

Martijn Demollin

Politechnika Warszawska

Katarzyna Budzynska

Politechnika Warszawska

Konrad Kiljan

Polish Debating Foundation

Ethos and Trust in Explainable AI

A central endeavour within explainable AI is facilitating transparency and trust in AI systems. However, in order to evaluate the trustworthiness of AI systems, there is a need for determining the characteristics that make an entity trustworthy.

In argumentation theory and philosophy, Ethos defines the trustworthiness of an individual in terms of their wisdom, virtue and good will. We employ this distinction to assess the fundamental qualities that are needed in AI systems in various contexts to be considered worthy of trust.