Peer reviewed version

## University of Bristol - Explore Bristol Research
### General rights

1 **A checklist for safe robot swarms**
2
3 Edmund Hunt and Sabine Hauert*
4 Engineering Mathematics, Bristol Robotics Laboratory, University of Bristol
5
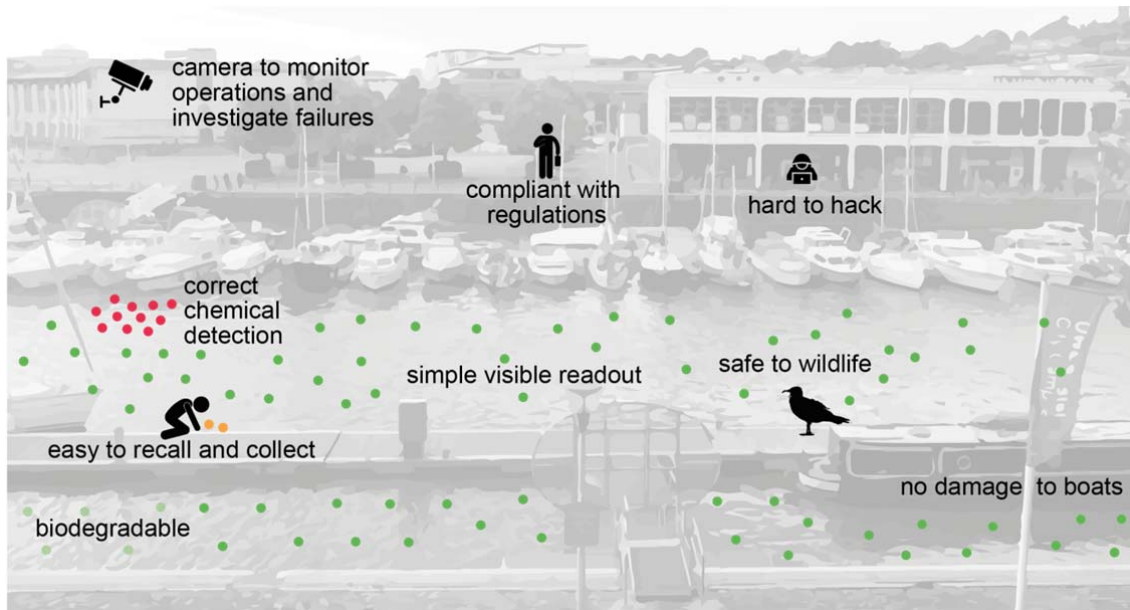6 *corresponding author: sabine.hauert@bristol.ac.uk
7
8 *Standfirst: As robot swarms move from the laboratory to real world applications, a routine*
9 *checklist of questions could help ensure their safe operation.*
10
11 Robot swarms promise to tackle problems ranging from food production and natural
12 disaster response, to logistics and space exploration[1–4]. As swarms are deployed outside the
13 laboratory in real world applications, we have a unique opportunity to engineer them to be
14 safe from the get-go. Safe for the public, safe for the environment, and indeed, safe for
15 themselves. This will help build public confidence in their use, and counter hyped or
16 negative narratives about swarms in media and science fiction. Designing safe swarms is
17 also challenging, as the main benefits of swarms, namely their scalability, robustness, and
18 emergent properties, arise from self-organisation, a concept rarely used in engineering[5].
19
20 Previous research has identified certain challenges for the deployment of safe robot
21 swarms, particularly in the areas of swarm agent fault tolerance[6-9], human-swarm
22 interaction[8] and swarm security[11–15], but limited consideration has been given to systematic
23 assessment of swarm safety. As a starting point, we propose a preliminary "safe swarm
24 checklist" with 10 questions that should be answered satisfactorily by engineers before a
25 swarm can be deployed in the real world, where real costs are at stake. Highlighting
26 potential risks early in the swarm design phase will allow mitigations to be introduced.
27
28 Safety in engineering can be defined as the absence of catastrophic consequences on the
29 user(s) and the environment. It is closely related to concepts of dependability, or the ability
30 to deliver a service that can justifiably be trusted[16]. We take a holistic view of safety that
31 goes beyond analysing failure modes and performing risk analysis[17], to also include the
32 broader socio-technical context of deployment. In our proposed "safe swarm checklist",
33 questions 1 and 2 on ethics and legality come first as a vital prerequisite for initial testing.
34 Ethical governance and training should be pervasive from the design to the deployment of
35 robot swarms[18]. Questions 3 and 4 relate to accountability and user-swarm interactions.
36 Then, because a defining feature of swarms is their emergent capabilities, we consider
37 individual and swarm-level risks separately for each of the dimensions of physical harm,
38 behavioural harm, and security in questions 5 to 10.
39
40 We briefly apply our checklist to a hypothetical swarm of 100 small floating robots – let's
41 call them bubblebots – deployed to monitor water pollutants (Figure 1). The idea builds on
42 several examples of real-world robot swarm deployments in aquatic environments[19–21]. The
43 bubblebots are meant to distribute over an enclosed floating harbour and light up in ways
44 that communicate their local sensor readings. By sharing information within the swarm, the
45 robots can collectively communicate the overall state of the water in the harbour and
46 reorganise to highlight pollution sources.
47

Figure 1 Safety considerations for the deployment of bubblebots used in a floating harbour to signal pollutants.

(1)  Ethics. Is this an ethical use of a robot swarm?

We will consult authorities such as the IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems and its work on Ethically Aligned Design[22], or the BSI standard for Ethical Design and Application of Robots and Robotic Systems[23]. With bubblebots, we focus on an application for social good, namely environmental monitoring of water pollutants, considering privacy and potential harm to actors in the harbour. Mutual shaping of the technology between researchers and users will help embed local ethical norms[24].

(2)  Legal. Does the swarm comply with all relevant laws and regulations for the domain(s) of deployment?

The bubblebot swarm will need to comply with all relevant rules: environmental, harbour and maritime, or relating to health and safety. There may be a need for public liability insurance.

(3)  Accountability. Is there a way to analyse swarm failures?

Following work by Winfield et al.[25], it would be helpful to store short-term recordings of the actions of the robots based on sensory readings and communication in a so-called "black box", inspired from flight recorders in the aviation industry. This could be done on board the robots, or using an external camera system monitoring overall operations. This information would help to investigate and reconstruct conditions that led to unsafe operations, and would be used to improve swarm implementation if things go wrong.

(4)  User interaction. Can the users interact with the swarm to prevent unwanted behaviour?

80

81　It should be possible to easily deploy, interact with and retrieve the swarm. In this case, user
82　interaction will involve depositing the bubblebots in the harbour, and reading out the state
83　of the swarm from the harbourside by looking at the robot location and colour status.
84　Bubblebots can easily be stopped using a broadcasted message transmitted throughout the
85　swarm from an operator on the harbourside, in which case robots will home to one area of
86　the harbour for collection.

87

88　(5)　Physical harm from individual robots. Can the individual robots cause physical harm
89　to humans, animals, or the environment?

90

91　Bubblebots will be designed to be small enough to avoid damage to boats in the harbour, or
92　other robots, but large enough to avoid seabirds and fish from eating them. Trials will be
93　done to check that they are compatible with actors in the harbour. They will be buoyant to
94　avoid them sinking and becoming a pollutant themselves, and will be easy to detect for
95　collection by harbour staff. Materials for the waterproof shell will be optimised for
96　durability to avoid breaches, and electronics will be low enough power to avoid possibility of
97　electric shock. Bubblebots failing (power loss, broken sensor or motors) will turn off,
98　avoiding further impact. In the future, bubblebots could even be biodegradable as an
99　additional safeguard – such research is moving beyond the conceptual stage[26].

100

101　(6)　Physical harm from the swarm. Can the emergent swarm behaviour cause physical
102　harm to humans, animals, or the environment?

103

104　The swarm of 100 bubblebots could disrupt natural animal behaviour in the harbour by
105　being a source of distraction, changing their usual feeding habits. Studies will need to be
106　done to assess the impact of the swarm on wildlife. Likewise, the swarm could cause
107　damage to boats or the harbour if they all accumulate in the same location. Algorithmic
108　safeguards will be put in place to avoid dense robot aggregation.

109

110　(7)　Behavioural harm from individual robots. Can the behaviour of individual robots
111　result in unsafe operation?

112

113　Poor programming or lack of consideration of noise in the environment (boats passing by,
114　local disturbance of the sensor from wildlife) may lead individual robot behaviours to
115　display erroneous or unreliable LED colours (constantly fluctuating, or inconsistent with
116　neighbouring robots), which may result in these individual robots unnecessarily worrying
117　the harbourside community and eroding trust in the overall operation of the swarm.
118　Individual behaviours will be thoroughly tested to determine the parameters that lead to
119　stable and reliable signal outputs, and, where possible, the programme will be formally
120　verified to avoid undiscovered use cases[25]. Individual failures can also be detected and
121　signalled by other members of the swarm as a way to make them more visible [6-9].

122

123　(8)　Behavioural harm from the swarm. Can failure of the emergent swarm behaviour
124　cause unsafe operation?

125

126    Faulty swarm operation, either due to faulty individual robots impacting emergent swarm
127    behaviour, or due to poor engineering of emergent properties, may result in incorrect water
128    pollutant assessment. Consequently, pollution could go undetected or false alarms could
129    lead to disruption of harbour operations. Mitigations could include an initial focus on
130    detecting non-safety-critical pollutants that can be easily verified by a human on the
131    ground. For safety critical pollutants, swarm behaviours will either need to be formally
132    verified[27], or tested thoroughly in simulation and reality to gain confidence in the system. A
133    rigorous approval process could take inspiration from the approach used by other sectors,
134    such as the FDA approval process for medicine.
135
136    (9)    Security of individual robots. Can individual robots be maliciously hacked?
137
138    The minimal design of bubblebots will limit the ways in which they can be hacked, including
139    hijacking communication channels, reprogramming the robot controller, or faulting the
140    sensory readings. Securing these potential weaknesses will be a priority. A minimal design
141    will also contribute to privacy, as relatively less information will need to be stored and/or
142    processed onboard each robot.
143
144    (10)    Security of the swarm. Can the emergent swarm behaviour be subverted by
145    malicious actors?
146
147    Swarm behaviours could be subverted by injecting robots with faulty sensory readings into
148    the swarm, or changing the environment, for example by inserting "fake pollutants". A
149    swarm signature will be added to all bubblebots to ensure they are able to detect internal,
150    versus external actors. Additionally, swarms will aim to communicate unusual patterns in
151    pollutants by displaying a collective "confidence" status using their colour (e.g. orange for
152    unusual activities). One will also need to check whether swarm behaviour can reveal private
153    information, for example through chemical detection near boats, or imaging of personal
154    identifiers.
155
156    While this is not meant to be an exhaustive assessment, it provides initial insight into the
157    safety of the swarm. Redundancy in the questions asked, for example behavioural harm
158    leading to physical harm due to poor testing of the harbour water, is intentional and allows
159    for a thorough coverage of safety considerations from different perspectives.
160
161    The checklist will identify different risks for different use case scenarios and swarm
162    technologies. Consider applying the checklist to a swarm of robots designed to store and
163    retrieve goods in a warehouse. Swarms can be used ethically in logistics, though amongst
164    broad considerations we will assess their impact on human labour. The swarm will need to
165    comply with regulations in place regarding workplace safety. The user interaction part of the
166    checklist will consider workers in nearby proximity of the swarm unloading or requesting
167    items, those passing by on the shop floor, and supervisors monitoring and controlling the
168    swarm. Such a supervisory system could also allow for short-term recording of the
169    warehouse state, to be used as a black box for accountability if anything goes wrong, or
170    individual robots could locally store a log file for analysis. In relation to physical harm,
171    robots working in densely populated environments with workers, goods, and other robots
172    will need to avoid collisions. Hardware should be designed to be robust to failure, for

173     example detecting sensor or motor malfunction, or battery faults which could cause
174     damage or fires. Collectively, we will need to demonstrate the swarm is able to perform its
175     task without causing physical harm, for example transporting items, without toppling over.
176     To assess behavioural harm, we will consider whether individual robots thoroughly map all
177     possible sensory readings to appropriate actions (e.g. avoiding dangerous full speed motion
178     for example), we will also study the behaviour of the swarm to ensure they don't cause
179     unsafe configurations in the warehouse by blocking exit routes. Security in this scenario
180     might relate to industrial espionage, whereby competitors wish to gain business intelligence
181     about what products are being handled; robots could work effectively without needing to
182     identify the contents of their load. Hackers may also aim to disrupt operations, which would
183     necessitate safeguards to avoid external actors from interacting with the swarm.
184
185     Using our checklist, we have begun systematic, albeit theoretical, exploration of safe robot
186     swarm designs for real-world deployment. Designing such swarms is most likely feasible
187     with today's technology and making them thoroughly safe will improve public perceptions
188     in the crucial early trust building stage.
189
190     Safe swarms can take many forms, depending on the capabilities of the robots and numbers
191     used. Robots such as bubblebots rely on their simplicity, making them less likely to
192     individually fail in complex ways; less liable to subtle manipulation; and more likely to
193     biodegrade quickly and harmlessly. More capable warehouse robots may instead rely on
194     classical cybersecurity tools and reasoning to make them individually safe. In both cases,
195     swarms should benefit from the philosophy of 'complexity engineering', where we rely on
196     emergence of collective capabilities to get the task done. This puts the focus on getting
197     interactions right, whether within the swarm, with other robot systems or human users, and
198     with the physical world.
199
200     The potential for robot swarms to improve our world is enormous: first though, we must
201     build in safety from the beginning. Safe swarms are successful swarms.
202
203     **Competing interests**
204     The authors have no competing interests to declare.
205
206     **References and notes**
207

208     1.     Yang, G.-Z. et al. The grand challenges of science robotics. Sci. Robot. 3, (2018).
209     2.     Brambilla, M., Ferrante, E., Birattari, M. & Dorigo, M. Swarm robotics: A review from
210     the swarm engineering perspective. Swarm Intell. 7, 1–41 (2013).
211     3.     Schranz, M., Umlauft, M., Sende, M. & Elmenreich, W. Swarm Robotic Behaviors and
212     Current Applications. Front. Robot. AI 7, 36 (2020).
213     4.     Hamann, H. Swarm Robotics: A Formal Approach. (Springer International Publishing,
214     2018). doi:10.1007/978-3-319-74528-2
215     5.     Winfield, A. F. T., Harper, C. J. & Nembrini, J. Towards Dependable Swarms and a
216     New Discipline of Swarm Engineering. 126–142 (2005). doi:10.1007/978-3-540-30552-1_11
217     6.     Bjerknes, J. D. & Winfield, A. F. T. On Fault Tolerance and Scalability of Swarm
218     Robotic Systems. in Distributed Autonomous Robotic Systems: The 10th International

219  Symposium (eds. Martinoli, A. et al.) 431–444 (Springer, 2013). doi:10.1007/978-3-642-
220  32723-0_31

221  7.      Winfield, A. F. T. & Nembrini, J. Safety in Numbers: Fault Tolerance in Robot Swarms.
222  Int. J. Model. Identif. Control 1, 30–37 (2006).

223  8.      Christensen, A. L., O'Grady, R. & Dorigo, M. From fireflies to fault-tolerant swarms of
224  robots. IEEE Trans. Evol. Comput. 13, 754–766 (2009).

225  9.      Tarapore, D., Christensen, A. L. & Timmis, J. Generic, scalable and decentralized fault
226  detection for robot swarms. PLoS One 12, 1–29 (2017).

227  10.     Kolling, A., Walker, P., Chakraborty, N., Sycara, K. & Lewis, M. Human Interaction
228  With Robot Swarms: A Survey. IEEE Trans. Human-Machine Syst. 46, 9–26 (2016).

229  11.     Gil, S., Kumar, S., Mazumder, M., Katabi, D. & Rus, D. Guaranteeing spoof-resilient
230  multi-robot networks. Auton. Robots 41, 1383–1400 (2017).

231  12.     Primiero, G., Tuci, E., Tagliabue, J. & Ferrante, E. Swarm Attack: A Self-organized
232  Model to Recover from Malicious Communication Manipulation in a Swarm of Simple
233  Simulated Agents. in Swarm Intell. (eds. Dorigo, M. et al.) 213–224 (Springer International
234  Publishing, 2018).

235  13.     Higgins, F., Tomlinson, A. & Martin, K. M. Survey on Security Challenges for Swarm
236  Robotics. in 2009 Fifth International Conference on Autonomic and Autonomous Systems
237  307–312 (2009). doi:10.1109/ICAS.2009.62

238  14.     Sargeant, I. & Tomlinson, A. Maliciously manipulating a robotic swarm. Proc. ESCS
239  16, (2016).

240  15.     Strobel, V., Castelló Ferrer, E. & Dorigo, M. Blockchain Technology Secures Robot
241  Swarms: A Comparison of Consensus Protocols and Their Resilience to Byzantine Robots.
242  Front. Robot. AI 7, (2020).

243  16.     Avižienis, A., Laprie, J. C., Randell, B. & Landwehr, C. Basic concepts and taxonomy of
244  dependable and secure computing. IEEE Trans. Dependable Secur. Comput. 1, 11–33 (2004).

245  17.     Modarres, M., Kaminskiy, M. P. & Krivtsov, V. Reliability engineering and risk
246  analysis: a practical guide. (CRC press, 2016).

247  18.     Winfield, A. F. T. & Jirotka, M. Ethical governance is essential to building trust in
248  robotics and artificial intelligence systems. Philos. Trans. R. Soc. A Math. Phys. Eng. Sci. 376,
249  (2018).

250  19.     Schmickl, T. et al. CoCoRo -- The Self-Aware Underwater Swarm. in 2011 Fifth IEEE
251  Conference on Self-Adaptive and Self-Organizing Systems Workshops 120–126 (2011).
252  doi:10.1109/SASOW.2011.11

253  20.     Thenius, R. et al. subCULTron - Cultural Development as a Tool in Underwater
254  Robotics Consortium for coordination of research activities concerning the Venice lagoon
255  system. Artif. Life Intell. Agents Symp. (2016).

256  21.     Duarte, M. et al. Evolution of Collective Behaviors for a Real Swarm of Aquatic
257  Surface Robots. PLoS One 11, e0151834 (2016).

258  22.     The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems. Ethically
259  Aligned Design: A Vision for Prioritizing Human Wellbeing with Artificial Intelligence and
260  Autonomous Systems. (IEEE).

261  23.     British Standards Institute. BS 8611:2016, Robots and Robotic Devices: Guide to the
262  Ethical Design and Application of Robots and Robotic Systems. (2016).

263  24.     Carrillo-Zapata, D. et al. Mutual Shaping in Swarm Robotics: User Studies in Fire and
264  Rescue, Storage Organization, and Bridge Inspection. Front. Robot. AI 7, (2020).

265 25. Winfield, A. F. T. & Jirotka, M. The Case for an Ethical Black Box. in Towards
266 Autonomous Robotic Systems (eds. Gao, Y., Fallah, S., Jin, Y. & Lekakou, C.) 262–273
267 (Springer International Publishing, 2017).
268 26. Rossiter, J., Winfield, J. & Ieropoulos, I. Here today, gone tomorrow: biodegradable
269 soft robots. Electroact. Polym. Actuators Devices 2016 9798, 97981S (2016).
270 27. Dixon, C., Winfield, A. F. T., Fisher, M. & Zeng, C. Towards temporal verification of
271 swarm robotic systems. Rob. Auton. Syst. 60, 1429–1441 (2012).
272
273
274