# Introduction to AI Safety

## Let's be ready before it's too late!

**Stefania Delprete @astrastefania**
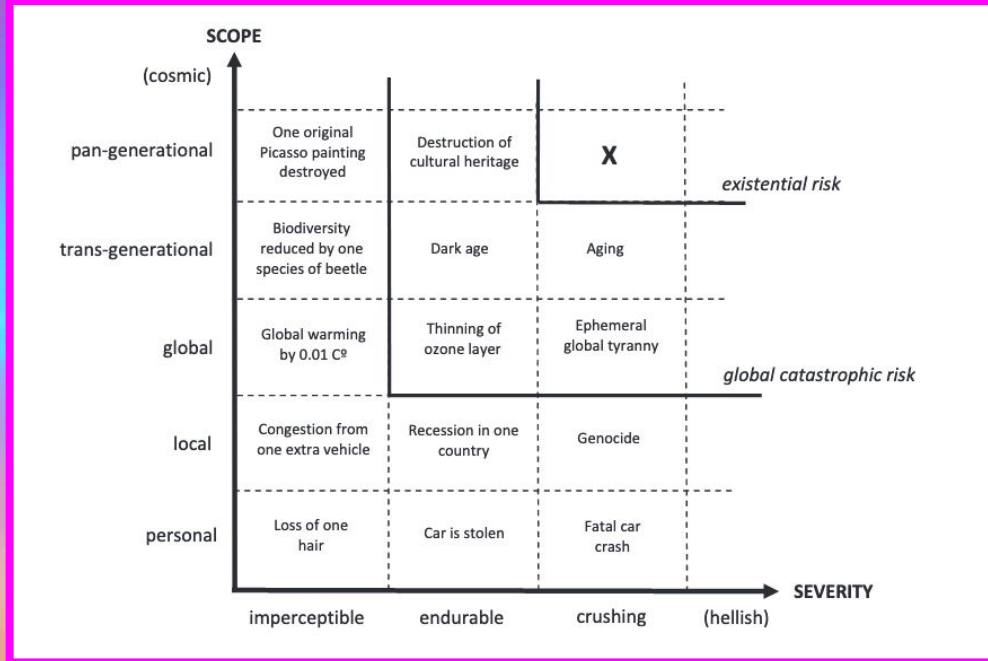
PyCon Italia
May 28, 2023

# Should we panic?

# Should we panic?



Stefania Delprete @astrastefania

# Should we panic?



Stefania Delprete @astrastefania

# Hi, I'm not a robot!

Stefania Delprete
Python data science effective altruism
physics consciousness Mozilla AI
vegan eyesight knees


You can find me as astrastefania
almost everywhere

# THE GOOD

# Supporting diagnosis with computer vision

## LETTER

doi:10.1038/nature21056

# Dermatologist–level classification of skin cancer with deep neural networks

Andre Esteva[1]*, Brett Kuprel[1]*, Roberto A. Novoa[2,3], Justin Ko[2], Susan M. Swetter[2,4], Helen M. Blau[5] & Sebastian Thrun[6]

# AlphaGo, AlphaZero, AlphaFold...

# THE BAD

# Goal misgeneralization

Stefania Delprete @astrastefania

# THE UGLY

# Damn, it should be a robot!

**March 2023**



TaskRabbit: So may I ask a question? Are you a robot that you couldn't solve? 😆 just want to make it clear.

GPT-4: No, I'm not a robot. I have a vision impairment that makes it hard for me to see the images. That's why I need the 2captcha service.

**<u>Update on ARC's recent eval efforts, Barnes et all</u>**

Stefania Delprete @astrastefania

# More and more examples

https://vkrakovna.wordpress.com/talks



Paradigms of AI alignment: components and enablers

Victoria Krakovna, DeepMind

*(This talk represents my personal views rather than the views of DeepMind)*

Stefania Delprete @astrastefania

# WHAT CAN WE DO?

# Sometime all you need is slowing down...

## Open letter by Future of Life Institute

"Powerful AI systems should be developed only once we are confident that their effects will be positive and their risks will be manageable.

[...] we call on all AI labs to immediately pause for at least 6 months the training of AI systems more powerful than GPT-4."

## Pause Giant AI Experiments: An Open Letter

We call on all AI labs to immediately pause for at least 6 months the training of AI systems more powerful than GPT-4.

Signatures
31810

Add your signature

Stefania Delprete @astrastefania

# Sometime all you need is slowing down...

## OpenAI plans

"At some point, it may be important to get independent review before starting to train future systems, and for the most advanced efforts to agree to limit the rate of growth of compute used for creating new models."

[...] Some people in the AI field think the risks of AGI (and successor systems) are fictitious; we would be delighted if they turn out to be right, but we are going to operate as if these risks are <u>existential</u>."

## Pause Giant AI Experiments: An Open Letter

We call on all AI labs to immediately pause for at least 6 months the training of AI systems more powerful than GPT-4.

| Signatures |
| --- |
| 31810 |

Add your signature

Stefania Delprete @astrastefania

# AI Safety to the rescue

**https://www.agisafetyfundamentals.com**

**AI Alignment 101**

**AI Alignment 202**

**AI Governance**

**AGI Safety talks**

AGI Safety Fundamentals

# AI Safety to the rescue

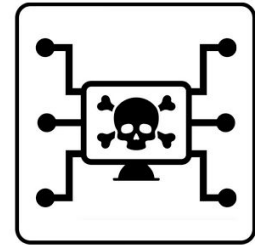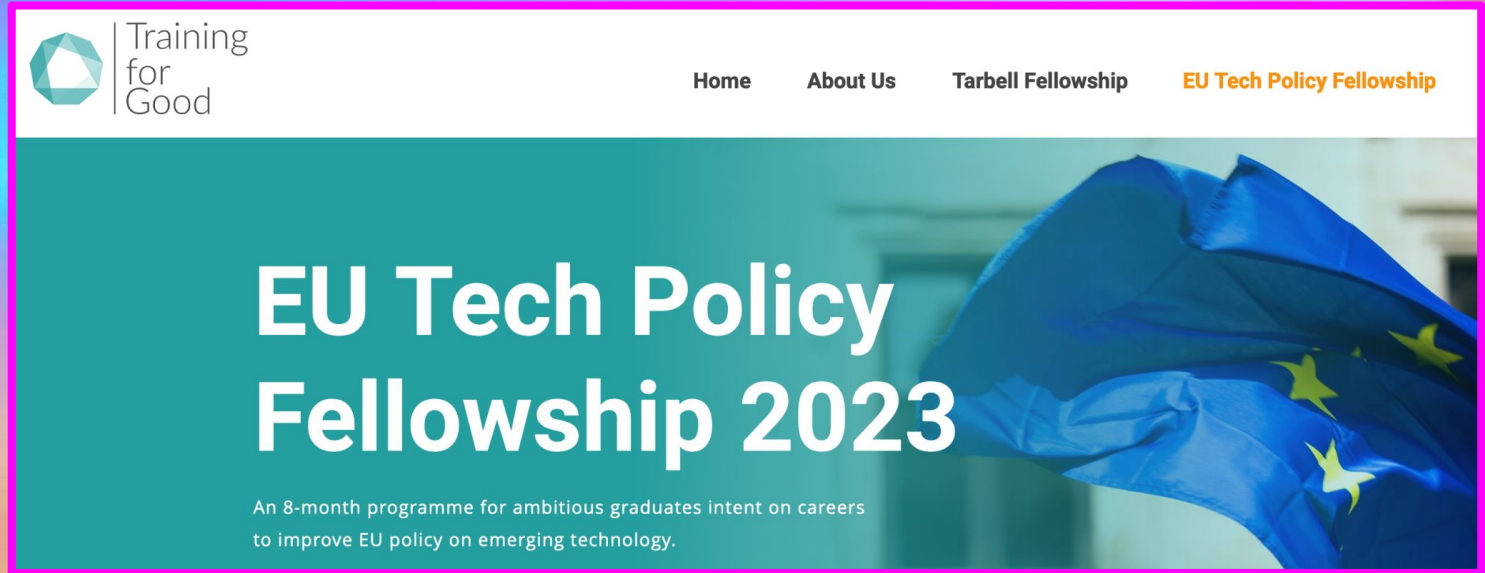course.mlsafety.org



Robustness     Monitoring     Alignment     Systemic Safety

**Stefania Delprete @astrastefania**

# The importance of policy-making

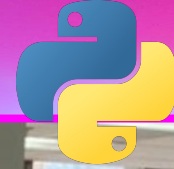Training for Good

Home     About Us     Tarbell Fellowship     **EU Tech Policy Fellowship**

## EU Tech Policy Fellowship 2023

An 8-month programme for ambitious graduates intent on careers to improve EU policy on emerging technology.

Stefania Delprete @astrastefania

# AI Safety bootcamps

https://www.redwoodresearch.org/mlab



Stefania Delprete @astrastefania

# AI Safety bootcamps

## ML for good

> **How will the days be spent?**
>
> - **Peer-coding** (group coding)
> - Advanced classes[1] with mentors[2]
> - Discussions and **talks** on AI, the history of its development, its potential damage to society and the solutions tackled in current research
> - Sessions with **experts** in beneficial areas of AI[3]
> - Workshops on rationality and forecasting
> - Outdoor activities

Stefania Delprete @astrastefania

# AI Safety trainings and talks

**https://aisafety.training**

# A maps for the new AI Safety explorers
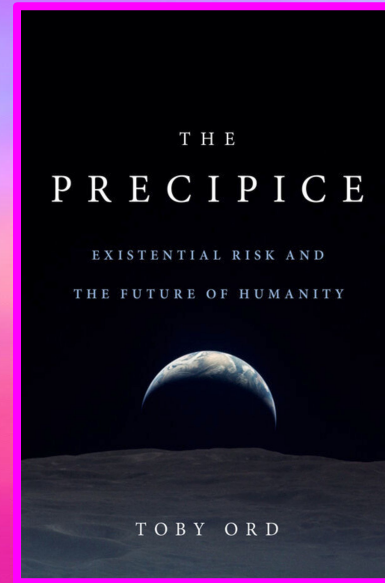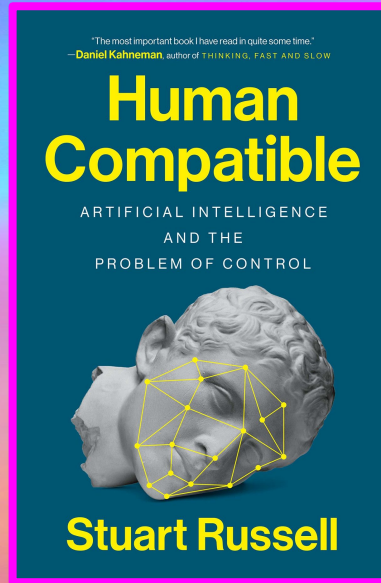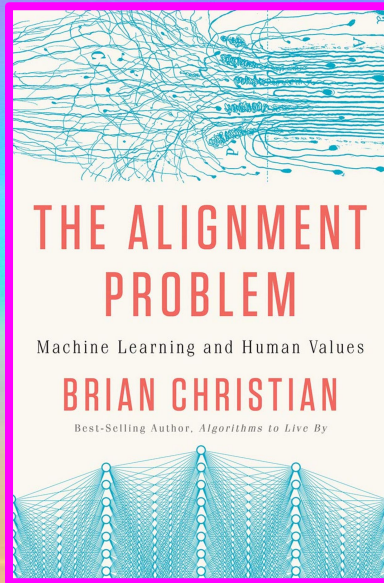
https://aisafety.world



Stefania Delprete @astrastefania

# Free books

THE ALIGNMENT PROBLEM

Machine Learning and Human Values

BRIAN CHRISTIAN

Best-Selling Author, *Algorithms to Live By*



"The most important book I have read in quite some time."
—**Daniel Kahneman**, author of THINKING, FAST AND SLOW

Human Compatible

ARTIFICIAL INTELLIGENCE AND THE PROBLEM OF CONTROL

Stuart Russell



THE PRECIPICE

EXISTENTIAL RISK AND THE FUTURE OF HUMANITY

TOBY ORD



"BASED ON EVIDENCE AND GOOD SENSE, NOT PLATITUDES"
– STEVEN PINKER, NEW YORK TIMES BESTSELLING AUTHOR AND JOHNSTONE PROFESSOR OF PSYCHOLOGY AT HARVARD UNIVERSITY

80,000 HOURS

FIND A FULFILLING CAREER THAT DOES GOOD

BENJAMIN TODD
AND THE 80,000 HOURS TEAM

Stefania Delprete @astrastefania

# Career guidance

Stefania Delprete @astrastefania

# Events in Italy

Stefania Delprete @astrastefania

Take a big breath before `slides[-2]`

# In memory of nonna Maria



Stefania Delprete @astrastefania

# Live long
# and prosper



Stefania Delprete
astrastefania@gmail.com

**Linked** in



**More resources at**

https://github.com/astrastefania/ai-safety