

Seton Hall University

eRepository @ Seton Hall

Law School Student Scholarship

Seton Hall Law

2021

Friction to Fight Misinformation: Content-Neutral Frictive Measures Under International Law

Robert J. Garcia

Follow this and additional works at: https://scholarship.shu.edu/student_scholarship



Part of the Law Commons

FRICION TO FIGHT MISINFORMATION:

CONTENT-NEUTRAL FRICTIVE MEASURES UNDER INTERNATIONAL LAW

INTRODUCTION

The Special Rapporteur on the Freedom of Opinion and Expression¹ (Special Rapporteur) is the position created by the U.N. Human Rights Council to address violations of the freedoms of opinion and expression, undertake fact-finding visits to countries, and publish annual reports relating to the freedom of opinion and expression.² In the Special Rapporteur's upcoming report on Disinformation and the freedom of expression,³ they should advocate the use of content neutral frictive measures in combating misinformation on social media.

The first section of this paper provides a synopsis of the Special Rapporteur's six reports on the freedom of opinion and expression in the Information, Communication, Technology sector, including Social Media Companies. The first section also discusses the Special Rapporteur's most recent report on disease pandemics. The second section details the international legal framework based on the International Covenant on Civil and Political Rights and the United Nations Guiding Principles on Business which is reiterated in the Special Rapporteur's reports. The third section defines frictive measures and what the Special Rapporteur has said about them in their reports. The fourth section argues that content neutral frictive measures comply international law, despite

¹United Nations Commission on Human Rights, *QUESTION OF THE HUMAN RIGHTS OF ALL PERSONS SUBJECTED TO ANY FORM OF DETENTION OR IMPRISONMENT*, U.N. Doc. E/CN.4/1993/L.48 (1993) (establishing the mandate of the Special Rapporteur on the and protection of the right to freedom of opinion and expression.).

² Office of the High Commissioner for Human Rights, <https://www.ohchr.org/EN/Issues/FreedomOpinion/Pages/mandate.aspx> (last visited May 2, 2021)

³ Irene Kahn (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *Disinformation and freedom of opinion and expression*, U.N. Doc. A/HRC/47/25 (Forthcoming June, 2021).

the potential for bias and present a viable method of combating disinformation which the Special Rapporteur should advocate in their upcoming report.

I. SYNOPSIS OF THE SPECIAL RAPPORTEUR REPORTS

Since 2015, the Special Rapporteur has published six reports on freedom of opinion and freedom of expression in the Information and Communication Technology sector.⁴ The Information and Communication Technology sector consists of the private actors involved in “organizing, accessing, populating and regulating the Internet.”⁵ Social media companies make up part of the ICT sector.⁶ The Special Rapporteur has written reports on the freedoms of opinion and expression in the ICT sector pertaining to: encryption and anonymity;⁷ role of states and the private sector in the digital age;⁸ the role of digital access providers;⁹ online content regulation;¹⁰ artificial intelligence;¹¹ and online hate speech.¹² In response to the COVID-19 pandemic, the Special

⁴ David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *Online Hate Speech*, ¶ 3, U.N. Doc. A/74/486 (Oct. 19, 2019).

⁵ David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *Freedom of Expression, States and the Private Sector in the Digital Age*, ¶ 15, U.N. Doc. A/HRC/32/38 (May 11, 2016) [hereinafter *States and the Private Sector in the Digital Age*].

⁶ Kaye, *States and the Private Sector in the Digital Age*, *supra* note 5, at ¶ 1.

⁷ David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *The Use of Encryption and Anonymity to Exercise the Rights to Freedom of Opinion and Expression in the Digital Age*, U.N. Doc. A/HRC/29/32 (May 22, 2015) [hereinafter *Encryption and Anonymity*].

⁸ Kaye, *States and the Private Sector in the Digital Age*, *supra* note 5.

⁹ David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *The Role of Digital Access Providers*, U.N. Doc. A/HRC/35/22 (May 30, 2017).

¹⁰ David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *Online Content Regulation*, U.N. Doc. A/HRC/38/35 (Apr. 6, 2018).

¹¹ David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *Artificial Intelligence Technologies and Implications for the Information Environment*, U.N. Doc. A/73/348 (Aug. 29, 2018) [hereinafter *Artificial Intelligence*].

¹² David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *Online Hate Speech*, U.N. Doc. A/74/486 (October 9, 2019).

Rapporteur published in April of 2020 on disease pandemics which discussed public health disinformation online.¹³

Encryption and Anonymity:

The Special Rapporteur's report on Encryption and Anonymity considers whether the rights to privacy, freedom of opinion and expression, include the ability to encrypt or anonymize.¹⁴ The Rapporteur defines encryption as “a mathematical ‘process of converting messages, information, or data into a form unreadable by anyone except the intended recipient...’”¹⁵ Anonymity is defined as the “condition of avoiding identification.”¹⁶ The report addressed concerns that encryption and anonymity would hide criminal activity online.¹⁷ While the report acknowledges that corporations have a duty to operate under the UN Guiding Principles on Business, the report specifically “is focused on State obligations.”¹⁸ Ultimately, encryption and anonymity are deemed “necessary for the exercise of the right to freedom of opinion and expression in the digital age.”¹⁹

States and the Private Sector in the Digital Age:

In A/HRC/32/38, the Special Rapporteur acknowledges the influence of private actors on freedom of expression through the internet and social media.²⁰ The report raises questions concerning the protection of freedom of opinion and expression as relating to the private sector

¹³ David Kaye (Special Rapporteur on the promotion and protection of freedom of opinion and expression), *Disease Pandemics and the Freedom of Opinion and Expression*, U.N. Doc. A/HRC/44/49 (Apr. 23, 2020) [hereinafter *Disease Pandemics*].

¹⁴ *Encryption and Anonymity*, *supra* note 7, at ¶ 3.

¹⁵ *Encryption and Anonymity*, *supra* note 7, at ¶ 7.

¹⁶ *Encryption and Anonymity*, *supra* note 7, at ¶ 9.

¹⁷ *Encryption and Anonymity*, *supra* note 7, at ¶ 13.

¹⁸ *Encryption and Anonymity*, *supra* note 7, at ¶ 17.

¹⁹ *Encryption and Anonymity*, *supra* note 7, at ¶ 56.

²⁰ *States and the Private Sector in the Digital Age*, *supra* note 5, at ¶ 3.

online.²¹ Finally, the report concludes with the Special Rapporteur concludes by offering some “normative guidance” in the most needed areas of the ICT sector.²²

The Special Rapporteur questions companies’ terms of service,²³ and raises derivative concerns from the lack of clarity in terms of service including “inconsistent enforcement”,²⁴ “overzealous censorship”,²⁵ “lack of an appeals process”,²⁶ and states’ opportunistically using such ambiguity to remove objectionable content.²⁷ For private actors, the primary source of regulation on their platforms is their terms of service.²⁸ These terms of service are formulated in such a way that makes it difficult for users to predict what content is restricted on the private actor’s platform.²⁹ Moving to content regulation, the Special Rapporteur examines the “design and engineering choices” of social media, recognizing that social media companies’ curation of content affects user’s access to information.³⁰ Such inconsistent enforcement and limitation of information poses threats to the freedom of opinion and expression.³¹

The Role of Digital Access Providers:

Beginning in 2016, the Special Rapporteur began the process of detailing the different facets of the “information and communications technology (ICT) sector.”³² The report outlines state obligations, under international law, to protect expression online.³³ The Rapporteur condemns

²¹ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 3.

²² *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 3.

²³ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 37.

²⁴ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 52.

²⁵ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 52.

²⁶ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 52.

²⁷ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 53.

²⁸ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 52.

²⁹ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 52.

³⁰ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 55.

³¹ *States and the Private Sector in the Digital Age*, supra note 5, at ¶ 52.

³² *The Role of Digital Access Providers*, supra note 9, at ¶ 4.

³³ *The Role of Digital Access Providers*, supra note 9, at ¶ 5.

internet and telecommunications shutdowns³⁴ and cautions that accessing online user information can interfere with privacy rights.³⁵ The role of digital access providers is also discussed, and the Rapporteur recognizes that access as critical to the freedom of expression.³⁶

Online Content Regulation:

The Special Rapporteur's report on online content regulation "...focuses on the regulation of user-generated content, principally by States and social media companies..."³⁷ Moderation describes "the process by which Internet companies determine whether user-generated content meets the standards articulated in their terms of service and other rules."³⁸ The report expresses concerns about national laws restricting speech across borders.³⁹ Because many social media companies have an international user base, "national laws are inappropriate for companies that seek common norms."⁴⁰ In lieu of national laws, the Special Rapporteur advocates compliance with the UN Guiding Principles on Business as an international framework for the content moderation policies of social media companies.⁴¹

The UN Guiding Principles on Business "and their accompanying body of 'soft law'",⁴² inform the Rapporteur's substantive standards for content moderation.⁴³ When developing standards, companies should seek policy that allows platforms "...for users to develop opinions, express themselves freely and access information of all kinds in a manner consistent with human

³⁴ *The Role of Digital Access Providers*, *supra* note 9, at ¶ 8.

³⁵ *The Role of Digital Access Providers*, *supra* note 9, at ¶ 17.

³⁶ *The Role of Digital Access Providers*, *supra* note 9, at ¶ 29.

³⁷ *Online Content Regulation*, *supra* note 10 at, ¶ 3.

³⁸ *Online Content Regulation*, *supra* note 10 at, ¶ 3.

³⁹ *Online Content Regulation*, *supra* note 10 at, ¶ 18.

⁴⁰ *Online Content Regulation*, *supra* note 10 at, ¶ 41.

⁴¹ *Online Content Regulation*, *supra* note 10 at, ¶ 6.

⁴² *Online Content Regulation*, *supra* note 10 at, ¶ 42.

⁴³ *Online Content Regulation*, *supra* note 10 at, ¶ 44,

rights law.”⁴⁴ Part of content moderation policy should include a standard of non-discrimination which requires companies to “transcend formalistic approaches” and “take into account the concerns of communities historically at risk of censorship and discrimination.”⁴⁵

The report ends with the recommendations of Special Rapporteur, including recognition that human rights law is “the authoritative global standard for ensuring freedom of expression” on social media,⁴⁶ and “smart” content moderation as opposed to “heavy-handed viewpoint-based regulation.”⁴⁷

AI and Free Speech:

The Special Rapporteur’s report on AI and Free Speech has three goals: “define key terms essential to a human rights discussion about AI; identify the human rights legal framework relevant to AI; and present some preliminary recommendations to ensure that, as the technologies comprising AI evolve, human rights considerations are baked into that process.”⁴⁸ Artificial Intelligence is defined as “a ‘constellation’ of processes and technologies enabling computers to complement or replace specific tasks otherwise performed by humans, such as making decisions and solving problems.”⁴⁹ The Rapporteur urges caution against removing human intervention from content moderation through the promotion of artificial intelligence.⁵⁰

⁴⁴ *Online Content Regulation*, *supra* note 10 at, ¶ 45.

⁴⁵ *Online Content Regulation*, *supra* note 10 at, ¶ 48.

⁴⁶ *Online Content Regulation*, *supra* note 10 at, ¶ 70.

⁴⁷ *Online Content Regulation*, *supra* note 10 at, ¶ 66.

⁴⁸ *Artificial Intelligence*, *supra* note 11 at, ¶ 2.

⁴⁹ *Artificial Intelligence*, *supra* note 11 at, ¶ 3.

⁵⁰ *Artificial Intelligence*, *supra* note 11 at, ¶ 6.

Concerning AI and Social Media, the Special Rapporteur notes areas of concern in content display, personalization,⁵¹ content moderation and removal.⁵² Because Social media companies tailor content based on user preferences, users might experience a total absence of “diverse views, interfering with individual agency to seek and share ideas and opinions across ideological, political or societal divisions.”⁵³ The Special Rapporteur reaffirms their advocacy for an international legal framework for Social Media companies to follow, calling attention to the right of freedom of opinion which “requires freedom from undue coercion.”⁵⁴

Online Hate Speech:

The Special Rapporteur’s report on Hate Speech Online discusses “Governments considering regulatory options and companies determining how to respect human rights online.”⁵⁵ While the report offers no concrete definition of Hate Speech, the Rapporteur cites ICCPR Article 20 Section, which prohibits “any advocacy of national, racial or religious hatred that constitutes incitement to discrimination, hostility or violence”.⁵⁶ Supplementing the definition, the report cites the 2013 Rabat Plan definitions of “Hatred’ and “hostility” refer to intense and irrational emotions of opprobrium, enmity and detestation towards the target group.”⁵⁷ The Rapporteur clarifies that hateful expression does not always “constitute advocacy or incitement.”⁵⁸

Social Media companies are criticized in the report for operating with a business model that values the spread of hateful content.⁵⁹ Simultaneously, the Special Rapporteur warns of

⁵¹ *Artificial Intelligence*, *supra* note 11 at, ¶ 10.

⁵² *Artificial Intelligence*, *supra* note 11 at, ¶ 13.

⁵³ *Artificial Intelligence*, *supra* note 11, at ¶ 12.

⁵⁴ *Artificial Intelligence*, *supra* note 11, at ¶ 23.

⁵⁵ *Online Hate Speech*, *supra* note 12, at ¶ 2.

⁵⁶ *Online Hate Speech*, *supra* note 12, at ¶ 8.

⁵⁷ *Online Hate Speech*, *supra* note 12, at ¶ 6.

⁵⁸ *Online Hate Speech*, *supra* note 12, at ¶ 20.

⁵⁹ *Online Hate Speech*, *supra* note 12, at ¶ 40.

ambiguous definitions of hate speech in Social Media policy and of automated content moderation that “is notoriously bad at evaluating context.”⁶⁰ To improve their moderation policies, the Rapporteur recommends “de-amplification, de-monetization, education, counter-speech.”⁶¹

Disease Pandemics:

In the wake of the COVID-19 Pandemic, the Special Rapporteur issued a report detailing five challenges to freedom of opinion and expression.⁶² One of the five challenges discussed is the spread of Public Health Disinformation.⁶³ The Special Rapporteur recognized the need to engage “rumors in order to correct them” whilst cautioning against disproportionate punishment for sharing disinformation.⁶⁴ The Special Rapporteur also acknowledged Social Media’s “enormous impact on public discourse and the rights of individuals on and off their platforms.”⁶⁵ To assist in stemming the flow of disinformation, Social Media companies “should aim towards maximum transparency of their policies and engage” with both public authorities and affected communities.⁶⁶

II.INTERNATIONAL LAW FRAMEWORK

The Special Rapporteur calls upon social media companies to tailor their content moderation policy based on an international law framework. National laws are an insufficient basis for content moderation policy because Social Media Companies involve a “geographically and culturally diverse user base.”⁶⁷ The Special Rapporteur presents an international law framework

⁶⁰ *Online Hate Speech*, *supra* note 12, at ¶ 50.

⁶¹ *Online Hate Speech*, *supra* note 12, at ¶ 58.

⁶² *Disease Pandemics*, *supra* note 13, at ¶ 6.

⁶³ *Disease Pandemics*, *supra* note 13, at ¶ 41.

⁶⁴ *Disease Pandemics*, *supra* note 13, at ¶ 42.

⁶⁵ *Disease Pandemics*, *supra* note 13, at ¶ 52.

⁶⁶ *Disease Pandemics*, *supra* note 13, at ¶ 52.

⁶⁷ *Online Content Regulation*, *supra* note 10 at, ¶ 41.

in their reports based on the International Covenant of Civil and Political Rights,⁶⁸ and the United Nations Guiding Principles on Business and Human Rights.⁶⁹ The Special Rapporteur uniformly cites ICCPR Art. 19, General Comment 34 to Article 19, and UN Guiding Principles on business for social media companies and the ICT sector in general.⁷⁰

International Covenant of Civil and Political Rights, Art. 19

The purpose of the ICCPR, explained in its preamble, is to promote “civil and political freedom...”⁷¹ The ICCPR was adopted by the United Nations General Assembly in 1966.⁷² Article 19 of the ICCPR protects the rights to freedom of expression and freedom of opinion.⁷³

As part of the rights to hold opinions and freedom of expression, Article 19 delineates important qualifications clarifying the breadth of those rights.⁷⁴ Not only does the Art. 19 protect the right to hold opinions, but to “hold opinions without interference.”⁷⁵ The right to freedom of expression under Art. 19, includes “freedom to seek, receive and impart information and ideas of all kinds...through any other media of his choice.”⁷⁶ However, Art. 19 also delineates when restrictions to the freedom of expression are acceptable.

Art. 19 allows derogation to protect the rights and reputation of others, “national security...public order...public health or morals.”⁷⁷ Under ICCPR Art. 4, states may derogate from

⁶⁸ UN General Assembly, *International Covenant on Civil and Political Rights*, 16 December 1966, United Nations, Treaty Series, vol. 999, p. 171 [hereinafter ICCPR].

⁶⁹ UN Office of High Commissioner on Human Rights, *Guiding Principles on Business and Human Rights*, U.N. Doc. HR/PUB/11/04 (June 16, 2011) [hereinafter Guiding Principles].

⁷⁰ See *States and the Private Sector in the Digital Age*, *supra* note 2, at ¶ 5 and ¶ 9; *The Role of Digital Access Providers*, *supra* note 6, at ¶ 5 and ¶ 45; *Online Hate Speech*, *supra* note 9 at ¶ 5 and ¶ 47.

⁷¹ ICCPR, *supra* note 68, at 1.

⁷² ICCPR, *supra* note 68, at 1.

⁷³ ICCPR, *supra* note 68, at art. 19.

⁷⁴ ICCPR, *supra* note 68, at art. 19.

⁷⁵ ICCPR, *supra* note 68, at art. 19.

⁷⁶ ICCPR, *supra* note 68, at art. 19.

⁷⁷ ICCPR, *supra* note 68, at art. 19(3)(a)-(b).

their obligations under the treaty excepting some articles.⁷⁸ Such derogation cannot be inconsistent with international law and cannot “involve discrimination solely on the ground of race, colour, sex, language, religion or social origin.”⁷⁹ Restrictions under Art. 19 must be “provided by law and...necessary.”⁸⁰ Restrictions can be applied to the freedom of expression, but not the freedom of opinion.⁸¹ To further clarify the extent of protections under Art. 19 and the requirements for restrictions on the freedom of expression, the Human Rights Committee of the United Nations issued General Comment No. 34.⁸²

General Comment No. 34

General Comment No. 34 was published in July 2011.⁸³ While the ICCPR binds state parties, the obligations under Art. 19 require states to protect people “from any acts by private persons or entities” infringing upon the freedoms of opinion and expression.⁸⁴

Comment No. 34 addresses the right to freedom of expression in media,⁸⁵ to access information,⁸⁶ to political rights,⁸⁷ and the scope of restrictions to expression.⁸⁸ “Internet-based modes of expression” are protected by Art. 19. The Comment reiterates that no restriction upon the holding of an opinion nor impairments on other rights because one holds an opinion are

⁷⁸ ICCPR, *supra* note 68, at art. 4(1).

⁷⁹ ICCPR, *supra* note 68, at art. 4(1).

⁸⁰ ICCPR, *supra* note 68, at art. 19.

⁸¹ See U.N. Human Rights Committee, *General Comment No. 34, Article 19 Freedom of Opinion and Expression*, ¶ 9, U.N. Doc. CCPR/C/GC/34 (Sept. 12, 2011) [hereinafter *General Comment No. 34*]; *States and the Private Sector in the Digital Age*, *supra* note 2 at ¶ 7.

⁸² *General Comment No. 34*, *supra* note 81, at ¶ 5.

⁸³ *General Comment No. 34*, *supra* note 81, at ¶ 5.

⁸⁴ *General Comment No. 34*, *supra* note 81, at ¶ 7..

⁸⁵ *General Comment No. 34*, *supra* note 81, at ¶ 13.

⁸⁶ *General Comment No. 34*, *supra* note 81, at ¶ 18.

⁸⁷ *General Comment No. 34*, *supra* note 81, at ¶ 20.

⁸⁸ *General Comment No. 34*, *supra* note 81, at ¶ 37.

permissible.⁸⁹ Part of the freedom to hold opinions is the prohibition on coercion to hold or not hold an opinion.⁹⁰

Comment No. 34 clarifies the requirements of Art. 19(3) restrictions of the freedom of expression.⁹¹ In addition to the requirements of legality, for the legitimate purposes listed in 19(3)(a-b), and the restrictions must be necessary and proportional.⁹²

The Comment goes into detail about what it means for a restriction to be required by law.⁹³ To be a law, means that a regulation is not merely discretionary, is sufficiently precise, and publicly available to allow people to tailor their conduct accordingly.⁹⁴

The extent of legitimate grounds—or legitimacy—for the restrictions to the freedom of expression are also clarified in the Comment.⁹⁵ Restrictions are permissible to protect rights under international law, both to individuals and members of a community.⁹⁶ While national security is a legitimate purpose for a restriction on expression, states must take care to ensure such restrictions take “extreme care” to conform with the requirements of Art. 19.⁹⁷ Finally, the comment explains that the legitimate grounds must not stem from any single social, philosophical or religious tradition.⁹⁸

⁸⁹ *General Comment No. 34, supra note 81, at ¶ 9.*

⁹⁰ *General Comment No. 34, supra note 81, at ¶ 10.*

⁹¹ *General Comment No. 34, supra note 81, at ¶ 35.*

⁹² *General Comment No. 34, supra note 81, at ¶ 13.*

⁹³ *General Comment No. 34, supra note 81, at ¶ 24.*

⁹⁴ *General Comment No. 34, supra note 81, at ¶ 25.*

⁹⁵ *General Comment No. 34, supra note 81, at ¶ 28.*

⁹⁶ *General Comment No. 34, supra note 81, at ¶ 28.*

⁹⁷ *General Comment No. 34, supra note 81, at ¶ 30.*

⁹⁸ *General Comment No. 34, supra note 81, at ¶ 32.*

For a restriction to be proportionate, it must not be overbroad.⁹⁹ Considering the form of expression restricted along with its method of dissemination, the restriction must be the least intrusive option to be proportionate.¹⁰⁰

To show the necessity of a restriction on the freedom on expression, there must be “a direct and immediate connection between the expression and the threat.”¹⁰¹

The Comment also outlines “certain specific areas” where restrictions should have a limited scope and are subject to greater scrutiny.¹⁰² The value placed on uninhibited expression relating to political discourse is particularly high by the ICCPR.¹⁰³ Accordingly, the Comment expresses the Committee’s concern for restrictions relating to political discourse.¹⁰⁴ Concerning websites and blogs, the Committee explains that restrictions by the state need to meet the Art. 19 requirements and cannot be based on prohibiting material critical of the government.¹⁰⁵ While Art.19 and its subsequent commentary detail state obligations including protections against private actors infringing upon the freedom of expression, the UN Guiding Principles directly address business’ role in protecting human rights.

UN Guiding Principles on Business and Human Rights

The UN Guiding Principles outline “The State Duty to Protect Human Rights”, “The Corporate Responsibility to Respect Human Rights”, and “Access to Remedy.”¹⁰⁶ In addition to reaffirming states’ duty to protect against human rights abuse,¹⁰⁷ when business enterprises are

⁹⁹ *General Comment No. 34, supra note 81, at ¶ 34.*

¹⁰⁰ *General Comment No. 34, supra note 81, at ¶ 34.*

¹⁰¹ *General Comment No. 34, supra note 81, at ¶ 34.*

¹⁰² *General Comment No. 34, supra note 81, at ¶ 37.*

¹⁰³ *General Comment No. 34, supra note 81, at ¶ 38.*

¹⁰⁴ *General Comment No. 34, supra note 81, at ¶ 38.*

¹⁰⁵ *General Comment No. 34, supra note 81, at ¶ 43.*

¹⁰⁶ *Guiding Principles, supra note 69 at iii.*

¹⁰⁷ *Guiding Principles, supra note 69 at 1.*

owned or controlled by the states, additional steps to protect against human rights abuse and oversight are encouraged.¹⁰⁸ Although the principles do not have the force of law, they provide a method of enhancing business practices to protect human rights.¹⁰⁹

The UN Guiding Principles state that “business enterprises should respect human rights.”¹¹⁰ For business to respect human rights at a minimum, they must not violate the International Covenant on Civil and Political Rights and the International Covenant on Economic, Social and Cultural Rights.¹¹¹ Respecting such rights includes preventing negative impacts to human rights through their business activities as well as from business relationships.¹¹² The responsibility to protect human rights regardless of the size of the business enterprise.¹¹³ However, businesses can craft their means of addressing and protecting human rights according to their size and circumstances.¹¹⁴

The Guiding Principles, recommend a policy commitment by businesses to protect human rights.¹¹⁵ The policy should be holistic, recognized at “the most senior level of the business”, informed by expertise, made publicly available and reflected through operational policies.¹¹⁶ Part of the policy should include a due diligence assessment to determine the potential adverse human rights impact of business practices.¹¹⁷ Upon assessment, business should implement the findings “across relevant and internal functions” affecting decision-making.¹¹⁸ The effectiveness of

¹⁰⁸ *Guiding Principles, supra note 69* at 4.

¹⁰⁹ *Guiding Principles, supra note 69* at 1.

¹¹⁰ *Guiding Principles, supra note 69* at princ. 11.

¹¹¹ *Guiding Principles, supra note 69* at princ. 12.

¹¹² *Guiding Principles, supra note 69* at princ. 13.

¹¹³ *Guiding Principles, supra note 69* at princ. 14.

¹¹⁴ *Guiding Principles, supra note 69* at princ. 15.

¹¹⁵ *Guiding Principles, supra note 69* at princ. 16.

¹¹⁶ *Guiding Principles, supra note 69* at princ. 16.

¹¹⁷ *Guiding Principles, supra note 69* at princ. 17.

¹¹⁸ *Guiding Principles, supra note 69* at princ. 19.

response should be tracked by the business.¹¹⁹ Communication between stakeholders effected by the human rights impact of businesses is expected.¹²⁰ When businesses come to learn that their practices have had adverse impacts on human rights, “they should provide for or cooperate in their remediation.”¹²¹

Finally, the UN Guiding Principles outline principles to provide access to remedy for persons experiencing an adverse impact on human rights from business practices.¹²² States have duty to provide remedy “through judicial, administrative, legislative, or other appropriate means.”¹²³ Businesses are recommended to have “effective operational-level grievance mechanisms” so effected persons and their communities can have direct access to remedy.¹²⁴ Operational-level grievance mechanisms help identify potential adverse impacts of business practices along with allowing grievances to be addressed and remedy before they escalate.¹²⁵ These mechanisms need not require a grievance that reaches the level of a human rights abuse, but rather should seek to act when legitimate concerns are identified.¹²⁶

Special Rapporteur on International Law

The Special Rapporteur identifies the freedoms protected in Article 19 as the basis for a wide range of human rights and are the foundation for free and democratic society.¹²⁷ The language

¹¹⁹ *Guiding Principles, supra note 69* at princ. 20.

¹²⁰ *Guiding Principles, supra note 69* at princ. 21.

¹²¹ *Guiding Principles, supra note 69* at princ. 22.

¹²² *Guiding Principles, supra note 69* at princ. 25.

¹²³ *Guiding Principles, supra note 69* at princ. 25.

¹²⁴ *Guiding Principles, supra note 69* at princ. 29.

¹²⁵ *Guiding Principles, supra note 69* at princ. 29.

¹²⁶ *Guiding Principles, supra note 69* at princ. 29.

¹²⁷ *See Online Hate Speech, supra note 4*, at ¶ 5; *Encryption and Anonymity, supra note 4*, at ¶ 22.

of Art. 19 is consciously devoid of a list of relevant media because the rights of freedom of opinion and expression accommodated future technological advances.¹²⁸

The Special Rapporteur uses the framework of Article 19 of the ICCPR to determine the conditions for restrictions on the freedom of expression to comply with international law.¹²⁹ The conditions are legality, necessity and proportionality, and Legitimacy.¹³⁰ The usefulness, reasonableness, and desirability of a restriction is neither dispositive nor sufficient to demonstrate any of requirements for the restriction under Article 19.¹³¹ The requirements for the restrictions are “to be applied strictly and in good faith, with robust and transparent oversight.”¹³² The burden of justifying the restriction falls on the authority imposing the restriction rather than the speakers subject to it.¹³³

A/HRC/29/32 includes a discussion on how the storage, transmission, and security of information in the digital age effect the right to freedom of opinion.¹³⁴ Because State and non-state actors have control over the storage, transmission, and security of information online, they both have the ability to interfere with rights in Art. 19 of the ICCPR. The right to freedom of opinion includes the ability to come to an opinion through reasoning, prohibiting undue coercion.¹³⁵ Offering preferential treatment to induce acceptance of an opinion could also constitute impermissible coercive conduct violating Art. 19.¹³⁶ The Special Rapporteur recognized in

¹²⁸ *Encryption and Anonymity*, *supra* note 4, at ¶ 26.

¹²⁹ *Online Content Regulation*, *supra* note 10 at ¶ 66; *Artificial Intelligence*, *supra* note 11, at ¶ 28; *Online Hate Speech*, *supra* note 12, at ¶ 6; *Disease Pandemics*, *supra* note 13, at ¶ 11.

¹³⁰ *Online Content Regulation*, *supra* note 10 at ¶ 7; *Artificial Intelligence*, *supra* note 11, at ¶ 28; *Online Hate Speech*, *supra* note 12, at ¶ 6; *Disease Pandemics*, *supra* note 13, at ¶ 11.

¹³¹ *Disease Pandemics*, *supra* note 13, at ¶ 15.

¹³² *Online Hate Speech*, *supra* note 12, at ¶ 7.

¹³³ *Online Hate Speech*, *supra* note 12, at ¶ 6.

¹³⁴ *Encryption and Anonymity*, *supra* note 7, at ¶ 12.

¹³⁵ *Artificial Intelligence*, *supra* note 11, at ¶ 23.

¹³⁶ *Artificial Intelligence*, *supra* note 11, at ¶ 23.

A/73/348 that content curation effects the ability to formulate opinions and “raises novel questions about the types of coercion or inducement that may be considered an interference.”¹³⁷

In recognizing the impact of ICT’s, the Special Rapporteur has consistently reiterated that companies should apply international law and “apply human rights principles in their operations.”¹³⁸ Rather than applying throughout business practices, most companies apply human rights principles in response to demands from the states.¹³⁹ In the A/HRC/38/35, the Special Rapporteur succinctly articulates the UN Guiding Principles’ minimum standards for business practices as an operative framework for companies to use:

(a) Avoid causing or contributing to adverse human rights impacts and seek to prevent or mitigate such impacts directly linked to their operations, products or services by their business relationships, even if they have not contributed to those impacts (principle 13);

(b) Make high-level policy commitments to respect the human rights of their users (principle 16);

(c) Conduct due diligence that identifies, addresses and accounts for actual and potential human rights impacts of their activities, including through regular risk and impact assessments, meaningful consultation with potentially affected groups and other stakeholders, and appropriate follow-up action that mitigates or prevents these impacts (principles 17–19);

(d) Engage in prevention and mitigation strategies that respect principles of internationally recognized human rights to the greatest extent possible when faced with conflicting local law requirements (principle 23);

(e) Conduct ongoing review of their efforts to respect rights, including through regular consultation with stakeholders, and frequent, accessible and effective communication with affected groups and the public (principles 20–21);

(f) Provide appropriate remediation, including through operational-level grievance mechanisms that users may access without aggravating their “sense of disempowerment” (principles 22, 29 and 31).¹⁴⁰

¹³⁷ *Artificial Intelligence*, *supra* note 11, at ¶ 24.

¹³⁸ *Online Content Regulation*, *supra* note 10 at, ¶ 9.

¹³⁹ *Online Content Regulation*, *supra* note 10 at, ¶ 10.

¹⁴⁰ *Online Content Regulation*, *supra* note 10 at, ¶ 11.

While this framework is not meant to carry the weight of law, it does provide an aspirational standard for businesses.¹⁴¹ One area that social media companies can assess their impact on the right to freedom of expression is the effect of frictionless sharing,¹⁴² on their platforms and how it leads to misinformation.¹⁴³

III. FRICTIVE MEASURES DEFINED

In media policy, there exists a concept called the “signal-to-noise ratio.”¹⁴⁴ Signal represents truthful information that supports democratic discourse.¹⁴⁵ Noise represents content that misinforms and undermines discourse.¹⁴⁶ As various kinds of content are amplified by algorithms on social media, it can become difficult to distinguish between signal and noise.¹⁴⁷ Friction can be introduced to help distinguish between the two.

Friction can be defined as the force that resists between two bodies.¹⁴⁸ To resist can be defined as “to exert force in opposition.”¹⁴⁹ Taken together, friction can be understood as the force exerted in opposition to the relative motion between two bodies in contact. In digital platforms seeking to streamline the flow of information, friction is an obstacle to be overcome.¹⁵⁰ However, whether friction is good or bad depends upon the force it is acting against.¹⁵¹ Friction can act “to disincentivize and disrupt practices that addict, surveil, and dull critical functions.”¹⁵² Frictive

¹⁴¹ *Guiding Principles*, *supra* note 69 at pg. 1.

¹⁴² William McGeeveran, *The Law of Friction*, 2013 U. CHI. LEGAL F. 15 at 19 (2013).

¹⁴³ Ellen P. Goodman, *Digital Information Fidelity and Friction*, Knight First Amendment Institute, 2, (2020) https://s3.amazonaws.com/kfai-documents/documents/c5cac43fec/2.27.2020_Goodman-FINAL.pdf.

¹⁴⁴ *Digital Information Fidelity and Friction*, *supra* 143 at 2.

¹⁴⁵ *Digital Information Fidelity and Friction*, *supra* 143 at 2.

¹⁴⁶ *Digital Information Fidelity and Friction*, *supra* 143 at 2.

¹⁴⁷ *Digital Information Fidelity and Friction*, *supra* 143 at 2.

¹⁴⁸ *The Law of Friction*, *supra* 142 at 15.

¹⁴⁹ Meriam-Webster, <https://www.merriam-webster.com/dictionary/resisting> (last visited May 5, 2021).

¹⁵⁰ *Digital Information Fidelity and Friction*, *supra* 143 at 20; *see also The Law of Friction*, *supra* 142 at 56.

¹⁵¹ *The Law of Friction*, *supra* 142 at 56.

¹⁵² *Digital Information Fidelity and Friction*, *supra* 143 at 3-4.

measures are forms of creating such friction and their goal is to “open pathways for reflection.”¹⁵³

While not exhaustive, the Knight Institute has identified three forms of friction relevant to digital platforms: “communication delays, virality disruptions, and taxes.”¹⁵⁴

Communication delays are a form of frictive measures that seek to systematize a pause.¹⁵⁵ There is some research to suggest that people are “more likely to resist manipulative communications” when they have the time to “raise cognitive defenses.”¹⁵⁶ In live broadcast media, short delays are already implemented for the sake of quality control.¹⁵⁷ In the stock market, the IEX ensure a degree of friction by running all trades through extra cable to avoid any trader from having an advantage by getting their information first.¹⁵⁸ A pause by way of communication delay gives people the time to think about the fidelity of information they are considering, to distinguish between signal and noise.¹⁵⁹

Virality Disruptors are a form of frictive measure that disrupt traffic once “a certain threshold of circulation” is reached.¹⁶⁰ Virality can be understood as the quality of “triggering quick, emotionally, intense responses.”¹⁶¹ Disrupting virality would involve a pause both for the user to process the incoming information and also, to allow content moderation through human review to ensure compliance of the disrupted communication.¹⁶² Moderation is defined as “the process by which internet companies determine whether user-generated content meets the

¹⁵³ *Digital Information Fidelity and Friction*, *supra* 143 at 3.

¹⁵⁴ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

¹⁵⁵ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

¹⁵⁶ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

¹⁵⁷ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

¹⁵⁸ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

¹⁵⁹ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

¹⁶⁰ *Digital Information Fidelity and Friction*, *supra* 143 at 22.

¹⁶¹ Anthony Nadler, Matthew Crain, and Joan Donovan, *Weaponizing the Digital Influence Machine: The Political Perils of Online Ad Tech*, *Data & Society*, 32 (2018) https://datasociety.net/wp-content/uploads/2018/10/DS_Digital_Influence_Machine.pdf

¹⁶² *Digital Information Fidelity and Friction*, *supra* 143 at 22.

standards articulated in their terms of service and other rules.¹⁶³ In the financial markets, a similar disruptive measure called a circuit-breaker is used to stop the flow of information that could overwhelm traders and contribute to instability.¹⁶⁴ When information flows result in a certain threshold of volatility in financial markets, regulators like the U.S. Securities and Exchange commission and the New York Stock Exchange can activate a circuit-breaker to cause a disruption. This disruption’s purpose is to allow time to process information and make informed decisions, “to create the space for the exercise of cognitive autonomy.”¹⁶⁵ Disruption to virality on digital platforms, allowing for time to process and, also for content moderation can limit the spread of noise.¹⁶⁶

Taxes can act as frictive measures to encourage businesses to avoid boosting noise over signal.¹⁶⁷ Taxes can act as friction by making companies suffer a financial cost for monetizing the virality of communication.¹⁶⁸ The revenue from such taxes could be put to use by supporting the production of signal.¹⁶⁹

Frictive Measures in the Special Rapporteur’s Reports

The Special Rapporteur recommends that companies in the ICT sector recognize the standard for freedom of expression is not based on the law of any state nor private interest, but rather human rights law.¹⁷⁰ Additionally, they recommend that companies engage in “smart regulation” focused on transparency and remediation, allowing people to choose whether to engage

¹⁶³ *Online Content Regulation*, *supra* note 10 at, ¶ 3.

¹⁶⁴ *Digital Information Fidelity and Friction*, *supra* 143 at 22.

¹⁶⁵ *Digital Information Fidelity and Friction*, *supra* 143 at 23.

¹⁶⁶ *Digital Information Fidelity and Friction*, *supra* 143 at 22.

¹⁶⁷ *Digital Information Fidelity and Friction*, *supra* 143 at 23.

¹⁶⁸ *Digital Information Fidelity and Friction*, *supra* 143 at 23.

¹⁶⁹ *Digital Information Fidelity and Friction*, *supra* 143 at 23.

¹⁷⁰ *Online Content Regulation*, *supra* note 10 at, ¶ 70.

in digital platforms.¹⁷¹ Companies should not engage in viewpoint-based regulation.¹⁷² At a minimum content moderation of digital platforms should be based on the UN Guiding Principles on business and restrictions to the freedom of expression should conform to the requirements of: Legality, Necessity and Proportionality, and legitimacy.¹⁷³ The Special Rapporteur also recognized a standard of non-discrimination in content moderation which required transcending formalistic approaches and considering “communities historically at risk of censorship and discrimination.”¹⁷⁴

The Special Rapporteur has acknowledged that digital platforms of Internet companies operate with a business model that benefits from attention and virality.¹⁷⁵ Artificial Intelligence and algorithmic personalization are “optimized for engagement and virality at scale.”¹⁷⁶ However, this preference for virality can threaten individual’s ability to find some content.¹⁷⁷ Because digital platforms value virality, content with lower levels of engagement can be deprioritized making some content obscure.¹⁷⁸

The Special Rapporteur recognizes that companies engage in content moderation and predicate access to their digital platforms upon compliance with user agreements and terms of service.¹⁷⁹ A major flaw in companies’ content moderation practices is a lack of transparency.¹⁸⁰ Automated content moderation creates risks that content moderation might violate human rights

¹⁷¹ *Online Content Regulation*, *supra* note 10 at, ¶ 66.

¹⁷² *Online Content Regulation*, *supra* note 10 at, ¶ 66.

¹⁷³ *Online Content Regulation*, *supra* note 10 at, ¶ 45-49; *Online Hate Speech*, *supra* note 12, at ¶ 18-20.

¹⁷⁴ *Online Content Regulation*, *supra* note 10 at, ¶ 48.

¹⁷⁵ *Online Hate Speech*, *supra* note 12, at ¶ 40.

¹⁷⁶ *Artificial Intelligence*, *supra* note 11, at ¶ 12.

¹⁷⁷ *Artificial Intelligence*, *supra* note 11, at ¶ 12.

¹⁷⁸ *Artificial Intelligence*, *supra* note 11, at ¶ 12.

¹⁷⁹ *Online Content Regulation*, *supra* note 10 at, ¶ 12.

¹⁸⁰ *Online Hate Speech*, *supra* note 12, at ¶ 45.

law.¹⁸¹ Some methods of company moderation that the Special Rapporteur has touched on include: automated flagging; automated removal; automated pre-publication filtering; user flagging of impermissible content; trusted flagging; human evaluation; action by the company against the account or content; notification; appeals and remedies.¹⁸² Additionally, the Special Rapporteur expressed concern at the prospect of promoting counter-narrative to user communications on digital platforms.¹⁸³ Use of such counter-narratives could transform digital platforms into propaganda carriers.¹⁸⁴ While the Special Rapporteur has acknowledged a myriad of content moderation tools, some which might even have the effect of creating friction, they have not directly advocated the use of frictive measures.

Despite not expressly using the term frictive measures, the Special rapporteur has recognized that companies have content moderation tools that can restrict the virality of communications and include a “range of options short of deletion.”¹⁸⁵ In countering the spread of disinformation—noise—on digital platforms, particularly in response to the COVID-19 pandemic, rumors must be addressed in order to correct them.¹⁸⁶ However, state penalization of disinformation is not proportionate and can deter communication of valuable information.¹⁸⁷

IV. FRICTIVE MEASURES UNDER THE INTERNATIONAL LAW FRAMEWORK

Although states, not business, are parties to the International Covenant of Civil and Political Rights,¹⁸⁸ the international law framework outlined in the special rapporteur’s reports can

¹⁸¹ *Online Content Regulation*, *supra* note 10 at, ¶ 56.

¹⁸² *Online Content Regulation*, *supra* note 10 at, ¶ 32-37.

¹⁸³ *Online Content Regulation*, *supra* note 10 at, ¶ 21.

¹⁸⁴ *Online Content Regulation*, *supra* note 10 at, ¶ 21.

¹⁸⁵ *Online Hate Speech*, *supra* note 12, at ¶ 51.

¹⁸⁶ *Disease Pandemics*, *supra* note 13, at ¶ 42.

¹⁸⁷ *Disease Pandemics*, *supra* note 13, at ¶ 42.

¹⁸⁸ *ICCPR*, *supra* note 65, at 1.

be applied to assess restrictions on the freedom of expression by social media companies. Frictive measures would have to satisfy the requirements of: Legality; Necessity and Proportionality; and Legitimacy.¹⁸⁹ Additionally, social media companies are covered under the UN Guiding Principles of Business and Human Rights.¹⁹⁰ Accordingly, frictive measures can be assessed under the framework of assessing restrictions to the freedom of expression under the ICCPR and recommendations of the UN Guiding Principles on Business.

Communication Delays

Systematizing a pause would be unlikely to violate either the freedoms under Art. 19 of the ICCPR or the Guiding Principles.

Legality

To be legal, a regulation must be publicly available and sufficiently precise for people to conform to the regulation.¹⁹¹ Nothing about a communication delay would be illegal under the law because it would be content neutral meaning that people subject to the restriction would not have to change their behavior in a meaningful way. So long as companies publicly inform users that they implement communication delays, such regulation would be publicly available.

Necessity and Proportionality

Communication delays comply with the requirement of necessity because of their benefit in reduction of noise on social media.¹⁹² To demonstrate necessity, the body issuing the restriction

¹⁸⁹ *Online Content Regulation*, *supra* note 10 at ¶ 7; *Artificial Intelligence*, *supra* note 11, at ¶ 28; *Online Hate Speech*, *supra* note 12, at ¶ 6; *Disease Pandemics*, *supra* note 13, at ¶ 11.

¹⁹⁰ *Guiding Principles*, *supra* note 69 at princ. 11.

¹⁹¹ *General Comment No. 34*, *supra* note 81, at ¶ 24.

¹⁹² *Digital Information Fidelity and Friction*, *supra* 143 at 22.

to freedom of expression needs to show a direct and immediate connection between the kind of expression and the threat.¹⁹³ In the context of digital platforms, the threat is disproportionate noise in the signal to noise ratio.¹⁹⁴ Frictionless sharing degrades the quality of information being shared.¹⁹⁵ The immediacy connection comes with the virality of the communication that is promoted through business models thriving on such virality.¹⁹⁶ On most social media, publication of communication is almost instantaneous.¹⁹⁷ The restriction of a communication delay allows for the creation of time for cognitive processing before information is shared.¹⁹⁸

Because delaying communication falls short of the most extreme restriction of deletion, it would likely qualify as proportionate. Considering the form of expression restricted along with its method of dissemination, the restriction must be the least intrusive option to be proportionate.¹⁹⁹ Because social media companies have the capacity to delete content from their platforms, any restriction short of deletion would be short of the most restrictive option.²⁰⁰ While it is likely untenable to create an exhaustive list of potential restrictions and assign a level to intrusiveness to each, simply creating a delay in the message does not compare to silencing the message in its entirety through deletion.

Legitimacy

¹⁹³ *General Comment No. 34*, *supra* note 81, at ¶ 34.

¹⁹⁴ *Digital Information Fidelity and Friction*, *supra* 143 at 13.

¹⁹⁵ *The Law of Friction*, *supra* 142 at 46.

¹⁹⁶ *The Law of Friction*, *supra* 142 at 48.

¹⁹⁷ Catherine O'Regan, *Hate Speech Online: an (Intractable) Contemporary Challenge?*, Vol. 71 *Current Legal Problems* 403, 416 (2018).

¹⁹⁸ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

¹⁹⁹ *General Comment No. 34*, *supra* note 81, at ¶ 34.

²⁰⁰ *Online Hate Speech*, *supra* note 12, at ¶ 51.

Communication delays qualify as legitimate restrictions to the freedom of expressions because they help protect rights under international law.²⁰¹ Some legitimate purposes for restrictions on the freedom of expression included in the language of Art. 19 are, “national security...public order...public health or morals.” (Art. 19, pg 11, 3(a)(b)). The purpose of communication delays is to boost signal, truthful information that supports democratic discourse. (Digital Information Fidelity and Friction, pg 2). The promotion of truthful, democratic discourse would likely be considered furtherance of public order and accordingly, legitimate.

Virality Disruptors

For much of the same reasons as Communication Delays, Virality Disruptors would likely be permissible under Art. 19’s international legal framework. However, because virality disruptors can be used to allow time for content moderation, there is increased likelihood for viewpoint discrimination. For this reason, it is valuable to distinguish between disrupting of the virality of content and disrupting virality for the purposes of determining whether a communication should be deleted. In this context, I choose to focus on virality disruptors “to create the space for the exercise of cognitive autonomy.”²⁰²

Legality

So long as the implementation of virality disruptors on social media platforms is made publicly available and is sufficiently precise to allow individuals to tailor their conduct,²⁰³ virality disruptors are legal under international law.

²⁰¹ *General Comment No. 34, supra* note 81, at ¶ 28.

²⁰² *Digital Information Fidelity and Friction, supra* 143 at 23.

²⁰³ *General Comment No. 34, supra* note 81, at ¶ 24.

Necessity and Proportionality

Virality disruptors seek the same goal as communication delays, to reduce noise.²⁰⁴ Whereas frictionless sharing increases noise,²⁰⁵ the introduction of a pause at a certain threshold of circulation allows for time to “exercise of cognitive autonomy.”²⁰⁶ The expression of frictionless threatens to drown out signal with noise,²⁰⁷ but a virality disruptor introduces friction in the hopes of increasing signal by allowing time to process. The immediate connection between the threat of noise and the expression of frictionless sharing should be sufficient to fulfill the necessity requirement of Art. 19.²⁰⁸

Because slowing the spread of information by way of virality disruptors falls short of deleting the content or censorship, there is a strong likelihood that such frictive measures would be proportional for the same reasons as communication delays.

Legitimacy

Virality disruptors help protect rights under international law by increasing the amount of signal, thereby supporting democratic discourse and public order. Instead of allow discourse on social media to be dominated by frictionless, quick, emotional responses, disrupting the virality of a communication allows more time for processing and increases the ability for cognitive choice. Cognitive choice helps democratic discourse by allowing people to make their intent a conscious

²⁰⁴ *Digital Information Fidelity and Friction*, *supra* 143 at 22.

²⁰⁵ *The Law of Friction*, *supra* 142 at 46.

²⁰⁶ *Digital Information Fidelity and Friction*, *supra* 143 at 23.

²⁰⁷ *The Law of Friction*, *supra* 142 at 46.

²⁰⁸ *General Comment No. 34*, *supra* note 81, at ¶ 34.

decision as opposed to an assumed intent based on a frictionless architecture.²⁰⁹ Supporting democratic discourse would benefit public order and sufficiently legitimize virality disruptors.²¹⁰

Frictive Measures under the UN Guiding Principles

Whereas Art. 19 has an analytical framework which can be readily applied to restrictions on freedom of expression to determine if they comply with international law, the UN Guiding Principles provide for business to better comply with human rights law.²¹¹ Frictive measures that are permissible under international law can be used to help businesses including social media companies comply with the UN Guiding Principles.

Frictive measures provide an opportunity for social media companies engage in “prevention and mitigation strategies” of the adverse impact their business can have on human rights²¹² Frictionless sharing is built into the architecture of social media companies.²¹³ Frictionless sharing results in noise,²¹⁴ and noise undermines the discursive potential of democratic discourse.²¹⁵ Frictive measures prevent and mitigate noise overwhelming signal on social media.²¹⁶ Specifically, Communication Delays and frictive measures can prevent and mitigate noise on social media,²¹⁷ without viewpoint discrimination.²¹⁸

²⁰⁹ *The Law of Friction*, *supra* 142 at 47.

²¹⁰ *Digital Information Fidelity and Friction*, *supra* 143 at 2.

²¹¹ *Guiding Principles*, *supra* note 69 at 1.

²¹² *Guiding Principles*, *supra* note 69 at princ. 13 and 23; *see also Online Content Regulation*, *supra* note 10 at ¶ 11.

²¹³ *The Law of Friction*, *supra* 142 at 19.

²¹⁴ *The Law of Friction*, *supra* 142 at 46.

²¹⁵ *Digital Information Fidelity and Friction*, *supra* 143 at 2.

²¹⁶ *Digital Information Fidelity and Friction*, *supra* 143 at 13.

²¹⁷ *Digital Information Fidelity and Friction*, *supra* 143 at 21-22.

²¹⁸ *Digital Information Fidelity and Friction*, *supra* 143 at 21.

While frictive measures will not satisfy all the UN Guiding Principles,²¹⁹ they provide an opportunity for social media companies to promote democratic discourse on their platforms and mitigate adverse impacts on human rights.

Concerns about AI

Despite the benefits of frictive measures, they are likely considered a form of artificial intelligence and are susceptible to the same kinds of discrimination posed by artificial intelligence. Artificial Intelligence is defined by the Special Rapporteur in, *A/73/348*, as a collection of technologies that allow computers to reinforce or replace tasks done by humans.²²⁰ The Special Rapporteur has identified the “potential for AI to embed and perpetuate bias and discrimination...in the exercise of freedom of opinion and expression.”²²¹ Discriminatory effects and bias in artificial intelligence are produced by the data sets in the design of the intelligence.²²² Additionally, the lack of transparency around the manner artificial intelligence effects the information environment prevents people “from understanding when and according to what metric information is disseminated, restricted or targeted.”²²³

On social media, the content users see and information personalized for their viewing is often dictated by artificial intelligence.²²⁴ Massive amounts of data including “browsing histories, semantic and sentiment analyses” are entered into algorithms to curate information displayed to users.²²⁵ Social media companies use subjective assessments to gauge how interesting content will

²¹⁹ *Online Content Regulation*, *supra* note 10 at, ¶ 11.

²²⁰ *Artificial Intelligence*, *supra* note 11, at ¶ 3.

²²¹ *Artificial Intelligence*, *supra* note 11, at ¶ 37.

²²² *Artificial Intelligence*, *supra* note 11, at ¶ 6.

²²³ *Artificial Intelligence*, *supra* note 11, at ¶ 31.

²²⁴ *Artificial Intelligence*, *supra* note 11, at ¶ 10.

²²⁵ *Artificial Intelligence*, *supra* note 11, at ¶ 10.

be to a user, thereby limiting exposure to different perspective across ideological or political lines.²²⁶ Artificial intelligence optimized for user engagement, promotes virality and demotes independent and user-generated content,²²⁷, thus promoting a noisy, frictionless information environment.²²⁸

In addition to determining what content is seen by users, artificial intelligence is also involved in content moderation and removal.²²⁹ The Special Rapporteur has raised concerns about AI-driven content moderation, expressly recognizing the limited ability for artificial intelligence to account for linguistic and cultural context.²³⁰ The exclusion of information by AI-driven content moderation and removal “increases the risk of manipulation of individual users through the spread of disinformation” by limiting diverse perspectives.²³¹

While nothing inherent to frictive measures requires that they be carried out through artificial intelligence, there is no reason computer technology cannot be used in implementing frictive measures. If measures like communication delays and virality disruptors are implemented using artificial intelligence, they are at risk for the kinds of discrimination and bias that threaten the freedoms of opinion and expression online.²³² While the potential for bias in implementing frictive measures is worth acknowledging, the potential bias does not necessarily change their legality under the international law framework reiterated by the Special Rapporteur.

CONCLUSION:

²²⁶ *Artificial Intelligence*, *supra* note 11, at ¶ 12.

²²⁷ *Artificial Intelligence*, *supra* note 11, at ¶ 12.

²²⁸ *Digital Information Fidelity and Friction*, *supra* 143 at 3; *see also The Law of Friction*, *supra* 142 at 46.

²²⁹ *Artificial Intelligence*, *supra* note 11, at ¶ 13.

²³⁰ *Artificial Intelligence*, *supra* note 11, at ¶ 15.

²³¹ *Artificial Intelligence*, *supra* note 11, at ¶ 18.

²³² *Artificial Intelligence*, *supra* note 11, at ¶ 37.

Frictionless sharing on social media degrades information quality on social media.²³³ Degradation of information quality increases noise which misinforms.²³⁴ One way for the social media companies to combat misinformation on social media is to introduce frictive measures which would result in boosting signal by slowing information flows and allowing time to raise cognitive defenses.²³⁵ Communication delays and virality disruptors present new frictions that can be implemented in a content neutral way.²³⁶ Content neutral frictive measures are especially valuable because they most readily comply with the international framework identified by the Special Rapporteur, based on the ICCPR and UN Guiding Principles on Business.

In the last seven special rapporteur reports discussing the freedom of opinion and expression in the ICT sector, frictive measures have not been mentioned by name. While concerns about states using their power to restrict speech through private actors like social media companies is a very legitimate concern, so too is limiting the spread of disinformation. In the Special Rapporteur's upcoming report on Disinformation and freedom of opinion and expression, the Special Rapporteur should analyze and advocate the measures social media companies have at their disposal to combat disinformation. Content neutral frictive measures comply with international law and present a viable bulwark against the proliferation against disinformation online.

²³³ *The Law of Friction*, *supra* 142 at 46.

²³⁴ *Digital Information Fidelity and Friction*, *supra* 143 at 2.

²³⁵ *Digital Information Fidelity and Friction*, *supra* 143 at 22.

²³⁶ *Digital Information Fidelity and Friction*, *supra* 143 at 21.