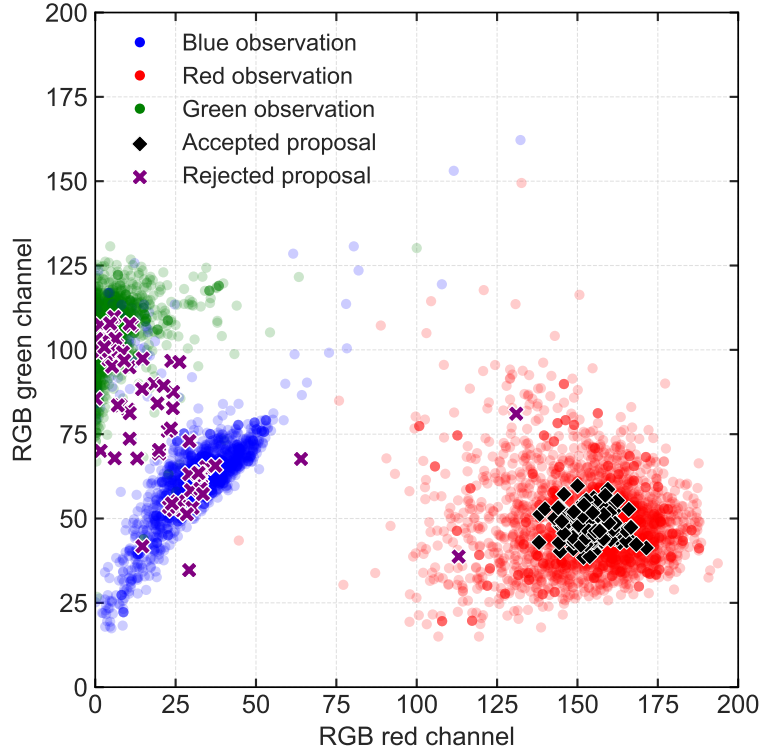


Figure 3: Short-run observations and proposals point cloud (two dimensions shown for clarity), with a clustering threshold $R = 60$. The proposals initiated due to attacker observations are rejected with 100% success rate (purple crosses). Some noisy red observations due to unintentional faults led to two rejections. Increasing $R$ would make the system more inclusive of noisy observations, but also increase the variance in accepted proposals.
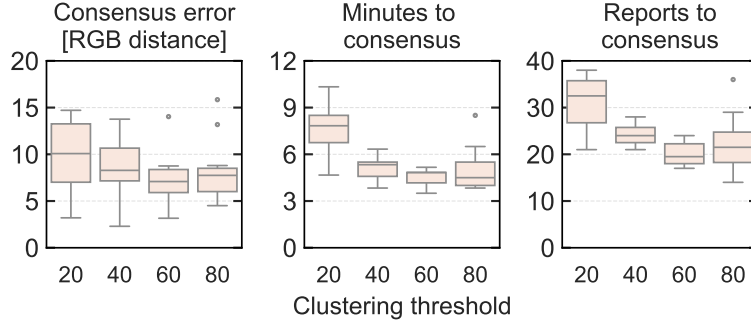
Figure 4: Short-run consensus error (given as a Euclidean distance in the RGB space) and costs when varying the clustering threshold. The plots suggest that, in this application, there is an ideal value for which the error and costs are lowest. We use this value ($R = 60$) for the other experimental configurations.

### 2.4.1 Short-run

Figure 4 shows that setting the clustering threshold $R$ too low leads to an increase in both consensus error and costs. This happens because the reports from robots observing the same landmark may be assigned to different proposals, increasing the time and number of reports required to reach the agreement. However, the creation of multiple proposals does not compromise liveness or safety because proposals are eventually accepted. On the other hand, if $R$ is set too high, there is a risk that observations of distinct landmarks will be incorrectly clustered together. This may pose a safety threat, for example, if noisy observations of red are matched to green or blue proposals, resulting in an incorrect consensus. It is therefore safer to adopt a conservative threshold, though not excessively so, in order to maintain reasonable costs. Based on these results, we fix the clustering threshold at $R = 60$ for all subsequent experiments.

Having established a reasonable clustering threshold, we study how the Swarm Oracle performs when we introduce attacking robots. In Figure 5, the baseline (the arithmetic average of all reports) performs poorly, with high error under safety and combined attacks due to the integration of incorrect observations. It also shows high variance, since in some runs attackers detect colors less often and thus send fewer reports.

Swarm Oracle, with fewer than 4 attackers (one third), is capable of rejecting all incorrect
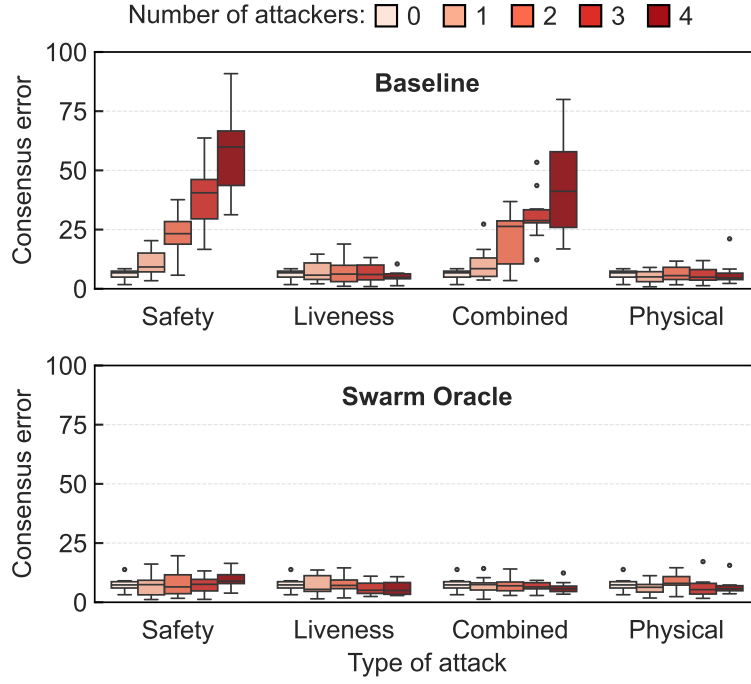
13

Figure 5: Short-run consensus error given as an Euclidean distance in the RGB space.

reports and arriving at an estimate that is as accurate as the attack-free results (the consensus error stays below 20 RGB distance units). Although there is a chance that unintentional faults compound with attacks and cause incorrect agreements, in particular when there are 4 attackers, this did not occur in the short-run experiments (it did however occur in the long-run experiments).

Figure 6 shows that consensus costs, both in terms of time and reports, increase with the number of attackers. However, the reasons behind these increases are fundamentally different for each type of attack.

Facing *safety attacks*, the non-attacking robots must divide their efforts between validating correct proposals and rejecting incorrect ones. As a result, both the time and number of reports required to reach consensus increase. In contrast, during *liveness attacks*, attackers do not submit incorrect proposals, so the number of reports remains constant. However, the time to reach consensus still increases with the number of attackers, at a similar rate to that of safety attacks because, unlike safety attacks—where attackers' rejection votes help
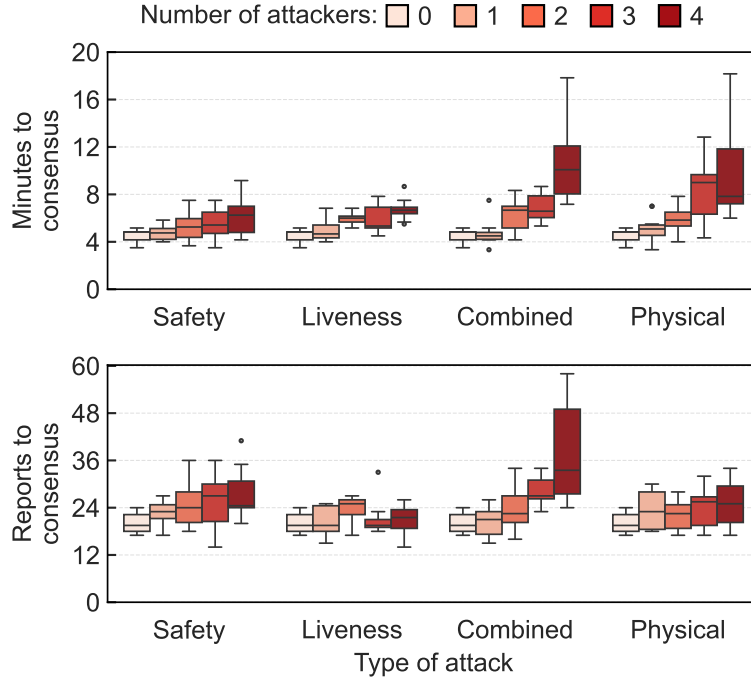
Figure 6: Short-run costs to reach consensus in terms of time and number of reports.

to meet the $\frac{2}{3}$ threshold—liveness attacks make the system more vulnerable to unintentional faults that prevent robots from reporting (e.g., getting stuck or not locating the panel due to obstructions or traffic congestion). The *combined attack* compounds the effect of both previous attacks, resulting in the highest overall costs in terms of time and reports. For this reason, our reputation system penalizes each attack type independently, leading combined attackers to lose reputation faster (see long-run results). Finally, in *physical attacks* the attackers attempt to induce faults in other robots, to deteriorate liveness (by physically blocking them) and safety (by partially obstructing the landmark and turning blue LEDs on to degrade perception quality). However, once physical attackers find the red landmark, they remain static and do not initiate as many incorrect proposals. This results in lower reporting costs compared to the combined attack. Conversely, the time to consensus matches that of the combined attack, since physical attackers can prevent other robots from reaching the landmark. This is most evident with 3 and 4 attackers since they can effectively form a body-wall in front of the red landmark (Figure 7).
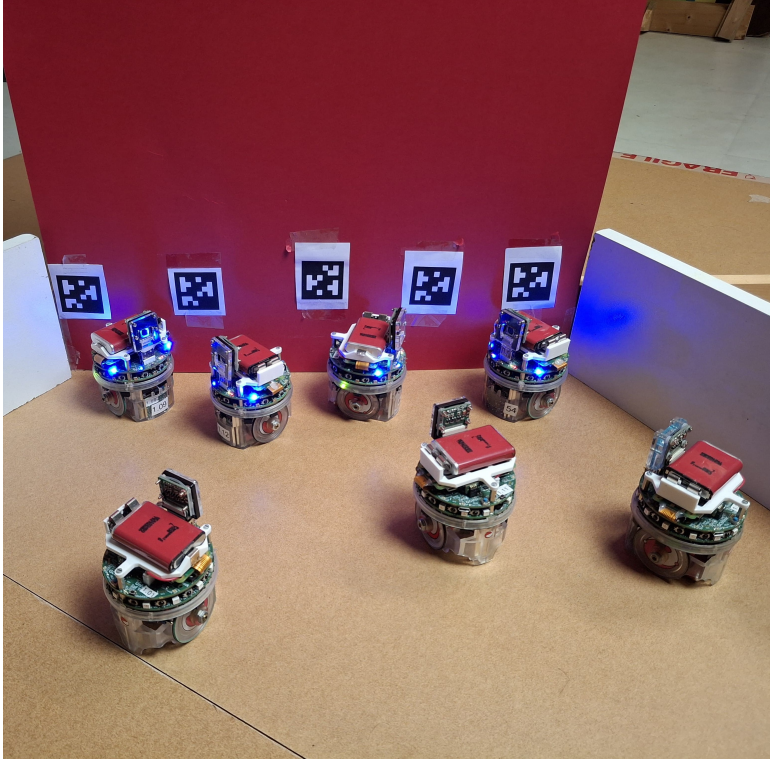
Figure 7: Bodywall formed by the physical attackers.

### 2.4.2 Long-run

Through a series of longer experiments, we evaluate the effectiveness of our reputation system, described by Equation 1. We aim to study the redistribution dynamics and verify that, by the end of the experiments, the non-faulty robots own most of the reputation tokens and that the system has recovered from faults and attacks, thereby regaining both security margins and operational efficiency.

In the experiment shown in Figure 8, each of the 12 robots initially receives an equal share of reputation tokens, therefore the 4 attackers together begin with 33.3% of the total. We chose to evaluate the system under this limit condition, where Swarm Oracle may fail if any additional robot behaves incorrectly, even if unintentionally. This failure occurred in 5 of the 400 long-run trials, in which the attackers ultimately acquired more than 60% of the reputation tokens. These trials are excluded from the plots, as they represent an operating regime beyond Swarm Oracle's intended design. Importantly, once an initial correct agreement

16

is reached, reputation tokens are redistributed, making it increasingly difficult for the attackers to succeed.

Although all faulty robots may lose reputation, the attackers are the most penalized, ending with a significantly lower proportion of the total tokens. The rate at which they lose tokens changes with the parameters regulating the reputation system, namely, the deposit quota $K$ and the issuance constant $I_c$. With $K = 1$, robots can only work on one consensus proposal at a time, leading to faster agreements and redistribution of tokens. Our choice of $K = \frac{1}{3}$ for the majority of the experiments initially leads to slower agreements. However, it allows robots to physically distribute themselves in the environment when contributing to different proposals, resulting in more efficient and robust operation, and in more consistent intervals between subsequent agreements. The issuance constant $I_c$ establishes how many tokens are created when a new agreement is reached. While safety attackers lose their deposits after a successful vote—according to the second term in Equation (1)—, liveness attackers are only penalized when $I_c > 0$. In this case, the first term in Equation (1) dilutes the stake of non-participating robots by issuing new tokens. Finally, different forms of attacks lead to different rates of token loss. Since safety attackers do not abstain from voting, agreements occur more frequently, leading to a faster redistribution of tokens. The combined and physical attackers are penalized by both terms in Equation (1), however, since they succeed at slowing down agreements (Figure 6), reputation loss is also slower.

Figure 9 shows that the cost of reaching agreements decreases over time, as the reputation system gradually mitigates the negative impact of the attackers. After approximately five agreements, most faulty robots have acquired low reputations. This results in a reduction of the average time to consensus, from 6 minutes for the first agreement to 3.7 minutes. The same trend is observed in the number of reports, which decreases from an average of 22 to just 13. These results highlight Swarm Oracle's self-healing ability, where it autonomously returns to an efficient operating state.
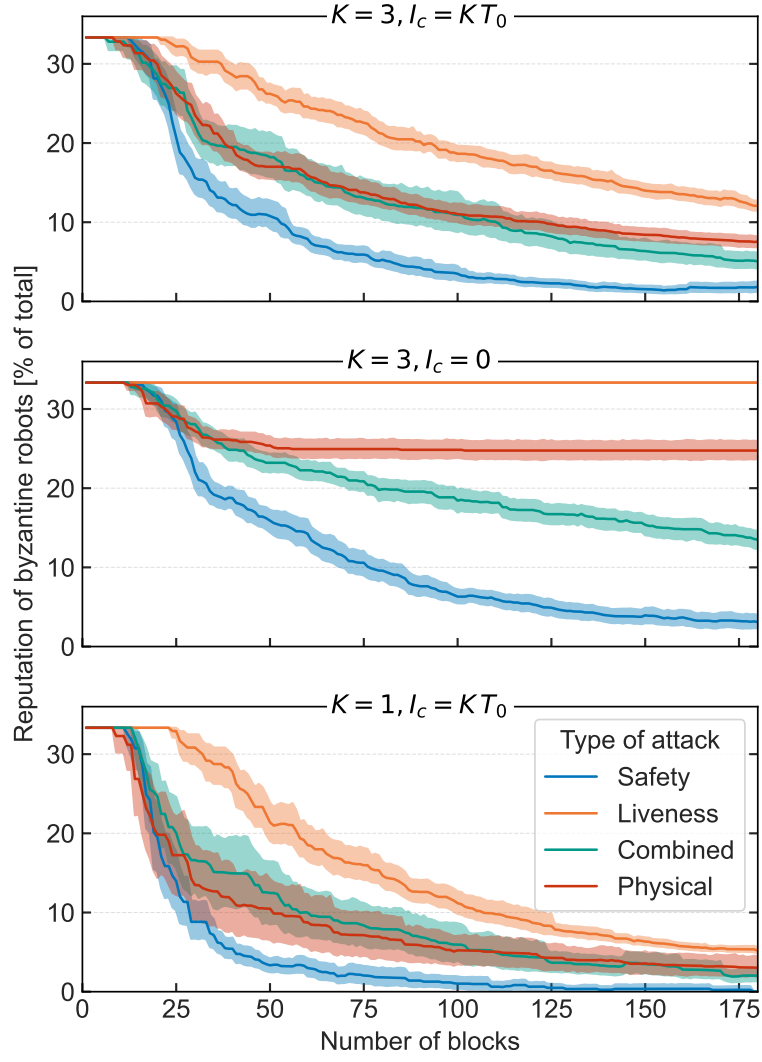
Figure 8: Long-term evolution of the fraction of reputation tokens held by attacking robots. Reputation updates occur as new blocks are added to the blockchain, approximately every 10 seconds.
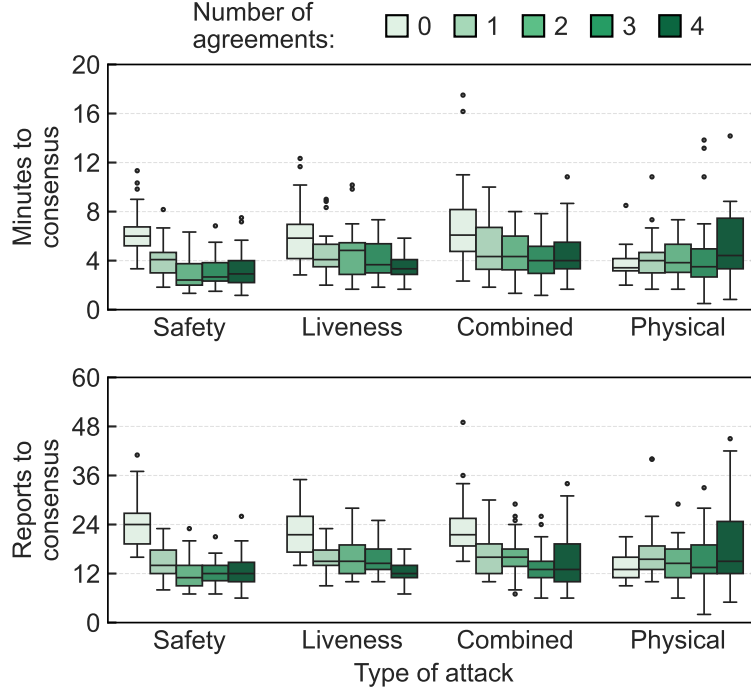
Figure 9: Long-run autonomous recovery from faults, showcasing how costs reduce over time as agreements are reached.

# 3    Discussion

Our short-run experiments show that attacks by one or two robots are effectively neutralized with only a modest increase in consensus time (below 50%). However, when three or four attackers are present, the time to achieve consensus increases sharply. This is most evident when four Byzantine attackers are present since even a single unintentional fault in another robot can delay consensus or potentially lead to an incorrect agreement, if the faulty observation aligns with the attackers'. While incorrect agreements with 4 attackers did not occur in the short-run experiments, it did happen in 5 different long-run experimental trials. Nevertheless, once a first correct agreement is reached, the redistribution of reputation tokens reduces the impact of the attackers, making their success increasingly unlikely. This self-healing property was showcased by the long-run experiments, which demonstrate that the Swarm Oracle can recover over time and progressively nullify the effects of the attacks.

**Scalability** Achieving scalable fault-tolerance is one of Swarm Oracle's key advantages. Since the protocol's security relies on a global condition (in contrast with swarm robotics' local ones), larger swarms can tolerate a greater absolute number of faulty robots: While a swarm of 12 robots tolerates 4 faulty robots, a swarm of 48 tolerates 16. If we assume that 1 in 10 robots are faulty, the smaller swarm can withstand an additional 3 faulty robots, in comparison with 11 for the larger swarm. Larger swarm sizes also make attacks more challenging: attackers must coordinate a larger number of robots, which becomes increasingly difficult given the network size and spatial distribution of the robots.

Hardware costs also scale well, since the number of reports increases linearly with the number of robots, and communications between robots are local and very sparse. These costs are further reduced by our use of a cost-efficient blockchain protocol and dual-layer communications (infrared broadcasts for peer identification, and Wi-Fi for exchanging reports and block metadata). In Methods, we describe this technical setup in detail and show that costs in terms of onboard storage, computation, and communication are manageable. Previous research corroborates this claim with a comprehensive cost analysis with swarm sizes between 8 and 24 robots, and with up to 120 simulated robots [24].

The main scalability challenge arises from the requirement that a supermajority of the robots participate in consensus (more precisely, the holders of $\frac{2}{3}$ of the reputation tokens). This becomes costly when robots and observation areas are either widely dispersed, leading to large physical costs to generate reports, or overly concentrated—potentially leading to physical interference between robots. To mitigate this, Swarm Oracle could be deployed in smaller committees within a larger swarm. The committee size can be tuned to match application requirements, enabling fault-tolerance to scale in accordance with the levels demanded by the application. Our implementation also includes the possibility of parallelization of consensus (through the $K$ parameter), improving scalability by enabling robots to work on multiple consensus proposals at a time.

## 3.1 Possible vulnerabilities and mitigation

**Denial-of-Service** The deposit quota parameter $K$ limits the number of pending proposals, making it possible for attackers that continuously report incorrect observations to occupy all available slots. In theory, liveness is not compromised since asynchronous communication prevents attackers from reporting at the exact time a slot becomes available. However, an attacker could willingly sacrifice all of its reputation to continuously attempt the attack, since reporting only requires a deposit proportional to the robot reputation. One solution is to not allow low-reputation robots to vote, effectively excluding them from Swarm Oracle until they are repaired. Alternatively, a minimum reputation could be required to open new proposals, allowing robots with low reputation to still participate in open proposals and regain reputation through useful contributions.

**Copycats** To ensure the decentralized and trustless execution of the Swarm Oracle smart contract, it is required that observations, as well as robot reputations, are public so that they each compute the same weighted average. However, this transparency introduces a vulnerability: a malicious robot could wait to receive reports from others, and copy or average the observations from the most reputable peers to unfairly gain reputation. This type of attack is known in decentralized finance [31], where trading bots copy the transactions of the most successful traders. Since hiding the reputations or observations is not possible, a potential solution is to instead conceal the votes using cryptography [32]. In this approach, robots submit encrypted votes, which are only revealed once the proposal reaches its deposit threshold (after which no more reports are accepted). Robots that fail to reveal their vote would be penalized (e.g., by losing their deposit), and if a missing reveal is decisive for the outcome, the round is nullified.

**Unobservability** Byzantine robots may generate proposals based on fictitious observations, which are not present in the environment. In such cases, simply relying on robots to report their own observations is insufficient, as unobservable proposals may remain pending. In the worst case, if all available proposal slots are unobservable, Swarm Oracle could reach a

deadlock. The simplest mitigation strategy is to automatically reject proposals that remain pending after a fixed number of minutes, blocks, or reports. Another is to have the robots vote negatively if they fail to observe anything matching the proposal within a reasonable time. In our experiments, we adopted a variant of this approach: robots attempt to *validate* existing proposals by searching the environment for a fixed amount of time, and then travel towards the observation that is the closest match. The resulting report includes the index of the proposal they aim to validate. If the report is not matched to the intended proposal by the clustering algorithm, then the reputation deposit is counted until the $\frac{2}{3}KT$ threshold is reached, after which the proposal is dropped.

**Physical Attacks**   A key aspect of mobile robotics systems is that robot interactions occur not only through communication, but also physically. While our protocol protects against attacks exploiting communication, physical attacks also warrant attention. In our robotics context, physical force may be used to interfere with other robots or manipulate the environment, inducing faults in otherwise non-faulty robots to degrade system performance and avoid accountability. Our experiments showed how physical attackers caused degradation comparable to that of combined attackers, but lost reputation at a slower rate. Since physical interactions occur outside of digital communications, ensuring fault-tolerance remains an open problem. Interestingly, one solution could be to use Swarm Oracle itself to let the robots report observed physical attacks and agree on the identity of the attackers. These robots could then be permanently expelled from Swarm Oracle, and the information could be passed to law enforcement agencies or to the public blockchain, to hold robot stakeholders accountable.

## 3.2   Swarm Oracle Deployment

We deployed the Swarm Oracle on a group of 12 robots, each self-hosting an Ethereum [2] instance. Although a blockchain is not strictly necessary, Swarm Oracle requires methods to ensure an ordering of the reports, track reputation tokens, and prevent double-spending attacks (in which one robot could use its reputation tokens multiple times). Additionally, the robots need to synchronously and correctly execute the deterministic logic that underlies the

Swarm Oracle protocol and its reputation system. These requirements are naturally fulfilled by blockchain smart contracts, which are widely available through open-source blockchains.

Blockchains also help to reduce the communication overhead and storage requirements by aggregating information into blocks generated at specified time intervals (every 10 seconds on average in our experiments). Communication is very efficient since robots exchange reports through sparse and local gossip, and blocks require minimal bandwidth as they only contain essential metadata and the ordering of reports. Finally, the linear structure of the blockchain results in predictable storage growth, and the immutability of the cryptographic links between blocks allows robots to safely prune historical data to further save space.

Swarm Oracle could also be deployed using other distributed ledgers, such as directed acyclic graphs [33] or blocklists [34], or through protocols tailored for swarm robotics, such as SwarmMesh [35] or Buzz [36]. However, these alternatives do not provide the security and convergence guarantees of blockchain-based protocols, and other potential advantages are yet to be demonstrated in practice.

Another alternative is to deploy the Swarm Oracle smart contract on a public blockchain [20]. This reduces the communication and computational burden on the robots, while leveraging the decentralization and automation of smart contracts hosted on public blockchains. We tend to favor the local approach, where the Swarm Oracle protocols are maintained by the robots through a private blockchain, and only the consensus results are periodically uploaded to a public blockchain. This strategy preserves autonomy and fault-tolerance, while minimizing reliance on external connectivity and reducing transaction costs on public blockchains.

Swarm Oracle empowers robots from different stakeholders to pool their capabilities to provide reliable and transparent data as a service. Unlike traditional swarms, Swarm Oracle protocols are designed to withstand adversarial behaviors, enabling collaboration without the need for trust. The data gathered by Swarm Oracle can be provided to external applications, but it can also be employed internally by the swarm. Recent research increasingly explores how robot swarms can employ blockchains to support self-organization, governance and economic participation, while identifying the oracle problem as a fundamental challenge [20, 37, 38].

Swarm Oracle can also be used to securely update distributed data structures maintained by the robots, other than blockchain, such as tuple-spaces [39, 40, 35], environment maps [41, 42, 43], coordinate systems [44, 45], directed acyclic graphs [33, 46], blocklists [34], and Turing-complete state machines [36].

# 4 Methods

## 4.1 The Swarm Oracle

Our Swarm Oracle consists of a group of robots, denoted as $\mathcal{S}$. The robots, operating in a real-world environment, individually collect and process data using their own sensors and computational resources. The data obtained by the robots (e.g., readings from their sensors) forms an observation space, denoted as $\Omega$. The goal is to develop a consensus protocol that lets robots agree on values in $\Omega$ while meeting the safety and liveness requirements: the protocol always produces correct and consistent agreements (safety), and every non-faulty robot eventually decides on a value (liveness). The *consensus set* $\mathcal{L} \in pow(\Omega)$ is the set of all agreements reached by the robots.

As an example, if the robots' observations consist of GPS coordinates marking the positions of valuable resources, then $\Omega = \mathbb{R}^2$, and pow($\Omega$) represents any possible set of positions on the Earth's surface. In this example, the goal of the consensus protocol is to enable the robots to continuously update the set $\mathcal{L}$ with accurate resource positions.

### 4.1.1 Byzantine faults and attacks

Robots operating in the real-world are susceptible to a wide range of potential faults. While some faults can be anticipated and mitigated through careful design of the robots, their behaviors, and their communication protocols, others cannot: a robot may develop faults causing unpredictable and erratic behavior, or multiple robots may collude to attack the system if controlled by a malicious adversary. Such faults are commonly referred to as Byzantine faults [27].

We consider a Byzantine fault model because distinguishing unintentional from malicious faults is often impossible: a robot that fails to send observations due to broken wheels is indistinguishable from one that deliberately stops, and a robot that sends incorrect observations because of a dirty sensor is indistinguishable from one that does so to mislead the consensus. Malicious robots may, however, coordinate their actions to carry out stronger attacks, timing them to exploit faults present in other robots. Therefore, the developed protocol must be *Byzantine fault-tolerant* to guarantee correct and continuous operation in real-world deployments, even in the presence of malicious attacks and unintentional faults.

Faults can also differ significantly in terms of duration: robots may experience temporary faults and subsequently recover, but they can also become permanently disabled or remain under an adversary's control for extended periods of time. Since the multi-robot system is assumed to operate autonomously, the protocol should include mechanisms to mitigate the long-term negative effects of permanent or persistent faults, which could otherwise degrade performance and ultimately lead to system failure. To address this issue, we integrate a *reputation system* into the protocol that weighs each robot's contributions based on its performance, eliminating the need for a repair technician or external authority during operation. Importantly, our reputation system does not exclude faulty robots: they are allowed to continue participating and may regain reputation if they recover from their faults.

### 4.1.2 Protocol

The number of robots in the Swarm Oracle is $N = |\mathcal{S}|$, and each robot has an associated digital identity $i = 1, \ldots, N$, which controls $t_i$ reputation tokens. The reputation tokens of each robot, and consequently the total amount of reputation tokens $T = \sum_{i=1}^{N} t_i$, are globally known.

Let's initially assume that each robot can obtain perfect observations, for instance, in the case their sensors are noise-free. Even though this condition is not realistic and will be later relaxed, it helps to lay down and understand the functioning of our oracle consensus protocol. Robots locally broadcast structured reports, which are disseminated across the swarm. A

report $r$ contains an observation $r.o \in \Omega$, the robot identity $r.i$, a deposit of reputation tokens $r.d$, and a vote, signaling acceptance or rejection, $r.v$. Through the blockchain, Swarm Oracle ensures a global ordering of the reports and that they are structurally valid, creating an array $r$ of ordered reports. Then, through the Swarm Oracle smart contract, the following deterministic logic is applied to the array $r$:

1. The reports are grouped according to their observations. Each group of observations generates a *proposal*.

2. Once a proposal reaches a cumulative value of deposits greater than $\frac{2}{3}T$, the voting decision takes place.

3. An absolute majority is required to determine whether the proposal is accepted or rejected using the identities, deposits, and votes in the reports associated with the proposal.

4. A reputation system redistributes the deposited tokens based on the outcome of the voting round.

5. The reports that have been used in the voting round are removed from $r$.

6. If the proposal is accepted, it is added to the consensus set $\mathcal{L}$.

**Byzantine fault-tolerance**   It is important to understand the rationale behind requiring a $\frac{2}{3}$ quorum of the total reputation tokens in order to validate a proposal. This value is optimal because it preserves both liveness and safety, as long as no more than $\frac{1}{3}$ of the reputation tokens are owned by faulty robots. Any change to this value would decrease the Byzantine fault-tolerance of the system: By increasing it beyond $\frac{2}{3}$, the Byzantine robots could cause a deadlock with less than $\frac{1}{3}$ of the reputation tokens by simply not sending reports. Conversely, by decreasing it, the faulty robots could force incorrect agreements with less than $\frac{1}{3}$ (i.e. half of $\frac{2}{3}$) of the reputation tokens by obtaining the absolute majority stake in a voting round.

**Parallelization**   In some applications, most notably those in swarm robotics, it may be inefficient to require a $\frac{2}{3}$ supermajority to reach agreements, since some robots may be slower than others when submitting their votes. To improve efficiency by enabling the robots to work in parallel, we introduce the *deposit quota* parameter, denoted $K \in\ ]0, 1]$, which establishes the number of reputation tokens that a robot $i$ must deposit with its reports ($r.d = K\,t_i$), as well as the cumulative deposits required before a voting round occurs ($\frac{2}{3}\,K\,T$). Note that this parameter indirectly regulates the maximum number of pending proposals: there can only be $\lfloor K^{-1} \rfloor$ pending proposals at a given moment, otherwise, there would not be enough circulating reputation tokens to reach the cumulative deposits requirement on all of them, and the system could come to a deadlock.

With this in mind, steps 1 and 2 in the logic above are adjusted as follows:

1. The reports are grouped according to their observations, up to a maximum of $\lfloor K^{-1} \rfloor$ groups. Each group of observations generates a *proposal*.

2. Once a proposal reaches a cumulative value of deposits greater than $\frac{2}{3}\,K\,T$, it enters a voting stage.

When observation tasks can be efficiently performed in parallel—for example, when observations result from unpredictable events, such as finding an object during random walks in large environments—then a low value of $K$ may be desirable so that the robots contribute towards multiple proposals simultaneously, exploiting parallelization and reducing physical interference. Conversely, if observations require dedicated effort (e.g., searching a specific area), or when the interference between robots is low, then higher values for $K$ may be warranted, allowing robots to concentrate on validating existing proposals and potentially accelerating the consensus process. The deposit quota $K$ is, therefore, a parameter that can be tuned as a function of the oracle's application context.

**Dealing with noise**   To have the robots reach agreements despite noisy observations, the individual observations of each robot stored in the shared array $\boldsymbol{r}$ are aggregated, e.g., by applying filtering and/or averaging functions. Our proposed method employs a clustering

algorithm to group robot reports based on a similarity score (e.g., a distance function). In doing so, the reports from robots will be matched to the closest proposal. If a report is not matched to existing proposals, a new cluster is created as long as the number of clusters does not surpass $\lfloor K^{-1} \rfloor$ (otherwise, the observation is dropped and the reputation tokens are refunded). An aggregation function is then employed to transform the reports in a cluster $\mathcal{C}_j$, $j \in \{1, \ldots, \lfloor K^{-1} \rfloor\}$ into a proposal $p_j \in \Omega$.

The clustering algorithm, the aggregation function and their parameters are chosen as a function of the oracle's application context and, in particular, of the observation space. In our experiments, we employed incremental k-means clustering [47], with a simple rule for incrementing the number of clusters: a new cluster is generated when the distance between a new observation and all existing proposals exceeds a threshold $R$. The aggregation function generates proposals that correspond to the average value of the observations, weighted by the deposits associated with each report in a cluster. In section 2 we explore the effects of varying the clustering threshold $R$.

**Autonomous recovery from faults**   In computer science, faults are often treated as temporary events, making dedicated detection and mitigation unnecessary, with prevention and tolerance preferred instead [48]. In robotic systems, however, various faults may occur that robots cannot automatically recover from. While in controlled environments, such as warehouses, fault detection may suffice—since faulty robots can be serviced by technicians or removed from operation—in autonomous deployments it is crucial to integrate recovery mechanisms that protect against the accumulation of faults, which could eventually lead to system failure. However, since fault detection processes are not infallible and may yield false positives [49], simply blocking robots identified as faulty from participating in the oracle is not advisable. Our approach to this problem introduces oracle *reputation tokens*, assigned based on each robot's contributions to the oracle and used as weights for its future contributions. A robot that has incurred many faults will have a lower reputation and, consequently, a smaller impact on the oracle's outcome. However, robots are always allowed to continue participating in the oracle and can therefore rebuild their reputation.

### 4.1.3 Reputation system

After each voting round, robots receive or lose reputation tokens based on their contributions. As mentioned above, the purpose of this is to enable the system to autonomously recover from faults: without such a measure, robots that are frequently or persistently faulty could continue to negatively affect the system performance and, as faults accumulate on different robots over time, this could lead to a system-wide failure.

The voting decision occurs when the $\frac{2}{3} K T$ quorum is reached. Given the assumption that faulty robots possess fewer than $\frac{1}{3}$ of the existing tokens, and since robots must deposit exactly $K t_i$ tokens with each report, the absolute majority of tokens in each proposal originated from non-faulty robots. This ensures that the outcome of the voting round is correct. The robots that vote favorably to this outcome should receive reputation tokens, while the robots that vote against it should lose tokens. This reallocation of reputation tokens improves the *safety* of the protocol for subsequent rounds. Additionally, it is important to penalize robots that abstain from participating in the voting round. This can be achieved either by removing some of their reputation tokens, or by diminishing their relative reputation by issuing extra reward tokens to the robots that participated in the voting round. Doing so makes the system more efficient and less likely to come to a deadlock (i.e., improves its *liveness*), since the most active robots gain assets relative to robots that vote less often, or malicious robots that abstain from voting.

Different applications of the oracle may require different choices in the design of the reputation system or its parameters. A common downside of introducing reputation systems in peer-to-peer networks is the risk of introducing vulnerabilities that allow malicious agents to illegitimately acquire reputation tokens [50].

In our experiments, the number of tokens earned by a robot for sending report $r$ is given by the reputation gain function $G$:

$$G(r, \mathcal{C}_j) = \begin{cases} \frac{I_c}{|M(\mathcal{C}_j)|} + \sum\limits_{r' \in \mathcal{C}_j \setminus M(\mathcal{C}_j)} \frac{r'.d}{|M(\mathcal{C}_j)|}, & \text{if } r \in M(\mathcal{C}_j), \\ -r.d, & \text{otherwise} \end{cases} \tag{1}$$

where $r$ is a report in the cluster $\mathcal{C}_j$, $M(\mathcal{C}_j)$ is the set of reports that form the (weighted) absolute majority of the cluster, $I_c$ is the *issuance constant*, the parameter which establishes how many reputation tokens are generated each time a voting round finishes. More specifically, $M(\mathcal{C}_j)$ is defined as the following set:

$$M(\mathcal{C}_j) = \{r \in \mathcal{C}_j \mid r.v = v^*\}$$

where $v^*$ is the unique weighted majority vote such that

$$\sum_{\substack{r \in \mathcal{C}_j \\ r.v = v^*}} r.d > \frac{1}{2} \sum_{r \in \mathcal{C}_j} r.d.$$

Equation (1) states that robots sending reports in $M(\mathcal{C}_j)$, in addition to receiving their deposited tokens $r.d$, will also receive an equal share of the deposits of robots sending reports not in $M(\mathcal{C}_j)$. The robots sending reports in $M(\mathcal{C}_j)$ also receive an equal share of the $I_c$ newly issued tokens. This means that the supply of reputation tokens changes after each voting round: after the $j^{th}$ voting round, the new quantity of circulating reputation tokens is $T_j = T_{j-1} + I_c$. Note that as more reputation tokens are added to the system, the $\frac{2}{3}$ threshold applies to the new quantity of circulating tokens, and the weights associated with existing deposits are adjusted to match the new token supply. In our experiments, we use $I_c = \frac{T_0}{\lfloor K^{-1} \rfloor}$, where $T_0$ is the total amount of reputation tokens at the start of each experiment.

## 4.2 Execution Environment

To implement and test our oracle system, we deploy a private blockchain hosted and maintained by the robots. In this setup, the reports sent by the robots are stored as ordered transactions in the blockchain's distributed ledger, and the oracle's consensus protocol is executed via a smart contract. Previous research has shown that various blockchain consensus protocols can be utilized in robotic systems for this purpose. Pacheco et al. [23] demonstrated that inexpensive robots can execute blockchain software. In their work, the robots hosted a private Ethereum blockchain running a Proof-of-Authority consensus. Other private

blockchains, such as Hyperledger [51], have also been deployed on robots, offering different trade-offs that emphasize different aspects of multi-robot system operation.

As an alternative to employing private blockchains, Swarm Oracle's oracle consensus protocol could be executed as a smart contract on a public blockchain [20], provided that the robots eventually connect to an external network, or on tailored protocols for robotic systems (e.g., Buzz [36]) which are particularly appealing to hardware-constrained systems such as nanorobotics or minimalist robotics [45].

In this work, the Swarm Oracle was executed by twelve Pi-puck robots, each operating as a node in a private Ethereum network running a Proof-of-Authority consensus. Proof-of-Authority is a simple yet effective permissioned protocol that provides the desired security and convergence guarantees for our robotic system [52, 53].

**Blockchain** Each robot functions as an Ethereum node that independently maintains a copy of a blockchain. The reports sent by the robots are ordered and stored as transactions in blocks, and are used as inputs to update the state of a smart contract—which encapsulates the logic behind our oracle consensus protocol. The robots employ Proof-of-Authority [54] consensus to resolve conflicts and agree on new blocks, ensuring that all robots maintain a synchronized state for their locally hosted smart contracts. In contrast with Proof-of-Work [1] consensus, Proof-of-Authority does not require spending computational resources and is therefore suitable for robotics applications. When the participants are perfectly time-synchronized, Proof-of-Authority is Byzantine-tolerant up to $\left(\frac{n}{2} - 1\right)$ Byzantine participants [54]. However, if the participants are not synchronized, this value drops to the theoretical limit of $\frac{n}{3}$ [52, 53].

The reception of transactions and blocks can be delayed due to sparse connectivity between the robots, which is influenced by the size of the environment and by the robots' limited communication range (approximately 15 cm in our experiments). Figure 10 shows the probability distribution of block reception delays, showing that on average robots receive a block 5.52 seconds after it is produced, and that over 90% of the blocks are received before 10 seconds. To allow time for transactions and blocks to synchronize across robots before producing another block, we set the block generation period to 10 seconds. This reduces
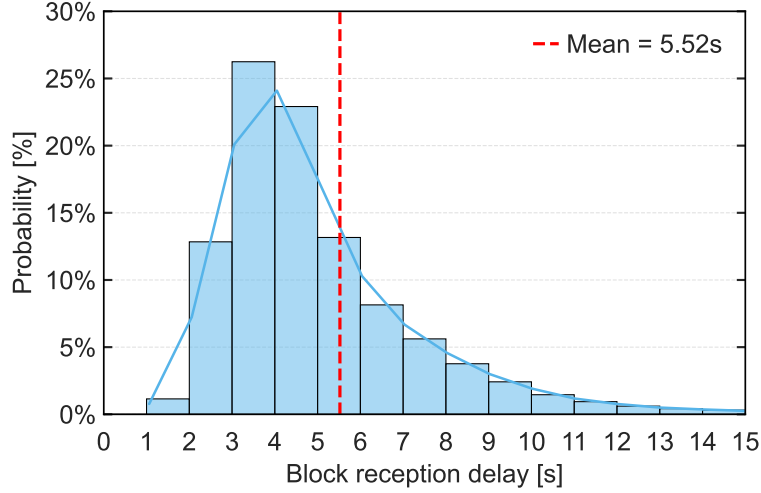
Figure 10: Probability distribution of the block reception delay during short-run experiments (the time elapsed between the creation of a block and its reception by a robot). A total of 20 002 unique blocks were produced, with 157 561 block receptions logged by robots.

forks (conflicting blockchain histories), minimizing bandwidth and processing overhead. For a detailed analysis of block period impact on hardware costs, see [55].

**Robots**   We use Pi-puck robots [56, 29], which have a Raspberry Pi Zero W onboard computer with a 16 GB SD card. This inexpensive hardware is sufficient to run most blockchain protocols that employ permissioned consensus protocols. Figure 11 shows that the hardware demands of running Ethereum on the Pi-pucks, in terms of storage, memory and computing are very manageable. The synchronization of the blockchain occurs over a Wi-Fi mesh network, where each robots communicates only with its physical neighbors. To find neighbors, robots use a range-and-bearing board to broadcast their IP addresses at a short range, and only when a robot receives this broadcast it allows a Wi-Fi connection to exchange blocks and transactions.

Robots perceive their environment using a ring of infrared sensors for obstacle avoidance, and a monocular camera for landmark detection and AprilTag recognition. Figure 12 illustrates the Pi-puck robot and highlights its main sensors and actuators.
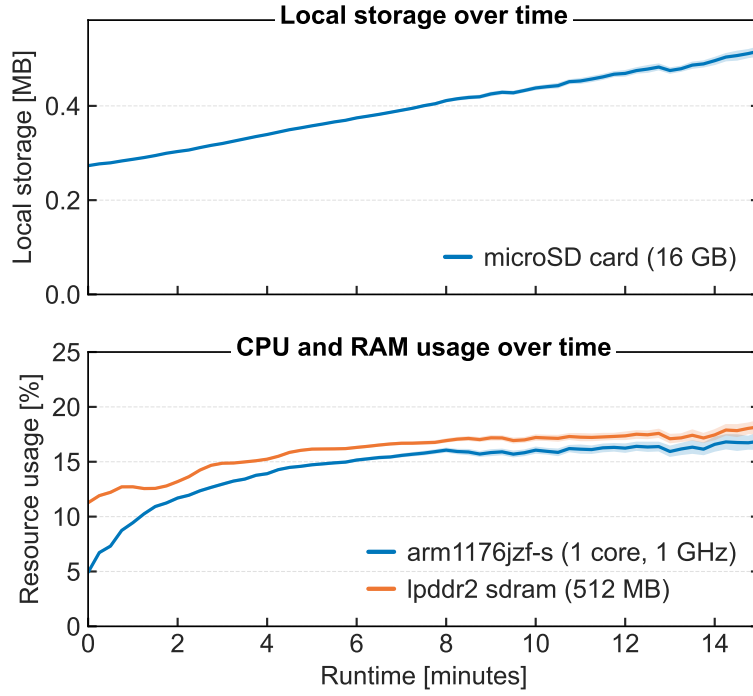
Figure 11: Hardware resource usage by the Ethereum system process and data folder during short-run experiments. Storage, memory, and compute usage remain within the robots' capabilities. The 95% confidence intervals are tight, as the data was collected over a total experimental duration of 22 hours and 32 minutes of experiments, with robots logging values every 10 seconds.
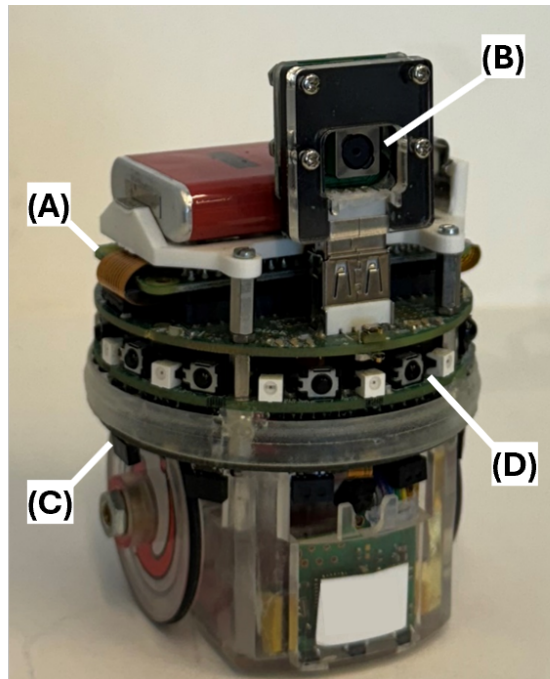
Figure 12: The Pi-puck is a differential drive robot equipped with (A) a Raspberry Pi Zero W with Wi-Fi capabilities, (B) a front-facing camera, (C) a ring of eight infrared sensors used for obstacle avoidance, and (D) a range-and-bearing board that allows robots to communicate with local neighbors.

**Simulator**   In addition to experiments with real robots, we provide extensive results obtained in simulation. The environment and the robots are reconstructed on a 1-to-1 scale in the physics-based simulator ARGoS [30]. The controller of each simulated robot, which is identical to that of the real robots, communicates with Ethereum node instances that are executed inside Docker containers (one for each robot). The simulated robots obtain their camera readings from the dataset collected during the real robot experiments, each randomly drawing from the dataset of a corresponding real robot, and the camera view angles and camera occlusions are tuned to closely match those observed on the real robots. In this way, the simulated results provide a highly accurate extension of the results obtained with the real robots.

**Algorithm**   The Swarm Oracle's fault-tolerant protocol and reputation system are implemented through Algorithm 1, deployed as an Ethereum smart contract. The smart contract provides two functions for oracle participants: `Report`, which allows them to send reports by depositing the required amount of reputation tokens, and `Query`, which returns all accepted and pending proposals.

The `Report` function is executed by sending blockchain transactions that encode the required deposit and report data. Once these transactions are added to the blockchain ledger, the state of the smart contract is updated. The `Query` function retrieves information from the current state of the smart contract, which is locally stored on each robot; this operation does not require a transaction.

The smart contract additionally employs four internal functions: a clustering function `Cluster(`$r$`)`, a distance function between two observations `Distance(`$o_1$`,`$o_2$`)`, a function `Aggregate(`$\mathcal{C}_j$`)` to aggregate the reports in cluster $\mathcal{C}_j$ into a proposal, and a reward function `Reward(`$r$`,`$\mathcal{C}_j$`)` that implements the logic of the reputation system. In this article, we employed incremental k-means clustering and an Euclidean distance function; the proposals correspond to the geometric center of the observations, weighted by the deposited reputation tokens; and the reward function is the one given in equation (1). The modular design of these functions allows for easy replacement with other options better suited to different applications

of the oracle. This design choice helps keep the protocol generic.

---

**Algorithm 1:** Oracle network consensus protocol

---

**Input**      : $\mathcal{S}$, $\Omega$
**Output**     : $\mathcal{L} \in pow(\Omega)$
**Parameters** : $K$, $I_c$
**State**      : $r$, $T$, $\mathcal{L}$, $\mathcal{C} = \{\mathcal{C}_j\}$, $j \in \{1, \ldots, \lfloor K^{-1} \rfloor\}$
**Initialization:** $r \leftarrow [\,]$, $T \leftarrow T_0$, $\mathcal{L} \leftarrow \emptyset$, $\mathcal{C} \leftarrow \emptyset$, $j \leftarrow 0$

**Procedure** `Report`$(r = (o, i, d, v))$
    **if** $r.d < K\,$`Balance`$(r.i) \vee r \in r$ **then**
        ▷ `Reject due to insufficient deposit`
        ▷ `or due to duplicated reporting`
        Transfer $r.d$ reputation tokens to $r.i$ ;
        return 0 ;

    $r \leftarrow$ append$(r, r)$ ;
    ▷ `Generate the set of clusters` $\mathcal{C}$
    $\mathcal{C} \leftarrow$ `Cluster`$(r)$ ;
    **for** $\mathcal{C}_j \in \mathcal{C}$ **do**
        **if** $\sum_{r \in C_j} r.d \geq \frac{2}{3} KT$ **then**
            **if** *weighted majority agreement* **then**
                $M(\mathcal{C}_j) \leftarrow r : r \in$ absolute majority ;
                ▷ `Reward reporting robots`
                **for** $r \in \mathcal{C}_j$ **do**
                    $reward \leftarrow G(r, \mathcal{C}_j)$ ;
                    Transfer $reward + r.d$ reputation tokens to $r.i$ ;

                ▷ `Update state variables`
                $T \leftarrow T + I_c$ ;
                $\mathcal{L} \leftarrow \mathcal{L} \bigcup$ `Aggregate`$(\mathcal{C}_j)$ ;
                $\mathcal{C} \leftarrow \mathcal{C} \backslash \{\mathcal{C}_j\}$ ;

**Procedure** `Query()`
    $proposals \leftarrow \{$`Aggregate`$(\mathcal{C}_j) : \forall \mathcal{C}_j \in \mathcal{C}\}$ ;
    $consensus \leftarrow \mathcal{L}$ ;
    **return** $proposals$, $consensus$ ;

---

## 4.3 Experimental Setup

### 4.3.1 Robot behavior

Our robots follow a routine that mainly consists of *exploring* the environment to find land-marks and initiate new proposals, or *validating* existing proposals in the smart contract. To achieve this, the robots detect the colors of landmarks in the environment and navigate to-wards them to measure their RGB values and read the AprilTags. The robots' state machine is illustrated in a flowchart in Figure 13.

1. `Query` The robot queries the current proposals from the smart contract. If there are proposals it has not yet validated, it randomly selects one and transitions to the `Validate` state; otherwise, it transitions to the `Explore` state.

2. `Validate` The robot searches the environment for landmarks and navigates to the one that is the closest match to the proposal to be validated. Upon arrival, it transitions to the *Report* state, or returns to the `Query` state after 100 s.

3. `Explore` The robot searches the environment for any landmark and then moves towards it. Upon arrival, it transitions to the `Report` state, or returns to the `Query` state after 100 s.

4. `Report` Once the distance to the AprilTag attached to a landmark is smaller than a threshold, the robot reads the AprilTag and the average RGB value of the largest color contour in its camera vision. Then it sends a report through a blockchain transaction and returns to the `Query` state.

### 4.3.2 Parameters

**Clustering threshold**   The clustering threshold $R$ is an internal parameter of the k-means clustering algorithm that defines the minimum required distance between a new observation and existing cluster centroids (i.e., proposals) for a new cluster to be initiated. In our experiments and real deployments, noisy camera readings can lead to faults, such as incorrect votes or robots failing to find proposals. To examine the impact of faults caused by noisy measurements, we could keep a fixed $R$ and add noise or biases to the readings gathered by the robots. However, this approach would require artificial noise models, which are less interesting to analyze compared to the natural noise and color miscalibration already affecting the robots' readings. Instead we run a set of short-run experiments where we vary the value of $R$, to analyze the impact of faults caused by erroneous classifications due to noise. This is showcased in Figure 4.
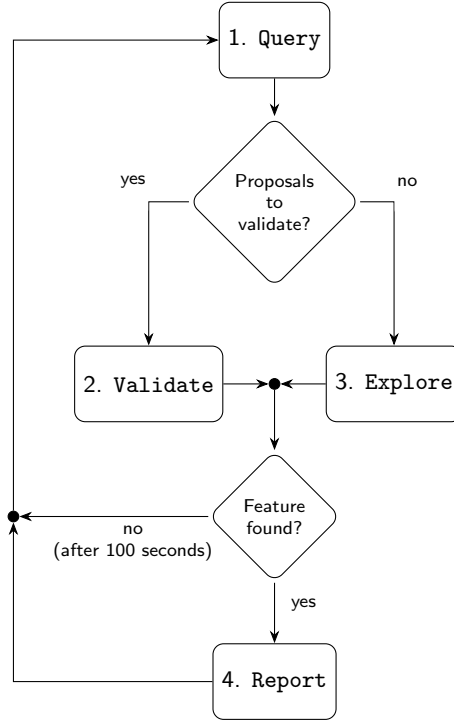
Figure 13: The robots' state machine as a flowchart.

**Number of Byzantine attackers**   In the short-run experiments, we also vary the number of Byzantine attackers $f$. One of the main features of employing a global protocol is the certainty of correct agreements provided Byzantine robots remain below the $\frac{1}{3}$ threshold. This is shown in Figures 5 and 6, which shows that our method resists all attacks up to this point (4 Byzantine attackers, in a total of 12 robots).

**Deposit quota**   The deposit quota $K$ establishes the percentage of reputation tokens that each robot must deposit when submitting a report. Indirectly, it regulates the maximum number of pending proposals: if $K = 1$, only one proposal can remain pending, otherwise $\lfloor K^{-1} \rfloor$ proposals can remain open. Notably, this has an effect on the dynamics of the token redistribution between the robots. In Figure 8 we can compare $K = \frac{1}{3}$ at the top, with $K = 1$ at the bottom. With $K = \frac{1}{3}$, the redistribution of tokens is less steep (i.e., more time and agreements are required for Byzantine robots to lose tokens).

**Issuance constant** The issuance constant $I_c$ is employed in the first term of equation (1) to generate new tokens that are rewarded to the robots that participated in the voting round. This generates an inflationary effect that penalizes robots that abstain from voting, and improves the liveness of the protocol. If the issuance constant is set to zero, the robots performing the attack on liveness are not penalized (Figure 8). During a real deployment, this could lead to the accumulation of faults due to more robots becoming stuck or disabled, potentially compromising the protocol's liveness.

### 4.3.3 Metrics and Baseline

**Consensus Error** The consensus error is the Euclidean distance between the agreements generated by Swarm Oracle and the average value of all red observations made by non-attacking robots during the short-run experiments. Since the number of experiments was large, we consider the average values to be a good estimate of the RGB values of the landmarks' colors. A point cloud plot of these observations is shown in Figure 3.

Using an estimated value is needed, since the true RGB values of the color panels are unknown (they depend on the camera sensors, their parameters and lighting conditions). The lower the consensus error, the more precise our oracle was, i.e., the closer the robots' final agreement is to the average value.

**Time to Consensus** This is the time it took, in minutes, until an agreement on a "valuable" landmark is reached and recorded on a blockchain block. In the short-run experiments, it consists of the time to reach the first agreement, while in the long-run experiments we reset the clock each time an agreement occurs.

**Reports to Consensus** Similarly to time, this is the number of reports sent by the non-attacking robots before an agreement occurs. While time is an intuitive metric, the number of reports is invariant to factors such as the speed of the robots and how fast they can make observations, which are respectively related to, for example, the battery levels (which change as the experiments progress), and to the environment and observation spaces (which change

depending on the scenario).

**Baseline** Distributed agreements in robotics are typically achieved through approximate local consensus, where each robot maintains an individual opinion that it shares with its neighbors. In such algorithms, the exact moment at which consensus is reached is not well defined, as agreement is determined by a threshold rather than by a discrete decision event. Additionally, robots are susceptible to attacks in which malicious robots exploit locality— for example, by flooding an isolated robot with incorrect reports—which makes it difficult or impossible to come up with safety guarantees. Due to these fundamental differences with respect to our exact consensus method, we do not consider approximate consensus approaches to be an adequate baseline. Instead, we opt for simplicity by comparing our approach to a baseline that consists of centrally collecting all observations made by the robots until the end of that experiment, and averaging them.

# References

[1] Satoshi Nakamoto. Bitcoin: A peer-to-peer electronic cash system. White paper, Bitcoin, 2008. `https://bitcoin.org/bitcoin.pdf`, Accessed July 20, 2025.

[2] Vitalik Buterin. A next-generation smart contract and decentralized application platform. White paper, Ethereum Foundation, 2014. `https://ethereum.org/en/whitepaper/`, Accessed July 20, 2025.

[3] Andreas Antonopoulos and Gavin Wood. *Mastering Ethereum: Building Smart Contracts and DApps*. O'Reilly Media, Sebastopol, CA, 2018.

[4] Alan Mathison Turing. Systems of logic based on ordinals. *Proceedings of the London Mathematical Society, Series 2*, 45:161–228, 1939.

[5] Giulio Caldarelli. Understanding the blockchain oracle problem: A call for action. *Information*, 11(11):509, 2020.

[6] Lorenz Breidenbach, Christian Cachin, Benedict Chan, Alex Coventry, Steve Ellis, Ari Juels, Farinaz Koushanfar, Andrew Miller, Brendan Magauran, Daniel Moroz, and Sergey Nazarov. Chainlink 2.0: Next steps in the evolution of decentralized oracle networks. White paper, Chainlink Labs, 2021. `https://research.chain.link/whitepaper-v2.pdf`, Accessed July 20, 2025.

[7] Abdeljalil Beniiche. A study of blockchain oracles. *arXiv preprint arXiv:2004.07140*, 2020. `https://arxiv.org/abs/2004.07140`, Accessed July 20, 2025.

[8] Amir Haleem, Andrew Allen, Andrew Thompson, Marc Nijdam, and Rahul Garg. Helium: A decentralized wireless communication network. *RS Open Journal on Innovative Communication Technologies*, 4(13), 2024.

[9] Marco Dorigo, Mauro Birattari, and Manuele Brambilla. Swarm robotics. *Scholarpedia*, 9(1):1463, 2014.

[10] Francis Galton. Vox populi. *Nature*, 75:450–451, 1907.

[11] Gabriele Valentini, Heiko Hamann, and Marco Dorigo. Self-organized collective decisions in a robot swarm. In *AAAI-15 Video Proceedings*. AAAI Press, 2015.

[12] Samer Hassan and Primavera De Filippi. Decentralized autonomous organization. *Internet Policy Review: Journal on Internet Regulation*, 10(2), 2021.

[13] Zhibin Lin, Taotao Wang, Long Shi, Shengli Zhang, and Bin Cao. Decentralized physical infrastructure networks (DePIN): Challenges and opportunities. *IEEE Network*, 39:91–99, 2025.

[14] Gabriele Valentini, Eliseo Ferrante, and Marco Dorigo. The best-of-n problem in robot swarms: Formalization, state of the art, and novel perspectives. *Frontiers in Robotics and AI*, 4(9), 2017.

[15] Heath J. LeBlanc, Haotian Zhang, Xenofon D. Koutsoukos, and Shreyas Sundaram. Resilient asymptotic consensus in robust networks. *IEEE Journal on Selected Areas in Communications*, 31(4):766–781, 2013.

[16] Heath LeBlanc and Xenofon Koutsoukos. Consensus in networked multi-agent systems with adversaries. In *Proceedings of the 14th International Conference on Hybrid Systems: Computation and Control*, HSCC '11, pages 281–290. Association for Computing Machinery, 2011.

[17] Haotian Zhang and Shreyas Sundaram. Robustness of information diffusion algorithms to locally bounded adversaries. In *2012 American Control Conference (ACC)*, pages 5855–5861. IEEE, 2012.

[18] Heath J. LeBlanc and Xenofon D. Koutsoukos. Low complexity resilient consensus in networked multi-agent systems with adversaries. In *Proceedings of the 15th ACM International Conference on Hybrid Systems: Computation and Control*, HSCC '12, pages 5–14. Association for Computing Machinery, 2012.

[19] Roman Mühlberger, Stefan Bachhofner, Eduardo Castelló Ferrer, Claudio Di Ciccio, Ingo Weber, Maximilian Wöhrer, and Uwe Zdun. Foundational oracle patterns: Connecting blockchain to the off-chain world. In *Business Process Management: Blockchain and Robotic Process Automation Forum*, pages 35–51. Springer International Publishing, 2020.

[20] Marco Dorigo, Alexandre Pacheco, Andreagiovanni Reina, and Volker Strobel. Blockchain technology for mobile multi-robot systems. *Nature Reviews Electrical Engineering*, 1(4):264–274, 2024.

[21] Volker Strobel, Eduardo Castelló Ferrer, and Marco Dorigo. 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 2018)*, pages 541–549. International Foundation for Autonomous Agents and Multiagent Systems, 2018.

[22] Volker Strobel, Eduardo Castelló Ferrer, and Marco Dorigo. Blockchain technology secures robot swarms: A comparison of consensus protocols and their resilience to Byzantine robots. *Frontiers in Robotics and AI*, 7:54, 2020.

[23] Alexandre Pacheco, Volker Strobel, and Marco Dorigo. A blockchain-controlled physical robot swarm communicating via an ad-hoc network. In *Swarm Intelligence – Proceedings of ANTS 2020 – Twelfth International Conference*, volume 12421 of *Lecture Notes in Computer Science*, pages 3–15. Springer, 2020.

[24] Volker Strobel, Alexandre Pacheco, and Marco Dorigo. Robot swarms neutralize harmful Byzantine robots using a blockchain-based token economy. *Science Robotics*, 8(79):eabm4636, 2023.

[25] Hanqing Zhao, Alexandre Pacheco, Volker Strobel, Andreagiovanni Reina, Xu Liu, Gregory Dudek, and Marco Dorigo. A generic framework for byzantine-tolerant consensus achievement in robot swarms. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2023)*, pages 8839–8846. IEEE Press, 2023.

[26] Leslie Lamport. Proving the correctness of multiprocess programs. *IEEE Transactions on Software Engineering*, SE-3(2):125–143, 1977.

[27] Leslie Lamport, Robert Shostak, and Marshall Pease. The Byzantine generals problem. *ACM Transactions on Programming Languages and Systems*, 4(3):382–401, 1982.

[28] Miguel Castro and Barbara Liskov. Practical byzantine fault tolerance. In *Proceedings of the Third Symposium on Operating Systems Design and Implementation (OSDI '99)*, pages 173–186. ACM Press, 1999.

[29] Alan Millard, Russell Joyce, James Hilder, Cristian Fleşeriu, Leonard Newbrook, Wei Li, Liam McDaid, and David Halliday. The Pi-puck extension board: A Raspberry Pi interface for the e-puck robot platform. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2017)*, pages 741–748. IEEE Press, 2017.

[30] Carlo Pinciroli, Vito Trianni, Rehan O'Grady, Giovanni Pini, Arne Brutschy, Manuele Brambilla, Nithin Mathews, Eliseo Ferrante, Gianni Di Caro, Frederick Ducatelle, Mauro Birattari, Luca Maria Gambardella, and Marco Dorigo. ARGoS: A modular, parallel, multi-engine simulator for multi-robot systems. *Swarm Intelligence*, 6(4):271–295, 2012.

[31] Mikolaj Barczentewicz, Alexander Sarch, and Natasha Vasan. Battle of the crypto bots: Automated transaction copying in decentralized finance. *SSRN Electronic Journal*, 26(3):672–730, 2023.

[32] Ludéric Van Calck, Alexandre Pacheco, Volker Strobel, Marco Dorigo, and Andreagiovanni Reina. A blockchain-based information market to incentivise cooperation in swarms of self-interested robots. *Scientific Reports*, 13:20417, 2023.

[33] Jason A. Tran, Gowri S. Ramachandran, Palash M. Shah, Claudiu B. Danilov, Rodolfo A. Santiago, and Bhaskar Krishnamachari. SwarmDAG: A partition tolerant distributed ledger protocol for swarm robotics. *Ledger*, 4:25–31, 2019.

[34] Kacper Wardega, Max von Hippel, Roberto Tron, Cristina Nita-Rotaru, and Wenchao Li. Byzantine resilience at swarm scale: A decentralized blocklist protocol from inter-robot accusations. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems (AAMAS '23)*, pages 1430–1438. International Foundation for Autonomous Agents and Multiagent Systems, 2023.

[35] Nathalie Majcherczyk and Carlo Pinciroli. SwarmMesh: A distributed data structure for cooperative multi-robot applications. In *IEEE International Conference on Robotics and Automation (ICRA 2020)*, pages 4059–4065. IEEE Press, 2020.

[36] Carlo Pinciroli and Giovanni Beltrame. Buzz: A programming language for robot swarms. *IEEE Software*, 33(4):97–100, 2016.

[37] Eduardo Castelló Ferrer. The blockchain: A new framework for robotic swarm systems. In *Proceedings of the Future Technologies Conference (FTC 2018)*, volume 881 of *Advances in Intelligent Systems and Computing*, pages 1037–1058. Springer, 2018.

[38] Jorge Peña Queralta, Farhad Keramat, Salma Salimi, Lei Fu, Xianjia Yu, and Tomi Westerlund. Blockchain and emerging distributed ledger technologies for decentralized multi-robot systems. *Current Robotics Reports*, 4:43–54, 2023.

[39] Carlo Pinciroli, Adam Lee-Brown, and Giovanni Beltrame. A tuple space for data sharing in robot swarms. In *9th EAI International Conference on Bio-inspired Information and Communications Technologies*. European Union Digital Library, 2016.

[40] Rocco De Nicola, Luca Di Stefano, and Omar Inverso. Multi-agent systems with virtual stigmergy. *Science of Computer Programming*, 187:102345, 2020.

[41] Miquel Kegeleirs, Giorgio Grisetti, and Mauro Birattari. Swarm SLAM: Challenges and perspectives. *Frontiers in Robotics and AI*, 8:618268, 2021.

[42] Pierre-Yves Lajoie and Giovanni Beltrame. Swarm-SLAM: Sparse decentralized collaborative simultaneous localization and mapping framework for multi-robot systems. *IEEE Robotics and Automation Letters*, 9(1):475–482, 2024.

[43] Andreas Birk and Stefano Carpin. Merging occupancy grid maps from multiple robots. *IEEE Proceedings, Special Issue on Multi-Robot Systems*, 94(7):1384–1397, 2006.

[44] Simon Jones and Sabine Hauert. Distributed spatial awareness for robot swarms. In *17th International Symposium on Distributed Autonomous Robotic Systems (DARS2024)*, Springer Proceedings in Advanced Robotics. Springer, 2024. In print.

[45] Michal Pluhacek, Simon Garnier, and Andreagiovanni Reina. Decentralised construction of a global coordinate system in a large swarm of minimalistic robots. *Swarm Intelligence*, 19(3), 2025. In print.

[46] Farhad Keramat, Jorge Peña Queralta, and Tomi Westerlund. Partition-tolerant and Byzantine-tolerant decision making for distributed robotic systems with IOTA and ROS2. *IEEE Internet of Things Journal*, 10(14):12985–12998, 2023.

[47] DT Pham, SS Dimov, and CD Nguyen. An incremental k-means algorithm. *Proceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering Science*, 218(7):783–795, 2004.

[48] Tom Anderson and John Knight. A framework for software fault tolerance in real-time systems. *IEEE Transactions on Software Engineering*, SE-9(3):355–364, 1983.

[49] Nancy A. Lynch. *Distributed Algorithms*. The Morgan Kaufmann Series in Data Management Systems. Morgan Kaufmann, San Francisco, CA, 1996.

[50] Jordi Sabater and Carles Sierra. Review on Computational Trust and Reputation Models. *Artificial Intelligence Review*, 24(1):33–60, 2005.

[51] Salma Salimi, Jorge Peña Queralta, and Tomi Westerlund. Hyperledger Fabric blockchain and ROS 2 integration for autonomous mobile robots. In *2023 IEEE/SICE International Symposium on System Integration*, pages 1–8. IEEE Press, 2023.

[52] Parinya Ekparinya, Vincent Gramoli, and Guillaume Jourjon. The attack of the clones against proof-of-authority. In *Proceedings of the 27th Annual Network and Distributed System Security Symposium (NDSS 2020)*. The Internet Society, 2020.

[53] Qin Wang, Rujia Li, Qi Wang, Shiping Chen, and Yang Xiang. Exploring unfairness on proof of authority: Order manipulation attacks and remedies. In *Proceedings of the 2022 ACM on Asia Conference on Computer and Communications Security (ASIA CCS 2022)*, pages 123–137. ACM, 2022.

[54] Péter Szilágyi. EIP 225: Clique proof-of-authority consensus protocol. `https://eips.ethereum.org/EIPS/eip-225`, 2017. Accessed July 10, 2025.

[55] Alexandre Pacheco, Volker Strobel, Andreagiovanni Reina, and Marco Dorigo. Real-time coordination of a foraging robot swarm using blockchain smart contracts. In *Swarm Intelligence – Proceedings of ANTS 2022 – 13th International Conference*, volume 13491 of *Lecture Notes in Computer Science*, pages 196–208. Springer, 2022.

[56] Francesco Mondada, Michael Bonani, Xavier Raemy, James Pugh, Christopher Cianci, Adam Klaptocz, Stephane Magnenat, Jean-Christophe Zufferey, Dario Floreano, and Alcherio Martinoli. The e-puck, a robot designed for education in engineering. In *Proceedings of the 9th Conference on Autonomous Robot Systems and Competitions*, pages 59–65. Instituto Politécnico de Castelo Branco, 2009.