

# SafeAI 2022



## The AAAI-22 Workshop on Artificial Intelligence Safety

Feb 28<sup>th</sup> - Mar 1<sup>st</sup>, 2022  
Virtual

Gabriel Pedroza, CEA LIST, France

José Hernández-Orallo, Universitat Politècnica de València, Spain

Xin Cynthia Chen, University of Hong Kong, China

Xiaowei Huang, University of Liverpool, UK

Huáscar Espinoza, ECSEL JU, Belgium

Mauricio Castillo-Effen, Lockheed Martin, USA

Seán Ó hÉigearthaigh, University of Cambridge, UK

Richard Mallah, Future of Life Institute, USA

John McDermid, University of York, UK

# Opening Remarks



*The main interest of SafeAI 2022 is to explore new ideas on **AI safety** by looking holistically at theoretical and practical, short-term and long-term, perspectives, jointly with the ethical and legal issues, to build trustable intelligent autonomous machines.*

- As SafeAI aims at bringing together **multiple perspectives**, it's probably **possible to harshly criticise any paper** here today... Most likely anyone has missed some important issue...
- So, please do be critical, but temper your criticism with **constructive discussions!**

# Program (See Website for Details)



**Day 1:** Scheduled on **Feb 28, 2022** from **13:00 to 18:10 UTC [5:00-10:10 PST]**

**Location:** **AAAI Virtual Venue, Red Building Room 5.**

Time (UTC)	Description
13:00-13:05	Welcome and Introduction – Chair: Gabriel Pedroza (CEA List)
13:05-13:50	<b>Keynote 1: Matthew Dwyer (University of Virginia), <i>Distribution-aware Test Adequacy for Neural Networks</i></b>
13:50-14:00	Short Break
14:00-15:30	<b>Special Session 1: <i>EnnCore</i> – Chair: Lucas Cordeiro (University of Manchester)</b> <b>EnnCore</b> addresses the fundamental problem of guaranteeing safety, transparency, and robustness in neural-based architectures. <a href="#">DETAILED ENNCORE PROGRAM...</a>
15:30-16:00	Coffee Break
16:00-16:25	<b>Invited Talk 1: Shiri Dori-Hacohen (University of Connecticut), <i>Quantifying Misalignment Between Agents</i></b>
	<b>Technical Session 1: <i>Bliss, Fairness and Value Alignment</i> – Chair: José Hernández</b>

**Day 2:** Scheduled on **Mar 1, 2022** from **8:00 to 17:15 UTC [0:00-9:15 PST]**

**Location:** **AAAI Virtual Venue, Red Building Room 5.**

Time (UTC)	Description
8:00-8:25	<b>Invited Talk 2: Roel Dobbe (TU Delft), <i>A System Safety Perspective for Developing and Governing Artificial Intelligence</i></b>
8:25-9:35	<b>Technical Session 3: <i>Robustness and Uncertainty</i> – Chair: Xin Cynthia Chen (University of Hong Kong)</b> <ul style="list-style-type: none"><li>– Efficient Adversarial Sequence Generation for RNN with Symbolic Weighted Finite Automata, Mingjun Ma, Dehui Du, Yuanhao Liu, Yanyun Wang and Yiyang Li.</li><li>– A Study on Mitigating Hard Boundaries of Decision-Tree-based Uncertainty Estimates for AI Models, Pascal Gerber, Lisa Jöckel and Michael Kläs.</li><li>– Quantifying the Importance of Latent Features in Neural Networks, Amany Alshareef, Nicolas Berthier, Sven Schewe and Xiaowei Huang.</li><li>– Maximum Likelihood Uncertainty Estimation: Robustness to Outliers, Deebul Nair, Nico Hochgeschwender and Miguel Olivares-Mendez.</li><li>– Debate Panel – Paper Discussants: Xiaowei Huang (University of Liverpool), Mauricio Castillo-Effen (Lockheed Martin)</li></ul>

# Some Additional Information



- SafeAI 2022 **Best Paper Award**.
- **Proceedings** is freely available at CEUR-WS: <http://ceur-ws.org/Vol-3087/>  
(URL is available at the SafeAI website)
- **Presentations** will be available on the website very soon
- We hope you enjoy **SafeAI 2022!**