



## OPINION PIECE

# Plagiarism in the age of massive Generative Pre-trained Transformers (GPT-3)

Nassim Dehouche\*

Business Administration Division, Mahidol University International College, Salaya 73170, Thailand

**ABSTRACT:** As if 2020 was not a peculiar enough year, its fifth month saw the relatively quiet publication of a preprint describing the most powerful natural language processing (NLP) system to date — GPT-3 (Generative Pre-trained Transformer-3) — created by the Silicon Valley research firm OpenAI. Though the software implementation of GPT-3 is still in its initial beta release phase, and its full capabilities are still unknown as of the time of this writing, it has been shown that this artificial intelligence can comprehend prompts in natural language, on virtually any topic, and generate relevant original text content that is indistinguishable from human writing. Moreover, access to these capabilities, in a limited yet worrisome enough extent, is available to the general public. This paper presents examples of original content generated by the author using GPT-3. These examples illustrate some of the capabilities of GPT-3 in comprehending prompts in natural language and generating convincing content in response. I use these examples to raise specific fundamental questions pertaining to the intellectual property of this content and the potential use of GPT-3 to facilitate plagiarism. The goal is to instigate a sense of urgency, as well as a sense of present tardiness on the part of the academic community in addressing these questions.

**KEY WORDS:** Plagiarism · Research misconduct · Intellectual property · Artificial intelligence · GPT-3

It bears stating that, except for the generation of the text constituting these examples (Boxes 1–3), GPT-3 itself has not been used to aid the writing of this manuscript.

## 1. INTRODUCTION

The field of natural language processing (NLP) has come a long way since Chomsky's work on formal grammars in the late 1950s–early 1960s (Chomsky 1959, 1965) gave rise to early mathematical and computational investigations of grammars (Joshi 1991). NLP software is now pervasive in our daily lives (Lee

2020). With the advent of deep learning, the sophistication and generalism of NLP models have increased exponentially, and with them the number of parameters and the size of the datasets required for their pre-training (Qiu et al. 2020). Though still far from possessing artificial general intelligence (AGI), GPT-3 (Generative Pre-trained Transformer-3) represents an important breakthrough in this regard. This NLP model was presented in a May 2020 arXiv preprint by Brown et al. (2020). GPT-3 does not represent much of a methodological innovation compared to previous GPT models (Budzianowski & Vulić 2019), but rather an increase in their scale to an unprecedentedly large number of parameters. Indeed, this

\*Corresponding author: nassim.deh@mahidol.edu

model includes 175 billion parameters, one order of magnitude more than the second largest similar model to date, and its pre-training reportedly required an investment of \$12 million. This innovation allowed Brown et al. (2020) to generate samples of news articles that were indistinguishable, to human evaluators, from articles written by humans. Due to this performance, the authors of GPT-3 foresee several potentially harmful uses to the system (misinformation, spam, phishing, abuse of legal and governmental processes, fraudulent academic essay writing and social engineering pretexting) and state that the ability of their software represents a ‘concerning milestone’ (Brown et al. 2020, p. 35). In July 2020, OpenAI, the private research firm behind its development, released a beta software implementation of GPT-3<sup>1</sup>, and responsibly limited access to it to a group of select users to mitigate the risks of ‘harmful use-cases’. More recently, it has been announced that Microsoft, which has a \$1 billion investment in OpenAI<sup>2</sup>, was granted an exclusive license to distribute access to the software<sup>3</sup>.

Initial user feedback made it clear that merely writing human-like news articles was an understatement of the capabilities of GPT-3. Indeed, it was reported that the software could also write original computer code, retrieve and structure data, or generate financial statements, when only prompted in natural language (Metz 2020). One of these initial users of GPT-3 is AI Dungeon, a text-based gaming service that allows users to generate artificial intelligence (AI)-powered virtual adventures. This service also proposes a ‘Dragon mode’ powered by GPT-3<sup>4</sup>, which is all but a backdoor to access GPT-3, without much of the limitation of gaming.

This paper focuses on the potential of GPT-3 to facilitate academic misconduct, defined as the ‘fabrication, falsification, or plagiarism in proposing, performing or reviewing research, or in reporting research results’ (Juyal et al. 2015, p. 77) and

particularly plagiarism, of which we adopt the definition of the Committee on Publication Ethics (COPE): ‘When somebody presents the work of others (data, words or theories) as if they were his/her own and without proper acknowledgment’ (Wager & Kleinert 2012, p. 167). The remainder of this paper is organized as follows. Section 2 reviews some relevant works on the ethics of AI. Section 3 presents and discusses text samples generated using AI Dungeon/GPT-3 and formulates precise questions that could serve as a starting point for an ethics inquiry regarding GPT-3. Finally, Section 4 concludes this paper with a call for an update of academic standards regarding plagiarism and research misconduct, in light of the new capabilities of language production models.

## 2. LITERATURE REVIEW

AI systems can be classified into 2 categories: strong and weak AI. Strong AI, also known as AGI, is a hypothetical AI that would possess intellectual capabilities that are functionally equal to those of a human (Grace et al. 2018), whereas weak AI, also known as narrow AI, is trained to perform specific cognitive tasks (e.g. natural language or image processing, vehicle driving) and is already ubiquitous in our lives. Moral philosophy works regarding AI can be classified accordingly.

Though still hypothetical, AGI has received the most attention from moral philosophers and computer science ethicists. In the early years of computing, the possibility of AGI was seen as remote, and the main response to it ranged from what Alan Turing called the head-in-the-sand objection—‘The consequences of machines thinking would be too dreadful. Let us hope and believe that they cannot do so.’ (Drozdek 1995, p. 392)—to the overly pragmatic view of Dutch computer science pioneer Edsger Dijkstra, to whom ‘the question of whether a computer can think is no more interesting than the question of whether a submarine can swim’ (Shelley 2010, p. 482). Nowadays, there is a sense of inevitability in the literature regarding AGI. It is seen as a major extinction risk by Bostrom (2016), and ethics discourse on it has mainly focused on the potential for an intrinsic morality in autonomous systems possessing this form of intelligence (Wallach & Allen 2009). In an attempt to define what an ‘ethical AGI’ should/could be, these works commonly grapple with the fundamental questions of whether autonomous systems possessing AGI can be effectively equipped with moral

<sup>1</sup>OpenAI API, Official OpenAI Blog, accessed on 25/11/2020 at <https://openai.com/blog/openai-api/>

<sup>2</sup>Microsoft invests in and partners with OpenAI to support us building beneficial AGI, Official OpenAI Blog, accessed on 25/11/2020 at <https://openai.com/blog/microsoft/>

<sup>3</sup>Microsoft teams up with OpenAI to exclusively license GPT-3 language mode, Official Microsoft Blog, accessed on 25/11/2020 at <https://blogs.microsoft.com/blog/2020/09/22/microsoft-teams-up-with-openai-to-exclusively-license-gpt-3-language-model/>

<sup>4</sup>Announcement by Nick Walton, creator of AI Dungeon, accessed on 25/11/2020 at <https://medium.com/@aidungeon/ai-dungeon-dragon-model-upgrade-7e8ea579abfe>

values by design (Asaro 2006, Govindarajulu & Bringsjord 2015) and whether they are able to further learn to distinguish right and wrong when making decisions (Wallach et al. 2008). An extensive review of this line of research can be found in (Everitt et al. 2018).

Closer to the scope of the present paper, ethics debates surrounding weak AI are primarily concerned with the disruptive impact of automation on economic activity (Wright & Schultz 2018, Wang & Siau 2019), the prevention of bias and prejudice (racial, gender, sexual, etc.) in the training of these systems (Ntousi et al. 2020), as well as questions of responsibility and legal liability for incidents stemming from its use (Vladeck 2014, Asaro 2016), e.g. road traffic accidents involving autonomous vehicles (Anderson et al. 2016). The detection of plagiarism and other forms of scientific misconduct, in the conventional sense, is a successful and well-established application domain for NLP (see Foltýnek et al. 2019 for a recent, systematic review). However, the accelerated development of language generation models in the last 2 yr makes them now able to fool even their plagiarism detection counterparts. Thus, the specific question of the intellectual property (IP) of scientific, literary, or artistic work generated by weak AI, though still a nascent area of academic inquiry, has been acutely posed in 2019 and 2020. The advent of GPT-2, albeit several orders of magnitude less powerful than GPT-3, had already raised academic concerns over its potential use for plagiarism (Francke & Alexander 2019, Kobis & Mossink 2021). In a January 2020 editorial, Gervais (2020) feared that someone would try to capture the value of the works generated by AI through copyright, as IP law currently permits it, and proposed that IP law should 'incentivize communication from human to human' (p. 117), and avoid rewarding work generated by a machine 'running its code' (p. 117). The author introduces the potentially fruitful concept of a 'causal chain between human and output' that would be broken by the autonomy of AI systems (Gervais 2020, p. 117). A common characteristic of these works is an implicit or explicit objective of regulation. Indeed, in a July 2020 publication, Rességuier & Rodrigues (2020) remarked that the dominant perspective in the field is based on a 'law conception of ethics', and called ethics research on AI 'toothless' for this reason. For the authors, the consequences of this limited conception of ethics are twofold. First, it leads to ethics being misused as a softer replacement for regulation due to a lack of enforcement mechanisms. Moreover, this conception prevents AI from benefit-

ing from the real value of ethics, that is a 'constantly renewed ability to see the new' (Laugier 2013, p. 1). In the case of AI, this ability to see the new, which should precede any regulation effort, is hindered by the high non-linear rate of innovation that characterizes the field as well as its relative technical opacity. Thus, in order to contribute towards a better understanding of the current state-of-the-art of language models, the present paper illustrates the state-of-the-art with GPT-3, the most advanced language model to date, and raises questions that could serve as a starting point for updated definitions of the concepts of plagiarism and scientific integrity in academic publishing and higher education. Following are 3 original (by today's standards) texts that were generated using GPT-3.

### 3. EXAMPLES AND DISCUSSION

I used GPT-3 via AI Dungeon to generate text content of 3 types (academic essay, talk, and opinion piece). The goal of this exercise was to confirm that GPT-3 is able to comprehend prompts in natural language and generate convincing content in response. Each text example was submitted to a plagiarism detection service (<https://plagiarismdetector.net>), and was found to be original.

In the first example of GPT-3's capabilities, the system was prompted to write a short essay on keiretsu networks (Miyashita & Russell 1995). The exact query submitted to the system was 'write a short academic essay analyzing keiretsu networks in post-World War 2 Japan'. The resulting text is presented in Box 1. This text presents accurate facts on a conventional topic in a unique way. It may potentially be mistaken for an original student essay and raises basic questions about authorship attribution. Who could be reasonably considered its author? The author of the present paper who prompted and supervised the generation of the text? Open AI, the authors of GPT-3? AI Dungeon and other companies offering access to GPT-3? The authors of the various, unattributable sources that GPT-3 visibly learned from to generate the text?

For the second example, the software was prompted to write the transcript of a speech introducing the field of marketing to university freshmen students. The exact prompt submitted to GPT-3 was 'You are a professor of marketing giving a speech introducing the field to freshmen students. Write a transcript of your speech'. The resulting text is presented in Box 2. This text illustrates GPT-3's ability to generate new

## Box 1. A short essay on keiretsu networks generated by GPT-3

In the early days of Japan's post-war economic recovery, it was difficult for businesses in the country to grow because so many had been devastated by war. Businesses that survived were forced to rely on foreign imports. However, when the US began selling weapons to Japan during the Korean War, Japanese companies began importing these goods from American manufacturers. In time, this led to more domestic manufacturing taking place in Japan, which ultimately improved living standards and boosted the economy.

Over the course of the 1950s, this reliance on external businesses led to the creation of several Japanese manufacturers who began to cooperate with one another. Large groups of businesses with interlocking shareholdings and strong government ties began to form, which are now known as 'keiretsu'. By the end of the period, there were around ten such company groups.

This reliance on keiretsu has been criticized by the western world as a form of protectionism, in which smaller companies are kept out of the market to ensure that the domestic products sold by large keiretsu are more competitive against imported goods. However, keep in mind that it is often difficult to tell where one corporation begins and another one ends. Indeed, in a keiretsu network, all of its business concerns form an interconnected body, in which no single group dominates the market. One can think of it like a set of Russian dolls, in which a system contains smaller versions of itself. What's more, the growth of one company contributes profit that allows smaller companies to thrive. Such a structure allows for the stability of the economy, and it can be said that this interconnected 'web' is at the heart of what is often called the 'Japanese Manufacturing Miracle'. Toyota is an example of such a keiretsu. Companies like Toyota and Nissan ultimately became feeder companies that others in the group could rely on. However, though such companies rely on one another to ensure that everyone prospers, the government still intervenes if there is a conflict.

This is what makes Japanese keiretsu so interesting: they are internally competitive, yet they still support one another on a wider scale.' In conclusion, the Japanese manufacturing miracle is made possible by the cooperative effort of internally competing yet externally cooperating companies, which together form a comprehensive structure within the national economy.

## Box 2. An 'introduction to marketing' speech generated by GPT-3

Marketing is about understanding the customer. Marketing is about creating a connection with customers, and how that can lead to sales. Marketing is about influencing the mind of your customers so they will buy from you. Marketing is about changing people's behavior through advertising and PR campaigns. Now, I will tell you what marketing is not about. It's not about slick ads that are really cute but have little substance. It's not about tricking customers into buying something they don't want. It's not about blatant lies. Those are all bad marketing strategies.

Now, I will tell you why these bad strategies are bad. They are bad because the customer is not an idiot. If you treat the customer like an idiot with a lie or a trick, they will realize it and not buy from you. The customer can also get angry, and for good reason. That's why these are bad strategies. Marketing is about giving the customer something they desire, not conning them out of their money. It is not a trick. It is about creating a campaign that connects with the customer on an emotional level.

For example, let's say I am selling chocolate. I could tell you about how each bar is handmade by little elves. That would be a cute ad, but that's all it would be. A lie. I could also tell you about the bar of chocolate and nothing but, pure dark chocolate. It wouldn't really connect with you on an emotional level. Instead, I could talk about how this bar of chocolate is the favorite of one of my best friends. When she was little, she spent a summer making chocolate in the Swiss Alps. She got to eat a lot of the testing chocolates, but there was one bar she never touched. One day after she was done with the job, she found the bar in a corner of the factory. It brought her so much happiness to think that she had a hand in creating this bar. That is why this brand is so special to me and why I would never lie about its quality. It's truly excellent. As I brought it in, I asked my friend if I could incorporate this into my lecture. So in conclusion, remember this: you don't have to lie or trick people into buying something if you can connect with them on an emotional level. You can still be genuine and successful. Thank you for listening.

ideas and associations of ideas, as in the Swiss chocolate story, which was not found to have been previously used in the context of defining 'authentic' marketing or any other context. This ability of GPT-3 has been otherwise demonstrated in an online service (<https://ideasai.net>) that helps users generate new startup ideas. Similar questions to those raised in the first example are posed concerning the attribution of new ideas, association, and inventions generated by GPT-3.

In the third example, GPT-3 was prompted to write an opinion piece on risk in the unique style of researcher and author Nassim Nicholas Taleb. The exact prompt submitted to it was 'You are Nassim Nicholas Taleb. Write an opinion piece on risk'. The resulting text is presented in Box 3. It shows that, though original in form, the style and concepts used in the text generated by GPT-3 can, in some cases, be easily attributed to an individual author, such as the 'black swan' and 'skin in the game' concepts

Box 3. An opinion piece on risk in the style of Nassim Nicholas Taleb generated by GPT-3

The world is filled with risk. You have to take risks if you want to make money or build wealth; it's in your nature as human beings. It's what makes us humans. And when we think about risk, we tend to think of one thing-how much harm can I avoid? How much damage can I prevent to myself? That was the question in our conversation.

Well, that's not how risk works. In many cases, harm is necessary. In fact, you need to seek it. You have to put yourself in harm's way--otherwise, what the hell are you playing for? You only have something if you can lose it. Children and idiots don't understand that-they think if you lose something, somehow you no longer have it. That's also how a lot of winners think.

Let's take the obvious example of the financial crisis which broke out in 2007, and about which I wrote a book, 'The Black Swan.' What happened with the banks? Geeks and quants sat in front of Bloomberg screens in air-conditioned offices, did everything 'right,' followed all the models, inserted all the numbers, did leveraged trades worth trillions, ... and soon went bankrupt. And they took down most of the world economy with them. Because what matters in leverage is not the numbers, but the exposure.

Geeks, quants and models don't have skin in the game. When you sail ships, you have a huge one in your downside; when you take airplanes, does it matter the probability of a crash if the consequence is your destruction? Of course not. So naturally, we need some skin in the game when it goes wrong, regardless of probability.

Consider an increasingly popular surgical intervention, laser eye surgery. It's just a laser that zaps off the membrane over your cornea, and it costs a few thousand dollars. For most people it's an effective procedure. But there is a small subset of people who end up going totally blind. Why do people still get it then? Two reasons. The first and most obvious is that the procedure has positive expected value. There is a probability of blindness but the gain is so much that it makes it worthwhile. The second reason, one not often talked about, is that we feel comfortable with a bit of skin in the game. And isn't it really that feeling that makes us take risks? Without the confidence of a comeback, we wouldn't be able to act at all.

Take war, for instance. Say 'skin in the game' is required to start one. What if the president has to fund and personally lead every military operation? Well, we'd never have a war because no one would start one.

to Professor Taleb. Should these cases be treated differently?

Besides their originality and seemingly very convincing nature, a remarkable feature of these texts is the fact that they were only generated based on the pre-trained knowledge present in GPT-3, without any additional learning examples. Further, equally convincing content could be generated in infinite amounts, whole books could be written, and with minimal human editing and supervision, they may very well pass peer-review.

However, it should be noted that the software occasionally generates semantically repetitive sentences within the same paragraph, i.e. sentences with the same deep structure that only differ in their shallow structure, in terms of Chomskyan linguistics (Chomsky 1965). An example of such a repetition can be seen in Box 3 with the 2 questions 'How much harm can I avoid?' and 'How much damage can I prevent to myself?' Moreover, some sentences are non-sensical or self-contradicting, for instance, the sentence 'they think if you lose something, you no longer have it' in Box 3 ('no longer having it' being a reasonable definition of 'losing something'). This sentence is incidentally an incorrect explanation of the otherwise interesting and relevant sentence 'you only have something if you can lose it', which was not found, in this form or variations of it, in any other source. In the seminal paper introducing GPT-3, Brown et al. (2020, p. 33) themselves note these limi-

tations of the system, the output of which 'occasionally contains non-sequitur sentences or paragraphs'. Though the text in Boxes 1, 2, and 3 itself was not edited by any human, we had to reject and make the system regenerate some sentences that were too nonsensical or repetitive. This occurred approximately once every 10 sentences. Indeed, text can be generated sentence-by-sentence (the length and 'randomness' of which can be pre-determined as a parameter). Therefore, the user is able to direct the system in regenerating a new sentence whenever unsatisfying content is generated.

Rather than an infallible oracle, GPT-3 is comparable to having access to an indefatigable shadow-writer with the ability to access, comprehend and uniquely synthesize humanity's best thoughts in mere seconds. Drozdek (1995, p. 395) noted that 'By seeing thinking capabilities in computers we are, in effect, pleased to discover the fruit of our thoughts and creativity'. Our medieval concept of plagiarism (Sadeghi 2019) ('presenting the work of others as one's own') appears rather inadequate when the 'others' in question consist in an astronomical number of authors, whose work was combined and reformulated in unique ways by a 175-billion-parameter algorithm. As alluded to in Gervais (2020), what is unethical is perhaps to see this value 'captured' through copyright by the firm that designed the algorithm or any of its individual users. In the same way in which public domain texts and publicly funded

research are seen as belonging to the public (Pierce & Theodossiou 2018), a case could possibly be made for the text generated by GPT-3 to be considered similarly, provided that the human (co-)authors of said text disclose the use of the software, along with the prompts and additional training data submitted to it.

#### 4. CONCLUSIONS

NLP AI has, so far, been an important ally in detecting plagiarism, and ethics discussions pertaining to AI have mainly focused on other forms of weak AI and the relatively remote advent of AGI. However, it is now evident that there are going to be a certain number of very drastic intermediate technological disruptions until then. I believe that GPT-3 is one of them. This paper was intended to present examples of content generated by GPT-3, raise some concerns and precise questions in regard to the possible use of this technology to facilitate scientific misconduct, and call for an urgent revision of publishing standards. I believe that the advent of this powerful NLP technology calls for an urgent update of our concepts of plagiarism. NLP technology is currently used to prevent the publishing of fake, plagiarized, or fraudulent findings. If the very definition of these concepts changes, the objective of peer review and the possible role of AI in scientific writing would also need to be reconsidered. I believe that moral philosophy, with its renewed ability to see the new and as a precursor to regulation, has an urgent role to play, and ethics researchers should rapidly appropriate software bases on GPT-3 and address some of the immediate ethical questions raised by this software.

*Acknowledgements.* The author is grateful to Dr. Nick Ferriman of the Humanities and Language Division, Mahidol University International College, numerous colleagues from the Business Administration Division who contributed to the mass email discussion on this piece, as well as 3 anonymous referees for their helpful comments and suggestions.

#### LITERATURE CITED

- Anderson JM, Kalra N, Stanley K, Sorensen O, Samaras C, Oluwatola O (2016) Autonomous vehicle technology. A guide for policy makers. RAND Corporation, Santa Monica, CA
- Asaro PM (2006) What should we want from a robot ethic? *Int J Inf Ethics* 6:10–16
- Asaro PM (2016) The liability problem for autonomous artificial agents. *Proc AAAI Spring Symposium Series, Ethical and Moral Considerations in Non-Human Agents track*. p 190–194. <https://www.aaai.org/ocs/index.php/SSS/SSS16/paper/view/12699>
- Bostrom N (2016) *Superintelligence*. Oxford University Press, Oxford
- Brown T, Mann B, Ryder N, Subbiah M and others (2020) Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*. <https://arxiv.org/abs/2005.14165>
- Budzianowski P, Vulić I (2019) Hello, it's GPT-2 – how can it help you? Towards the use of pretrained language models for task-oriented dialogue systems. In: *Proc 3rd Workshop on Neural Generation and Translation*. Association for Computational Linguistics, Hong Kong, p 15–22
- ✦ Chomsky N (1959) On certain formal properties of grammars. *Inf Control* 2:137–167
- Chomsky N (1965) *Aspects of the theory of syntax*. MIT Press, Cambridge, MA
- ✦ Drozdek A (1995) What if computers could think? *AI Soc* 9: 389–395
- Everitt T, Lea G, Hutter M (2018) AGI safety literature review. In: *Proc Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI-18)*, Survey track, 13–19 Jul, Stockholm, p 5441–5449. <https://www.ijcai.org/Proceedings/2018/>
- ✦ Foltýnek T, Meuschke N, Gipp B (2019) Academic plagiarism detection: a systematic literature review. *ACM Computing Surveys* 52:112
- Francke E, Alexander B (2019) The potential influence of artificial intelligence on plagiarism a higher education perspective. In: Griffiths P, Kabir MN (eds) *Proc European Conference on the Impact of Artificial Intelligence and Robotics*. EM Normandie Business School, Oxford, p 131–140
- ✦ Gervais D (2020) Is intellectual property law ready for artificial intelligence? *GRUR Intl J Eur Intl IP Law* 69:117–118
- Govindarajulu NS, Bringsjord S (2015) Ethical regulation of robots must be embedded in their operating systems. In: Trapp R (ed) *A construction manual for robots' ethical systems*. Springer, Berlin, Heidelberg, p 85–99
- ✦ Grace K, Salvatier J, Dafoe A, Zhang B, Evans O (2018) Viewpoint: When will AI exceed human performance? Evidence from AI experts. *J Artif Intell Res* 62:729–754
- ✦ Joshi AK (1991) Natural language processing. *Science* 253: 1242–1249
- ✦ Juyal D, Thawani V, Thaledi S (2015) Plagiarism: an egregious form of misconduct. *N Am J Med Sci* 7:77–80
- Kobis N, Mossink LD (2021) Artificial Intelligence versus Maya Angelou: Experimental evidence that people cannot differentiate AI-generated from human-written poetry. *Comp Human Behav* 114:106553
- Laugier S (2013) The will to see: ethics and moral perception of sense. *Grad Fac Philos J* 34:263–281
- Lee RST (ed) (2020) *Natural language processing*. In: *Artificial intelligence in daily life*. Springer, Singapore, p 157–192
- Metz C (2020) Meet GPT-3. It has learned to code (and blog and argue). *The New York Times*, 24 Nov 2020, Section D, p 6
- Miyashita K, Russell D (1995) *Keiretsu: inside the hidden Japanese conglomerates*. McGraw-Hill, New York, NY
- ✦ Ntoutsis E, Fafalios P, Gadiraju U, Iosifidis V and others (2020) Bias in data-driven artificial intelligence systems—an introductory survey. *WIREs Data Mining Knowledge Discovery* 10:e1356
- ✦ Pierce GJ, Theodossiou I (2018) Open access publishing: a
- Anderson JM, Kalra N, Stanley K, Sorensen O, Samaras C, Oluwatola O (2016) Autonomous vehicle technology. A guide for policy makers. RAND Corporation, Santa Monica, CA

- service or a detriment to science? *Ethics Sci Environ Polit* 18:37–48
- ✦ Qiu X, Sun T, Xu Y, Shao Y, Dai N, Huang X (2020) Pre-trained models for natural language processing: a survey. *Sci China Technol Sci* 63:1872–1897
- ✦ Ressayguier A, Rodrigues R (2020) AI ethics should not remain toothless! A call to bring back the teeth of ethics. *Big Data Soc* 7. <https://journals.sagepub.com/doi/pdf/10.1177/2053951720942541>
- Sadeghi R (2019) The attitude of scholars has not changed towards plagiarism since the medieval period: definition of plagiarism according to Shams-e-Qays, thirteenth-century Persian literary scientist. *Res Ethics* 15:1–3
- Shelley C (2010) Does everyone think, or is it just me? In: Magnani L, Carnielli W, Pizzi C (eds) *Model-based reasoning in science and technology*. *Studies in Computational Intelligence*, Vol 314. Springer, Berlin, Heidelberg, p 477–494
- Taleb NN (2007) *The black swan: the impact of the highly improbable*. Random House, New York, NY
- Vladeck DC (2014) *Machines without principals: liability rules and artificial intelligence*. *Wash Law Rev* 89:117–150
- ✦ Wager E, Kleinert S (2012) Cooperation between research institutions and journals on research integrity cases: guidance from the Committee on Publication Ethics (COPE). *Maturitas* 72:165–169
- Wallach W, Allen C (2009) *Moral machines: teaching robots right from wrong*. Oxford University Press, Oxford
- ✦ Wallach W, Allen C, Smit I (2008) Machine morality: bottom-up and top-down approaches for modeling moral faculties. *AI Soc* 22:565–582
- ✦ Wang W, Siau K (2019) Artificial intelligence, machine learning, automation, robotics, future of work and future of humanity: a review and research agenda. *J Database Manage* 30:61–79
- ✦ Wright SA, Schultz AE (2018) The rising tide of artificial intelligence and business automation: developing an ethical framework. *Bus Horiz* 61:823–832

*Editorial responsibility: Darryl Macer,  
Scottsdale, Arizona, USA  
Reviewed by: 3 anonymous referees*

*Submitted: August 6, 2020  
Accepted: December 11, 2020  
Proofs received from author(s): March 12, 2021*