

Assessing Trustworthiness of Autonomous Systems

Greg Chance¹, Kerstin Eder¹

Trustworthy Systems Lab, University of Bristol, Bristol, UK

Abstract—As autonomous systems (AS) become more ubiquitous in society, more responsible for our safety and our interaction with them more frequent, it is essential that they are trustworthy. Assessing the trustworthiness of AS is a mandatory challenge for the verification and development community. This will require appropriate standards and suitable metrics that may serve to objectively and comparatively judge trustworthiness of AS across the broad range of current and future applications. The meta-expression ‘trustworthiness’ is examined in the context of AS capturing the relevant qualities that comprise this term in the literature. A list of challenges are presented in the form of a process that can be used as a trustworthiness assessment framework for AS.

1 Introduction

Autonomous systems (AS) are pervasive in current society and set to become even more so with current technological growth trends and adoption rates. Systems with embedded artificial intelligence (AI) and machine learning (ML) algorithms can be found in numerous applications from mobile phones [mediumaiphones], insurance pricing [kuo2020towards] and self-driving vehicles to medical diagnostics [kononenko2001machine], detecting structural damage to buildings [avci2021review] and predicting the shape of protein molecules [alpha-fold] to name a few. For successful adoption of these systems then there needs to be demonstrable assurance of their trustworthy operation which becomes increasingly difficult in complex and changing environments. There is also growing use of machine learning in a range of safety-critical systems (SCSs), for example in the aerospace and automotive industry where low reliability of these systems could result in catastrophic failure and potentially loss of life or damage to property and the environment. These safety-critical autonomous systems present a complex but essential challenge to the safety assurance and verification community.

Verification and validation (V&V) is the process to gain confidence in the correctness of a system relative to its requirements. Prior to, and separate from verification, a specification must clearly define the trustworthy operational behaviour of the system and many chal-

lenges are associated with this task for autonomous systems [Abeywickrama2022]. If autonomous systems are to be fully accepted into society, there must be acknowledgement of, and evidence to show, compliance with a broad range of *trustworthiness qualities*. A major challenge in verifying this broad category of trustworthiness qualities, are the lack of standards and regulations against which they can be evaluated. And whereas verification methods of assessing, for example, functional correctness are relatively mature, there also exists the challenge of developing robust assessment methodologies and metrics for these more nuanced trustworthiness qualities. This paper focuses on reviewing what trustworthiness means in the fields of AS, robotics, HRI and Cyber-Physical Systems (CPS). An assessment process is presented along with the challenges associated with practical application.

1.1 V&V for AS in complex environments

A widely held tenet is that there can never be a suitable amount of verification that gives complete assurance for complex, safety-critical autonomous systems, although limits on reliability rates have been proposed [Butler1993]. Corroborative V&V [webster2020corroborative] attempts to improve confidence through combining mutually consistent evidence from multiple and diverse assessment methods, e.g. formal, testing. But even this may not be enough for the diverse operational domains of some AS, e.g. automated vehicles in high-density urban areas, and thinking should move beyond verification at only the system design stage, to a more continuous operational evaluation such as runtime verification. Runtime verification brings other currently unresolved issues, such as suitable oracle design [Leucker2009], but some authors propose valid ideas to this using edge computing and cloud-based verification [CyRes20, eder2021cyres]. The use of *serious games* can be another interesting opportunity for building trustworthiness in complex autonomous systems and has been used in the context of mission planning for NASA [Allen2018] and in a self-driving vehicle controller leveraging the power of crowdsourcing for test generation [ref test gen game, need to make github page, add code and youtube videos].

Further to this issue, are the lack of *standards* against which some trustworthy qualities should be appraised and the *methods* by which they should be evaluated. For example, there are standards for correct road driving conduct [highwayCode] but no ethical standards by which those driving decisions should be made [ref like trolley problem?]. Although headway is being made into developing standards for non-functional properties, such as guidelines for ethical AI

Statements about authorship contribution. Greg Chance (e-mail: greg.chance@bristol.ac.uk), and Kerstin Eder (e-mail: kerstin.eder@bristol.ac.uk) are with the Trustworthy Systems Lab, Department of Computer Science, University of Bristol, Merchant Ventures Building, Woodland Road, Bristol, BS8 1UQ, United Kingdom.

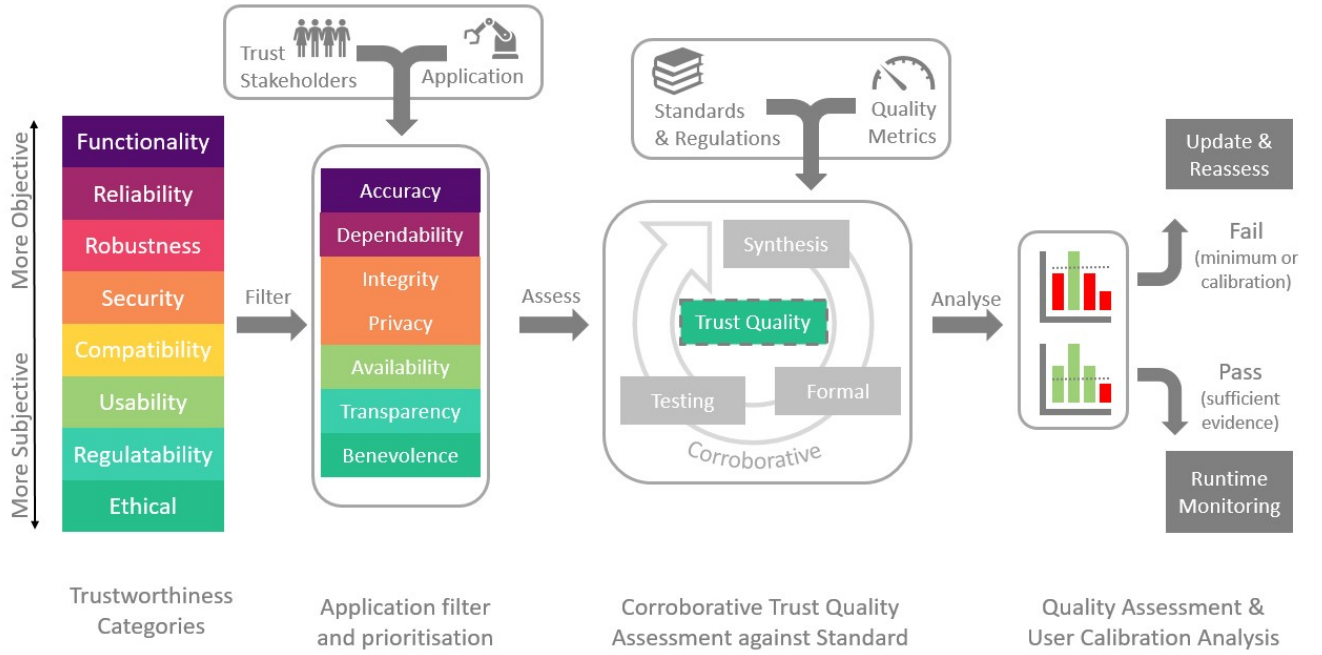


Fig. 1: AS trustworthiness assessment process

[ref EU AI high level expert group], checklists for HRI best practice [kraus2022trustworthy] and transparency [winfield2021ieee], there are still areas that need attention, such as standards for adaptability, cooperation and fairness [Abeywickrama2022]. Where standards are lacking or immature will require engagement with *trust stakeholders*, expert steering groups that can define and prioritise the necessary trustworthiness qualities for each subject domain or application.

Additionally, there is more that can be done at the design stage to improve *verifiability* [add ref]. Evidence for functional correctness is essential, but this must be supported with decision explanation [koopman2018toward] whilst maintaining intellectual property rights around, for example, sensitive software algorithms and trade secrets [ref].

In addition to assessing the AS trustworthiness, there must also be consideration to gain, calibrate and maintain user trust in the system [kok2020trust, Chiou2021], as miscalibration of trust between system and user can have serious consequences [kok2020trust].

1.2 Trustworthiness Qualities

As computing and automation has developed, systems are now both more capable and users more reliant on them. This extension of capability has resulted in a broadening of the terms which encompass trustworthiness, as, for example, the important trustworthiness qualities of a calculator may be less numerous than those of a medical decision support system. Advancement in automation then, has led us to question and challenge these new capabilities, or, with more automation comes more responsibility [Yazdanpanah2021].

Trust can be expressed in a number of ways and directions; trust the user has in the system, the

objective trustworthiness of the system and the context in which the interaction between the two takes place [Hancock2021]. Trustworthiness can also be described as the probability that a system holds some established property or quality, and that greater trustworthiness begets greater likelihood that the system may exhibit that quality. In this research we consider the trustworthiness of the system and the specific qualities that must be demonstrated, but we acknowledge the importance of the other mechanisms where human-system trust can be gained or lost in which there has been much contribution from the HRI, psychology and human factors community [Floridi2019, Lee2004, kok2020trust, Chiou2021, Kohn2021, kraus2022trustworthy]. Trustworthiness of autonomous systems in the context of this work then, results from objective assessment of the system with respect to a set of appropriate standards. There has been much academic deliberation on the specific qualities that comprise trustworthiness of AS, specifically for AI [Thiebes2021, Wing2021] and HRI [kraus2022trustworthy, atkinson2012trust]. Devitt argues that reliability and accuracy are the two central pillars of trustworthiness of AS and that all other properties stem from these, for example, stating that adaptability and redundancy are higher-order properties of reliability [devitt2018trustworthiness].

[ts'foundation] state 5 facets of trustworthy software: Safety: The ability of the software to operate without causing harm to anything or anyone. Reliability: The ability of the software to operate correctly. Availability: The ability of the software to operate when required. Resilience: The ability of the software to recover from errors quickly and completely. Security: The ability of the software to remain protected against the hazards posed by malware, hackers or accidental misuse.

There is, in fact, a spectrum of qualities that comprise user trust in an autonomous system...(leads into categories)

A summary of AS trustworthy qualities in the literature can be found at [https://github.com/TSL-UOB/TAS-Verif]

1.3 Ontology of AS Trustworthiness Qualities

An ontology of trustworthiness qualities is presented. This is an independent set of characteristics, where one category is not necessary influenced or related to its neighbours, that supports clarity of the assessment [connett 2018]

(Lee & See 2004) performance, process, purpose, which although are broadly capturing of all trustworthy qualities miss some nuance to give a complete and expressive account of trustworthiness.

avizenis2004basic a set of general concepts are required for dependable and secure computing, which may cover a wide range of applications and system failures, which comprise; availability (readiness for correct service), reliability (continuity of correct service), safety (absence of catastrophic consequences on the user(s) and the environment), integrity (absence of improper system alterations) and maintainability (ability to undergo modifications and repairs).

Thiebes et al. argue for five foundational principles of trustworthy AS: beneficence (doing good), non-maleficence (not harming), autonomy (preserving human decision making), justice (fair and reasonable), and explicability (easily understood) [Thiebes2021]. These are based on and related to other discussion on ethically principled foundations of trustworthiness, Floridi 2018 [Floridi2018], and should be weighted appropriate to the strong international contribution to this area Jobin 2019 [jobin2019global] but fail to capture the full trustworthiness spectrum of qualities presented here.

2 Assessment Framework Vision

Below is a checklist for assessemnt of AS trustworthiness: **Metrics** Floridi suggests the need for agreed upon metrics for trustworthiness of AI systems and suggests an AI Trust comparison index, metrics are needed for benchmarking AI suitability to the public. Rudas and Haidegger also supports the idea of agreed upon metrics from the verification community that can be used to ensure reliability of complex autonomous systems [Rudas2020]. Wang et al. go further and propose a theoretical framework of *tripartite trustworthiness* covering; *to-be trust* (trustfulness of an entity or structure), *to-do trust* (trust in an action or behaviour) and *system trust* (a statistical runtime evaluation of performance) and set out 18 formal definitions [Wang2020]. Garbuk presents the idea of *applied intellimetry* to assess the quality of AI systems by formulating a list quality characteristics in a functional characteristic vector [garbuk2018intellimetry]. Kaur et al. suggest explainability metrics based on the euclidean distance between the system output compared to

a panel of experts [kaur2021trustworthy]. trustworthiness of computer systems using metrics designed to assess security, trust, resilience and agility [cho2019stram]

Bolster and Marshall proposes the idea of *multi-vector trust metrics* for networks of autonomous systems, indicating that the use of *grey relational analysis*, a theory to describe and model uncertainty, could be beneficial for combining temporally sparse, low fidelity metrics with unknown statistical distributions [Bolster2014].

Verification Methods: Kress-Gazit et al. state that assessment in the correctness of AS can be broken down into four approaches: synthesis of correct-by-construction systems, formal verification at design time, runtime verification or monitoring, and test-based methods [kress2021formalizing].

2.1 Existing Standards for AS

Existing standards on verification (from Rudas 2020): P1872.1, P2817, P7000 and P7007.

Safety of autonomous systems (from Hawkins 2022): UL4000 [Underwriters Laboratories. Standard for evaluation of autonomous products, 2020] or SCSC-153B [Safety of Autonomous Systems Working Group. Safety assurance objectives for autonomous systems, 2022. URL: https://scsc.uk/scsc-153B]

Ethical framework for AI [Floridi2018] “offer 20 concrete recommendations to assess, to develop, to incentivise, and to support good AI”

Porter2022 presents an ethical assurance argument for AS, extending the assurance case considered for safety to include ethical standards

devitt2018trustworthiness

What human factors are considered eg. ISO29119

Dhaminda to contribute here?

2.2 Future Challenges in Standards

Riaz et al. [Riaz2018] suggest the idea of using social norms and human emotions as a standard by which better self-driving controllers may be developed. This idea sets the way for not just development of higher functioning AS, but also standards of trustworthiness by which they can be judged. Although there is much scholarly work on the theory and modelling of social norms, e.g. [hechter2001social], there is yet to be published a standard that could be used to objectively assess an autonomous system.

In some cases, e.g. driving, legislation on appropriate conduct is presented to society in the form of guidelines such as the UKHC in the UK [highwayCode] but must be translated to a computer readable format to act as an appropriate standard, or set of assertions [harper2021safety], if these guidelines can be used to assess AS trustworthiness. A similar process will have to be undertaken for other standards which have yet to be defined, e.g. cooperation, fairness or verifiability, to ensure all aspects of trustworthiness can be assessed.

-
1. protection against intentional subversion or forced failure, malicious access, use, modification, destruction, or disclosure
 2. defining,
-

Table 1: A description of the test scenarios showing the test number, the actors included, whether a collision occurred and if so then between which actors. n , the number of repeats is set to 1000 and $\max \sigma$ is the maximum simulation deviation. The term *unrestricted* refers to an unrestricted account of the results including results of any resource utilisation. To understand the impact of collisions and high resource utilisation, the *restricted* column shows a subset of the results where post-collision data and experiments above 75% resource utilisation have been removed.

2.3 Assessment Methods & Corroborative Evidence

Gaining reliability assurance of SCASs using testing alone is unfeasible given the often high-dimensional operational state space. Multiple testing methodologies should be employed where appropriate, e.g. verification, falsification and testing, [Harper Corroborative 2022] combining mutually consistent evidence from multiple and diverse assessment methods will raise the confidence in system trustworthiness.

Knowledge of the internal state of the system is often hidden, e.g. blackbox, due to IP and commercial sen-

sitivity, but whitebox access will be essential for certain aspects of trustworthiness assessment. This may not need to reveal sensitive algorithms but just enough information through observability points in the software architecture could go a long way to understanding if automated decisions are made for the right reason [koopman2018toward].

3 Conclusion