

From Explanation Correctness to Explanation Goodness: Only Provably Correct Explanations can Save the World

Maike Schwammbberger¹[0000-0002-3344-6282]

Karlsruhe Institute of Technology, Karlsruhe 76131, Germany
schwammbberger@kit.edu mase.kastel.kit.edu

Abstract. Explainability Engineering gets evermore important in the era of self-learning and automated systems. We motivate the necessity for interdisciplinary research to engineer verifiably correct and good explanations: Systems engineering research must ensure that correct and machine-understandable explanations can be derived from system specifications and social sciences research must ensure that a context-dependent and stakeholder-tailored explanation can be provided in a fitting manner. We describe our first steps in the direction of a holistic and interdisciplinary explainability engineering process for tackling these challenges.

Keywords: self-explainable software systems · explanation correctness · explanation goodness · eXplainable Artificial Intelligence · trustworthy systems.

1 Introduction

Automated and self-learning software systems are increasingly used in a variety of domains and in people’s everyday life: from driving assistance systems and products manufactured in smart factories to smart home technologies and applications on our smartphones. Often, the level of automation and system functionality might be known to stakeholders interacting with the automated system to some degree. For instance, an owner of a semi-automated vehicle will have a certain degree of knowledge about the automated distance keeping functionality of their car. However, there still might be features that they do not understand, e.g. the car’s behaviour in some special outlier situations (e.g. how the car operates in exceptionally bad weather).

We perceive two key reasons why the design and functionality of such automated and self-learning software systems must be made explainable. (a) Relevant stakeholder groups that interact with such systems need to be sufficiently informed about the systems’ functionality. E.g., to be enabled to safely interact with the system or for trusting the automated system. Secondly, (b) for an explainable system, analysis and verification of correct system behaviour can be aided. We postulate that, without a certain level of system explainability, a system should not be launched into our markets and with that be integrated into our

societies (cf. IEEE Standard for Transparency of Autonomous Systems [2,27]). We develop and investigate a system’s *self-explainability* capabilities. With self-explainability, a system can explain its decision making process without the help of an external explainer.

A challenge for engineering self-explainability is the sheer number and complexity of components that a system comprises: Even a (seemingly) simple system like a robot vacuum cleaner comprises a collection of different software components (see Fig. 1). These could, e.g., be a sensor processing unit for avoiding collisions, a communication unit for interacting with a connected smart home system, a behaviour model for specifying when the robot must return to its charging station and a sub-symbolic Artificial Intelligence (AI) component that learns a map about the cleaning area. Due to this level of complexity, having a system engineer write explanations manually cannot be desirable nor feasible and would certainly lead to human errors. This would result in unreliable, even incorrect, explanations. Further on, with an increase in sub-symbolic AI components that learn new phenomena during run-time, not all explanations can be built during design time. Due to this, we focus on automatically extracting provably *correct explanations* from system models. In our understanding, a correct explanation is one where the validity of explanation content can be proven through system analysis techniques. Such extracted explanations would be in an internal, machine-readable, format, allowing for formal reasoning about explanation correctness by the system itself. A benefit of this is model-reuse: the necessary system models are built during the software system design process, with examples for system models including communication diagrams, automata models and architecture models[3].

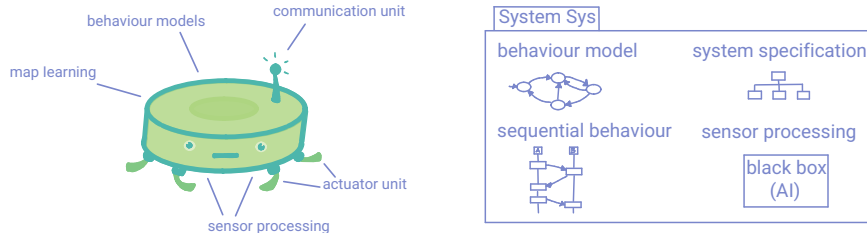


Fig. 1: A robot vacuum cleaner with diverse functionality, defined as a system Sys, comprising several software system components.

However, purely considering the formal correctness of an explanation is not enough: *explanation goodness* [12] is a factor that must be considered. Explanation goodness means that explanations must be actually helpful in increasing properties like understanding or trust into the system for specific explanation recipients. One explainability research consensus comes to the fore in explainability research discussions: only *verifiably good explanations* can be successful in increasing system understandability [10,16,28]. We argue in this paper, that,

to gain such a good explanation, both explanation goodness and explanations correctness must be analysed and validated.

For proving both the goodness and correctness of an explanation, a variety of research areas must meet and join forces (see Fig. 2); On the one side, a technical explanation must be derived in the right moment [4,21] and from adequate sources (e.g. system models, environment models, mental models) [6,23]. Such a technical explanation will be in some internal, machine-readable, format (e.g. a logical expression), allowing the autonomous system itself to reason about explanations. This step might be provided through work from the engineering sciences, e.g. computing science. On the other side, this technical explanation must be presented adequately to fit the needs of specific types of recipient stakeholders, in varying contexts [14]. For such an adequate presentation of an explanation, an interdisciplinary viewpoint on explainability is the key, especially taking knowledge from social sciences into account [15,14,18]. Note that the graphic presented in Fig. 2 is a very coarse segmentation of research directions necessary for engineering correct and good explanations. It suits the needs of this paper, but a finer segmentation is of interest for a more detailed topic discussion in future work.

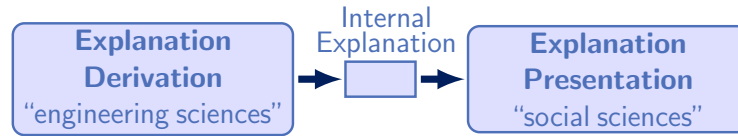


Fig. 2: To engineer correct and good explanations, technical and interdisciplinary research on explanations must be brought together.

In this paper, we summarise our roadmap for a *holistic explainability engineering process* from explanation correctness to explanation goodness in Sect. 2. In our current research, the key argument is that we need automatically derived and formalised explanations to enable a formal validation of explanation correctness (cf. Sect. 3). However, it would be presumptuous to believe that we can fully validate explanations on a purely logical level. Hence, we point out necessary interdisciplinary steps towards explanation goodness (cf. Sect. 4). We summarise our key findings towards a *holistic explainability engineering process* in Sect. 5.

2 Holistic Explainability Engineering

Our overall vision is to enable autonomous systems to self-explain their behaviour and functionality. Generally, the systems that need explaining comprise several *system components* as we discussed for the robot vacuum cleaner from Fig. 1. We postulate that explanations must be providable for the entire system,

and not only for a singular system component. For instance, an end-user might want to know why their robot vacuum cleaner did not clean a specific room. The reasons for this must be found in the entire system and not just within a singular component: It might be that due to a defective sensor, the map learning was not fully successful. Another explanation could be that the robot was not capable of moving over a cable laying in the threshold of the room. Such a need for holistic system explanations is also discussed in [19], where the authors argue that a majority of approaches within the eXplainable Artificial Intelligence (XAI) community only considers explanations for specific AI algorithms and for the community itself, meaning experts. While this can help said experts in debugging a system, it does not meet the explainability requirements for most non-expert stakeholders (e.g. end-users, lawyers, regulation bodies, ...) [8,5]. The authors of [9] go even further than [19] by showing that most XAI approaches focus on “low-level”, narrow, explanations, while instead we need to go towards “high-level”, strong, explanations. They call their approach “Broad eXplainable Artificial Intelligence (Broad-XAI)”. Their approach entails to map explanations derived by XAI approaches to human models of explanation, thereby also arguing for the necessity of connecting the technical explanation derivation with explanation presentation (cf. Fig. 2).

Our own research towards holistic explainability engineering does not only focus on sub-symbolic AI models, but instead takes the entire system into account, with both symbolic and sub-symbolic components. Our work can be boiled down to two key hypotheses, that can be associated with each one of the research areas we motivate in Fig. 2:

Explanation Derivation To verify correctness of an explanation, we must consider the logical core of explanations and their formal source.

Explanation Presentation To connect our formal notion of explanation correctness to actual goodness of explanations, we must enrich our research with interdisciplinary expertise, e.g. from social sciences.

We give more details on both hypotheses in the following sections.

3 Explanation Derivation

A precondition for engineering a self-explainable system is to consider how the system reasons about its behaviour; namely in some formalised machine code, following its system description. There is a striking benefit of first considering machine code, instead of natural language, for explanations: the inherent ambiguities of natural language are not existing in the machine code counterpart [22].

Following this reasoning, the logical core E of an explanation that we consider is a machine-readable and machine-producible intermediate format of an explanation. Its formal source are artefacts from system development processes: system models Sys (e.g. architecture diagrams, communication protocols, ...) and environment models Env (describing the operating context). We assume that Env also includes mental models for human behaviour (e.g. derived from cognitive architectures [1]). We define explanation correctness as follows.

Definition 1 (Explanation Correctness). *An internal explanation E is correct, if it can be deduced for an explanandum X and from provably correct system models Sys and environment descriptions Env .*

Here, the *explanandum* X describes the system phenomenon that needs explaining. This could, e.g., be the driving assistance function of an automated car. It is not the goal of this paper, to provide a full formalisation of explainability concepts, but we refer to [13] for this. Through definition 1, we can conclude explanation correctness from correctness of specification models. Considering this formal core of an explanation comes with a striking benefit: We can use known techniques from formal verification to analyse explanations and prove their correctness.

Deriving the logical core of an explanation from system descriptions is not a trivial step. Instead, we must develop means on how to extract explanations from system models. We describe a reference framework on how to automatically derive *explanation models* from system models in previous work [23] and give a simplified visualisation for it in Fig. 3. An explanation model could, for instance, be a causal behaviour tree [11]. The starting point for the framework is a system specification Sys . By analysing potential causes for system actions within Sys , we extract an initial explanation model. We extract the explanation model from the different system components contained within Sys . This leads to several explanation models, which are combined to one explanation model in this phase to avoid redundancy and computational overhead. We give details on this process of merging explanation models in our previous work [17]. Through information provided by environment Env , we refine the initial explanation model in the tailoring phase. The output of the tailoring phase is an explanation model which is tailored toward specific recipient stakeholders. The assumption for this is that an end-user stakeholder will need different information from an explanation than an expert stakeholder. In the era of self-learning adaptive systems [26], we must also take updates of the explanation model into account. This is needed, if, for instance, the vacuum cleaner encounters a new type of obstacle to avoid or is confronted with a new type of surface to clean.

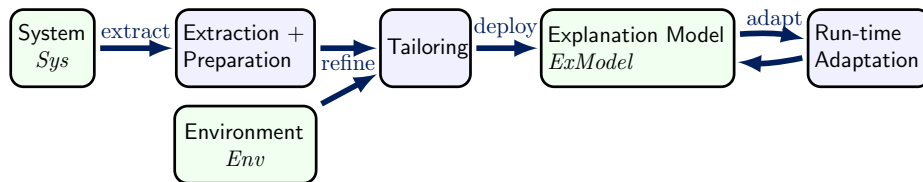


Fig. 3: Simplified visualisation of the reference framework from [23] on the derivation of an explanation model from a formalised system model. The green elements are the artefacts used and created within the blue-coloured phases.

We depict an exemplary explanation model as of [23] in Fig. 4. A technical explanation extracted from this model can be an *explanation path* within the tree structure. An example for such an explanation path could be $\textit{because}(C, \textit{and}(x_2, y_2))$. In case of the vacuum cleaner example, instances of actions could be “emergency brake” or “drive to station” with potential reasons “out of energy”, “obstacle in the way” and “sensor obscured”. The technical explanation $\textit{because}(C, \textit{and}(x_2, y_2))$ could thus translate to “An emergency brake (C) was done because of an obstacle in the way (x_2) and an obscured sensor (y_2)”.

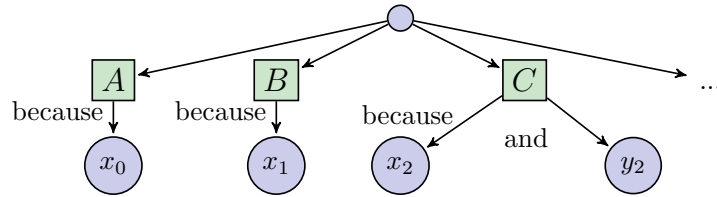


Fig. 4: An exemplary explanation model, comprising reasons for actions A , B and C .

The idea of explanation models as explanation source is embedded into the MAB-EX framework for self-explainable systems that we developed in previous work [6]. MAB-EX suggests to identify the need for an explanation through system monitoring and analysis. Whenever the need for an explanation is detected, an explanation path must be extracted from our explanation model. We refine MAB-EX in [25] by integrating *levels of explainability* that were suggested in [5]. With this, explanations can be provided by a singular or multiple systems. This enables that explanations can be provided even if the system itself has an incomplete explanation model, e.g. through missing environmental information. The missing information is then provided by another system. Equally, we must take the explanation context (i.e. where and what to explain) [5] and timing (i.e. when to explain) [4,21] into account to provide holistic explanations.

Hitherto, the framework from [23] is only conceptual and has been instantiated exemplarily for an autonomous driving controller in [23] and for digital twins in [17]. Our next steps are to examine different types of system and environment models from which explanation models must be derived, and how to formalise an automated explanation model extraction process for a variety of system models.

4 Explanation Presentation

For our goal of a holistic explainability engineering process, we must now tackle the challenge of “informalising” our logical explanation path [22]. From it, we must achieve an explanation presentation that follows explanation requirements

for diverse explanation recipient groups G [14,8]. This means that now, from a *correct* explanation, we must derive a *good* explanation.

Definition 2 (Explanation Goodness). *An explanation E can be labelled as good, if it stems from a correct explanation E_c and measurably helps a recipient group G in understanding the explanandum X .*

So, explanation goodness depends on different types of recipient stakeholders. To measure explanation goodness, *explanation validation* through user-studies must be considered. As it would be out of the scope of this visionary paper, we do not go into details on this, but instead refer to two review papers on explanation validation through user studies [14,20]. We tailor our explanation model towards different recipient groups in the respective phase of our framework (cf. Fig. 3). Situation- and stakeholder-dependent explanations have various benefits to different stakeholders, for which we provide examples here:

- a system engineer can improve and debug the system during design time and
- an end-user is enabled to use the system safely and to justifiably trust in its automated decisions,
- political and societal bodies can decide whether to allow a system to be launched into the markets, and
- lawyers can decide who is to be blamed in tort claim cases that involve automated and self-learning systems.

To allow for explanations to be used on such central societal levels, a holistic explanation validation process is of the utmost importance. Such an explanation validation must contain the formal verification part that we discussed in the previous section, but also needs a user-dependent validation to assess explanation goodness [12]. One of the shortcomings that [14] discover in an extensive literature review is the lack of empirical evidence for explanation goodness.

In our research, we aim to approach this problem from two sides: We need to investigate what information explanations for different recipient stakeholders must contain, to allow for a correct tailoring of our explanation model towards recipients. This is necessary as an explanation model containing too much information leads to larger computation times for deriving explanation paths from the model, and an explanation model missing information would lead to phenomena that cannot be explained at all. We started this endeavour by investigating explanations needed for lawyers in [7]. In sum only verifiable explanations can help stakeholders like lawyers, courts and regulation bodies to assess liabilities and to admit systems into the markets. Further research includes a translation of our internal, technical, explanations into adequate presentation formats to investigate and formalise explanation goodness on top of explanation correctness.

5 Conclusion

We motivate the need for a connection of formal reasoning about explanations with research from social sciences to validate both explanation correctness and

explanation goodness. We argue that trustworthiness of explanations can only be reached through a holistic explainability engineering process. Such a holistic explainability engineering process entails that the entire system, with all its components and together with environmental influence factors, must be explained. Moreover, different recipient stakeholders and their specific explanation requirements must be taken into account. This makes reasoning about explanation correctness and goodness a complex endeavour, where a challenging amount of requirements and factors must be taken into account.

For tackling this challenge, we recently suggest requirements for explainability levels [24], dividing explanations into local and global explanations. We also discuss a notion of *explanation quality* in that paper. Equally, through tailoring our explanation model from [23] towards different stakeholders, we do not intend to formalise one model that must comprise explanations for all possible stakeholders. Instead, each one explanation model exists for each group of stakeholders. The motivation for this is that different stakeholders might need very different types and degrees of information within an explanation model. Even with such steps for decreasing complexity in explainability engineering, one should not assume to be able to formalise and formally prove each aspect of explanation correctness and goodness right away, for each stakeholder group and varying application domains. Instead a focus on specific application domains and a fixed amount of stakeholder groups certainly makes sense as a starting point.

We follow the argumentation of [19,9,14] and emphasise that a focus of the explainability community must continue to shift from expert explanations for narrow and highly isolated system functions towards holistic explainability research for complex systems of systems. Several interdisciplinary challenges and starting points for doing so have been summarised in [15]. One challenge particularly to be overcome is the derivation of explanations for probabilistic AI systems into approaches of explainability engineering, and we sketch some steps for that in [24].

Acknowledgments. This research was supported by the Innovation Campus for Future Mobility (www.icm-bw.de) and by the Helmholtz Association within the Core Informatics project.

References

1. The Cambridge Handbook of Computational Psychology. Cambridge Handbooks in Psychology, Cambridge University Press (2008)
2. IEEE standard for transparency of autonomous systems. IEEE Std 7001-2021 pp. 1–54 (2022). <https://doi.org/10.1109/IEEESTD.2022.9726144>
3. Iso/iec/ieee 42010:2022 software, systems and enterprise architecture description 2, 1–62 (2022), <https://www.iso.org/standard/74393.html>
4. Bairy, A., Hagemann, W., Rakow, A., Schwammberger, M.: Towards formal concepts for explanation timing and justifications. In: 30th IEEE International Requirements Engineering Conference Workshops, RE 2022 - Workshops, Melbourne, Australia, August 15-19, 2022. pp.

- 98–102. IEEE (2022). <https://doi.org/10.1109/REW56159.2022.00025>, <https://doi.org/10.1109/REW56159.2022.00025>
5. Bersani, M.M., Camilli, M., Lestingi, L., Mirandola, R., Rossi, M., Scandurra, P.: A conceptual framework for explainability requirements in software-intensive systems. In: 2023 IEEE 31st International Requirements Engineering Conference Workshops (REW). pp. 309–315 (2023). <https://doi.org/10.1109/REW57809.2023.00059>
 6. Blumreiter, M., Greenyer, J., Garcia, F.J.C., Klös, V., Schwammberger, M., Sommer, C., Vogelsang, A., Wortmann, A.: Towards self-explainable cyber-physical systems. In: 22nd ACM/IEEE International Conference on Model Driven Engineering Languages and Systems Companion. pp. 543–548 (2019). <https://doi.org/10.1109/MODELS-C.2019.00084>
 7. Buiten, M.C., Dennis, L.A., Schwammberger, M.: A vision on what explanations of autonomous systems are of interest to lawyers. In: Schneider, K., Dalpiaz, F., Horkoff, J. (eds.) 31st IEEE International Requirements Engineering Conference, RE 2023 – Workshops, Hannover, Germany, September 4-5, 2023. pp. 332–336. IEEE (2023). <https://doi.org/10.1109/REW57809.2023.00062>, <https://doi.org/10.1109/REW57809.2023.00062>
 8. Chazette, L., Brunotte, W., Speith, T.: Explainable software systems: from requirements analysis to system evaluation. *Requir. Eng.* **27**(4), 457–487 (2022). <https://doi.org/10.1007/s00766-022-00393-5>, <https://doi.org/10.1007/s00766-022-00393-5>
 9. Dazeley, R., Vamplew, P., Foale, C., Young, C., Aryal, S., Cruz, F.: Levels of explainable artificial intelligence for human-aligned conversational explanations. *Artificial Intelligence* **299**, 103525 (2021). <https://doi.org/https://doi.org/10.1016/j.artint.2021.103525>, <https://www.sciencedirect.com/science/article/pii/S000437022100076X>
 10. de Bruijn, H., Warnier, M., Janssen, M.: The perils and pitfalls of explainable ai: Strategies for explaining algorithmic decision-making. *Government Information Quarterly* **39**(2), 101666 (2022). <https://doi.org/https://doi.org/10.1016/j.giq.2021.101666>, <https://www.sciencedirect.com/science/article/pii/S0740624X21001027>
 11. Garcia, F.J.C., Robb, D.A., Liu, X., Laskov, A., Patrón, P., Hastie, H.F.: Explain yourself: A natural language interface for scrutable autonomous robots. *CoRR abs/1803.02088* (2018), <http://arxiv.org/abs/1803.02088>
 12. Hoffman, R.R., Mueller, S.T., Klein, G., Litman, J.: Measures for explainable ai: Explanation goodness, user satisfaction, mental models, curiosity, trust, and human-ai performance. *Frontiers in Computing Science* **5** (2023). <https://doi.org/https://doi.org/10.3389/fcomp.2023.1096257>, <https://doi.org/10.3389/fcomp.2023.1096257>
 13. Köhl, M.A., Baum, K., Langer, M., Oster, D., Speith, T., Bohlender, D.: Explainability as a non-functional requirement. In: RE. pp. 363–368. IEEE (2019)
 14. Langer, M., Oster, D., Speith, T., Hermanns, H., Kästner, L., Schmidt, E., Sesing, A., Baum, K.: What do we want from explainable artificial intelligence (xai)? – a stakeholder perspective on xai and a conceptual model guiding interdisciplinary xai research. *Artificial Intelligence* **296**, 103473 (2021). <https://doi.org/https://doi.org/10.1016/j.artint.2021.103473>, <https://www.sciencedirect.com/science/article/pii/S0004370221000242>
 15. Longo, L., Brcic, M., Cabitza, F., Choi, J., Confalonieri, R., Ser, J.D., Guidotti, R., Hayashi, Y., Herrera, F., Holzinger, A., Jiang, R., Khosravi,

- H., Lecue, F., Malgieri, G., Páez, A., Samek, W., Schneider, J., Speith, T., Stumpf, S.: Explainable artificial intelligence (xai) 2.0: A manifesto of open challenges and interdisciplinary research directions. *Information Fusion* **106**, 102301 (2024). <https://doi.org/https://doi.org/10.1016/j.inffus.2024.102301>, <https://www.sciencedirect.com/science/article/pii/S1566253524000794>
16. Markus, A.F., Kors, J.A., Rijnbeek, P.R.: The role of explainability in creating trustworthy artificial intelligence for health care: A comprehensive survey of the terminology, design choices, and evaluation strategies. *Journal of Biomedical Informatics* **113**, 103655 (2021). <https://doi.org/https://doi.org/10.1016/j.jbi.2020.103655>, <https://www.sciencedirect.com/science/article/pii/S1532046420302835>
 17. Michael, J., Schwammberger, M., Wortmann, A.: Explaining cyberphysical system behavior with digital twins. *IEEE Softw.* **41**(1), 55–63 (2024). <https://doi.org/10.1109/MS.2023.3319580>, <https://doi.org/10.1109/MS.2023.3319580>
 18. Miller, T.: Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence* **267**, 1–38 (2019). <https://doi.org/https://doi.org/10.1016/j.artint.2018.07.007>, <https://www.sciencedirect.com/science/article/pii/S0004370218305988>
 19. Miller, T., Howe, P., Sonenberg, L.: Explainable ai: Beware of inmates running the asylum or: How i learnt to stop worrying and love the social and behavioural sciences (2017)
 20. Nauta, M., Trienes, J., Pathak, S., Nguyen, E., Peters, M., Schmitt, Y., Schlötterer, J., van Keulen, M., Seifert, C.: From anecdotal evidence to quantitative evaluation methods: A systematic review on evaluating explainable ai. *ACM Comput. Surv.* **55**(13s) (jul 2023). <https://doi.org/10.1145/3583558>, <https://doi.org/10.1145/3583558>
 21. Rakow, A., Hajnorouzi, M., Bairy, A.: What to tell when? - information provision as a game. In: Farrell, M., Luckcuck, M., Gleirscher, M., Schwammberger, M. (eds.) *Proceedings Fifth International Workshop on Formal Methods for Autonomous Systems, FMAS@iFM 2023, Leiden, The Netherlands, 15th and 16th of November 2023. EPTCS, vol. 395, pp. 1–9* (2023). <https://doi.org/10.4204/EPTCS.395.1>, <https://doi.org/10.4204/EPTCS.395.1>
 22. Ranta, A.: Translating between language and logic: What is easy and what is difficult. In: Bjørner, N.S., Sofronie-Stokkermans, V. (eds.) *Automated Deduction - CADE-23 - 23rd International Conference on Automated Deduction, Wrocław, Poland, July 31 - August 5, 2011. Proceedings. Lecture Notes in Computer Science, vol. 6803, pp. 5–25. Springer* (2011). https://doi.org/10.1007/978-3-642-22438-6_3, https://doi.org/10.1007/978-3-642-22438-6_3
 23. Schwammberger, M., Klös, V.: From specification models to explanation models: An extraction and refinement process for timed automata. In: Luckcuck, M., Farrell, M. (eds.) *Proceedings Fourth International Workshop on Formal Methods for Autonomous Systems (FMAS) and Fourth International Workshop on Automated and verifiable Software sYstem DEvelopment (ASYDE), FMAS/ASYDE@SEFM 2022, and Fourth International Workshop on Automated and verifiable Software sYstem DEvelopment (ASYDE)Berlin, Germany, 26th and 27th of September 2022. EPTCS, vol. 371, pp. 20–37* (2022). <https://doi.org/10.4204/EPTCS.371.2>, <https://doi.org/10.4204/EPTCS.371.2>
 24. Schwammberger, M., Mirandola, R., Wenninghoff, N.: Explainability engineering challenges: Connecting explainability levels to run-time explainability. In: *Proceed-*

- ings of 2nd World Conference on Explainable Artificial Intelligence Conference (XAI2024). Springer (2024), 46.23.01; LK 01
25. Schwammberger, M., Mirandola, R., Wenninghoff, N.: Explainability engineering challenges: from requirement definition to run-time explainability (2024), submitted to 2nd World Conference on eXplainable Artificial Intelligence (XAI)
 26. Weyns, D., Iftikhar, M.U., de la Iglesia, D.G., Ahmad, T.: A survey of formal methods in self-adaptive systems. In: Proceedings of the Fifth International C* Conference on Computer Science and Software Engineering. p. 67–79. C3S2E '12, Association for Computing Machinery, New York, NY, USA (2012). <https://doi.org/10.1145/2347583.2347592>, <https://doi.org/10.1145/2347583.2347592>
 27. Winfield, A.F.T., Booth, S., Dennis, L.A., Egawa, T., Hastie, H., Jacobs, N., Muttram, R.I., Olszewska, J.I., Rajabiyazdi, F., Theodorou, A., Underwood, M.A., Wortham, R.H., Watson, E.: IEEE p7001: A proposed standard on transparency. *Frontiers in Robotics and AI* **8**, 225 (2021). <https://doi.org/10.3389/frobt.2021.665729>, <https://www.frontiersin.org/article/10.3389/frobt.2021.665729>
 28. Wing, J.M.: Trustworthy ai. *Commun. ACM* **64**(10), 64–71 (sep 2021). <https://doi.org/10.1145/3448248>, <https://doi.org/10.1145/3448248>