*Q:*
 "If human level general AI is developed, then what are likely outcomes?"

# Lee Sedol v. AlphaGO

**Oct 2015**   "Based on its level seen in the match (against Fan), I think I will win the game by a near landslide"

**Feb 2016**   "I have heard that Google DeepMind's AI is surprisingly strong and getting stronger, but I am confident that I can win at least this time"

**Mar 9, 2016**   "I was very surprised because I didn't think I would lose"

**Mar 10, 2016**   "I'm quite speechless ... I am in shock. I can admit that ... the third game is not going to be easy for me"
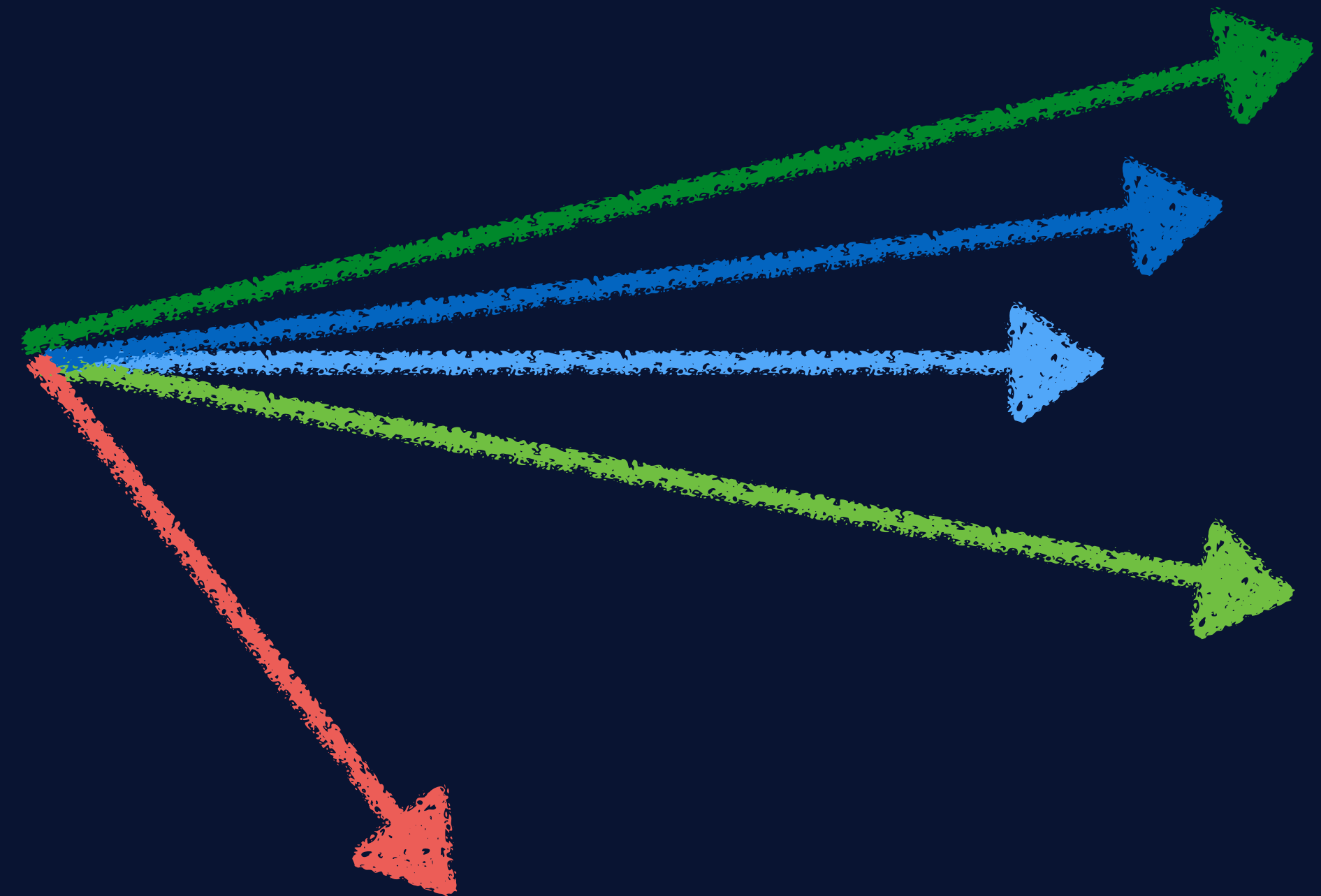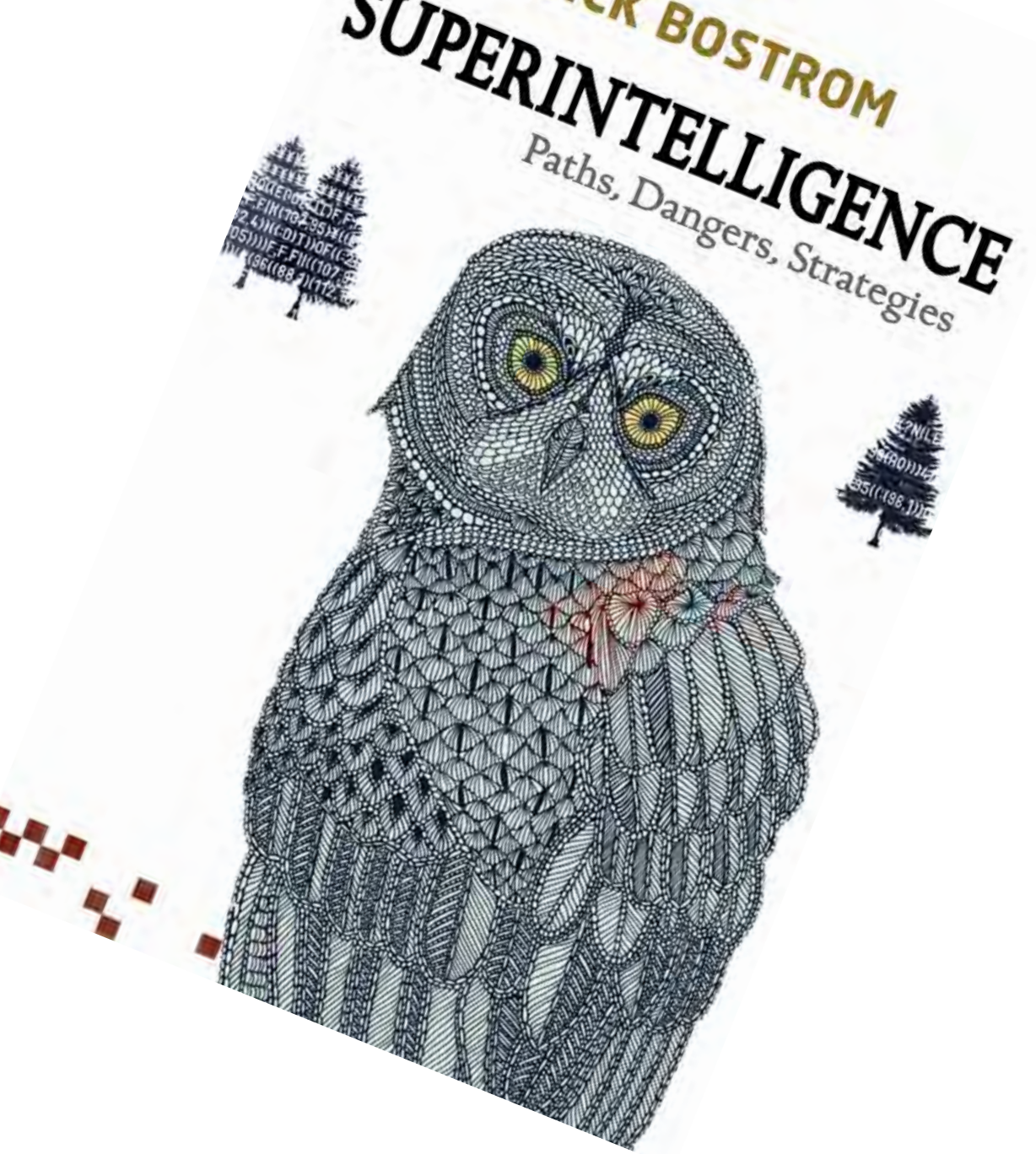
**Mar 12, 2016**   "I kind of felt powerless."

**SUPERINTELLIGENCE**
Paths, Dangers, Strategies

BOSTROM

0 of 1 people found the following review helpful

★☆☆☆☆ **One Star**, September 10, 2016

By **Amazon Customer**

Verified Purchase (What's this?)

This review is from: **Superintelligence: Paths, Dangers, Strategies (Paperback)**

I read it in 3 days and I'm profoundly depressed.

Help other customers find the most helpful reviews

Was this review helpful to you? [ Yes ] [ No ]

Report abuse | Permalink

*Q:*
  "What can we do now to maximize the probability of a positive outcome?"

*A:*
- solve intelligence
- solve scalable control
- solve AI governance problem

ICML 2016
THE 33rd INTERNATIONAL CON

nature
LEARNING CURVE

Science
$10 17 JULY 20
sciencemag.org

SPECIAL ISSUE
ARTIFICIAL
INTELLIGENCE

AREA CHAIRS

www.icml.cc/2016/

TWEETS 18.4K  FOLLOWING 1,182  FOLLOWERS 3,520  LIKES 17.5K

Miles Brundage @Miles_Brundage · Jan 1
"A Joint Speaker-Listener-
for Referring Expressions,"
arxiv.org/abs/1612.09542

Miles Brundage @Miles_Brundage · Jan 1
"Feedback Networks," Zar
arxiv.org/abs/1612.09508

Miles Brundage @Miles_Brundage · Jan 1
"Automatic Discoveries of
Semantic Concepts via As
Neuron Groups," Li et al.:
arxiv.org/abs/1612.09438

undage @Miles_Brundage · Jan 1
on Recognition Based on Joint Trajectory
s with Convolutional Neural Networks,"
g et al.: arxiv.org/abs/1612.09401

Hierarchical Multiscale Recurrent Neural Networks
Junyoung Chung, Sungjin Ahn, Yoshua Bengio
12/14/2016  (v1: 9/6/2016)  cs.LG
1609.
show similar papers

Learning both hierarchical and temporal representation has been among the long-standing challenges of recurrent neural networks. Multiscale neural networks have been considered as a promising approach to resolve this issue, yet there has been a lack of empirical evidence show type of models can actually capture the temporal dependencies by discovering the latent hierarchical structure of the sequence. In this paper, a novel multiscale approach, called the hierarchical multiscale recurrent neural networks, which can capture the latent hierarchical stru sequence by encoding the temporal dependencies with different timescales using a novel update mechanism. We show some eviden proposed multiscale architecture can discover underlying hierarchical structure in the sequences without using explicit boundary information. our proposed model on character-level language modelling and handwriting sequence modelling.

Generating images with recurrent adversarial networks
Daniel Jiwoong Im, Chris Dongjoo Kim, Hui Jiang, Roland Memisevic
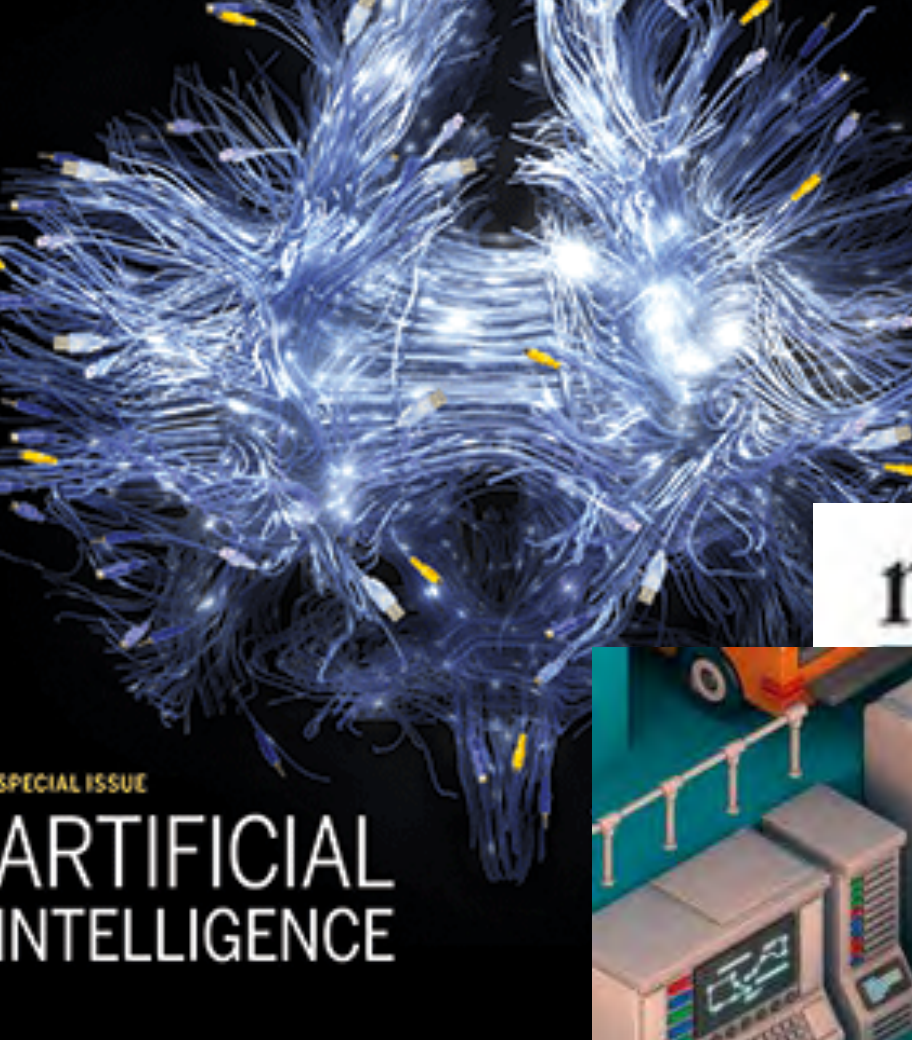12/13/2016  (v1: 2/16/2016)  cs.LG | cs.CV
1602.
show similar pap

Gatys et al. (2015) showed that optimizing pixels to match features in a convolutional network with respect reference image features is a wa images of high visual quality. We show that unrolling this gradient-based optimization yields a recurrent computation that creates images by i adding onto a visual "canvas". We propose a recurrent generative model inspired by this view, and show that it can be trained using adversa to generate very good image samples. We also propose a way to quantitatively compare adversarial networks by having the gene discriminators of these networks compete against each other.

SeqGAN: Sequence Generative Adversarial Nets with Policy Gradient
Lantao Yu, Weinan Zhang, Jun Wang, Yong Yu
12/9/2016  (v1: 9/18/2016)  cs.LG | cs.AI
The Thirty-First AAAI Conference on Artificial Intelligence (AAAI 2017)
1609.
show similar pap

NIPS 2016
BARCELONA · SPAIN · DECEMBER 5 - 10, 2015  |  http://nips.cc/

# technical research agendas

- inverse reinforcement learning
- adversarial examples
- models of control failure
- approval-maximizing agents
- imitation agents
- architectural composition
- corrigibility
- foundations of reflective agents
- detecting context change
- interpretability and explanation
- control diversification

# the AI governance problem

# Openness



- safety measures ✓
- values ✓
- (capability)
- source code, platforms
- science
- training data, environments, benchmarks

} ?

# Observation

Openness reduces the gap between the leading developer and the nearest follower.

- a couple of years in a low openness scenario?
- a few months in a high openness scenario?
- zero in the limiting case of maximal openness

This could help reduce the risk that a small group monopolizes all the benefits.

# Suppose that…



- safety requires some significant extra work after AI is completed

OR ➡️ doom

- safe operation initially incurs a significant performance penalty

OR ➡️ doom

- the Vulnerable World Hypothesis is true in the post-AI-transition world

# Vulnerable world hypothesis

There is some level of technology at which offense strongly dominates defense, in the sense that any small group of reasonably competent people with access to the technology would be able to take some action that would lead to the destruction of the world (independently of what other people did after the action was taken).

biotechnology?
nanotechnology?
doomsday device?

# Suppose that…



- safety requires some significant extra work after AI is completed ➡️ doom

OR

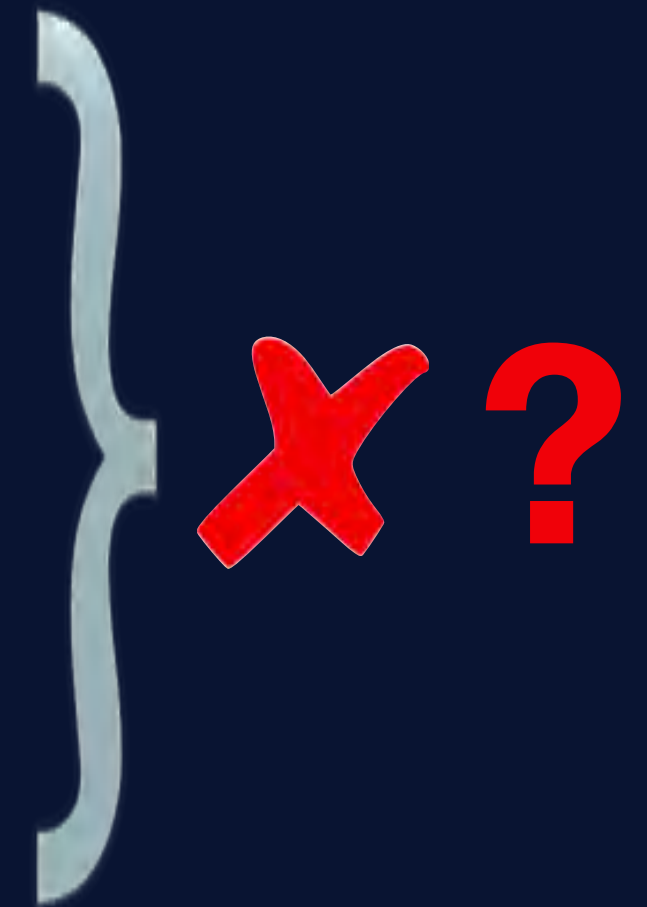- safe operation initially incurs a significant performance penalty ➡️ doom

OR

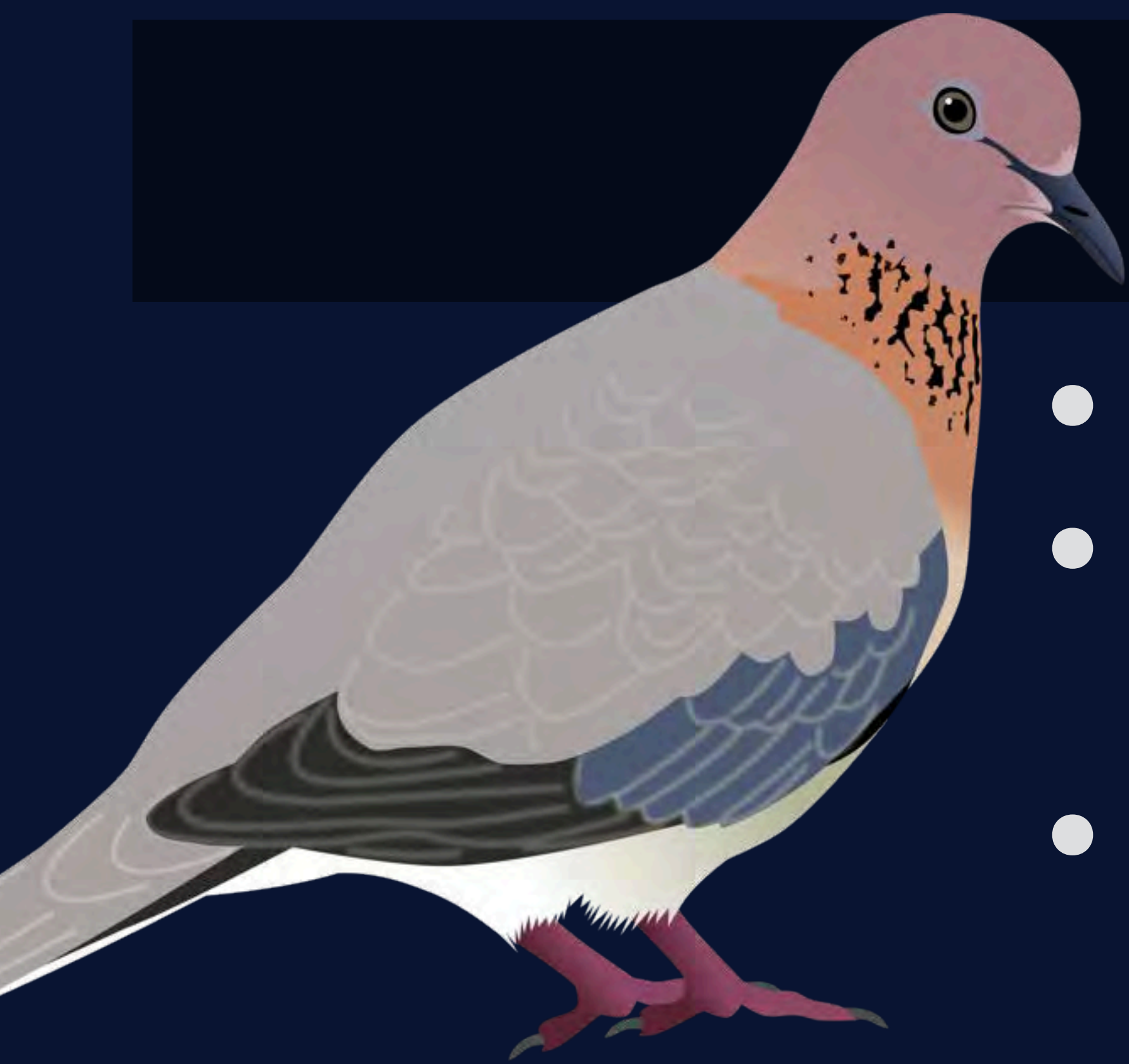- the Vulnerable World Hypothesis is true in the post-AI-transition world

# Openness



- safety measures ✓
- values ✓
- (capability)
- source code, platforms
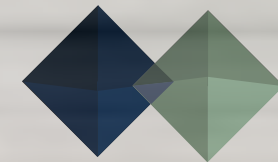- science
- training data, environments, benchmarks } ✗ ?

# What to do?

- openness for now
- desired property: **conditional stabilization…**
- lay the foundations for a collaborative approach later:
  - coordinate (or ideally pool) research among trusted leading groups
  - create ability **not** to share science and algorithms until it is safe to do so
  - credibly commit to sharing benefits and influence