



This document is downloaded from the
VTT's Research Information Portal
<https://cris.vtt.fi>

VTT Technical Research Centre of Finland

Safety Challenges of AI in Autonomous Systems: A Human Factors Perspective

Karvonen, Hannu; Heikkilä, Eetu; Wahlström, Mikael

Published: 26/11/2019

Document Version
Publisher's final version

[Link to publication](#)

Please cite the original version:

Karvonen, H., Heikkilä, E., & Wahlström, M. (2019). *Safety Challenges of AI in Autonomous Systems: A Human Factors Perspective*. Poster session presented at AI Day 2019, Espoo, Finland.



VTT
<http://www.vtt.fi>
P.O. box 1000FI-02044 VTT
Finland

By using VTT's Research Information Portal you are bound by the following Terms & Conditions.

I have read and I understand the following statement:

This document is protected by copyright and other intellectual property rights, and duplication or sale of all or part of any of this document is not permitted, except duplication for research use or educational purposes in electronic or print form. You must obtain permission for any other use. Electronic or print copies may not be offered for sale.

Safety Challenges of AI in Autonomous Systems: A Human Factors Perspective

Hannu Karvonen, Eetu Heikkilä, Mikael Wahlström

VTT Technical Research Centre of Finland Ltd

AI in autonomous systems

Autonomous systems (AS) are among the most potential application areas for AI technologies. AI, especially machine learning, is currently used in AS in limited areas such as object detection, but more advanced decision-making and adaptation in operations are also targeted.

AI safety challenges

AI-enabled AS hold big promises for increases in productivity and safety. However, the application of AI also introduces new safety and security risks. In the literature, several concrete safety challenges of AI applications and AS have been identified. For an overview, a list with numbers 1–5 in the left side of Figure 1 is provided. The safety challenges addressed here are adopted from the widely-cited work of Amodei et al. (2016). These issues are relevant for all AI systems, but become especially crucial with systems where AI is embedded as a part of a physical AS that interacts with humans.

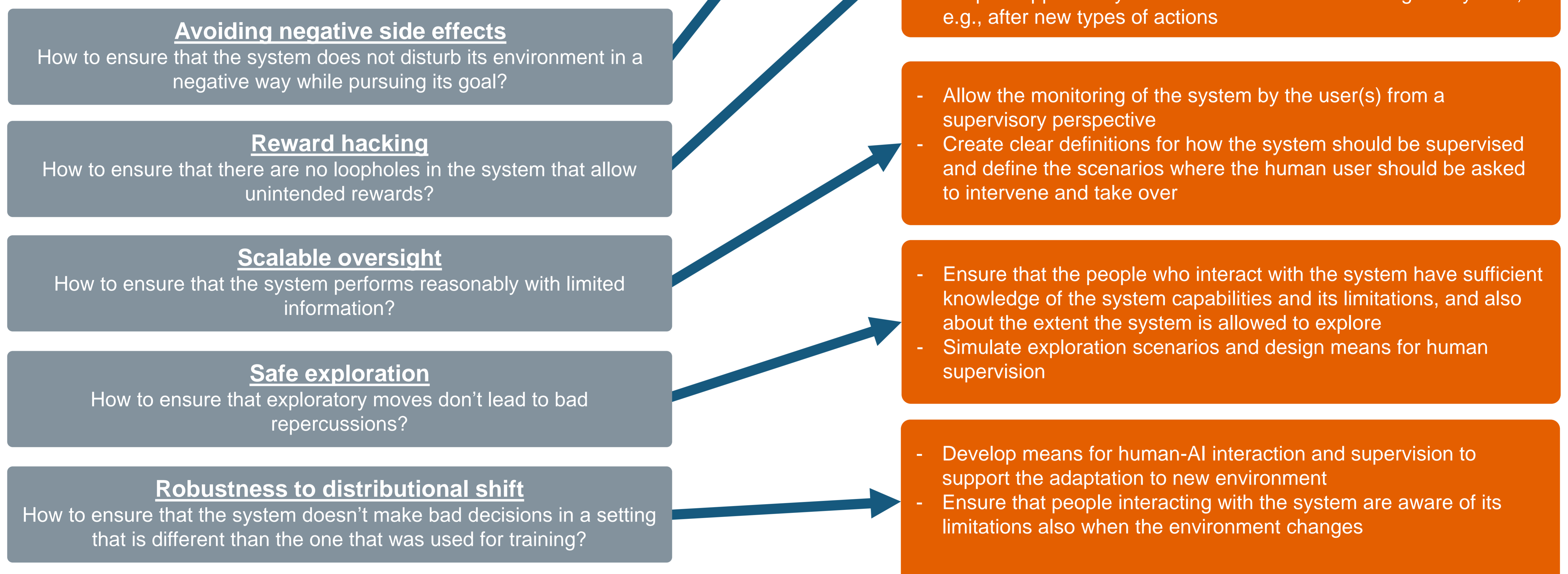


Figure 1. AI safety challenges and proposals for approaches to address them from the human factors perspective

Conclusions

- Autonomous systems are a major application area for AI technologies
- Although the level of autonomy increases, the systems will still be in interaction with human users
- Several concrete AI safety issues have been raised in literature
 - Technical assurance of AI systems is important, but also a broader systemic view is needed
 - Human factors is one part of this consideration
- The key safety issues of AI should be addressed from a human factors point of view to ensure the safety of users and other people interacting with autonomous systems that employ AI