**Supporting Document to C4C's Response to the EC Consultation on the White Paper on 'Artificial Intelligence – A European Approach'**

## INTRODUCTION

The Copyright for Creativity (C4C) Coalition – more details below – welcomes the opportunity to respond to the European Commission's consultation on the White Paper on 'Artificial Intelligence - A European Approach'. For further questions on our submission, please do not hesitate to contact Caroline De Cock, C4C's Coordinator, at cdc@n-square.eu.

## ABOUT C4C

C4C is a broad-based coalition that seeks an informed debate on how copyright can more effectively promote innovation, access, and creativity. C4C brings together libraries, scientific and research institutions, digital rights groups, technology businesses, and educational and cultural heritage institutions that share a common view on copyright. See our members here: http://copyright4creativity.eu/about-us/.

---

**Section 1 - An ecosystem of excellence**

**To build an ecosystem of excellence that can support the development and uptake of AI across the EU economy, the White Paper proposes a series of actions.**

**Are there other actions that should be considered?**

---

## OUR RECOMMENDATIONS

The use of artificial intelligence (AI) is starting to become more widespread, and with it the complaints about and examples of bias and lack of fairness of AI systems are also rapidly growing.[1]

Therefore, we would like to start our submission by emphasising an important point raised by the International Federation of Library Associations and Institutions (IFLA) in their submission to the draft WIPO issues paper on AI and intellectual property (IP),[2] namely that:

> "(…) regardless of the type of intellectual property right used to protect algorithms, it is important to address the question of the use that others may make of them. IP rules should not be used to create secrecy, to prevent the testing of algorithms and consultation of their source code in the interests of exploring their workings, or to take preservation copies in order to allow for archiving and future access."

With regards to the interplay between AI and IP, we observe that the Commission's AI White Paper neglects an important facet in the training of AI systems, namely the role of text and data mining (TDM) technologies and the surrounding legislative copyright framework. This is of crucial importance, because **limitations to what can be freely mined risks amplifying AI's bias and unfairness, and could negatively impact the usefulness of AI applications and projects.**

The 'BlueDot project', which uncovered the COVID-19 outbreak, is an interesting example of how TDM and AI projects need to be able to rely on the use of copyright-protected work. This project analysed "a variety of information sources, including chomping through 100,000 news reports in 65 languages a day", data which is then "compared with flight records to help predict virus outbreak patterns".[3]

Therefore, it is important to ensure that European AI actors can obtain adequate access to in-copyright material and this without impediments at the legal and/or technical level. This aspect and other issues that we have identified are set out in more details below.

### THE IMPORTANCE OF COPYRIGHT LEGISLATION TO HELP AVOIDING THE BIASED, LOW-FRICTION DATA TRAP

Amanda Levendowski, Associate Professor of Law at Georgetown University, observes that "**copyright law has the power to bias AI systems, but copyright law also has the profound power to unbias them**".[4]

This is an important element to take into account in the EU's approach to AI and for EU Member States in the implementation of the TDM provisions – Articles 3 and 4 – of the recently adopted Directive on copyright in the Digital Single Market[5] ('the Copyright Directive'), which should be transposed by 7 June 2021.

A restrictive framework to engage in TDM activities is bound to intensify negative externalities surrounding AI, impacting on users' fundamental rights,[6] and this to the detriment of European AI development. Therefore, we recommend that the Commission pays close attention to TDM technologies and encourages EU Member States to take an open and non-restrictive approach in implementing the Copyright Directive's TDM provisions.

Levendowski explains that "AI systems are commonly 'taught' by reading, viewing, and listening to copies of works created by humans", adding that "many of those works are protectable by copyright law".[7] She warns that even "impeccable algorithms will nevertheless generate biased results if reliant on biased data", the so-called 'garbage in, garbage out' problem.[8] As a result, Levendowski observes that:[9]

> "Copyright law causes friction that limits access to training data and restricts who can use certain data. This friction is a significant contributor to biased AI. The friction caused by copyright law encourages AI creators to use biased, low-friction data (BLFD) for training AI systems (…) despite those demonstrable biases. (…)"

The 'BLFD'-trap implies that "the biases encoded in BLFD (…) are picked up by AI systems trained using those data".[10] In Levendowski's view, public domain works and Creative Commons-licensed works, which could be used as alternative training data to overcome copyright hurdles, should be included in the BLFD category because these also contain a lot of embedded biases.[11]

In the context of public domain works, Levendowski warns that "this flavor of BLFD can easily erase entire perspectives and replicate the biases of a more homogenous authorship and less tolerant society".[12] She explains that: "A dataset reliant on works published before 1923 would reflect the biases of that time, as would any AI system trained with using that dataset."[13]

This exemplifies how restricted access to in-copyright works can hinder, and even negatively impact, the training of AI applications in fields such as natural language processing, if the basis for the training data is limited to only older out-of-copyright works instead of focussing today's broad repertoire of modern (progressive) works that shape today's society.

In order to avoid this 'BLFD'-trap, Levendowski concludes that: "**If we hope to create less biased commercial AI systems, using copyright-protected works as AI training data will be key**."[14]

Finally, Levendowski also observes how copyright can stifle innovation and facilitate anti-competitive behaviour in this area, which is especially detrimental towards startups, including in Europe: "Without the resources to get the vast amounts data easily acquired by major AI players, meaningful competition becomes all but non-existent."[15]

## TECHNICAL PROTECTION MEASURES (TPMs) SHOULD BE REMOVED WITHIN 72 HOURS WHEN THEY HINDER LAWFUL ACCESS TO CONTENT

Next to any legal impediments to TDM activities, there are also potential technical pitfalls that can stop or severely hinder the use of content for AI training purposes. The Copyright Directive introduces protections against contractual override for newly introduced exceptions, including for the TDM exceptions. However, as long as this protection can still be bypassed through the introduction of technical protection measures (TPM), such as digital right management systems (DRMs), it will be rendered useless.[16] Therefore, **the Commission and EU Member States need to pay careful attention to ensure that in the transposition of the Copyright Directive mechanisms are put in place to guarantee that TPMs are removed within 72 hours when they hinder lawful access to content**.

According to a survey by LIBER, the Association of European Research Libraries, resolving such issues now takes-up from 1 week to over 2.5 months.[17] Currently, all negative consequences are to the end-user's detriment. As a result, clear mechanisms are needed to guarantee the removal of technical protection measures within 72hrs when hindering lawful access to content. Furthermore, financial compensations should apply if access to purchased content is not restored within the set timeframe. Moreover, a reasonable cut-off timeframe should be imposed, after which institutions and users should be empowered to legally circumvent technical protection measures when they have lawful access to content.

Such mechanism can be implemented through Article 6 of the InfoSoc Directive.[18] This would help overcome the fact that existing mechanisms to deal with technical protection measures are often unclear, non-transparent and time-consuming. These protective measures will be especially important in the context of the text and data mining (TDM) exceptions in Articles 3 and 4 of the Copyright Directive.

## AVOID THE (PREMATURE) DELETION OF DATASETS

It is important to avoid the (premature) deletion of the datasets used to train AI algorithms, in order to:

1. foster more transparency as regards these algorithms, notably to ensure they do not comprise undue biasses.
2. ensure the possibility to verify the  scientific soundness of the algorithm used and/or the conclusions it may reach, where appropriate.

When these datasets stem from TDM activities, there is a potential risk that such (premature) deletion could be mandated at national level. This will depend on the national implementations of the TDM provisions in the Copyright Directive, and in particular Article 4(2).[19] Therefore, we consider that the retention period of the datasets should be allowed for the entire process, including preservation for evidence and reverse engineering purposes. It would hence be best that national implementations do not define specific retention periods, and remains as close as possible to the original language of Article 4(2).

Such an approach would allow the retention of datasets used for AI training, and as such ensure their availability for in-depth assessments of AI algorithms' fairness, be it internally or through an auditing body. It would also make it possible to fine-tune or further refine the training data if bias issues are uncovered.

### OVERCOMING THE LIMITATIONS OF ARTICLE 4 OF THE COPYRIGHT DIRECTIVE

In order to facilitate the use of the general TDM exception under Article 4 of the Copyright Directive, a number of other pitfalls, next to avoiding the (premature) deletion of datasets – see above, should also be tackled in the national transposition process, namely the need to:

– **Ensure that beneficiaries of Article 3, the 'scientific research' TDM exception, can benefit from Article 4 when engaging in activities not meeting the threshold of the 'for the purposes of scientific research' criterion:** Article 3 beneficiaries (*i.e.* research organisations and cultural heritage institutions) should also be covered by Article 4 when facing the limitations of the 'for the purposes of scientific research' criterion, which could support them in their efforts to engage in the development of AI solutions.

– **Enforce a machine-readable opt-out mechanism:** The legal text should clarify that a machine-readable opt-out mechanism should be used for explicit rights reservation under Article 4(3). Such a mechanism should:

1) Be mandatory for all types of digital content, both online and offline: Mandatory machine-readable opt-out mechanisms should be the norm for all types of minable content, and not be limited to online content. This would:

   a) enable TDM algorithms to check rights reservations before engaging in mining activities; and,
   b) avoid the need for manual checks that entail a substantial time and resource investment, which increases overhead costs and defeats the purpose of engaging in TDM.

   If such mechanism cannot be used and content is licenced, TDM activities should be deemed as allowed by default, unless the licencing terms unambiguously stipulate that TDM activities are prohibited.

2) Not lead to blanket/automatic rights reservation: Schemes to try and licence TDM of online content are likely to prove complicated and inefficient due to the huge number of rightholders. Therefore, it has to be avoided at all cost that such a mechanism opens the door for abuses from rightholders through, for example, blanket or automatic rights reservation, which is not a hypothetical threat.[20]

3) Not impact other uses, including TDM activities covered under Article 3: Recital 18 DCDSM clearly states that: "Other uses should not be affected by the reservation of rights for the purposes of text and data mining."

4) Preferably be uniform across the EU: A uniform opt-out mechanism across the EU would be favourable, especially in cross-border contexts: one coherent regime could not only benefit users, but also rightholders.

– **Ensure an inclusive stakeholder dialogue including Article 4 beneficiaries:** Article 3(4) stipulates that rightholders and the beneficiaries of the exception should define commonly agreed best practices. Discussions regarding best practices could be best facilitated through an all-inclusive stakeholder dialogue. Article 4 beneficiaries should also be included, to ensure all possible perspectives are met. This could facilitate a broader dialogue, and also help to (un)cover other aspects that require attention, such as rightholders' expressly reserving their rights (cf. Article 4(3) – see above)

– **Reject any attempts to make this exception subject to remuneration:** It is important to clarify that uses under the Article 4 exception are not subject to remuneration.[21]

---

[1] Murray, J. (2019, 24 April). Racist Data? Human Bias is Infecting AI Development. Towards Data Science. Available at, https://towardsdatascience.com/racist-data-human-bias-is-infecting-ai-development-8110c1ec50c?gi=b625fa921eca.

[2] IFLA, (2020, February). IFLA Response to the Draft WIPO Issues Paper - Artificial Intelligence and Intellectual Property. Available at, https://www.ifla.org/publications/node/92895.

[3] See Prosser, Marc, (2020, 5 February). How AI Helped Predict the Coronavirus Outbreak Before It Happened. Singularity Hub. Available at, https://singularityhub.com/2020/02/05/how-ai-helped-predict-the-coronavirus-outbreak-before-it-happened/; and, Flynn, Sean, Geiger, Christophe, & Quintais, João Pedro, (2020, 21 April). Implementing User Rights for Research in the Field of Artificial Intelligence: A Call for Action at International Level. Kluwer Copyright Blog. Available at, http://copyrightblog.kluweriplaw.com/2020/04/21/implementing-user-rights-for-research-in-the-field-of-artificial-intelligence-a-call-for-action-at-international-level/.

[4] Levendowski, A. (2018). How Copyright Law Can Fix Artificial Intelligence's Implicit Bias Problem. *93 Wash. L. Rev. 579*. Available at, https://digitalcommons.law.uw.edu/wlr/vol93/iss2/2/. p. 630.

[5] Directive (EU) 2019/790 of the European Parliament and of the Council of 17 April 2019 on copyright and related rights in the Digital Single Market and amending Directives 96/9/EC and 2001/29/EC (2019), Official Journal of the European Union L130, p. 92ff. Available at, http://data.europa.eu/eli/dir/2019/790/oj.

[6] The European Union Agency for Fundamental Rights (FRA) highlights that: "Algorithms used in machine learning systems and artificial intelligence (AI) can only be as good as the data used for their development. (…) AI systems based on incomplete or biased data can lead to inaccurate outcomes that infringe on people's fundamental rights, including discrimination." – See FRA. (2019). Data quality and artificial intelligence – mitigating bias and error to protect fundamental rights. Available at, https://fra.europa.eu/en/publication/2019/artificial-intelligence-data-quality. p. 1.

[7] Levendowski, A. (2018). *Ibid.* p. 582.

[8] Levendowski, A. (2018). *Ibid.* p. 585.

The 2016 report "Preparing for the Future of Artificial Intelligence", authored by the US National Science and Technology Council (NSTC), also observed that: "AI needs good data. If the data is incomplete or biased, AI can exacerbate problems of bias." See Executive Office of the President. (2016). Preparing for the Future of Artificial Intelligence, Washington DC. Available at, https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/microsites/ostp/NSTC/preparing_for_the_future_of_ai.pdf. p. 30.

[9] Levendowski, A. (2018). *Ibid.* p. 589.

[10] Levendowski, A. (2018). *Ibid.* p. 613.

[11] Levendowski, A. (2018). *Ibid.* p. 613-614.

[12] Levendowski, A. (2018). *Ibid.* p. 616.

Similarly, with regards to Creative Commons content, Levendowski observes that: "Because Wikipedia is easily discoverable, fully machine readable, and CC licensed, its articles are especially appealing as training data for AI systems. (…) Wikipedia also has a significant gender imbalance: in 2011, only 8.5% of Wikipedia editors were women. The editorship gender gap has measurable effects on the content of Wikipedia articles" – see p. 618.

[13] Levendowski, A. (2018). *Ibid.* p. 615.

[14] Levendowski, A. (2018). *Ibid.* p. 621.

[15] Levendowski, A. (2018). *Ibid.* p. 609.

[16] Christophe Geiger, Giancarlo Frosio, and Oleksandr Bulayenko, from the Centre for International Intellectual Property (CEIPI), observe in the context of TDM that: "(…) given that dominant market players customarily override exceptions by imposing both contractual and technological measures—depriving users of the enjoyment of exceptions and lawful uses—limitations to technological blocking should have been introduced as well by clearly spelling out that both TPMs and network security and integrity measures should not undermine the effective application of the exception. Accordingly, protection against contractual and technological override should also be clearly extended to TDM mining materials not protected by IPRs, including those made available in a database." See: Geiger, C., Frosio, G., & Bulayenko, O., (2019). Text and Data Mining: Articles 3 and 4 of the Directive 2019/790/EU. *Center for International Intellectual Property Studies Research Paper No. 2019-08*. Available at, https://ssrn.com/abstract=3470653. pp. 36-37.

[17] For more details, see LIBER, (2020, 10 March). Europe's TDM Exception for Research: Will It Be Undermined By Technical Blocking From Publishers?. Available at, https://libereurope.eu/blog/2020/03/10/tdm-technical-protection-measures/.

[18] Directive 2001/29/EC of the European Parliament and of the Council of 22 May 2001 on the harmonisation of certain aspects of copyright and related rights in the information society (2001), Official Journal of the European union L167, p. 10ff. Available at, http://data.europa.eu/eli/dir/2001/29/oj.

[19] Article 4(2) of the copyright Directive (EU 2019/790): "Reproductions and extractions made pursuant to paragraph 1 may be retained for as long as is necessary for the purposes of text and data mining."

[20] Rossana Ducato and Alain M. Strowel, from the Interdisciplinary Research Center Jean Renauld Law Enterprise and Society at UCLouvain, analysed in 2018 the terms and conditions (T&C) of twenty-one online platforms, equally distributed among three sectors: mobility, accommodation and food. They observed a trend toward a general contractual ban of TDM, often very broadly worded: "The prohibition is broad and refers to all the website's contents and services, thus including the informative pages containing the legal conditions." More specifically, they found that "20 out of 21 platforms published the T&C on their website and

---

14 of them contained specific intellectual property clauses, directly or indirectly, related to TDM activities" – see page 22: Ducato, R. & Strowel, A., (2018, 31 October). Limitations to Text and Data Mining and Consumer Empowerment: Making the Case for a Right to "Machine Legibility", CRIDES Working Paper Series. Available at, https://ssrn.com/abstract=3278901.

[21] Recital 17 clearly states that "Member States should (…) not provide for compensation for rightholders as regards uses under the text and data mining exceptions introduced by this Directive", which covers both Articles 3 and 4.