

## Response to the Royal Society Call for Evidence:

### Technologies for Spreading and Detecting Misinformation

Submitted by Prof. Derek McAuley, Dr. Ansgar Koene and Dr. Jiahong Chen of  
Horizon Digital Economy Research Institute, University of Nottingham

15 August 2020

1. Horizon<sup>1</sup> is a Research Institute centred at The University of Nottingham and a Research Hub within the UKRI Digital Economy programme<sup>2</sup>. Horizon brings together researchers from a broad range of disciplines to investigate the opportunities and challenges arising from the increased use of digital technology in our everyday lives. Prof. McAuley is Director of Horizon and Principal Investigator of the EPSRC-funded DADA<sup>3</sup> (Defence Against Dark Artefacts) project, addressing smart home IoT network security, and its acceptability and usability issues, the ESRC-funded CaSMa<sup>4</sup> (Citizen-centric approaches to Social Media analysis) project to promote ways for individuals to control their data and online privacy, and the EPSRC-funded UnBias<sup>5</sup> (Emancipating Users Against Algorithmic Biases for a Trusted Digital Economy) project for raising user awareness and agency when using algorithmic services. Dr Koene was a lead researcher of the CaSMa and UnBias projects, is Research co-Investigator on the EPSRC-funded ReEnTrust<sup>6</sup> (Rebuilding and Enhancing Trust in Algorithms) project and chairs the working group for developing the IEEE P7003 Standard for Algorithm Bias Considerations. Dr Jiahong Chen is a Researcher Fellow of Horizon, working on the DADA project and a book project based on his doctoral research on regulating online advertising.

#### ***Propagation and impacts of misinformation***

##### ***Q1. What impact have digital technologies had on patterns of information consumption? What evidence exists on their wider social impact?***

2. The prevalent business model of major sources of online information, including news outlets, blogs, and social media, depends heavily on monetisation of granular user profiles. The “ad tech” industry has developed advanced technologies, such as programmatic trading and universal IDs, to target internet users with highly personalised content. This may change the patterns of information consumption in several ways: First, internet users may find themselves trapped in their own “echo chambers” with a feedback loop of information from like-minded groups, which may intensify the polarisation of the society.<sup>7</sup> Second, political campaigners may gain unfair advantage by exploiting

---

<sup>1</sup> <http://www.horizon.ac.uk>

<sup>2</sup> <https://epsrc.ukri.org/research/ourportfolio/themes/digitaleconomy/>

<sup>3</sup> <https://www.horizon.ac.uk/project/defence-against-dark-artefacts/>

<sup>4</sup> <http://casma.wp.horizon.ac.uk>

<sup>5</sup> <http://unbias.wp.horizon.ac.uk>

<sup>6</sup> <https://ReEnTrust.org>

<sup>7</sup> Tien T. Nguyen et al. "Exploring the filter bubble: the effect of using recommender systems on content diversity." *Proceedings of the 23rd international conference on World wide web*. 2014; Haim, Mario, Andreas Graefe, and Hans-Bernd Brosius. "Burst of the filter bubble? Effects of personalization on the diversity of Google News." *Digital journalism* 6.3 (2018): 330-343.

online personal data, as shown by the ICO's investigation into the Cambridge Analytica scandal.<sup>8</sup> Third, conspiracy and pseudoscientific theorists may find it easier to channel their messages to potentially more susceptible audiences, which is evidenced by, for example, the revelation about the possibility to target social media users based on a "vaccine controversies" category.<sup>9</sup>

### **Q2. How do digital technologies contribute to the spread of misinformation?**

3. Many of the new technological phenomena are "neutral" in the sense that they can facilitate the dissemination of information regardless of the nature of such information. Social media, for example, have significantly augmented individuals' ability to create and share information, a large part of which, however, can be misinformation. The rise of marketing by online influencers has further contributed to the spread of misleading or mistaken information.<sup>10</sup> Deepfake is another controversial area where synthetic videos can be created to help circulate false news.<sup>11</sup>

### **Q3. What tools exist to create synthetic text, audio or visual media, and what are the likely near-term future developments of these technologies?**

4. Machine learning has been applied to create machine-generated content, including misinformation. Text-based applications, such as natural language generation, for example, have been proved capable of fabricating convincing news stories.<sup>12</sup> Image- and video-based applications are also widely used in context less associated with the spreading of misinformation, such as smartphone camera filters, but have also raised concerns about user privacy, dignity and mental health.<sup>13</sup> These technologies can be easily repurposed for conducting political or personal attacks.<sup>14</sup>

### **Detection and tracing of misinformation**

#### **Q4. Which technologies can currently be used to identify or trace misinformation? What are the strengths and weaknesses of these technologies?**

5. Current approaches to detecting misinformation can be largely categorised as text-based, image-based and profile-based. Text-based strategies identify misinformation by analysing lexical, semantic and statistical features.<sup>15</sup> Image-based solutions identify deepfakes by examining biological signals<sup>16</sup> or inter-frame dissimilarities.<sup>17</sup> Profile-based approaches identify misinformation spreaders

---

<sup>8</sup> <https://ico.org.uk/media/action-weve-taken/2260271/investigation-into-the-use-of-data-analytics-in-political-campaigns-final-20181105.pdf>

<sup>9</sup> <https://www.theguardian.com/technology/2019/feb/15/facebook-anti-vaccination-advertising-targeting-controversy>

<sup>10</sup> Catalina Goanta and Gerasimos Spanakis. "Influencers and Social Media Recommender Systems: Unfair Commercial Practices in EU and US Law." (2020). [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3592000](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3592000)

<sup>11</sup> <https://www.bbc.co.uk/news/business-51204954>

<sup>12</sup> <https://www.theguardian.com/technology/2019/feb/14/elon-musk-backed-ai-writes-convincing-news-fiction>

<sup>13</sup> <https://www.nbcnews.com/tech/security/face-swapping-app-takes-china-making-ai-powered-deepfakes-everyone-n1049501>; <https://www.bbc.co.uk/news/business-50152053>;

<https://www.theguardian.com/commentisfree/2019/jun/29/deepnude-app-week-in-patriarchy-women>

<sup>14</sup> <https://www.fireeye.com/blog/threat-research/2020/02/information-operations-fabricated-personas-to-promote-iranian-interests.html>; <https://www.vox.com/2020/6/8/21284005/urgent-threat-deepfakes-politics-porn-kristen-bell>

<sup>15</sup> Dinesh Kumar Vishwakarma and Chhavi Jain. "Recent State-of-the-art of Fake News Detection: A Review." 2020 *International Conference for Emerging Technology (INCET)*. IEEE, 2020.

<sup>16</sup> Umur Aybars Ciftci, Ilke Demir and Lijun Yin. "Fakecatcher: Detection of synthetic portrait videos using biological signals." *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020).

<sup>17</sup> Irene Amerini et al. "Deepfake video detection through optical flow based CNN." *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2019.

by distinguishing abnormal user behavioural patterns.<sup>18</sup>

**Q5. What technological advancements in the next ten years could improve the ability to identify or trace misinformation?**

6. Provenance has been held out as a solution to this problem as a means to clearly identify trustworthy information, and hence, by implication, everything else should be untrusted. The World Wide Web Consortium had a programme of work on Provenance culminating in 2013<sup>19</sup> which is widely ignored – one issue is that it tries to revert the unstructured web to the mentality of databases and predefined schemas. Rather, for provenance, we must draw the lessons from web search and adopt statistical means to define probabilistic provenance graphs. The challenge will be that a significant number of the original sources of misinformation are not openly available on the web, so a global system would require the collaboration of many platform providers and the federation of provenance information.

**Q6. What role could these technologies play in building a trustworthy information environment?**

7. The algorithmic editorial processes that are at the heart of much social media are currently statistical optimisations with a primary goal to drive profit. It would not be surprising to find that this is the antithesis of provenance, so challenges will include whether citizens have any faith that asking for “provenance ordered search results” or equivalent has not been tampered with based on commercial considerations.

**Q7. Are there any current regulatory or policy barriers to the successful development or deployment of detection technologies?**

8. Currently in the UK, there is no primary legislation prohibiting the publication or circulation of misinformation. Nor is there a general obligation for platforms to monitor or remove misinformation. Equally, however, the law does not stop platforms from taking measures to address misinformation, although such measures must be fully in line with human rights standards, especially with respect to freedom of speech. The adoption of detection technologies by online platforms are facing two major regulatory barriers: The lack of economic incentives<sup>20</sup> and the legal uncertainty of what counts as misinformation<sup>21</sup>. The Government’s Online Harms White Paper proposes to address these issues by introducing a statutory duty of care,<sup>22</sup> but the approach is subject to criticisms about the unclear definition of “harm”.<sup>23</sup>

---

<sup>18</sup> Liang Wu et al. "Misinformation in social media: definition, manipulation, and detection." *ACM SIGKDD Explorations Newsletter* 21.2 (2019): 80-90.

<sup>19</sup> <https://www.w3.org/TR/prov-overview/>

<sup>20</sup> <https://www.theguardian.com/technology/2016/nov/15/facebook-fake-news-us-election-trump-clinton>

<sup>21</sup> <https://www.pewresearch.org/internet/2017/10/19/the-future-of-truth-and-misinformation-online/>

<sup>22</sup> <https://www.gov.uk/government/consultations/online-harms-white-paper/online-harms-white-paper>

<sup>23</sup> <https://uhra.herts.ac.uk/handle/2299/21431>