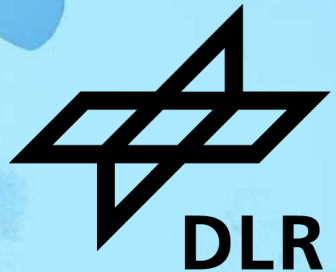


THE NECESSITY AND POSSIBILITY OF TRUSTWORTHY AI

Hagen Braun & Lukas Albrecht



Aims of the paper



De-buzzwordify Trustworthy AI

Trustworthy AI seems like a “magic concept” if ever there was one.

Give a theory-driven account of the concept

The EU Assessment List for Trustworthy AI is eclectic. We aim to take the existing scholarship on trust seriously and use it to unify the concept of trustworthy AI.

Explain how trust in the AI context helps bring about better ethics

If the ALTAI list is an ethics guideline, why is this aim expressed in terms of trust?

Necessity of Trustworthy AI



Why do we need trustworthy AI?

- Relying on AI is risky in the same characteristic way that relying on other (human) agents is risky.
- This risk cannot be removed entirely without losing what makes AI use desirable in the first place.
- The paradigmatic way to manage low, but non-zero amounts of risk in the context of cooperation between agents is through trust.
- But this necessitates a trustworthy cooperative partner.
- **Therefore, trustworthy AI is necessary.**

What is Trust?



Trust (relation)

Trust is a three-place relation between two agents and a task:

$$T(A_1, A_2, \Phi)$$

The first agent A_1 (the trustor) trusts the second agent A_2 (the trustee) with the task Φ .

With regard to its agential components, this relation is not symmetric, reflexive or transitive.

Important takeaways

- The basic conception of trust is *between agents*.
- Trust is specific to a *task*.

What is Trustworthiness?



Trustworthiness

Trustworthiness is a complex property of agents. It reduces to the properties of *reliability*, *competence* with regard to a certain task, and the possession of a *properly normatively grounded motivation* to perform that task.

Trust is semi-voluntary: An agent can decide to trust anyone they want.

But *well-placed trust* requires identifying what agents *deserve* your trust (i.e.: who is trustworthy).

Only well-placed trust actually manages risk in cooperative relationships.

Preconditions of Trustworthiness

Breaking Down Trustworthiness

Reliability



Competence



Normative Motivation



- Reliability is a general precondition for well-placed trust
- Competence makes a person trustworthy with regard to a specific task
- Normative motivation provides an explanatory basis of how trust can manage risk

Are AI Agents Real Agents?



Basic Theory of Action

An action requires a belief-desire pair that constitutes the reason for the action.

It both *causally explains* the action and *justifies* it from the perspective of the agent.

- For AI agents, we can defend the claim that they possess functionally equivalent states to beliefs and desires that can play their causal role.
- But to the best of our knowledge, **AI agents do not possess a perspective from which they justify their own actions in normative terms.**

Are AI Agents Real Agents?

If we take action theory seriously, we cannot categorize AI agents as full agents.

However, from the outside, their behaviour is indistinguishable from real action in many cases.

We therefore want to categorize them as **Pseudo-Agents**: Their actions have the causal component of a real action, without having any internal justifying component. **This means they are still a source of agential risk.**



An AI agent

The Possibility of Trustworthy AI



However full normative agency is not a precondition for trustworthiness.

The possession of a certain kind of motivation for one's actions that can be adequately explained in normative terms is.

AI agents can possess motivational equivalent states that are explainable in this way, even though they cannot provide these explanations themselves.

But AI developers can.

If we view AI actions as an exercise of joint agency, then the concept of trustworthy AI does make sense. Notably, it does not *just* reduce to trust in AI developers. The agent who acts, and whose agential risk is managed through the trust relationship is still the AI agent.