

Modelling and regulating safety compliance: Game theory lessons from AI development races analyses

The Anh Han

AI Safety @ IJCAI, August 19-20, 2021



Tom Lenaerts,
ULB Brussels



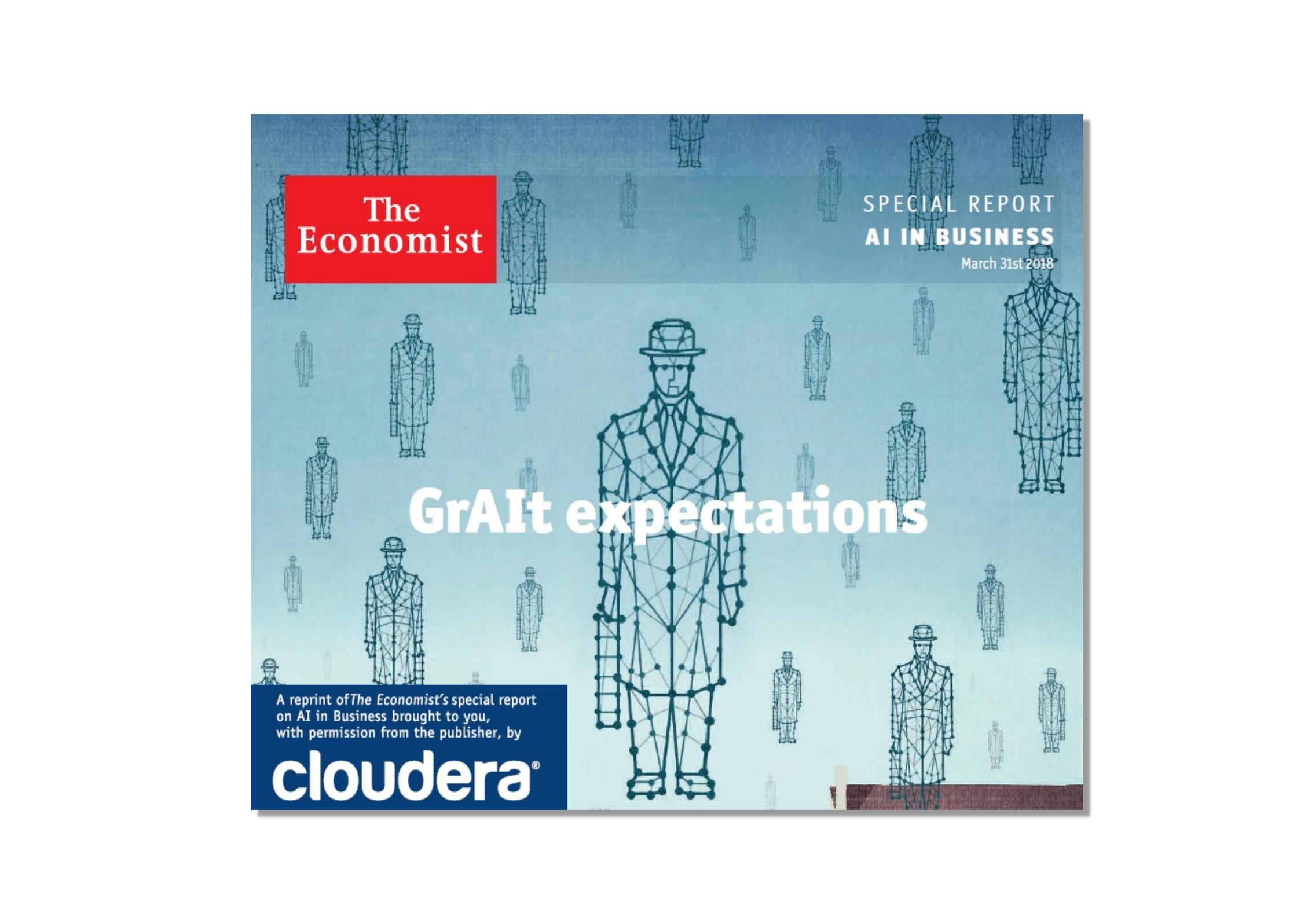
Luis Moniz Pereira,
UNL Lisbon



Francisco C. Santos,
IST Lisbon



t.han@tees.ac.uk



The
Economist

SPECIAL REPORT
AI IN BUSINESS

March 31st 2018

GrAIIt expectations

A reprint of *The Economist's* special report
on AI in Business brought to you,
with permission from the publisher, by

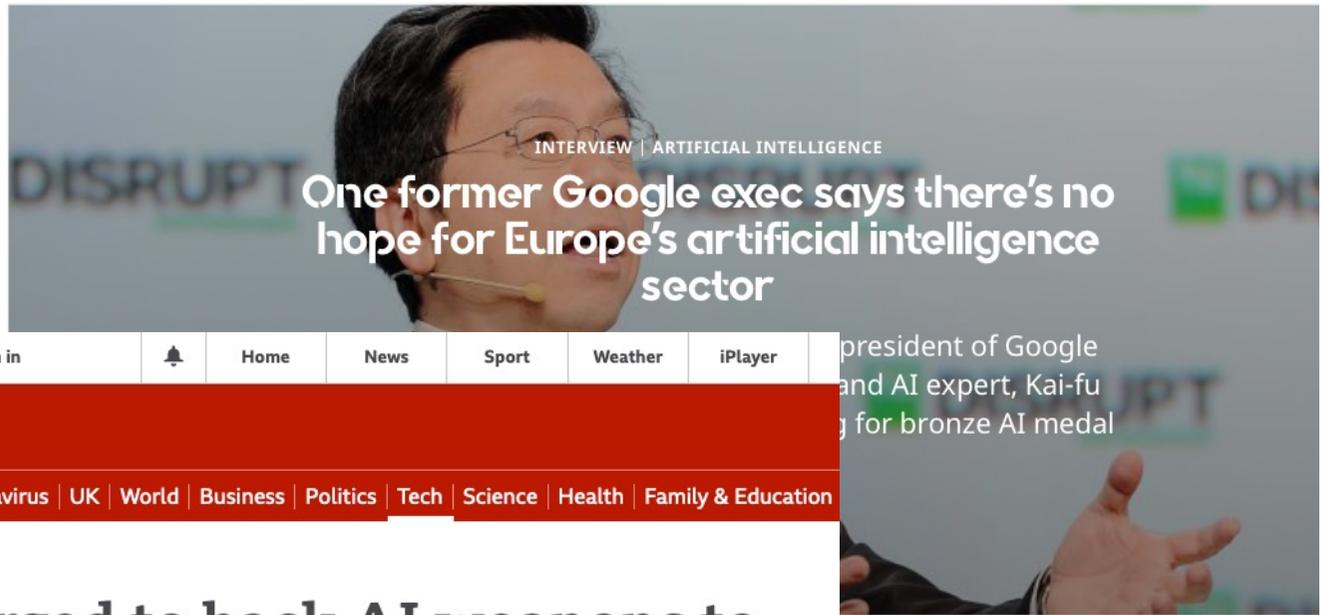
cloudera[®]

Commentary: Are China, Russia winning AI arms race?

Peter Apps

13 MIN READ

In October 31 Chinese teenagers reported to the Beijing Institute of Technology the country's premier military research establishments. Selected from millions of applicants, Chinese authorities hope they will design a new generation of



One former Google exec says there's no hope for Europe's artificial intelligence sector

INTERVIEW | ARTIFICIAL INTELLIGENCE

president of Google and AI expert, Kai-fu Ma, says there's no hope for Europe's artificial intelligence sector.

29.06.2018, TRYING TO KEEP UP

Europe's



Sign in



Home

News

Sport

Weather

iPlayer

NEWS

Home | Brexit | Coronavirus | UK | World | Business | Politics | Tech | Science | Health | Family & Education

Being a leading tech innovator is the future. But can Europe catch up?

Long Reads

From imitation to... China became a

Technology

Biden urged to back AI weapons to counter China and Russia threats

By Leo Kelion
Technology desk editor

9 hours ago



Figure 2 – Expected gains from AI



US WINNING AI RACE:
OR THE UNITED STATES?

MCLAUGHLIN, AND ELINE CHIVOT | AUGUST 2019



AI
Innovation/
Research

Question 1: What do we actually know about the dynamics of this (alleged) AI race?

Question 2: What do we know about the impact of regulatory actions on such race dynamics?

Question 3: Is regulation always useful? Useful in what sense?

Question 1: What do we actually know about the dynamics of this (alleged) AI race?

To Regulate or Not: A Social Dynamics Analysis of an Idealised AI Race

The Anh Han

*School of Computing, Engineering and Digital Technologies,
Teesside University, Middlesbrough, UK TS1 3BA*

T.HAN@TEES.AC.UK

Luís Moniz Pereira

*NOVA Laboratory for Computer Science and Informatics
(NOVA LINCS), Faculdade de Ciências e Tecnologia,
Universidade Nova de Lisboa, 2829-516 Caparica, Portugal*

LMP@FCT.UNL.PT

Francisco C. Santos

*INESC-ID and Instituto Superior Técnico,
Universidade de Lisboa, IST-Taguspark, 2744-016,
Porto Salvo, Portugal
§ Machine Learning Group, Université Libre de Bruxelles,
Boulevard du Triomphe CP212, Brussels, Belgium*

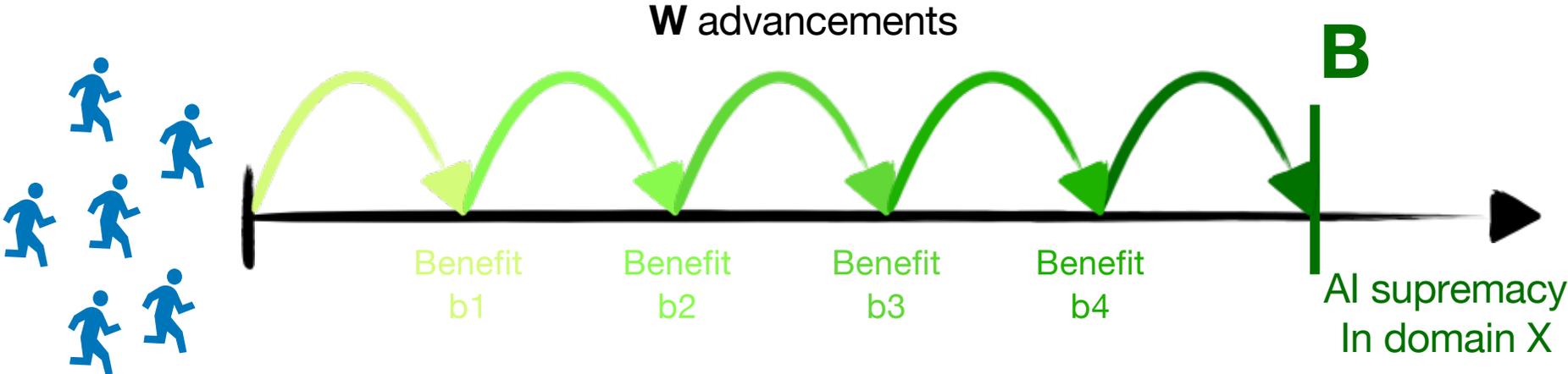
FRANCISCOCSANTOS@TECNICO.ULISBOA.PT

Tom Lenaerts

*Machine Learning Group, Université Libre de Bruxelles,
Boulevard du Triomphe CP212, 1050 Brussels, Belgium
§ Artificial Intelligence Lab, Vrije Universiteit Brussel,
Pleinlaan 2, 1050 Brussels, Belgium*

TOM.LENAERTS@ULB.AC.BE

AI Supremacy Race (AISR) model



Each advancement requires **a choice**: play **SAFE** / play **UNSAFE**

SAFE is costly (c) and takes more time than **UNSAFE**

s=1

s>1

Simplifying assumptions:

- Each advancement requires the same time
- Each advancement generates the same benefit **b**

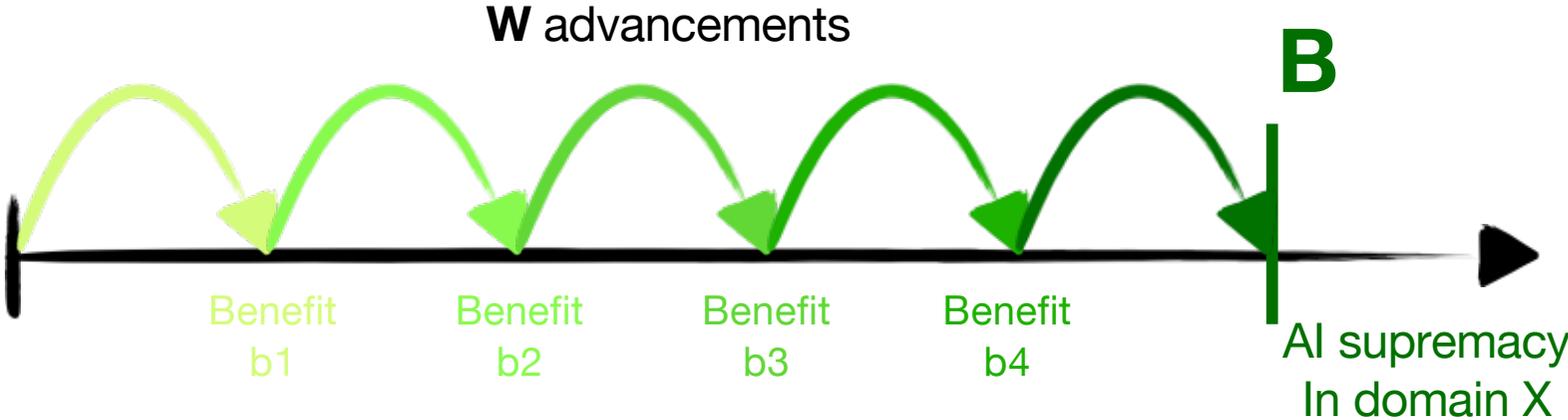
Payoffs in one round

	 SAFE	 UNSAFE
 SAFE	$-c + b/2$ 	$-c + b/(s+1)$ 
 UNSAFE	$bs/(s+1)$ 	$b/2$ 

We have also expanded this to $N > 2$ players

Assume simultaneous moves for each advancement

Strategies in AISR



A
S
Always safe



A
U
Always unsafe

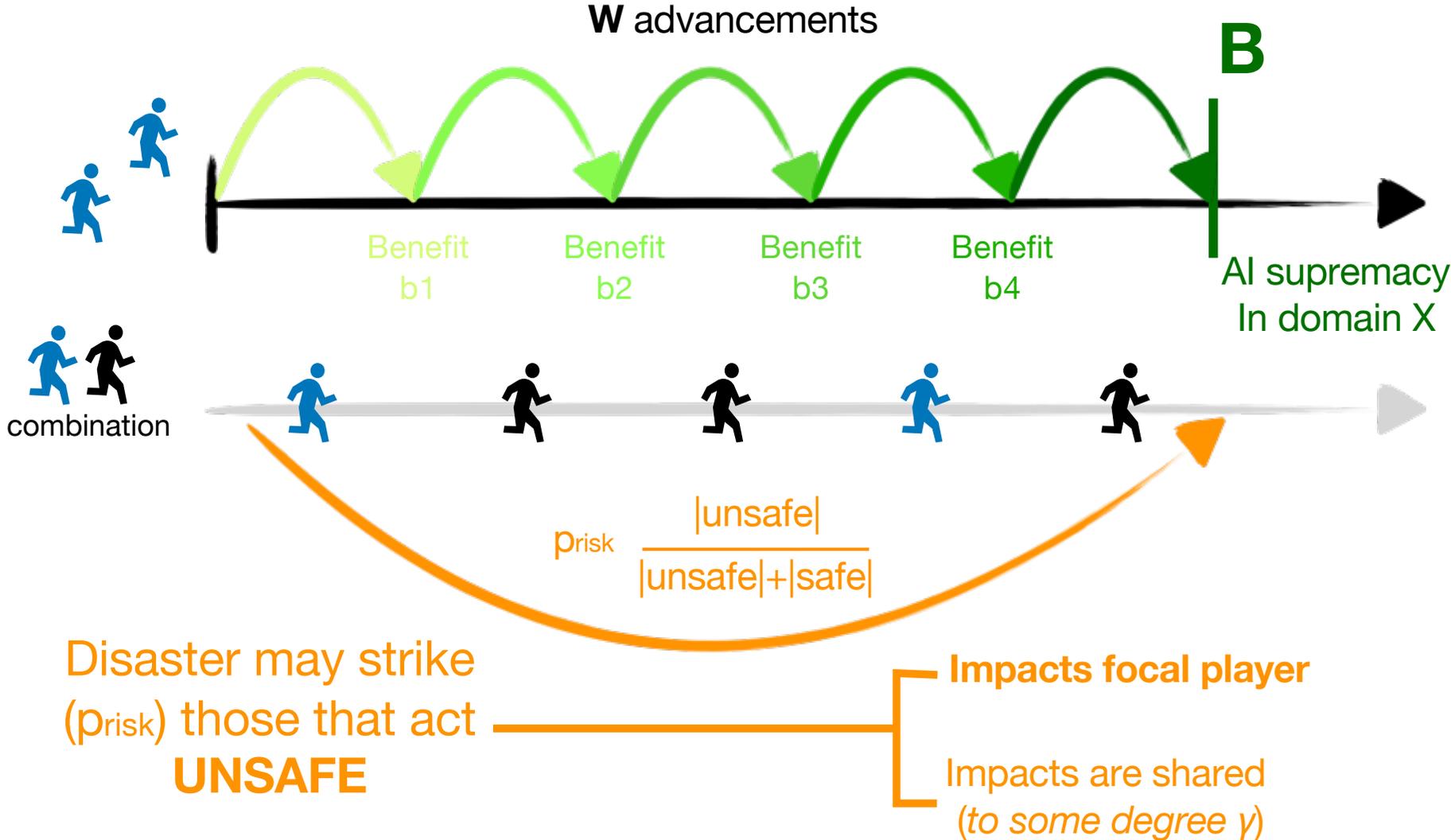


Risk of disaster

A
S
A
U
combination



Strategies in AISR

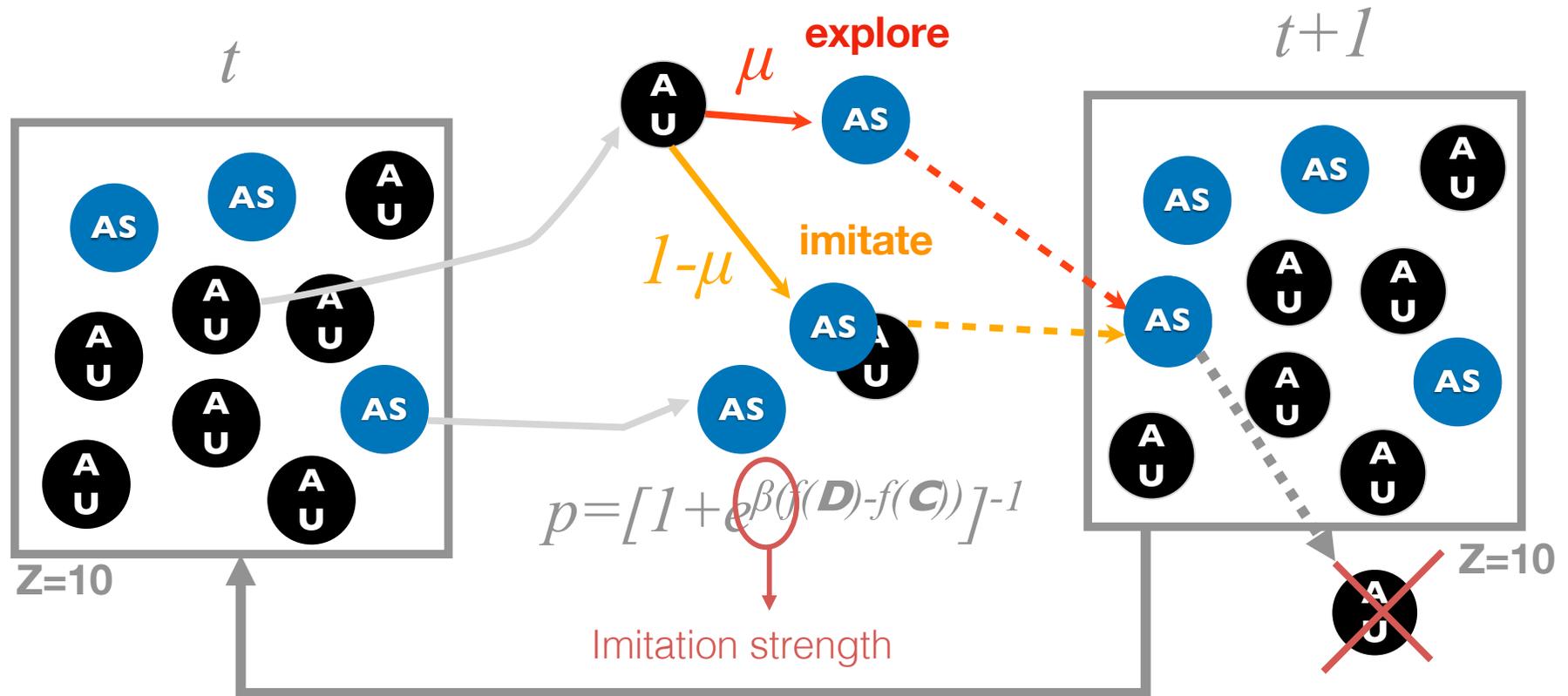


Average Gains
(fitness / social success)

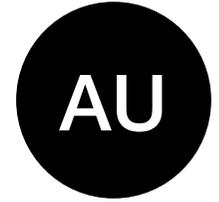
	AU	AS
AU	$(1 - p_{risk})$ [ + $sB/(2W)$]	$(1 - p_{risk})$ [ + (sB/W)]
AS		 + $B/(2W)$

Evolutionary Game Theory

Defining the social learning process

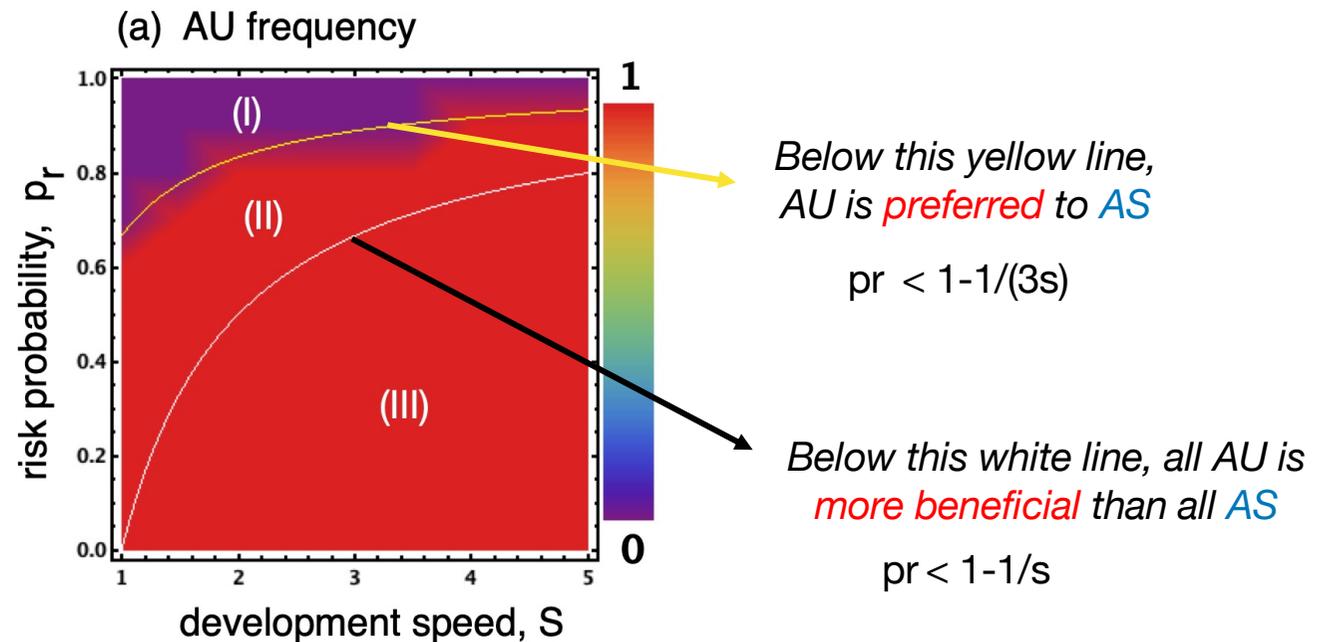


Short-term AI (small W)



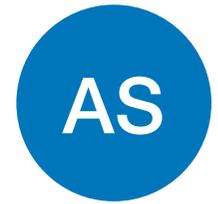
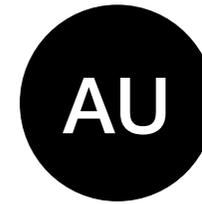
Q1: What behaviour dominates?

Q2: Which parameters influence dominance?

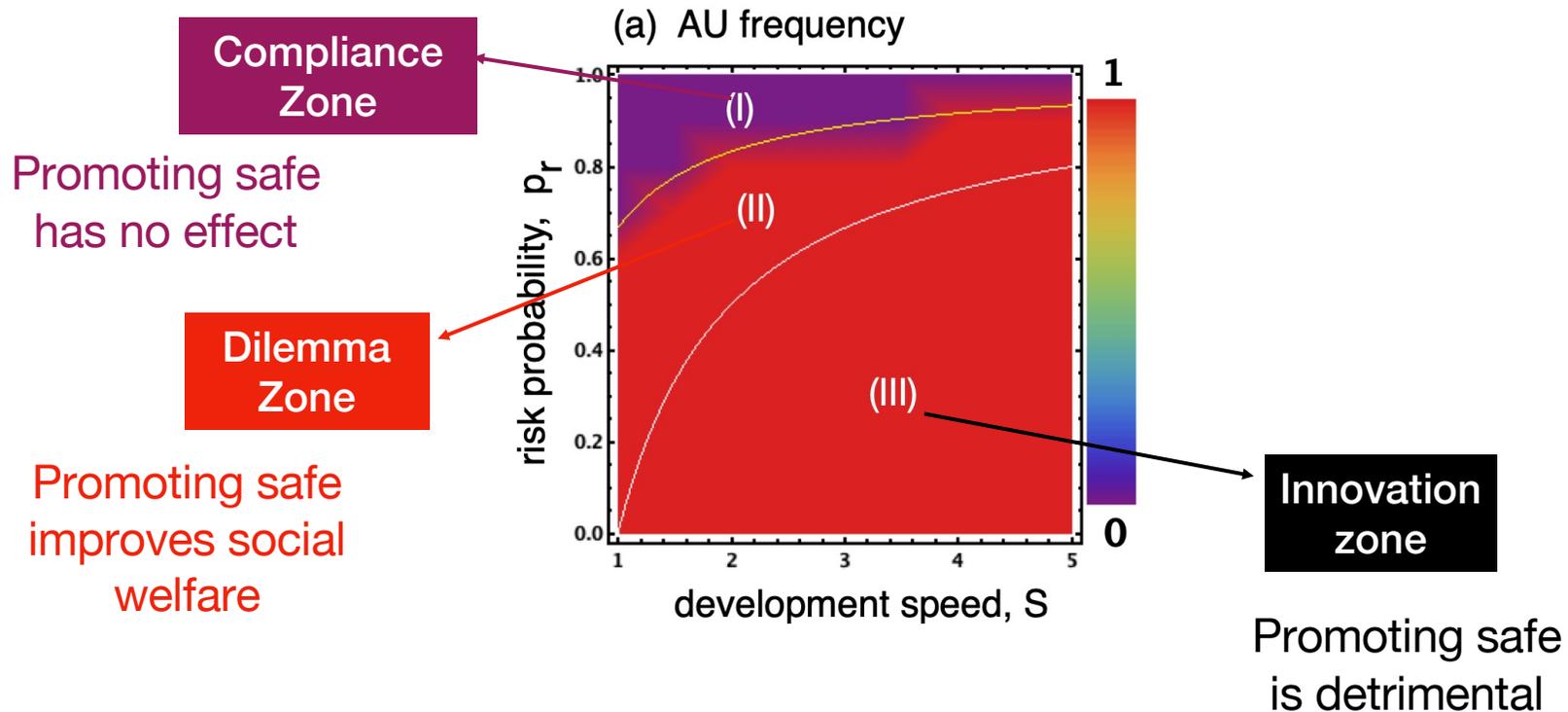


$c=1, b=4, W=100, B=10000, Beta=0.1, Z=100$

Short-term AI requires regulation of unsafe development



Q3: When are regulatory actions required? When not?

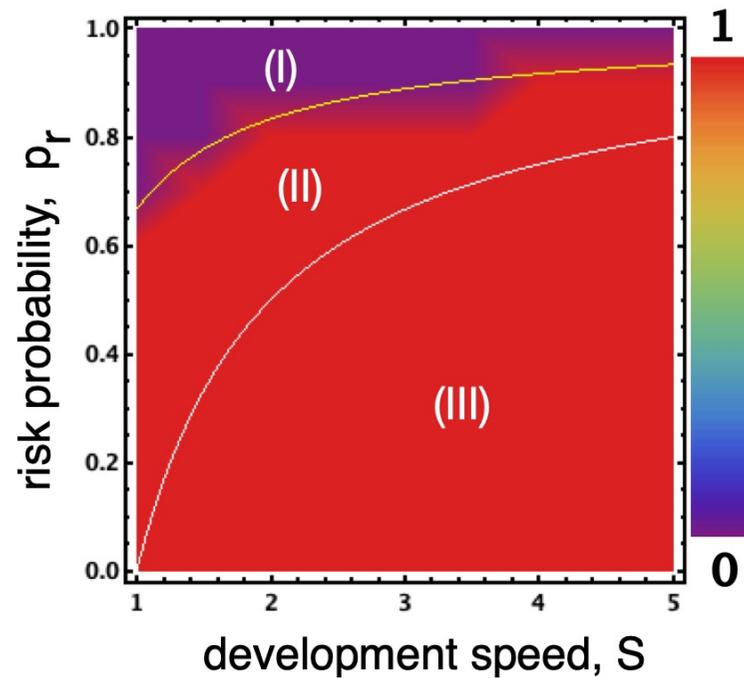


$c=1, b=4, W=100, B=10000, \text{Beta}=0.1, Z=100$

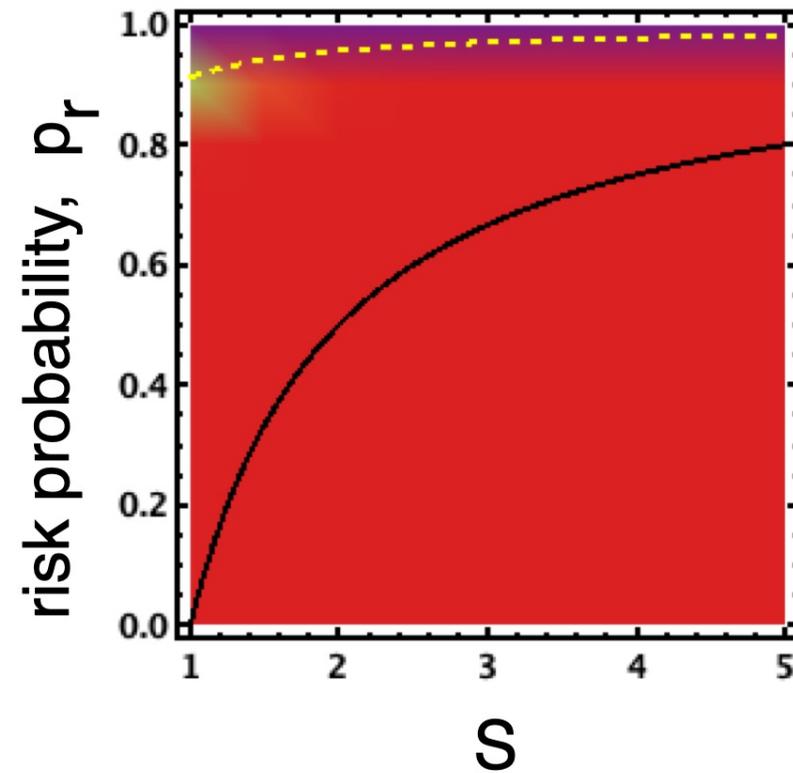
Bigger groups lead to larger dilemma zones

N = 2

(a) AU frequency



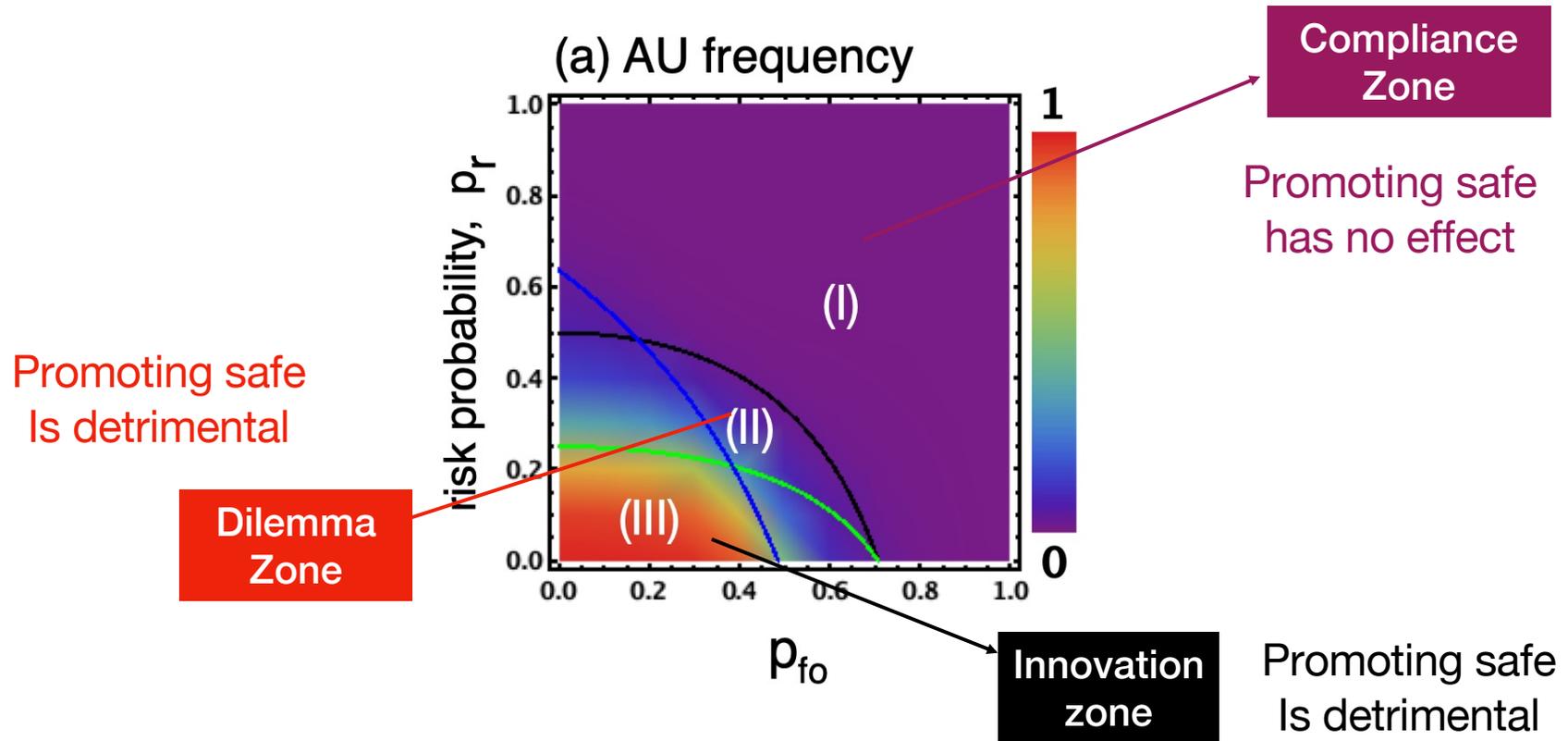
N = 5



$c=1, b=4, W=100, B=10000, \text{Beta}=0.1, Z=100$

Long-term AI (large W)

**Long-term AI requires
promotion of risk-taking!**



$W = 100000$, $c=1$, $b=4$, $s=1.5$, $B=10000$,
 $Beta=0.1$, $Z=100$

DUDE, WHERE'S MY **DATA?**

FEATURING
KIP AND GARY



BY DIANE ALBER

Question 2: What do we know about the impact of regulatory actions on such race dynamics?

Question 3: Is regulation always useful?

PLOS ONE

RESEARCH ARTICLE

Mediating artificial intelligence developments through negative and positive incentives

The Anh Han ^{1*}, Luís Moniz Pereira ², Tom Lenaerts^{3,4}, Francisco C. Santos ⁵

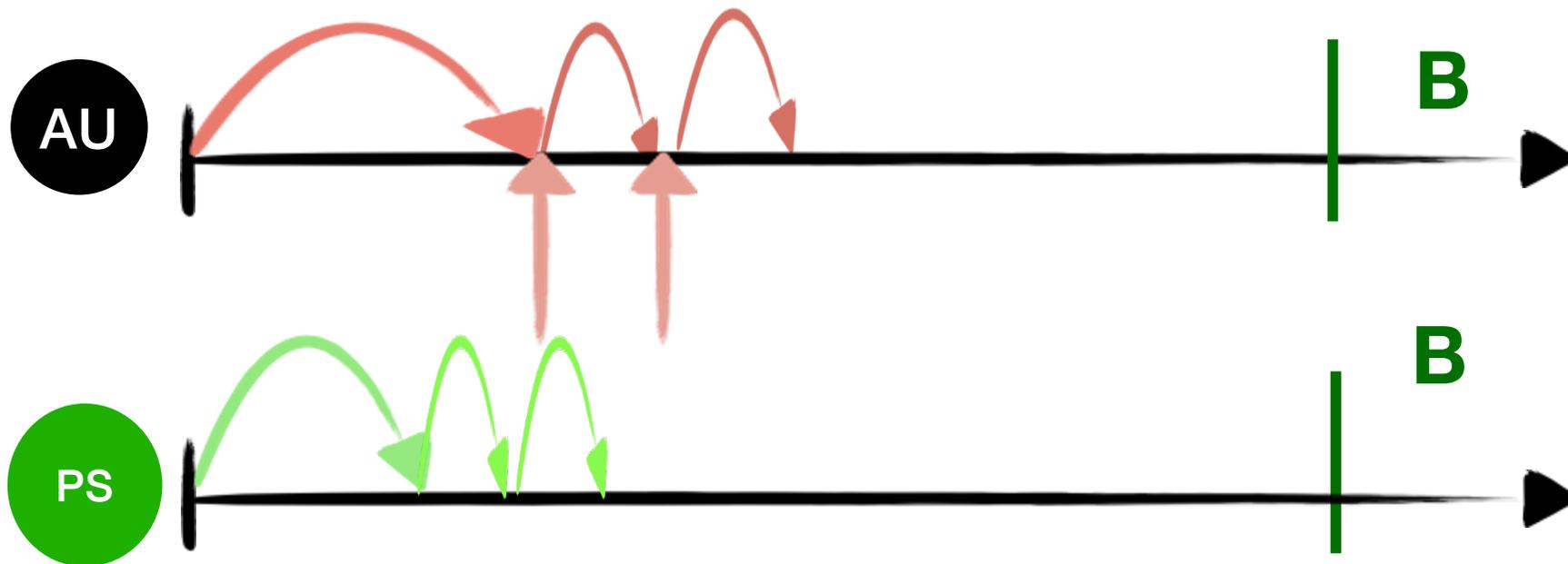
Citation: Han TA, Pereira LM, Lenaerts T, Santos FC (2021) Mediating artificial intelligence developments through negative and positive incentives. PLoS ONE 16(1): e0244592. <https://doi.org/10.1371/journal.pone.0244592>

Abstract

The field of Artificial Intelligence (AI) is going through a period of great expectations, introducing a certain level of anxiety in research, business and also policy. This anxiety is further energised by an AI race narrative that makes people believe they might be missing out. Whether real or not, a belief in this narrative may be detrimental as some stake-holders will feel obliged to cut corners on safety precautions, or ignore societal consequences just to “win”. Starting from a baseline model that describes a broad class of technology races where winners draw a significant benefit compared to others (such as AI advances, patent race, pharmaceutical technologies), we investigate here how positive (rewards) and negative (punishments) incentives may beneficially influence the outcomes. We uncover conditions in which punishment is either capable of reducing the development speed of unsafe participants or has the capacity to reduce innovation through over-regulation. Alternatively, we show that, in several scenarios, rewarding those that follow safety measures may increase the development speed while ensuring safe choices. Moreover, in the latter regimes, rewards do not suffer from the issue of over-regulation as is the case for punishment. Overall, our findings provide valuable insights into the nature and kinds of regulatory actions most suitable to improve safety compliance in the contexts of both smooth and sudden technological shifts.

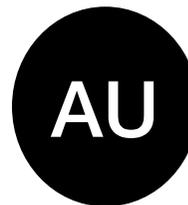
Negative Incentives in AISR

Peer Punishment (PS): reduces speed of an UNSAFE opponent, whereas PS's speed is reduced

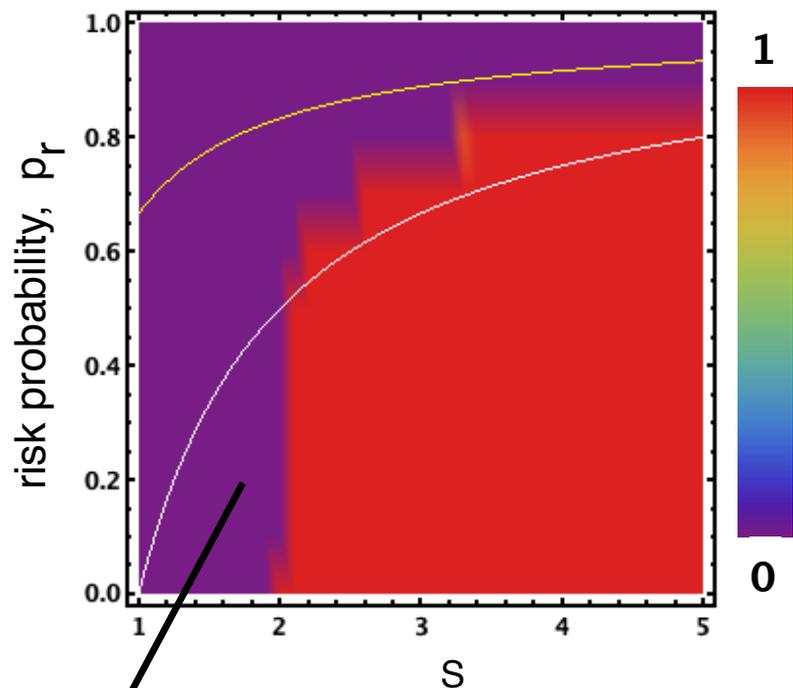


Institutional Punishment: an institution, not part of the system, reduces speed of UNSAFE players

Always regulate unsafe players leads to over-regulation

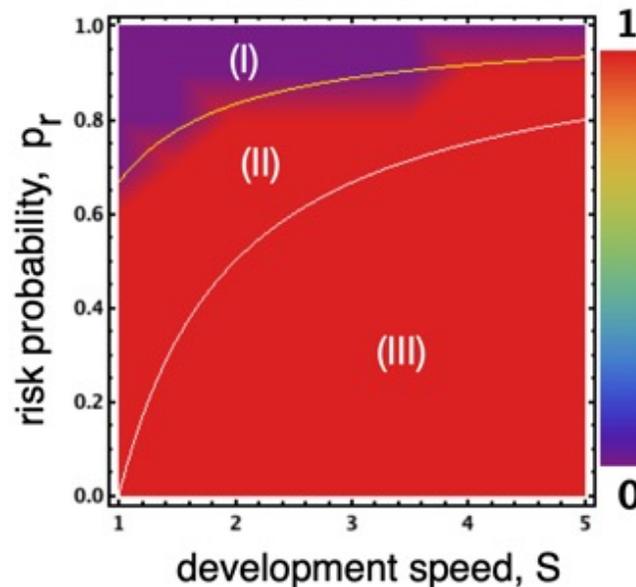


Effect of incentive
 $s' = 1$



Over-regulation

(a) AU frequency



Short-term AI

$c=1, b=4, W=100, B=10000, \text{Beta}=0.1, Z=100$

Conclusions

We describe a **plausible model** that can be useful when thinking about AI governance policies and regulations

Time-scale to reach AI supremacy strongly influences regulation

Incentives can improve safety outcome in the dilemma zones but can lead to **over-regulation**

Han et al. *To Regulate or Not: A Social Dynamics Analysis of an Idealised AI Race*. Journal of Artificial Intelligence Research, 69: 881-921, 2020

Han et al. *Mediating Artificial Intelligence development through positive and negative incentives*. PLoS ONE 16(1): e0244592, 2021.

Talk and poster presentation - Journal track @IJCAI

#J21 Time-scale Differences will Influence the Regulation Required in an Idealised AI Race Game

*The Anh Han (Teesside University), Luis Moniz Pereira (Universidade Nova de Lisboa), Francisco Santos (INESC
Lenaerts (Université Libre de Bruxelles)*

Video streams:

Aug 23th at 13:30 EDT – Room Red 1 **and**

Aug 26th at 10:00 EDT – Room Green 2

Posters:

Aug 23th at 14:00 EDT – Room Red **and**

Aug 26th at 10:30 EDT – Room Green

Home Who We Are Activities Existential Risk Get Involved Contact

future of life INSTITUTE

Technology is giving life the potential to flourish like never before... ..or to self-destruct. Let's make a difference!

THIS HOLIDAY SEASON, GIVE TO THE FUTURE. CLICK HERE TO DONATE.

AI Biotech Nuclear Climate Podcasts

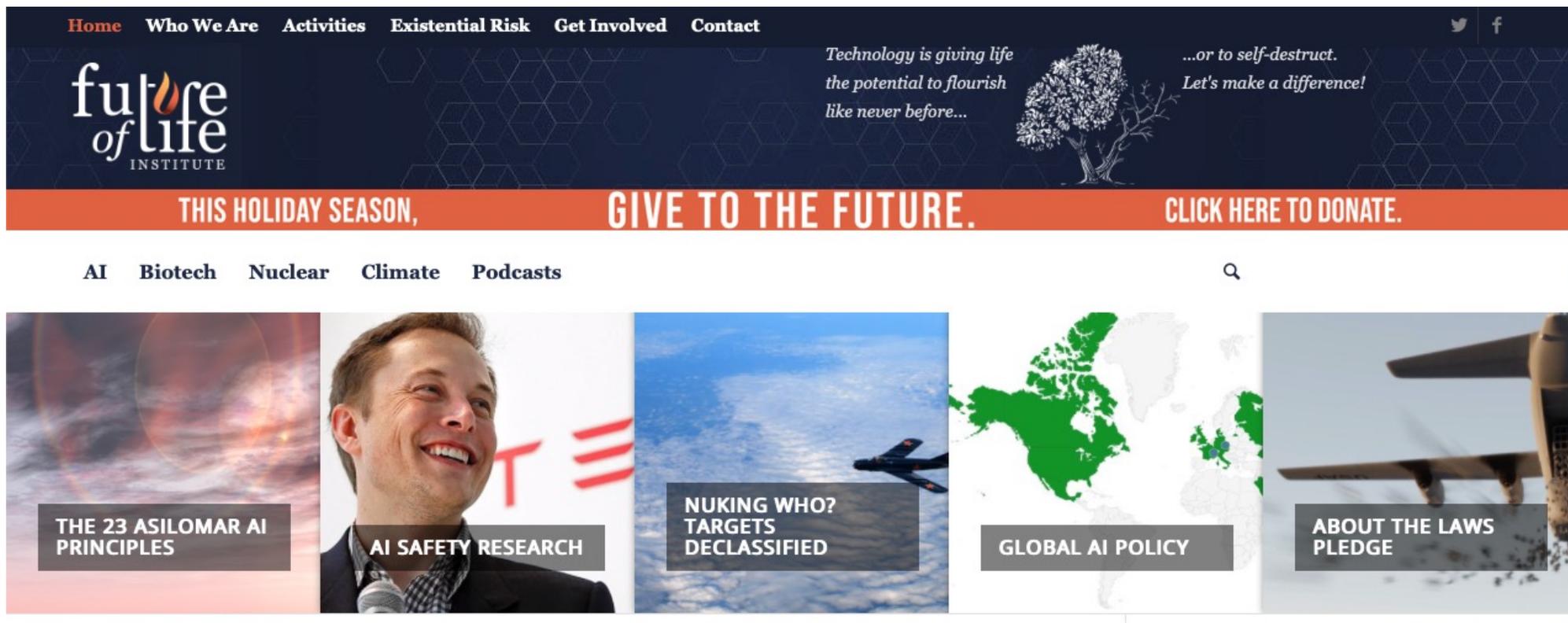
THE 23 ASILOMAR AI PRINCIPLES

AI SAFETY RESEARCH

NUKING WHO? TARGETS DECLASSIFIED

GLOBAL AI POLICY

ABOUT THE LAWS PLEDGE



“Incentives for Safety Agreement Compliance in AI Race ” (2019-2021)

<http://futureoflife.org/2018-ai-grant-recipients#Han>