



HAL
open science

Artificial Intelligence Awareness in Work Environments

Hannu Karvonen, Eetu Heikkilä, Mikael Wahlström

► **To cite this version:**

Hannu Karvonen, Eetu Heikkilä, Mikael Wahlström. Artificial Intelligence Awareness in Work Environments. 5th IFIP Working Conference on Human Work Interaction Design (HWID), Aug 2018, Espoo, Finland. pp.175-185, 10.1007/978-3-030-05297-3_12 . hal-02264611

HAL Id: hal-02264611

<https://inria.hal.science/hal-02264611>

Submitted on 7 Aug 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Artificial Intelligence Awareness in Work Environments

Hannu Karvonen¹, Eetu Heikkilä¹ and Mikael Wahlström¹

¹ VTT Technical Research Centre of Finland Ltd., Finland
firstname.lastname@vtt.fi

Abstract. Based on the concepts of situation and automation awareness, we present a new concept called “artificial intelligence awareness”. We also examine in detail how these three phenomena relate to each other particularly in work environments. To open this discussion, we shortly go through the ideas behind these concepts and focus especially on artificial intelligence (AI) from the machine-learning perspective and on some related human-AI interaction issues. In addition, we present an illustration of a theoretical taxonomy where our understanding of the relationships between the three key awareness concepts is visualized. We conclude by giving pointers for further research and design regarding how to support automation and AI awareness of intelligent systems users in work environments.

Keywords: Artificial Intelligence, Work, Automation, Situation Awareness, AI Transparency, Intelligent Systems, Trust in Automation, Human Factors

1 Introduction

We are living in an era of work transformation: much of routine work tasks have already been replaced with different kinds of automated solutions. It is foreseeable that in the future increasingly complex work assignments will be taken care by automated systems that use artificial intelligence (AI). On the way to that change, new kinds of interaction issues in work environments between automation, AI and humans need consideration. Already now it has been noticed that in automated work environments the workers might not be in some situations able to understand what the complex automation is doing, why it is doing that, and what it is going to do next.

Related to these issues, the concept of automation awareness (AA) has been recently discussed [1-4] as a human factors concept to be taken into account with modern automation solutions. Similarly to what good situation awareness (SA) allows in complex work environments, AA enables the workers to stay in the loop with the automated system and operate it in an appropriate way in highly automated environments. An example of a complex work environment here could be ship operation from the command bridge of the ship. In contrast, as a highly automated environment, the remote monitoring (and possible intervening in case of exception situations) of the operations carried out by an autonomous ship from an office-type control room could serve as an example. We see that AA becomes important especially in the latter case, but it nevertheless does

not exclude SA as the remote operators still need to comprehend and take into account also the demands of the work situation at hand in the object environment.

Currently, novel technologies for work environments are being rapidly introduced to address the increasing demand for systems with adaptive and autonomous capabilities. Many new advances are based on different AI technologies, such as machine learning (ML), which includes also deep learning (DL). An example of these technologies in the work context would be predictive maintenance approaches, which can detect anomalies in temperatures indicating a possible upcoming engine failure. We see that the increase of these AI technologies brings along further awareness-related interaction challenges for the users of these systems.

The focus of this paper is to consider the concepts of SA, AA and AI awareness (AIA) and discuss how these phenomena relate to each other especially in work environments. To open this discussion, in this paper we suggest a theoretical taxonomy where these concepts are visually presented to illustrate their relationships between each other. However, before going into the details of this taxonomy, the definition of the related key concepts and issues here is essential.

2 Definition of Key Concepts and Issues and Their Relevance to the Topic

Key concepts related to this paper are situation and automation awareness, artificial intelligence, machine learning and artificial intelligence awareness. Next, each of these concepts and their relevance to the topic of this paper will be discussed in detail.

2.1 Situation and Automation Awareness

The most cited definition of situation awareness is the one by Endsley [5], who stated that SA is ‘*the perception of elements in the environment within a volume of time and space, the comprehension of their meaning, and the projection of their status in the near future*’. Therefore, in detail, this definition presumes three separate stages referring to different information processing events: 1) *the perception* of relevant data and elements in the environment, 2) the interpretation of these perceptions and *the comprehension* of their meaning for the situation at hand, and 3) the ability to *project* how these elements will change in the future [5]. In this paper, we utilize these three stages to consider also other types of awareness concepts (i.e., automation and AI awareness), which we see are relevant in highly automated and artificial intelligence-supported work contexts.

Automation awareness has been mentioned in previous human factors literature (e.g., [6]) already, but it is not yet an established concept. In line with the definition of SA stages by Endsley [5], the development and maintenance of AA has been defined as a continuous process that comprises of perceiving the current status of the automation, comprehending this status and its meaning to the system behaviour, as well as projecting its future status and meaning [1]. From a work contexts’ operational per-

spective, AA can be seen to be the worker's conception of the utilized automation system's state in such a manner that enables the operator to observe, control, and anticipate the events initiated by the automation [7]. Therefore, the difference of AA to SA is that in AA, the focus of awareness is specifically in the automation system. Therefore, we see that AA can be counted to be as a part of SA, which is a much broader concept concerning all the relevant matters in a situation.

2.2 Artificial Intelligence and Machine Learning

Artificial intelligence is another concept without a single agreed definition. Some of its definitions focus on the capability of an AI-based technical system to understand its environment and make rational decisions accordingly (e.g., [8]), whereas others emphasize the technologies used to develop such capabilities. In this paper, we focus specifically on the group of AI technologies with machine-learning characteristics.

Much of the current excitement around AI technologies is based on the advances in ML methods, which are increasingly founded on the utilization of artificial neural networks (ANNs). ML technologies have produced impressive results, especially when trained with large data sets – the approach now known as deep learning [9]. DL includes a network of mathematical connections (or neural networks) that are initially random. However, through trial and error, for example, by using manually pre-labelled data, this network modifies its connections so that it is capable of generating robust predictions. Simply put, the network's "routes" towards successes are strengthened, while the error-producing connections are weakened in the learning process.

In industrial work environments, ML technologies can be used, for example, in various image processing and object recognition tasks. These tasks are needed especially when building advanced sensor systems for increasingly autonomous machines. An example of this kind of machine could be an unmanned and autonomous forest harvester that has been trained with machine learning (e.g., by going through a large number of pictures where certain types of trees have been pre-labeled) so that the harvester's environmental sensors are able to detect the trees that are supposed to be cut from the forest.

Additionally, ML technologies are planned to be used in work environments in different decision-making tasks, such as path planning of autonomous vehicles. This kind of decision-making is typically based on complex algorithms that process data, for example, from the system's environmental sensors and their sensor fusion. For instance, an autonomous car should be able to detect objects (e.g., the road, people, traffic signs, and other cars) from the environment to be able to adjust its speed, plan its path and navigate accordingly. In addition, for example, vehicle-to-vehicle communication may be utilized, which can also potentially increase system complexity from the perspective of the users and bystanders.

2.3 Underlying Machine Learning Paradigms

While a majority of ML applications utilize ANN algorithms in their implementation, there are also underlying fundamental ML paradigms that need to be considered when

moving towards increasing machine system autonomy and in supporting AIA. These paradigms can roughly be divided as follows [10]:

1. Supervised learning is used in situations where training data includes both the input data and the corresponding desired results. Typical use cases are classification and regression problems. In classification, the input data is placed in pre-defined classes, such as in the harvester image classification example discussed above where a pre-labeled data set of trees may be used for training. Regression problems, on the other hand, involve the prediction of continuous variables. For example, a prediction of a machine system performance can be performed based on sensor data if the historical data of relevant parameters is available. Due to their versatility, supervised learning methods are currently the most widely used ML method.
2. Unsupervised learning is used when only unstructured data is available. In other words, this refers to a situation where no labeled training data is available and the aim is to describe the dataset in a useful way. This description can be achieved, for example, by clustering the data or by detecting anomalies from the data [11].
3. Reinforcement learning is based on the interaction between the AI system and its operating environment. A reinforcement learning system observes responses of the environment when performing various actions and the system is “rewarded” for desired responses. While supervised and unsupervised learning paradigms are powerful for various perception and classification tasks, reinforcement learning provides a potential solution for implementing dynamic decision-making capabilities [12].

Implications of ML paradigms for AIA

Each of the above general-level ML paradigms contain a number of different methods and algorithms that can be used for the actual implementation. From the system user’s point of view, it is not typically evident which paradigm (or, even on a more detailed level, which algorithm) was used to implement the system. However, the implications of these paradigms for ML system design and AIA are important: the user should be provided with sufficient information about the limitations there might be associated with the ML system, because of the paradigm it is based on.

For example, a supervised learning system should be capable of providing the user a sufficient overview of the data it has been trained with, and to demonstrate the limitations arising from the used data. For reinforcement learning systems, on the other hand, it could be beneficial for the user to be aware, for example, of the logic how the ML system is rewarded for its various actions.

2.4 Unsolved Issues in AI-Worker Interaction

While the use of AI technologies can be seen as beneficial from a purely technological point of view, there are still major unsolved issues regarding their use in work environments involving the systems' interaction with workers. Some of these key issues will be next discussed in separate subsections.

System Transparency

System transparency refers to how transparent the functioning of the AI is (see, e.g., [13] and [14]). In this context, AI awareness means that the human is to some extent able to understand the process and estimate the results of ML and can therefore decipher the decision-making rationale of the AI system. This comprehension can be supported, for example, with understandable explanations of the process, reasons behind certain decisions and results, and simplifications or illustrative visualizations of the used algorithms to tell the human about the basis of how they are working and why they ended up in a certain solution. In this way, the worker can anticipate the functioning of the intelligent system in the work context the system is utilized. Consequently, the worker should also be able to better understand the grounds of the possible decisions made by the AI system and choose whether they are made on a suitable basis or not considering the work situation at hand.

Currently, many ML systems can be described to some degree, as "black boxes", which process input data and produce predictive outputs without providing a clear logic or justification for the results. The deep neural networks include a myriad of connections that have been produced through trial and error training and the particularities of these networks are therefore not easily understandable for the user. To support AI transparency, the black box of ML should be replaced by a more transparent one, which tells simply and clearly the intentions, capabilities, and limitations of the systems in an easily understandable format for the human. We know this is a difficult task, but at least the probabilities of the predictions in tasks, such as image recognition, could be presented for the user to enable the estimation of their reliability.

In addition to the probabilities involved in the output parameters, a worker using and monitoring an AI-based system should have a generic understanding and view of the input parameters. For example, the worker should understand what variables produce the predictions of the AI system. By understanding this prediction basis, the worker could have a better understanding on the limits of AI system; for example, the system might not be comprehensive enough and might not consider all the possible elements involved in the generic situational awareness. For instance, a fire-detection system AI's object detection functioning might or might not involve data from a smoke-detection system. With smoke detection-based input variables, an AI that would try to detect, for example, wild fires, could provide results different from object detection based on thermal camera technologies only.

Furthermore, a skilled worker's understanding could cover the data used for AI-based learning and its generic mechanisms. We now know that AI can involve certain biases. For instance, it has been claimed that AI used as a risk management tool in profiling people involves indicators related to race and age, and identifies as risk those

who have never committed a serious crime [15]. A worker should therefore understand that the predictions (and the functioning that results out of these predictions) are, indeed, predictions and not a certain truth. Hence, work system user interfaces should also support this kind of understanding with their design.

Computer-Human Communication

Computer-human communication refers to both the ways in which the AI system communicates its functioning, intentions, capabilities, and limitations to the user and also to the possibilities of the AI to understand human communication. Many of the ways to communicate the state of the AI in modern intelligent systems are far too complicated to be understood by non-expert users (e.g., log-type textual representations that include all the previous command actions of the process). To compensate this problem, for example, simplifications and clear visualizations are needed. We see that, for instance, explicit information visualization techniques, multimodal output user interfaces and interaction technologies that are natural for humans are in a key role here. For example, modern multimodal output user interfaces, such as subtle ambient sound feedback or augmented reality may be useful in some situations.

Vice versa, the communication of humans understood by the computer is an important human-computer interaction question. Natural ways of humans to communicate include, for example, voice or gesture input. However, regarding input communication methods from human to the computer it should not be forgotten that often the traditional solutions, like mouse and a keyboard, are actually rather efficient ways of commanding computers. Often, in the longer term fancy input systems may become tiring, tedious and frustrating especially in routine tasks.

Regarding more implicit human to computer communication issues, for example, Medina-Mora et al. [16] present The Coordinator tool, which allows creating and managing reports of conversations based on the universal vocabulary of speech acts. The goal of this tool is to enable a structure of interactions that is effective for coordination within work organizations [17]. Related to the possibilities of AI, the Coordinator type of approach might be beneficial in trying to understand how computers could be able to know from implicit information what humans are doing (e.g., praising vs. insulting) or aiming at when we, for example, utter something on a certain (e.g., cheerful vs. aggressive) tone of voice.

Appropriate Trust

Appropriate trust in this context refers to how the worker can trust the AI, based for example, on gained knowledge that the AI system has been developed using a sufficient amount of relevant data, or that the system has learned the required skills without becoming biased. According to Lee & See [18], trust can be defined as the “*attitude that an agent will help achieve an individual’s goals in a situation characterized by uncertainty and vulnerability*”. In this paper’s context, we see that the ‘agent’ is AI, which is used by human beings for a variety of different work-related purposes.

In case of ML, to calibrate the trust in AI to an appropriate level (see [18]) the user has to have a clear idea of the capabilities of the used algorithms and also what kind of

data has been used in the machine learning process. This would allow mitigating the possible problematic effects of overtrust or distrust in the results of the AI.

In case of overtrust, for example, accidents may happen in safety-critical environments as the worker trust the system in situations it should not be trusted. In a distrust situation, possible inefficiencies from the functioning of the human-AI team operation perspective may occur as the human does not trust the system and tries to do everything manually by herself. We also see that here, intuitive visualizations with possible simplifications of the ideas of how the AI algorithms work may help in calibrating the trust to an appropriate level. In this way, the AI-based system in question can be made more acceptable and successful among the workers using the system.

2.5 Artificial Intelligence Awareness and Affecting Matters

Our proposed concept of artificial intelligence awareness is related to all of the above-mentioned three types of theoretical human-AI interaction issues. AIA in work environments can be defined based on the previous definitions of SA and AA as *the worker's perception of the current decision made by the AI, her comprehension of this decision and her estimate of the decision(s) by AI in the future*.

In contrast to theoretical issues, on the practical level AIA is very much based on the ways for the worker to gain understanding about the AI's functioning. This understanding is affected by concrete matters, such as 1) the AI system's user interfaces, 2) the provided AI-specific training of the work organization, 3) the worker's general knowledge of the principles of computer systems and AI, 4), momentary high level of cognitive workload of the worker, and 5) the worker's subjectively experienced level of complexity of the AI system. Next, each of these factors will be discussed shortly in detail.

Firstly, if the system's user interfaces are badly designed and the usability of the system meant to provide a view to the functioning of the AI is poor, the workers will be able to follow the AI's functioning less. In the worst case, the system does not provide a possibility to follow the functioning of the AI at all. This may ultimately result in human overtrust in the AI system, as the workers have no other choice than to trust the system even in situations where it was not originally designed to cope with. Especially in exceptional situations such as AI-related faults, it is essential to provide information in the user interface about where the fault originates and what are the possible options for next steps.

Secondly, the training provided by the work organization is in a crucial role. Both the increase of theoretical knowledge and practical skills development related to the AI system is essential in training. Theoretical training may focus, for example, on the functioning principles of the AI and the limitations of the system. The training should therefore also consider situations in which the AI may make wrong decisions, the reasons behind these decisions and what to do when this type of a situation happens. Practical skills can be trained, for example, in simulator environments where different scenarios are practiced in hands-on situations. In addition, the AIA of the worker in different simulated situations may be evaluated with various human factors oriented methods, such as interviews and surveys.

Thirdly, the worker's general knowledge of the principles of computer systems and AI is relevant. Even though the worker does not necessarily have to have a background education in, say, computer science, in the future, the understanding of the functioning of algorithms and intelligent systems becomes increasingly important in many work domains. Basic understanding of different logic operations and Boolean algebra may help in deciphering the basic functioning of computer logic. In some cases, similar logic symbols are used even in the user interfaces, operating procedures and manuals of complex intelligent systems, which makes understanding them even more relevant.

Fourthly, momentary high level of cognitive workload of the worker may affect the level of AIA in specific situations. The worker is less likely to follow the functioning of the AI if other tasks cause too much cognitive workload. Therefore, smart and adaptive solutions are required to take into account the user's current task so that AI-related information is not given in situations where the user is not able to focus on the given information because of being busy with other relevant tasks.

Finally, a high level of subjectively experienced complexity of an AI system makes the development and maintenance of AIA harder. The objective structure and couplings of the AI system work as a basis for the subjective experience of the complexity of the system. However, also, for example, the views of colleagues affect subjectively experienced complexity of the system as well. In a teamwork setting, the experience of subjective complexity can also be lowered by joint sense making among the team regarding the AI system and its functioning.

With the influence of the above-mentioned factors, along with many others as well, the worker forms a mental model about the AI system. This model works as the basis for creating AI awareness in different work tasks conducted with the system. However, the mental model does not necessarily have to include details about the AI system's functioning, because a general-level understanding of the system's functioning is often enough to form an understanding about is the system is working correctly and the end-result credible. Ultimately, a clear overview mental model of the system helps in developing and maintaining an artificial intelligence awareness that forms the basis of safe and efficient work activity.

3 Taxonomy of Situation, Automation and Artificial Intelligence Awareness

We suggest that the awareness-related phenomena presented in the previous chapter enclose each other in the following way: AA encloses AIA, while SA encloses both AA and AIA. To illustrate this theoretical taxonomy, we present AIA inside of AA and AA inside of SA in Figure 1. The stages of awareness of all the three concepts in the figure are based on the three-level SA model by Endsley [5]. However, instead of the linear cognitive information processing approach behind Endsley's original SA model, the circular form of this illustration refers to the formation of SA, AA and AIA as a continuous process, along the lines of Ulric Neisser's [19] idea of action-perception

cycles in human cognitive activity. This view of SA is also similar to the situated cognition perspective, which emphasizes how the current awareness of a situation effects the process of acquiring and interpreting new awareness in an ongoing cycle [20].

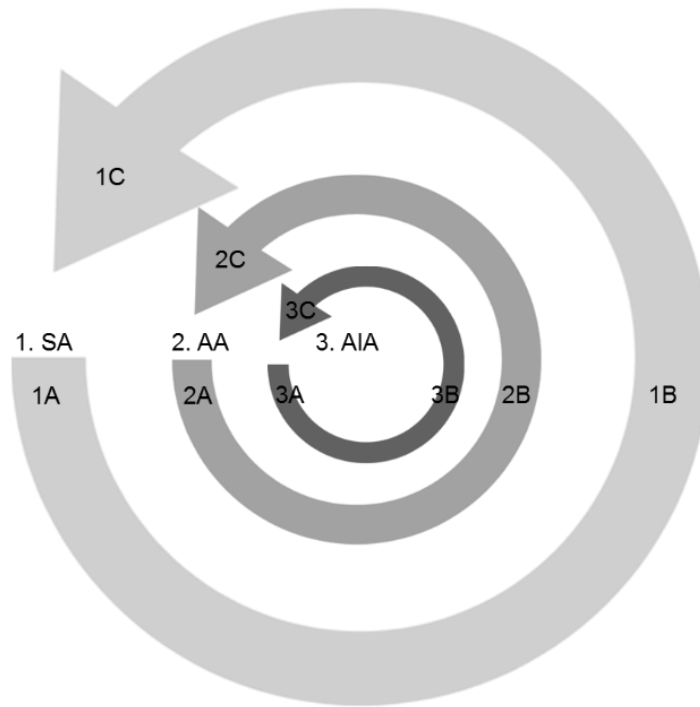


Figure 1. The Awareness Circles of SA, AA and AIA. The Figure is meant to reflect the relationship between 1. Situation Awareness, 2. Automation Awareness, and 3. AI Awareness. Explanations for the used codes are, 1A: *Perception of the current situation*, 1B: *Comprehension of the current situation*, 1C: *Estimate of the future situation*; 2A: *Perception of the current status of the automation*, 2B: *Comprehension of the current status of the automation*, 2C: *Estimate of the future status of the automation*; 3A: *Perception of the current AI-based decision(s) by the system*, 3B: *Comprehension of the current AI-based decision(s) of the system and their basis*, 3C: *Estimate of the system's AI-based decision(s) in the future and their basis*.

The different shades of grey used in the circles of Figure 1 do not have any other function than just to more clearly separate the circles from each other. However, the thickness of the circles in Figure 1 indicates the broadness of the awareness in question: SA can be seen as the broadest awareness, as it includes also AA and AIA inside of it. Naturally, there is a lot of interaction between these different levels of awareness and in practice, the formation of a certain level of awareness is as clearly separated in reality.

The presented model also bears similarities to the semiotic theory of Charles Sanders Peirce [21]. We see that in Figure 1's awareness circles, the perception level refers to the sign, the comprehension level refers to the object, and the future estimate refers to the interpretant found in Peirce's triadic relation model [21]. The interpretant can also work as a new sign for something else and the semiosis may continue from there.

Theory-wise, this model opens up opportunities for conceptual discussions regarding the nature of human cognition and awareness. We see that instead of the linear information-processing model that is based on cognitive psychology, also cyclical models in which human awareness is considered to be formed in the interaction with one's environment are worth investigating. As suggested above, for example, Neisser's ecological psychology [19], Suchman's situated cognition thinking [20], and Peirce's pragmatism [21] may provide very relevant views to the concepts of situation, automation and artificial intelligence awareness.

4 Discussion and Conclusions

It is evident that the increasing the level of automation in work environments has many effects on work activity. For example, in highly automated environments, it is no longer enough for the workers to be sufficiently aware of the situation of the work-related process (i.e., to maintain a good SA), but they also need to follow the functioning of the automated systems related to this process in order to achieve a good automation awareness. Due to AI's important role in future automated systems, we suggest that also "AI awareness" is a key concept to be considered when studying and designing future socio-technical work environments, which rely on AI technologies.

Similar problems that have already nowadays been identified with automation in work environments (see e.g., [22]) can potentially be avoided by considering how to support the workers' AIA through, for example, design and organizational practices. This support may include, for instance, the appropriate design of the intelligent system's user interfaces and the provided training related to the AI system. In this way, better worker understanding and interaction with and also appropriate trust in AI may be possible with intelligent work systems.

In addition to appropriate trust, key AIA-related concepts that are relevant here are system transparency and computer-human communication. While system transparency refers to how transparent the functioning of the AI is, computer-human communication refers to both the way in which the AI system communicates its functioning, intentions, capabilities, and limitations to the user and to the possibilities of the AI to understand human communication.

Clearly, this paper is just an initial small step in this endeavor and presents only some preliminary thoughts as a discussion opening. Both theory and practice should be developed to consider AIA more holistically. A lot more research is needed, for example, on how automation and AI affect work activity in different situations and how worker AA and AIA could be supported to mitigate possible human out-of-the-loop problems with future systems. Undoubtedly, practical case studies both in naturalistic and experimental settings are needed to validate the concepts and design. Furthermore, guidelines

for the design of AI systems and their associated training from the perspective of supporting human AIA are needed. We hope this paper serves as a good beginning for research and development to consider the importance of AIA with future intelligent work systems.

References

1. Karvonen, H., Liinasuo, M., Lappalainen, J.: Assessment of Automation Awareness. In: Proceedings of Automaatio XXI Conference 2015, 44, Publication Series of the Finnish Society of Automation, Helsinki (2015).
2. Karvonen, H., Lappalainen, J., Liinasuo, M.: Automation Awareness User Interface Study – Preliminary Results. In: Proceedings of the Man-Technology-Organisation Sessions at the 2014 Enlarged Halden Project Group Meeting, Vol. 1, C2.4., OECD Halden Reactor Project, Norway (2014).
3. Karvonen, H., Aaltonen, I., Lappalainen, J., Liinasuo, M., Kuula, T., Wahlström, M., Aikala, M., Laitio, P., Savioja, P., Laarni, J.: Studying automation awareness in nuclear power plants. In: Hämäläinen J., Suolanen, V. (eds.), The Finnish Research Programme on Nuclear Power Plant Safety 2010–2014, Final Report, VTT Technology 213, pp. 92–102, VTT, Espoo (2015).
4. Laitio, P., Savioja, P., Lappalainen, J.: Exploring Automation Awareness in Nuclear Power Plant Control Rooms. In: Proceedings of the Man-Technology-Organisation Sessions at the 2013 Enlarged Halden Project Group Meeting, OECD Halden Reactor Project, Norway (2013).
5. Endsley, M. R. Toward a theory of situation awareness in dynamic systems. *Human factors* 37(1), 32–64 (1995).
6. Whitlow, S. D., Dorneich, M. C., Funk, H. B., Miller, C. A.: Providing appropriate situation awareness within a mixed-initiative control system. In: Proceedings of the 2002 IEEE International Conference on Systems, Man and Cybernetics, (2002).
7. Karvonen, H., Liinasuo, M., Lappalainen J.: Assessment of situation and automation awareness, VTT Technical Research Centre of Finland, VTT Research Report VTT-R-05997-14 (2014).
8. Nilsson, N. J.: *The Quest for Artificial Intelligence*, Cambridge University Press (2009).
9. Stone, P., Brooks, R., Brynjolfsson, E., Calo, R., Etzioni, O., Hager, G., ..., Leyton-Brown, K.: *Artificial Intelligence and Life in 2030. One Hundred Year Study on Artificial Intelligence: Report of the 2015-2016 Study Panel*, Stanford University, Stanford, CA (2016).
10. Jordan, M. I., & Mitchell, T. M.: Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245), 255–260 (2015).
11. Lee, J. H., Shin, J., & Realf, M. J.: Machine learning: Overview of the recent progresses and implications for the process systems engineering field. *Computers & Chemical Engineering*, 114(1), 111–121 (2018).
12. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... Hasabis, D.: Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533 (2015).
13. Theodorou, A., Wortham, R. H., Bryson, J. J.: Designing and implementing transparency for real time inspection of autonomous robots. *Connection Science* 29(3), 230–241 (2017).
14. Pynadath, D. V., Barnes, M. J., Wang, N., Chen, J. Y.: Transparency Communication for Machine Learning in Human-Automation Interaction. In: *Human and Machine Learning*, pp. 75–90. Springer, Cham (2018).

15. Amnesty International: Trapped in the matrix: Secrecy, stigma, and bias in the Met's Gangs Database (2018). Published online by Amnesty International at <https://www.amnesty.org.uk/files/reports/Trapped%20in%20the%20Matrix%20Amnesty%20report.pdf>
16. Medina-Mora, R., Winograd, T., Flores, R., Flores, F.: The action workflow approach to workflow management technology. *The Information Society* 9(4), 391–404 (1993).
17. Winograd, T.: Categories, disciplines, and social coordination. *Computer Supported Cooperative Work (CSCW)* 2(3), 191–197 (1993).
18. Lee, J. D., See, K. A.: Trust in automation: Designing for appropriate reliance. *Human factors* 46(1), 50–80 (2004).
19. Neisser, U.: *Cognition and reality: Principles and implications of cognitive psychology*. WH Freeman/Times Books/Henry Holt & Co (1976).
20. Sandom, C.: *Situation Awareness*. In: Noyes, J., Bransby, M. (eds.), *People in control: human factors in control room design*, IET, Herts, UK (2001).
21. Peirce C. S.: *The Essential Peirce, Selected Philosophical Writings*. Indiana University Press (1998).
22. Bainbridge, L.: Ironies of automation. *Automatica* 19(6), 775–779 (1983).