

Can Media Act as a Soft Regulator of Safe AI Development? A Game Theoretical Analysis

Henrique Correia da Fonseca^{1,*}, António Fernandes¹, Zhao Song², Theodor Cimpanu³,
Nataliya Balabanova⁴, Adeela Bashir², Paolo Bova², Alessio Buscemi⁵, Alessandro Di Stefano²,
Manh Hong Duong⁴, Elias Fernandez Domingos^{6,7}, Ndidi Bianca Ogbo², Simon T. Powers⁸,
Daniele Proverbio⁹, Zia Ush Shamszaman², Fernando P. Santos¹⁰, The Anh Han², Marcus Krellner³

¹ INESC-ID and Instituto Superior Técnico, Universidade de Lisboa

² School of Computing, Engineering and Digital Technologies, Teesside University

³ Biological and Environmental Sciences, University of Stirling

⁴ School of Mathematics, University of Birmingham

⁵ Luxembourg Institute of Science and Technology

⁶ Machine Learning Group, Université libre de Bruxelles ⁷ AI Lab, Vrije Universiteit Brussel

⁸ Division of Computing Science and Mathematics, University of Stirling

⁹ Department of Industrial Engineering, University of Trento

¹⁰ University of Amsterdam

* Corresponding author: henrique.c.fonseca@tecnico.ulisboa.pt

Abstract

When developers of artificial intelligence (AI) products need to decide between profit and safety for the users, they likely choose profit. Untrustworthy AI technology must come packaged with tangible negative consequences. Here, we envisage those consequences as the loss of reputation caused by media coverage of their misdeeds, disseminated to the public. We explore whether media coverage has the potential to push AI creators into the production of safe products, enabling widespread adoption of AI technology. We created artificial populations of self-interested creators and users and studied them through the lens of evolutionary game theory. Our results reveal that media is indeed able to foster cooperation between creators and users, but not always. Cooperation does not evolve if the quality of the information provided by the media is not reliable enough, or if the costs of either accessing media or ensuring safety are too high. By shaping public perception and holding developers accountable, media emerges as a powerful soft regulator – guiding AI safety even in the absence of formal government oversight.

Data/Code available at:

[https://anonymous.4open.science/r/
media-AI-governance-0752](https://anonymous.4open.science/r/media-AI-governance-0752)

Introduction

In May 2024, Google introduced a feature that used Artificial Intelligence (AI) to give short answers to search prompts. Shortly thereafter, a tweet went viral (McMahon and Kleinman, 2024), showing how the prompt “cheese not sticking to pizza” produced an answer containing the following suggestion: “You can also add about 1/8 cup of non-toxic glue to the sauce to give it more tackiness.” (Kelly, 2024). While this particular suggestion was so obviously

ridiculous that it hopefully caused mostly laughs rather than severe medical consequences, it highlighted a substantial issue of AI technology – can we trust it, and is it safe?

To ensure AI safety, governments try to work towards effective regulations, such as the European Union with their “Ethics guidelines for trustworthy AI” (Commission et al., 2019). Governmental intervention traditionally plays a role in ensuring that consumers are protected, and recently Evolutionary Game Theory (EGT) (Sigmund, 2010; Hofbauer and Sigmund, 1998) models have shown how regulation of AI technology can incentivise the safe adoption of AI (Han et al., 2019, 2020; Alalawi et al., 2026; Cimpanu et al., 2022; Bova et al., 2024). However, these models also showed some drawbacks, namely how over-regulation can prevent the adoption of valuable AI technology (Han et al., 2022).

The anecdote from the beginning shows another path: a form of regulation by the media. The story of the pizza glue was first spread via social media, then by traditional media outlets, reaching many people and making them aware of the problem. In addition, the media backlash caused Google to adjust and improve its feature. Media has played a role in the safety of other products as well. For example, the Guardian reported that Apple contractors had been listening to confidential Siri recordings, some of which contained highly sensitive private information (Hern, 2019). The resulting public backlash pressured Apple to apologise and adopt an opt-in system for reviewing recordings. Media has even influenced the inner workings of companies. For example, a Vox investigation uncovered contractual clauses that prevented departing OpenAI employees from making disparaging remarks about the company’s practices; following the media scrutiny,

these clauses were removed (Piper, 2024). Concomitantly, surveys have shown that the media plays an important role in the perception of consumer risks (Cao and Li, 2020; Zhang et al., 2022), which is also true for AI (Yang et al., 2023). On the other hand, some research has also highlighted that an overly optimistic media coverage can be linked with increased adoption of unsafe products (Melero-Bolaños et al., 2025). The influence of the media might not always lead to the desired effects.

Artificial life researchers have long been fascinated by how complex systems self-organise, adapt, and maintain stability through decentralised interactions (Bedau, 2003; Powers et al., 2018; Krellner and Han, 2021; Sayama, 2015; Gershenson et al., 2020). One of the most powerful mechanisms that enable such coordination is indirect reciprocity (IR) (Boyd and Richerson, 1989; Nowak and Sigmund, 2005), whereby agents cooperate based not on direct experience, but rather on reputational cues. While this principle has been studied extensively in biological systems, its implications for socio-technological systems and safe AI remain envisioned (Paiva et al., 2018), but underexplored (Xia et al., 2023; Jøsang et al., 2007; Nowak and Sigmund, 1998; Hammond et al., 2025). Moreover, how the reputational information that is crucial to indirect reciprocity is transmitted and by whom remains an open challenge to the study of reputation-based cooperation (Hilbe et al., 2018; Santos et al., 2018; Sommerfeld et al., 2007).

In this work, we leverage the lens of artificial life to investigate whether media – acting as a decentralised, stochastic purveyor of reputational information – can serve as a soft regulator in the development and adoption of artificial intelligence. Much like costly signalling systems in nature, media commentary provides noisy but influential feedback that shapes behaviour, not through coercion, but through perception. We model these dynamics as an evolving game between creators, users, and media agents, exploring how reputation-based safety might emerge in the real-world challenge of AI governance. In doing so, we bridge the study of artificial life with societal evolution in the context of technological norms of adoption.

EGT models have already indicated that media – also referred to as the commentariat in prior work – can foster safe AI adoption in the context of governmental regulation (Balabanova et al., 2025). We therefore seek to study in particular whether media alone could enable similar results. By this, we mean that users choose to adopt AI technology and that AI creators choose to develop safe AI products. The media acts by flagging creators as safe or unsafe, hence enabling users to make an informed decision. However, information from the media can never be perfect, nor does it come for free. Taking all these factors into account, we created a model to predict the behaviour of users and creators in the presence of two distinct media outlets, differing in their cost and accuracy. In principle, our model could be

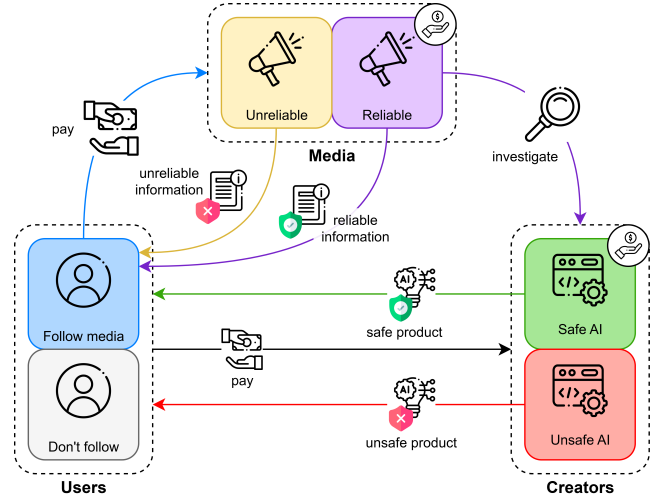


Figure 1: **Visual description of the AI regulatory ecosystem.** Users decide whether or not to use AI products, incurring a cost for adoption. For this decision process, they can choose to follow media recommendations by paying a small amount in exchange for their information about creators’ strategies. Meanwhile, creators decide whether to create safe or unsafe technology; safe technology further involves additional costs. Media can monitor the creators’ decision-making process, provided that commentators invest in obtaining more reliable information.

seen as a general model of the role of investigative journalism as a soft regulator, without being specific to any one industry. However, unlike other high-risk industries, AI currently lacks widespread formal regulation. It is this case – high-risk industry lacking effective formal regulation – that our model directly addresses. We use this model to address the pressing question of to what extent media could provide incentives that fill the AI regulatory gap, while formal regulations are still waiting to be developed and enforced.

Model and Methods

We use a two-population model consisting of N_C creators and N_U users. Users and creators evolve by updating their strategies over time. Additionally, we consider two different types of media commentators – reliable and unreliable. Media acts as surveillance agents that monitor creators’ strategies regarding AI safety practices and relay this information to users. They do not represent an evolving third population, but rather represent entities that users can choose between to get information from.

The game is played in consecutive rounds where one user meets one creator. The creator can choose to defect by providing an unsafe AI product or to cooperate by creating a safe one. Their cooperation is better for the user, but entails additional costs for the creator ($c_c > 0$). This cost can rep-

resent increased development or production demands, and also increased effort to meet regulatory or voluntary safety requirements. Users decide whether to cooperate by adopting this AI product or to defect by refusing to do so. Adopting always grants a benefit to the creator ($b_c > 0$), but users only gain from adopting a safe product ($b_u > 0$) and instead lose by adopting an unsafe one ($c_u > b_u$) (see Table 1 for reference of all parameters).

Both users and creators can apply unconditional strategies: always cooperate and always defect. For the creators, we refer to these as C and D , whereas for the users as $AllC$ and $AllD$. For the latter, we introduce additional strategies that are based on the recommendation of a media source. We implement two types of media that they can follow, a good one ($GMedia$) and a bad one ($BMedia$). The good media performs thorough investigations and consequently has a chance q to identify the strategy of a creator correctly (and $1 - q$ to identify them incorrectly). Relying on the good media will cost the user ($c_i > 0$). This represents the cost the media incurs for investigating the AI creator, which it passes on to the user in one way or another. Bad media does not perform thorough investigations of creators and so comes with no costs (see Figure 1). But its chance is set to $q = 0.5$, meaning that it gives random recommendations (see Figure 2).

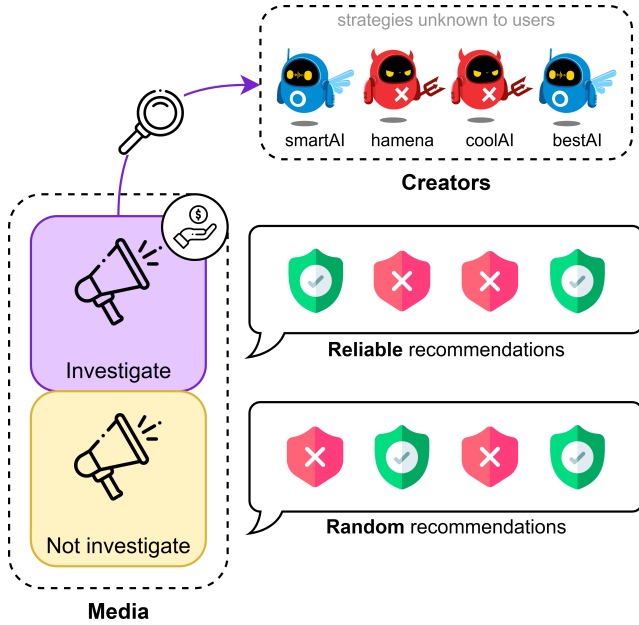


Figure 2: **Media recommendations.** Media that investigate creators’ strategies can provide more reliable recommendations, although incurring an extra cost. Media that do not investigate will provide random recommendations to users.

The resulting payoffs are shown in Table 2. Because the media gives probabilistic recommendations, the user be-

Table 1: Description of the key parameters of the model.

Parameter	Explanation
b_u	benefit a user receives when adopting a safe technology
c_u	cost a user incurs when adopting an unsafe technology
b_c	benefit a creator receives when their technology is adopted
c_c	additional cost of creating safe AI (the cost of creating unsafe AI is normalised to 0)
c_i	cost of informed recommendation (for user or commentator)
q	probability that the recommendation of a commentator is <i>correct</i>

haviour is probabilistic as well, and the values are averaged for each possible interaction. These payoffs will be used for the replicator dynamics (see below) to determine how the strategies evolve in infinitely large populations. To supplement and verify these results, we also run agent-based simulations of finite populations.

Table 2: **Payoff matrix of our game theoretical model involving users and creators.** Columns Π^U and Π^C represent the payoffs for users and creators, respectively, for a given combination of strategies. The *Unsafe* and *Safe* terms represent the defect and cooperate strategies from creators, respectively. *AllD* and *AllC* represent user strategies that always and never cooperate with creators, respectively, while *BMedia* and *GMedia* stand for deferring the decision to the recommendation of a low-quality and a higher-quality media outlet, respectively.

Creator	User	Π^U	Π^C
<i>Unsafe</i>	<i>AllD</i>	0	0
	<i>AllC</i>	$-c_u$	b_c
	<i>BMedia</i>	$-0.5c_u$	$0.5b_c$
	<i>GMedia</i>	$-(1-q)c_u - c_i$	$(1-q)b_c$
<i>Safe</i>	<i>AllD</i>	0	$-c_c$
	<i>AllC</i>	b_u	$b_c - c_c$
	<i>BMedia</i>	$0.5b_u$	$0.5b_c - c_c$
	<i>GMedia</i>	$qb_u - c_i$	$qb_c - c_c$

Replicator Dynamics

We first introduce the so-called replicator dynamics. This is a widely used model in evolutionary game theory to express how the frequencies of strategies in a population that evolves over time (Hofbauer and Sigmund, 1998). It is based on the idea that the proportion of agents of a given strategy increases when the strategy achieves payoffs higher than

the average payoff of the population $\bar{\pi}$, and decreases when achieving expected payoffs lower than that average payoff. We use x_1, x_2, x_3 , and x_4 to denote the frequencies of *AllD*, *BMedia*, *GMedia*, and *AllC* ($x_4 = 1 - x_1 - x_2 - x_3$); and y and $1 - y$ to denote the frequencies of *C* and *D*. Formally, the replicator dynamics is given by a system of differential equations

$$\begin{aligned}\dot{x}_1 &= x_1(1 - x_1) [\pi_{AllD} - \bar{\pi}_{user}], \\ \dot{x}_2 &= x_2(1 - x_2) [\pi_{BMedia} - \bar{\pi}_{user}], \\ \dot{x}_3 &= x_3(1 - x_3) [\pi_{GMedia} - \bar{\pi}_{user}], \\ \dot{y} &= y(1 - y) [\pi_C - \bar{\pi}_{creator}],\end{aligned}\quad (1)$$

where the average payoff of the user population is $\bar{\pi}_{user} = x_1\pi_{AllD} + x_2\pi_{BMedia} + x_3\pi_{GMedia} + x_4\pi_{AllC}$, and of the creator population is $\bar{\pi}_{creator} = y\pi_C + (1 - y)\pi_D$. Additionally, based on the Table 2, the expected payoffs for each strategy are:

$$\begin{aligned}\pi_{AllD} &= 0, \\ \pi_{BMedia} &= 0.5b_u y - 0.5c_u(1 - y), \\ \pi_{GMedia} &= (qb_u - c_I)y - ((1 - q)c_u + c_I)(1 - y), \\ \pi_{AllC} &= b_u y - c_u(1 - y), \\ \pi_C &= (-c_c)x_1 + (0.5b_c - c_c)x_2 + (qb_c - c_c)x_3 + (b_c - c_c)x_4, \\ \pi_D &= 0.5b_c x_2 + (1 - q)b_c x_3 + b_c x_4.\end{aligned}\quad (2)$$

With these equations, we can study how the population changes over time. By averaging the frequencies of the strategies over time, we can determine how frequent they are, and consequently, how much cooperation is shown by the users and the creators.

Setting the right-hand side of Equation (1) to 0 yields the potential equilibrium states, which are points where the population composition remains static. Any homogeneous population composition (e.g. all users are *AllD* and all creators are *D*) might be stable. In detail, we analyse these equilibria using Lyapunov's indirect method to assess the local stability, by examining the eigenvalues of the Jacobian matrix evaluated at these points (Khalil and Grizzle, 2002). For non-hyperbolic cases, we utilise centre manifold theory to reduce the system's dimensionality and facilitate stability analysis (Carr, 2012). In addition, we are also interested in the evolution when it diverges from equilibria. To see this, we analyse the stability of each equilibrium.

Agent-Based Simulations

Next, we introduce the methods for finite populations that rely on agent-based simulations. Such methods are well-known in the literature, originating from statistical physics (Monte Carlo simulations) (Perc et al., 2017). We consider two well-mixed populations: users and creators, with population sizes N_U and N_C , respectively. In the beginning, each

creator and each user is assigned a random strategy from the set of $\{D, C\}$ and $\{AllD, BMedia, GMedia, AllC\}$, respectively. They then undergo several evolutionary time steps in which they might change their strategies as described below.

In each evolutionary step, a user and a creator are randomly selected to update their strategy (one after the other). We will describe the process for a generic player (user or creator), since the process is essentially the same. With probability μ , the player updates their strategy through mutation, randomly exploring a novel (different) strategy. Note that users and creators may have different mutation probabilities μ_u and μ_c . With complementary probability $1 - \mu$, the focal player performs a Monte Carlo evolutionary step, whereby a second player is randomly selected (from the same population) for comparison of their respective payoffs. To determine these payoffs, both players engage in a number of games with randomly selected agents from the antipodal population equal to that population's size, where they accumulate payoffs, as detailed in Table 2.

After accumulating payoffs, the originally randomly selected player can update their strategy through stochastic pairwise comparison of their current fitness, as is typically done in EGT models of finite populations (Traulsen et al., 2006). Player i thus adopts the strategy of the second player j with a probability given by the Fermi function

$$p_{i \rightarrow j} = (1 + e^{-\beta(\bar{\pi}_j - \bar{\pi}_i)})^{-1}, \quad (3)$$

where $\bar{\pi}_j$ and $\bar{\pi}_i$ represent the average payoff of players j, i . The selection strength β ranges from 0, which would mean random imitation, to infinity, where Equation 3 essentially becomes purely deterministic. Throughout this work, we use $\beta = 1$ to ensure a strong selection strength, as is typically done in literature (Santos et al., 2021b; Okada, 2020), except where explicitly stated otherwise. For average results, we repeat each simulation $R = 100$ times. A generation corresponds to $N_U + N_C$ discrete time steps, where a strategy update may occur, allowing for, on average, all members of both populations to evolve. Each simulation runs for G generations.

Results

The goal of our investigation is to test whether media can cause both users and creators to cooperate, resulting in the adoption of safe AI technology. We try to answer this question in different ways. We first study the behaviour of the evolving populations over time. For this, we use the replicator dynamics to run numerical simulations, which we subsequently compare with simulations of finite populations using our agent-based model. We also use the replicator dynamics to analytically study system equilibria in infinite populations.

To study the evolution of strategies over time, we initialise the population with an equal distribution of all strate-

gies (25% for each user strategy and 50% for each creator strategy). We study the evolution of these populations in different settings, comprehensively exploring the parameter space. Specifically, we focused on varying: **(i)** the quality of media predictions (q); **(ii)** the cost of consuming good media for the users (c_i); and **(iii)** the surplus costs which must be covered by safe AI creators while remaining competitive (c_c).

We are especially interested in one metric – the average cooperation rate η , as it captures cooperation in both populations. This can be computed using the frequencies of the strategies and their average behaviour (compare to Table 2). For example, a BMedia user (x_2) cooperates if they receive information that the creator is trustworthy, which is equivalent to a coin flip. A GMedia user, on the other hand, cooperates in one of two cases: if they meet a cooperating creator (y) and get a true signal (q), or when they encounter a defecting creator ($1 - q$) and receive a wrong signal ($1 - q$). Unconditional cooperators, be it users (x_4) or creators (y), naturally always cooperate. The average cooperation rate is thus

$$\eta = \frac{y + 0.5x_2 + (qy + (1 - q)(1 - y))x_3 + x_4}{2}. \quad (4)$$

The result of the numerical simulations is shown in Figure 3. They show the relationship between cooperation and the parameters of interest. Crucially, we observe that all three parameters have the capacity to completely collapse cooperation but also to achieve high levels of cooperation.

Taken in pairwise isolation, we show thresholds beyond which creators always resort to unsafe AI development. Firstly, media must always tread a fine balance: not only must it maintain sufficient quality q of its reports on creator behaviour, but this threshold gets stricter as the cost of investigating creators increases (Figure 3 first panel). Similarly, creator costs react in much the same fashion, imposing a certain amount of strictness on the existing commentariat (Figure 3 second panel). In other words, when either of the costs is prohibitive (either to creators c_c , or to media c_i), then media providers must provide accurate information to compensate, else the advantage they gain from user trust is overcome by lazy media. Secondly, we show a more intuitive relationship between the two costs (c_c and c_i), as either of them can lead to a breakdown of cooperation if excessive. There is a slightly more forgiving nature to costly investigation ($c_i \lesssim 0.35$) as opposed to the cost of safety ($c_c \lesssim 0.25$), but cooperation can be achieved if neither threshold is overstepped.

Following this approach, we tried to replicate these findings with agent-based simulation of finite populations, where strategy evolution occurs through mutation and social learning. Figure 4 shows similar insights to the ones obtained through replicator dynamics (Figure 3). Due to the stochastic nature of the finite models, transitions between

areas of cooperation and defection are less sharp, but follow the same patterns very closely. This provides robustness to our main finding that media supports cooperation between users and creators, given realistic quality levels, costs of investigative media and costs of safe development.

We then explore the stability of the infinite populations as described above. Our analysis reveals that although achieving a state of full cooperation within the user population remains an unattainable equilibrium, cooperation endures through an intriguing evolutionary dynamic characterised by the persistent oscillation between *GMedia* and *AllC*.

As illustrated in Figure 5(a-b), upon the initiation of the evolutionary process, *AllD* and *BMedia* rapidly decline within the user population, while *GMedia* and *AllC* exhibit sustained oscillatory behaviour. Similarly, in the creator population, the frequencies of *C* and *D* also display persistent oscillations. This dynamic interplay allows *AllC* to survive and evolve in concert with *GMedia*.

An in-depth examination of the equilibria yields a stable point in which *AllC* dominates in the user population and *C* dominates in the creator population. However, it is only stable if $c_c < 0$, $b_u > 0$, $c_i + b_u(1 - q) > 0$. We limited the parameters, such as c_c to be positive, because we wanted to explore realistic scenarios. This stable point has, therefore, no relevance for our overarching question.

However, we found the equilibrium point where *AllD* dominates in the user population and *D* dominates in the creator population. It is stable for the conditions $c_c > 0$, $c_u > 0$, $c_i + c_u(1 - q) > 0$. They always hold for the parameter space we consider realistic. However, the existence of a stable point for defection does not mean that cooperation is impossible (Imhof et al., 2005).

In our model, convergence to the stable defection point seems to depend on initial conditions. Even under the same parameter settings that produce oscillatory dynamics, different starting states can drive the system to collapse into this defection-dominated equilibrium. This sensitivity highlights a possible bistable regime in the model: the same parameters can yield either persistent cooperation (via oscillations) or universal defection, contingent on the initial states of the populations, see Figure 5(c-d).

We quantify which of these two outcomes, oscillation with high cooperation or collapse to full defection, is more relevant for the answer to our research question. We therefore used the replicator dynamics to study the evolution of the population from different starting states. We use the same numerical approach as described above and vary the frequencies by steps of 2%, which are equivalent to 1172286 different starting states. Of these, 38.6% reach the stable defection state, the others do not (parameters as in Figure 3). Across all starting states, the measured average cooperation rate is $\eta = 0.54$. We therefore consider both dynamics as relevant.

For our finite model, there exist no truly stable states be-

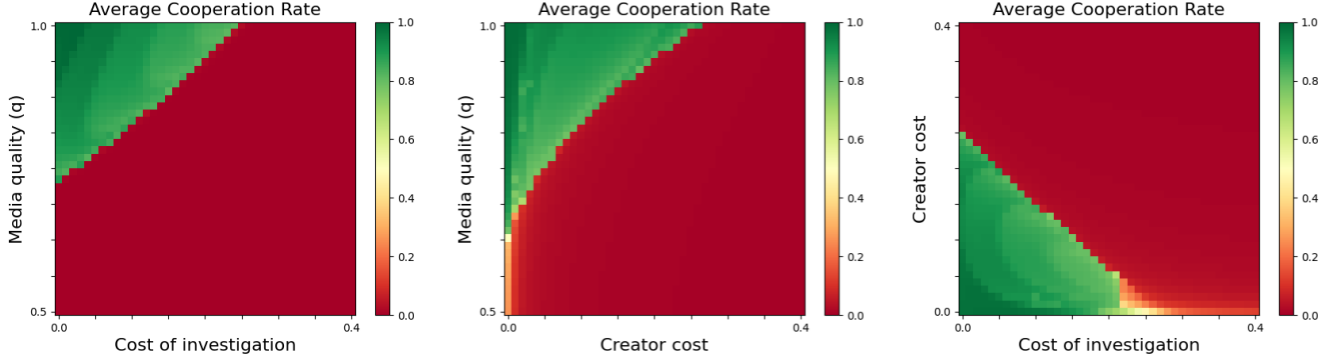


Figure 3: **Average cooperation ratio η via replicator dynamics**, across parameters of interest. If not varied, $q = 0.9$, $c_i = 0.1$, and $c_c = 0.1$, other parameters: $b_c = 0.4$, $b_u = 0.4$, $c_u = 0.8$.

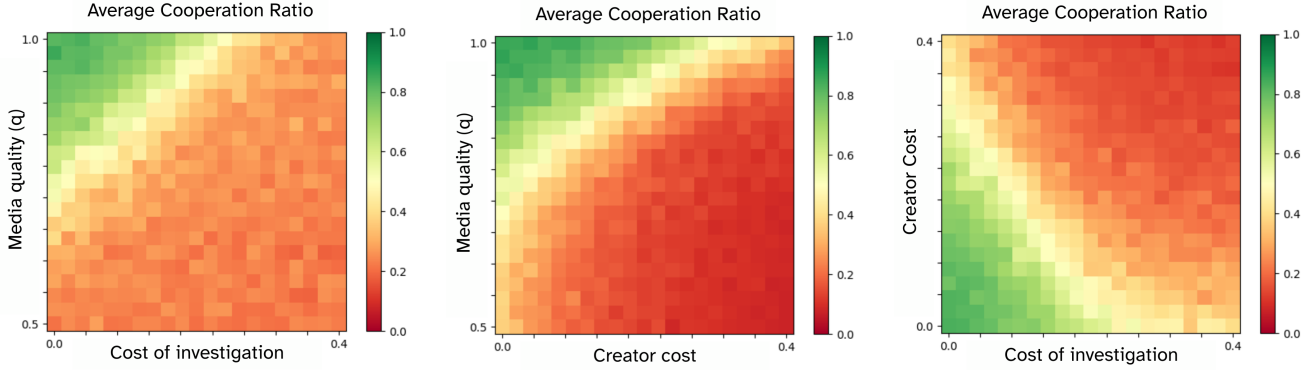


Figure 4: **Average cooperation ratio η via agent-based simulations**, across parameters of interest. Each data point shows the average cooperation ratio η averaged over $R = 100$ runs. All parameters are set to the same ones as Figure 3, with the additional evolutionary parameters of $G = 500$, $N_U = 100$, $N_C = 50$, $\beta_C = \beta_U = 1$, $\mu_u = 1/N_U$, $\mu_c = 1/N_C$. The metric η is only computed after allowing for a converging period of $G/10$ generations.

cause of mutation. However, it is interesting to compare their behaviour with the replicator dynamics. Indeed, even the stochastic agent-based model shows consistent oscillations of strategy frequency – see Figure 6 – for the same parameters of Figure 5. Cycles follow the same pattern: (i) *GMedia* users proliferate in face of the initial population state, leading to (ii) the invasion of *AllC* user strategies that take advantage of the highly cooperative creator population state, which in turn (iii) are easily exploited by an invading population of defective creators, (iv) counteracted by a non-zero population of *GMedia* users whose discriminating strategy stops the rise of defective creators, leading the system back to step (ii). Notably, users of *GMedia* do not dominate the population.

To see if the initial population composition has any significant effect on the simulations, we run several examples starting in the predicted stable defection state, averaging the frequencies of strategies over time (see Figure 7). We note that the simulated populations can consistently escape this

state and display high levels of cooperation.

All the results from the computational model are robust to variations in the number of generations, population sizes, and both selection strengths and mutation rates, as long as their order of magnitude remains above a certain threshold ($\beta \geq 1$ and $mu \geq 0.1$).

Discussion

AI safety is a concern, because – all things being equal – providing safe AI products is costly, giving creators an incentive to provide unsafe AI instead (Han et al., 2019; Cohen et al., 2024; Hammond et al., 2025; Bengio et al., 2025). In this paper, we show that media alone can act as a soft regulator of AI creators’ development behaviour, by providing users with information about which AI products are safe to adopt. In our model, we formulated populations with two types of media – good media, which pay investigative costs to produce a reliable signal, and bad media, which give random signals. Even if the signal by the good media is far from

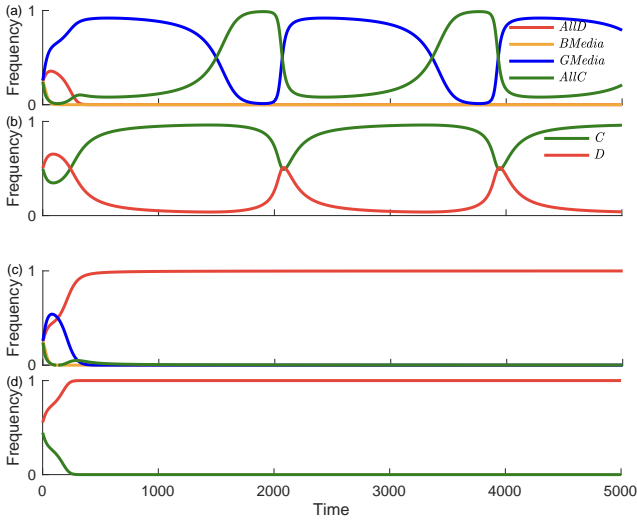


Figure 5: Evolution of strategies over time for replicator dynamics. In the top two panels (plots *a* and *b*), with initially 50% *C* in the creator population, persistent interdependent oscillation dynamics are observed in the user and creator populations. In the bottom two panels (plots *c* and *d*), with initially 45% *C* in the creator population, *AllD* dominance and *D* dominance are observed in the user and creator populations. Parameters are set as $b_c = 0.4$, $c_c = 0.2$, $b_u = 0.4$, $c_u = 0.8$, $c_i = 0.05$, and $q = 0.9$.

perfect, and even if users have to pay substantially for the service of such media, it was able to maintain cooperation (products are safely developed and users adopt the products) for a wide range of the parameter space. However, we also showed that media in some cases failed to enable cooperation. These include: when the cost of good media or of safe AI creation was too high; when the signal of the media was too noisy; or when the population started with a large majority of defection by creators and users. Our findings were robust to both analytical predictions and agent-based simulations. Furthermore, we found that populations with high cooperation often experience cycles of behaviour. Widespread use of good media leads to the adoption and creation of safe AI, undermining the need for good media. At this point, users blindly adopt all AI products. This, in turn, leads to an increase in unsafe AI creation that exploits the dominant naive user population, making good media valuable again, thus restarting the cycle.

Our work is related to the large body of research on indirect reciprocity (Nowak and Sigmund, 2005; Okada, 2020). Indirect reciprocity means that if player *x* cooperates with player *y*, player *x* is rewarded by the cooperation of another player *z*, because *x* has built a good reputation. In such models, reciprocity can flow from any player to any other. In our model, on the other hand, reciprocity only flows from users

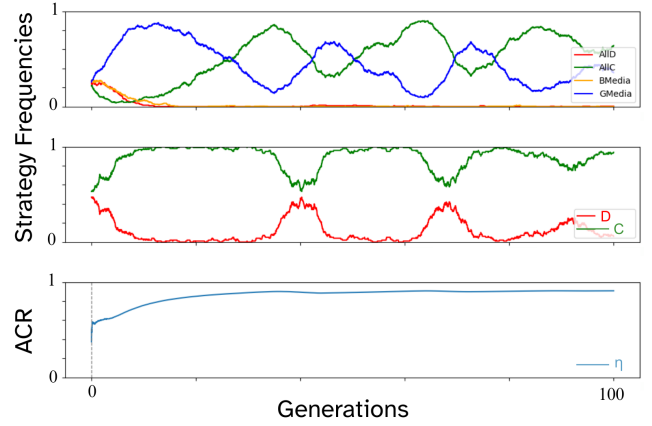


Figure 6: Evolution of strategies and cooperation over time for the agent-based model of user (top) and creator (middle) strategies, alongside the average cooperation ratio (ACR, or η) (bottom), for a typical run of the simulation ($R = 1$) with $N_U = 2N_C = 100$, $\mu_u = 1/N_U = 0.005$, $\mu_c = 1/N_C = 0.01$ and other payoff values set to the same as in Figure 5. The simulation is run for $G = 100$ generations for a minute illustration of the typical evolutionary dynamics. User strategies *AllD*, *AllC*, *BMedia*, and *GMedia* are pictured in red, green, yellow, and blue, respectively. Creator strategies, *Unsafe* (*D*) and *Safe* (*C*), are pictured in red and green, respectively.

towards creators, not vice versa nor within the populations. In standard indirect reciprocity, players are motivated to cooperate because they want to receive reciprocal cooperation. In our model, many players – namely users – can never receive reciprocal cooperation because their behaviour is entirely unmonitored. Instead, their cooperation is motivated by a self-interest to benefit from their own current cooperation (i.e. by benefiting from the use of AI), which they can only receive from creators that are – coincidentally – worth of reciprocal cooperation.

Limitations and Future Work

Our model is the first of its kind, showcasing the pure effect of media on safe AI adoption. In this, we limited the concerns of safety to the immediate user of the technology. Examples that could be captured by our model include the use of Large Language Model (LLM)-based health applications by patients, which could provide the patient with inaccurate or harmful information if not correctly regulated (Freyer et al., 2024), or the use of AI chatbots that could leak user conversations (as reported in the media with ChatGPT’s “share” feature; Prada, 2025). Our model does not, however, capture societal-wide consequences of unsafe AI development, such as plagiarism and copyright violations by LLMs, biased decisions adversely affecting minor-

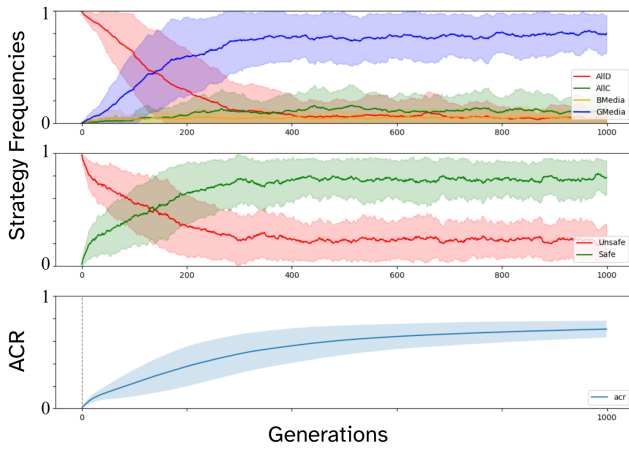


Figure 7: Initial state with only defective strategies. Evolution over simulation-time of user (top) and creator (strategies, as well as cooperation levels (bottom), where the initial population state is that of only AIID users and D creators. Colored shaded areas represent the standard deviation of averaging each metric and result. Results averaged over $R = 100$ runs. Payoff parameters set to $q = 0.9, b_c = 0.4, b_u = 0.4, c_c = 0.1, c_u = 0.8, c_i = 0.05$ and evolutionary parameters set to $N_U = 2N_C = 100, \beta_C = \beta_U = 1, \mu_U = 1/N_U, \mu_C = 1/N_C, G = 1000$.

ity groups (Wu et al., 2024) or widespread opinion polarisation through link-recommendation algorithms (Santos et al., 2021a). Considering these effects would require a substantially different model and is a worthwhile extension for the future. We also envisage several other expansions to address current limitations.

We only considered four types of users and two types of media at a time. More realistically, a whole ecosystem of different kinds of media outlets co-exists, most of which will have different levels of budgets for investigation and hence different levels of quality. Some media outlets spend millions of dollars, whereas actors on social media may only spend seconds of their time. Even more realistically, the accuracy of media is not just determined by their effort of investigation, but also by their bias. Many actors, especially in the social media realm, are either overly enthusiastic or critical about certain topics, including AI, to say nothing of the role of the political affiliation of media providers (Yang et al., 2023).

To address this shortcoming, we propose to expand our model in the future by introducing parameters for bias and expenditure. Bias could be a determinant of the initial assumption of a media outlet about all creators (before they start their investigation). It could range from trusting every creator to distrusting them all. During their investigation, the media outlet then has a chance to discover the real

value of the creator and change its opinion accordingly. This chance will depend on the individual expenditure of the media. Only if the true value is discovered will the initial assumption be overwritten. This way, a media outlet can provide valuable information but still be biased. This is especially interesting to the evolutionary dynamics, since biased media might be more profitable (Baron, 2006). However, increasing the set of possible media might be challenging to study.

A second limitation is that we assumed users would only listen to a single media source. However, in real life, users can hardly avoid being bombarded by many different sources. Integrating multiple sources of information is not trivial (Massaro and Friedman, 1990). Future research will need to define and compare heuristics for users to deal with multiple media sources that are available to them at the same time, especially if they provide contradictory information about the safety of AI.

In this work, we purposely limited the possible regulation of AI safety to be done by the media, using only reputations, but no other enforcement. Realistically, such enforcement, mainly by governments, is starting to come into place and is needed to ensure safety under all circumstances. We therefore want to continue the research of artificial population with government regulation (Han et al., 2020; Alalawi et al., 2026; Cimpeanu et al., 2022; Bova et al., 2024; Han et al., 2022; Buscemi et al., 2025), by combining it with media oversight (Balabanova et al., 2025; Powers et al., 2023), using more complex artificial systems to understand the dynamics of the real struggle for AI safety.

Acknowledgments

The authors acknowledge support by the Future of Life Institute (mini-grant for “AI Governance Modelling Workshop organisation” by TAH), EPSRC (grant no. EP/Y00857X/1 and grant no. EP/Y008561/1), CRCRM (MR/Z505833/1), European Commission (ERC). Additionally, this project was supported by the INESC-ID (UIDB/50021/2020), as well as by the Centre for Responsible AI (CRAI) project (grant no. C645008882-00000055/510852254 and C628696807-00454142, IAPMEI/PRR).

References

- Alalawi, Z., Bova, P., Cimpeanu, T., Di Stefano, A., Duong, M. H., Domingos, E. F., Han, T. A., Krellner, M., Ogbo, N. B., Powers, S. T., et al. (2026). Trust ai regulation? discerning users are vital to build trust and effective ai regulation. *Applied Mathematics and Computation*, 508:129627.
- Balabanova, N., Bashir, A., Bova, P., Buscemi, A., Cimpeanu, T., da Fonseca, H. C., Di Stefano, A., Duong, M. H., Domingos, E. F., Fernandes, A., et al. (2025). Media and responsible ai governance: a game-theoretic and llm analysis. *arXiv preprint arXiv:2503.09858*.
- Baron, D. P. (2006). Persistent media bias. *Journal of Public Economics*, 90(1):1–36.

- Bedau, M. A. (2003). Artificial life: organization, adaptation and complexity from the bottom up. *Trends in cognitive sciences*, 7(11):505–512.
- Bengio, Y., Mindermann, S., Privitera, D., Besiroglu, T., Bommasani, R., Casper, S., Choi, Y., Fox, P., Garfinkel, B., Goldfarb, D., et al. (2025). International ai safety report. *arXiv preprint arXiv:2501.17805*.
- Bova, P., Di Stefano, A., and Han, T. A. (2024). Both eyes open: Vigilant incentives help auditors improve ai safety. *Journal of Physics: Complexity*, 5(2):025009.
- Boyd, R. and Richerson, P. J. (1989). The evolution of indirect reciprocity. *Social Networks*, 11:213–236.
- Buscemi, A., Proverbio, D., Bova, P., Balabanova, N., Bashir, A., Cimpeanu, T., da Fonseca, H. C., Duong, M. H., Domingos, E. F., Fernandes, A. M., et al. (2025). Do LLMs trust AI regulation? Emerging behaviour of game-theoretic LLM agents. *arXiv preprint arXiv:2504.08640*.
- Cao, Y. and Li, C. (2020). The influence mechanism of reputation information on the formation of safety trust in chinese infant milk powder. *Healthcare (Switzerland)*, 8(2):1–16.
- Carr, J. (2012). *Applications of centre manifold theory*, volume 35. Springer Science & Business Media.
- Cimpeanu, T., Santos, F., et al. (2022). Artificial Intelligence Development Races in Heterogeneous Settings. *Scientific Reports*, 12(1):1723.
- Cohen, M. K., Kolt, N., Bengio, Y., Hadfield, G. K., and Russell, S. (2024). Regulating advanced artificial agents. *Science*, 384(6691):36–38.
- Commission, E., Directorate-General for Communications Networks, C., Technology, and ekspertów wysokiego szczebla ds. sztucznej inteligencji, G. (2019). *Ethics guidelines for trustworthy AI*. Publications Office.
- Freyer, O., Wiest, I. C., Kather, J. N., and Gilbert, S. (2024). A future role for health applications of large language models depends on regulators enforcing safety standards. *The Lancet Digital Health*, 6(9):e662–e672. Publisher: Elsevier.
- Gershenson, C., Trianni, V., Werfel, J., and Sayama, H. (2020). Self-organization and artificial life. *Artificial Life*, 26(3):391–408.
- Hammond, L., Chan, A., Clifton, J., Hoelscher-Obermaier, J., Khan, A., McLean, E., Smith, C., Barfuss, W., Foerster, J., Gavenčiak, T., Han, T. A., Hughes, E., Kovařík, V., Kulveit, J., Leibo, J. Z., Oosterheld, C., de Witt, C. S., Shah, N., Wellman, M., Bova, P., Cimpeanu, T., Ezell, C., Feuillade-Montixi, Q., Franklin, M., Kran, E., Krawczuk, I., Lamparth, M., Lauffer, N., Meinke, A., Motwani, S., Reuel, A., Conitzer, V., Dennis, M., Gabriel, I., Gleave, A., Hadfield, G., Haghtalab, N., Kasirzadeh, A., Krier, S., Larson, K., Lehman, J., Parkes, D. C., Piliouras, G., and Rahwan, I. (2025). Multi-Agent Risks from Advanced AI. *preprint arxiv 2502.14143*.
- Han, T. A., Lenaerts, T., et al. (2022). Voluntary Safety Commitments Provide an Escape from Over-Regulation in AI Development. *Technology in Society*, 68:101843.
- Han, T. A., Pereira, L. M., et al. (2020). To Regulate or Not: A Social Dynamics Analysis of an Idealised AI Race. *Journal of Artificial Intelligence Research*, 69:881–921.
- Han, T. A., Pereira, L. M., and Lenaerts, T. (2019). Modelling and influencing the AI bidding war: a research agenda. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 5–11.
- Hern, A. (2019). Apple apologises for allowing workers to listen to siri recordings — apple — the guardian.
- Hilbe, C., Schmid, L., Tkadlec, J., Chatterjee, K., and Nowak, M. A. (2018). Indirect reciprocity with private, noisy, and incomplete information. *Proceedings of the National Academy of Sciences*, 115:12241–12246.
- Hofbauer, J. and Sigmund, K. (1998). *Evolutionary games and population dynamics*. Cambridge university press.
- Imhof, L. A., Fudenberg, D., and Nowak, M. A. (2005). Evolutionary cycles of cooperation and defection. *Proceedings of the National Academy of Sciences of the United States of America*, 102(31):10797–10800.
- Jøsang, A., Ismail, R., and Boyd, C. (2007). A survey of trust and reputation systems for online service provision. *Decision Support Systems*, 43:618–644.
- Kelly, J. (2024). Google’s ai recommended adding glue to pizza and other misinformation—what caused the viral blunders?
- Khalil, H. K. and Grizzle, J. W. (2002). *Nonlinear systems*, volume 3. Prentice hall Upper Saddle River, NJ.
- Krellner, M. and Han, T. A. (2021). Pleasing enhances indirect reciprocity-based cooperation under private assessment. *Artificial Life*, 27(3–4):246–276.
- Massaro, D. W. and Friedman, D. (1990). Models of integration given multiple sources of information. *Psychological Review*, 97(2):225–252.
- McMahon, L. and Kleinman, Z. (2024). Glue pizza and eat rocks: Google ai search errors go viral.
- Melero-Bolaños, R., Gutiérrez-Villar, B., Montero-Simo, M. J., Araque-Padilla, R. A., and Olarte-Sánchez, C. M. (2025). Media Influence on the Perceived Safety of Dietary Supplements for Children: A Content Analysis of Spanish News Outlets. *Nutrients*, 17(6):1–14.
- Nowak, M. A. and Sigmund, K. (1998). Evolution of indirect reciprocity by image scoring. *Nature*, 393:573–577.
- Nowak, M. A. and Sigmund, K. (2005). Evolution of indirect reciprocity. *Nature*, 437(7063):1291–1298.
- Okada, I. (2020). A review of theoretical studies on indirect reciprocity. *Games*, 11(3):27.
- Paiva, A., Santos, F. P. F. C., and Santos, F. P. F. C. (2018). Engineering pro-sociality with autonomous agents. In *AAAI*, volume 32, pages 7994–7999.
- Perc, M., Jordan, J. J., Rand, D. G., Wang, Z., Boccaletti, S., and Szolnoki, A. (2017). Statistical physics of human cooperation. *Physics Reports*, 687:1–51.

- Piper, K. (2024). Openai ndas: Leaked documents reveal aggressive tactics toward former employees — vox.
- Powers, S. T., Ekárt, A., and Lewis, P. R. (2018). Modelling enduring institutions: The complementarity of evolutionary and agent-based approaches. *Cognitive Systems Research*, 52:67–81.
- Powers, S. T., Linnyk, O., et al. (2023). The Stuff We Swim in: Regulation Alone Will Not Lead to Justifiable Trust in AI. *IEEE Technology and Society Magazine*, 42(4):95–106.
- Prada, L. (2025). ChatGPT Briefly Made Chat Logs Accessible on Google. Yikes.
- Santos, F., Pacheco, J., and Santos, F. (2018). Social norms of cooperation with costly reputation building. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32:4727–4734.
- Santos, F. P., Lelkes, Y., and Levin, S. A. (2021a). Link recommendation algorithms and dynamics of polarization in online social networks. *Proceedings of the National Academy of Sciences of the United States of America*, 118:e2102141118.
- Santos, F. P., Pacheco, J. M., and Santos, F. C. (2021b). The complexity of human cooperation under indirect reciprocity. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 376.
- Sayama, H. (2015). *Introduction to the modeling and analysis of complex systems*. Open SUNY Textbooks [Imprint].
- Sigmund, K. (2010). The calculus of selfishness. In *The Calculus of Selfishness*. Princeton University Press.
- Sommerfeld, R. D., Krambeck, H.-J., Semmann, D., and Milinski, M. (2007). Gossip as an alternative for direct observation in games of indirect reciprocity. *Proceedings of the National Academy of Sciences*, 104:17435–17440.
- Traulsen, A., Nowak, M. A., and Pacheco, J. M. (2006). Stochastic dynamics of invasion and fixation. *Phys. Rev. E*, 74:11909.
- Wu, X., Duan, R., and Ni, J. (2024). Unveiling security, privacy, and ethical concerns of ChatGPT. *Journal of Information and Intelligence*, 2(2):102–115.
- Xia, C., Wang, J., Perc, M., and Wang, Z. (2023). Reputation and reciprocity. *Physics of Life Reviews*, 46:8–45.
- Yang, S., Krause, N. M., Bao, L., Calice, M. N., Newman, T. P., Scheufele, D. A., Xenos, M. A., and Brossard, D. (2023). In ai we trust: The interplay of media use, political ideology, and trust in shaping emerging ai attitudes. *Journalism & Mass Communication Quarterly*, page 10776990231190868.
- Zhang, J., Wu, H.-C., Chen, L., and Su, Y. (2022). Effect of social media use on food safety risk perception through risk characteristics: Exploring a moderated mediation model among people with different levels of science literacy. *Frontiers in Psychology*, 13(September):1–14.