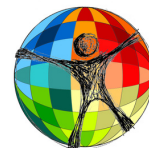


VOL. 10 | NO. 1 | SUMMER 2021

Regulating the ~~Marketplace~~ Metaverse of Ideas? Free speech as a Triangle, New School Moderation, and the Algorithm in Canada

Christoph Ivancic

McGill Centre for
Human Rights
and Legal Pluralism



Centre sur les droits de la
personne et le pluralisme
juridique de McGill



McGill FACULTY OF
Law

ABOUT CHRLP

Established in September 2005, the Centre for Human Rights and Legal Pluralism (CHRLP) was formed to provide students, professors and the larger community with a locus of intellectual and physical resources for engaging critically with the ways in which law affects some of the most compelling social problems of our modern era, most notably human rights issues. Since then, the Centre has distinguished itself by its innovative legal and interdisciplinary approach, and its diverse and vibrant community of scholars, students and practitioners working at the intersection of human rights and legal pluralism.

CHRLP is a focal point for innovative legal and interdisciplinary research, dialogue and outreach on issues of human rights and legal pluralism. The Centre's mission is to provide students, professors and the wider community with a locus of intellectual and physical resources for engaging critically with how law impacts upon some of the compelling social problems of our modern era.

A key objective of the Centre is to deepen transdisciplinary collaboration on the complex social, ethical, political and philosophical dimensions of human rights. The current Centre initiative builds upon the human rights legacy and enormous scholarly engagement found in the Universal Declaration of Human Rights.

ABOUT THE SERIES

The Centre for Human Rights and Legal Pluralism (CHRLP) Working Paper Series enables the dissemination of papers by students who have participated in the Centre's International Human Rights Internship Program (IHRIP). Through the program, students complete placements with NGOs, government institutions, and tribunals where they gain practical work experience in human rights investigation, monitoring, and reporting. Students then write a research paper, supported by a peer review process, while participating in a seminar that critically engages with human rights discourses. In accordance with McGill University's Charter of Students' Rights, students in this course have the right to submit in English or in French any written work that is to be graded. Therefore, papers in this series may be published in either language.

The papers in this series are distributed free of charge and are available in PDF format on the CHRLP's website. Papers may be downloaded for personal use only. The opinions expressed in these papers remain solely those of the author(s). They should not be attributed to the CHRLP or McGill University. The papers in this series are intended to elicit feedback and to encourage debate on important public policy challenges. Copyright belongs to the author(s).

The WPS aims to meaningfully contribute to human rights discourses and encourage debate on important public policy challenges. To connect with the authors or to provide feedback, please contact human.rights@mcgill.ca.

ABSTRACT

The contours of freedom of expression have changed with the advent of social media platforms; extremist speech is now accessible by broad audiences at the click of a button and private companies are responsible for regulating and enforcing speech policies. Governments seeking to combat extremist speech must find a way to balance freedom of expression guarantees for both individuals and companies, like Twitter and Facebook, with the censorship of illegal content. Canada attempted to do just that with proposed online harms legislation, but the balance could be better.

This paper makes three overlapping arguments: first, that the landscape of freedom of expression has been fundamentally changed by the triangular nature of freedom of expression, new school regulation, and hybrid moderation. Second, that the Canadian online harms legislation is under inclusive and it will cause collateral censorship and prior restraint. And, third, that alternative approaches of intermediary immunity, information fiduciary duties, and the duty to act responsibly are better legislative responses to online harms.

CONTENTS

1. INTRODUCTION	6
2. SYSTEM-LEVEL CONSTRAINTS AND FREEDOM OF EXPRESSION AS A TRIANGLE	7
3. CANADA'S NEW HATE SPEECH REGIME	17
4. LESSONS FROM THE EUROPEAN APPROACH TO ONLINE HARMS	22
5. ALTERNATIVE REGULATORY REGIMES	25
CONCLUSION	31
BIBLIOGRAPHY	34
ANNEX A: TECHNICAL DISCUSSION PAPER ON ONLINE HARMS LEGISLATION	41

1. Introduction

The Trudeau liberals are taking a new stance on hate speech. With the possible return of section 13 of the *Human Rights Act*¹ and proposed Online Harms legislation,² Canada is signaling that it means to seriously regulate online expression. However, the government should look to evolutions in the literature of hate speech regulation in the online sphere before embarking on such a project. A growing body of literature has come to reshape the online hate speech debate and accommodate for the antiquated framing of the marketplace of ideas given the ubiquity of social networks and the advent of web 2.0. Applying this new conception of online expression, expression as a triangle,³ it becomes apparent that the proposed Online Harms framework is misguided.

This paper makes three overarching arguments, first that the landscape of freedom of expression has been fundamentally changed by the triangular nature of freedom of expression, new school regulation, and hybrid moderation. Second, that the Canadian Online Harms legislation is underinclusive and it will cause collateral censorship and prior restraint. And third, that alternative approaches of intermediary immunity, information fiduciary duties, and the duty to act responsibly are better legislative responses to Online Harms. The Government would be better served by expanding their definition of Online Harms to encompass the new challenges posed by web 2.0 and creating a flexible regulatory regime which will incentivize social media platforms (SMPs) to be responsible stewards of free speech, privacy rights, and online security. Section 2 will examine the evolving theorization of freedom of expression and the system-level constraints of online spaces. Section 3 will briefly set out the regulatory regime proposed by the federal government for regulating online harms. Section 4 will look at examples of intermediary liability regimes in the EU and Germany. And section

¹ See Anja Karadeglija, “New Hate Law could have Chilling Effect, Free Speech Advocates Say”, *National Post* (4 June 2021), online: nationalpost.com/news/politics/new-hate-law-could-have-chilling-effect-free-speech-advocates-say.

² The Technical Briefing for this legislation has been attached as Annex A.

³ See Jack M Balkin, “Free Speech is a Triangle” (2018) 118:7 *Colum L Rev* 2011 [Balkin, “Triangle”].

5 examines alternative regulatory regimes and their benefits relative to the approach being contemplated in Canada.

2. System-level constraints and freedom of expression as a triangle

a. Expression as a triangle

In the 21st century the marketplace of ideas is radically different than ever before. Its scale is grander, its content broader, and mediums more varied. What was once a quaint farmers' market with a couple of stands is now a super-shopping center complete with restaurants, movie theatres, escape rooms, and an indoor ski hill. As put by Evelyn Douek, "[i]f the 'marketplace of ideas' analogy was ever more than an evocative oversimplification, it surely does not apply to platform ecosystems that optimize for engagement rather than truth."⁴ The phenomenon of SMPs has completely changed the landscape where freedom of expression is negotiated. In their seminal 2018 essay, Jack Balkin posits that free speech is now a triangle with its three corners being "nation-states, private infrastructure, and speakers."⁵ The theory posits that each corner of the triangle has rights and obligations in relation to one another. Just as SMPs have proven an important platform for social mobilization in authoritarian states, governments must exercise a check on SMPs when their action harms the democratic structure through censorship, spreading misinformative propaganda, or allowing extremist threats to flourish. The pluralist, globalized, "algorithmic" society has rendered the dualist conception of expression, which theorizes speech rights in the context of only nation-states and speakers, unpersuasive and difficult to apply.

The importance of this theorization is heightened by considering the monopolization of the SMP market by certain players and the possible democratic deficit created by such a

⁴ Evelyn Douek, "Governing Online Speech" (2021) 121:3 Colum L Rev 759 at 777.

⁵ See Balkin, "Triangle", *supra* note 3 at 2055.

monopoly. There are 4.66 billion people with access to internet⁶ and 2.89 billion monthly active Facebook users.⁷ Over 3 billion including other social media companies owned by Facebook.⁸ While this should not be taken definitively, as users might have Facebook accounts but primarily use another social media service or single users might have multiple accounts, it is nonetheless indicative of the importance of regulating these online spaces. Speech is, after all, a definitive pillar of healthy democracy. The possibility of the expression of over 3 billion people being censored according to the whims of a profit-driven corporation should raise democratic alarms. It is worrying that competition law has done little to prevent the development of Facebook's chokehold on the social media market however the issue of Facebook's monopolistic hold⁹ is beyond the ambit of this paper.

The intersection between the press and social media is also worth considering, the number of people that get their news from social media is substantial and rising.¹⁰ The media, often referred to as the "fourth branch"¹¹ of government, plays an important role in keeping the executive, legislative, and judicial branches in check. SMPs wield significant power over what news citizens do and do not see therefore jeopardizing the ability of the press to hold the government to account. SMPs "[l]ike twentieth-century

⁶ See Joseph Johnson, "Worldwide digital population as of January 2021" Statista (10 September 2021), online: www.statista.com/statistics/617136/digital-population-worldwide/.

⁷ See Statista Research Department, "Facebook: Number of Monthly Active Users Worldwide 2008–2021" Statista (1 November 2021), online: www.statista.com/statistics/264810/number-of-monthly-active-facebook-users-worldwide/ [Statista, "Monthly Users"].

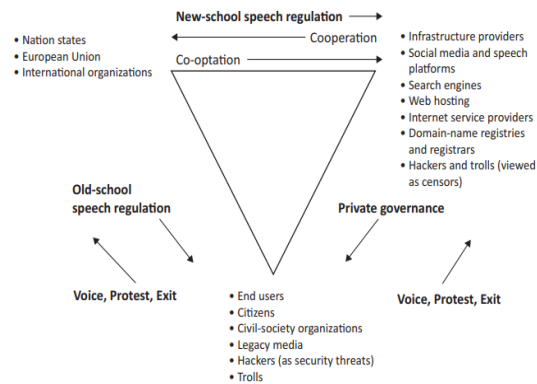
⁸ See Statista Research Department, "Distribution of Instagram users worldwide as of October 2021, by age group" Statista (23 November 2021), online: www.statista.com/statistics/325587/instagram-global-age-group/.

⁹ See Dipayan Ghosh, "How the Free Market Incentivized Facebook's Harmful Monopoly", *Centre for International Governance Innovation* (6 January 2021), online: www.cigionline.org/articles/how-free-market-incentivized-facebooks-harmful-monopoly/.

¹⁰ See Evidence for Democracy, "Misinformation in Canada: Research and Policy Options" (2021) at 9, online (pdf): evidencefordemocracy.ca/sites/default/files/reports/misinformation-in-canada-evidence-for-democracy-report_.pdf. 90% of Canadians used online sources to find information on COVID-19.

¹¹ See Rachel Luberd, "The Fourth Branch of Government: Evaluating the Media's Role in Overseeing the Independent Judiciary" (2014) 22:2 *Notre Dame JL Ethics & Pub Pol'y* 507 at 508.

mass media... have become important custodians of the public sphere and democratic self-government.”¹² There is a significant democratic imperative for states to adopt a view of free speech as a triangle. This has driven a desire for “regulation of regulation”¹³ or a judicial review of “the law of the platforms”¹⁴ in areas once thought to be beyond the purview of the state.



15

An issue with free speech as a triangle is that the lines between public and private become blurred. Why, after all, should private companies be part of the conversation regarding freedom of expression which is only guaranteed vis-à-vis the government and its citizens? One reason is increasing public-private cooperation and co-optation.¹⁶ The online harms legislation proposed in Canada is an excellent example; it essentially deputizes SMPs to regulate speech for the government with the threat of financial penalties, creating de facto government regulation of expression without the burdens of procedural fairness and substantive rights analysis that would otherwise be triggered. The deputization of private corporations to do the work of governments has been found to implicate constitutional

¹² See Balkin, “Triangle”, *supra* note 3 at 2041.

¹³ See Frank Pasquale, “Toward a Fourth Law of Robotics: Preserving Attribution, Responsibility, and Explainability, in an Algorithmic Society” (2017) 78:5 Ohio St LJ 1243 at 1244.

¹⁴ See Luca Belli, Pedro A Francisco and Nicolo Zingales, “Law of the Land or Law of the Platform? Beware of the Privatisation of Regulation and Police” in Luca Belli & Olga Cavalli, *Internet Governance and Regulations in Latin America*, 1st ed (Rio de Janeiro: FGV Direito Rio, 2019) 423.

¹⁵ Balkin’s pluralist model of speech regulation.

¹⁶ See Balkin, “Triangle”, *supra* note 3 at 2019.

guarantees in a limited set of circumstances¹⁷ though it is unclear at what point regulation of SMPs would trigger constitutional guarantees.¹⁸ This is sure to depend upon the nature of the “state action doctrine” and constitutional freedom of expression interpretation being applied.

The level of information sharing between governments and SMPs also gives the theory cogency. The Canada’s Online Harms regime once again illustrates this as it creates new methods for law enforcement to get ahold of data collected by SMPs about their users. The fragile divide between private and public in this context can make it difficult to determine at what point a citizen’s rights are being impacted by government or private action. The conception of free speech as a triangle does not necessarily advocate for the application of judicial free speech norms to the conduct of private corporations,¹⁹ but it recognizes that free speech protection will require the government to regulate these corporations given their democratic importance. This regulation must also respect constitutional divisions of public/private spheres and to do so governments must take a minimally impairing approach.

b. Old-School vs new-school speech regulation

Speech regulation has not only been revolutionized by the actors at play, as in the triangle, but also how regulation happens. While speech has traditionally been regulated using penalties, or disincentives, online space creates the opportunity for far more carrots and sticks to encourage or deter behaviour. SMPs don’t need to choose between censorship and non-censorship, they can tailor their regimes proportionately by responding creatively. They can force bots to identify themselves as such to prevent them from creating false consensus or disharmony.²⁰ They can create warnings on news articles that contain information that is

¹⁷ Exceptions to the state action doctrine are in the US. In Canada, situations where “extensive government control” leads to the application of the *Charter* as contemplated in *Eldridge v British Columbia (Attorney General)*, [1997] 3 SCR 624, 151 DLR (4th) 577.

¹⁸ See Balkin, “Triangle”, *supra* note 3 at 2046.

¹⁹ See *ibid.*

²⁰ See Madeline Lamo & Ryan Calo, “Regulating Bot Speech” (2019) 66:4 UCLA L Rev 988.

misleading or false.²¹ They can “quarantine” potentially harmful content by removing it from the recommendation algorithm forcing users to choose to view it.²² They can program the algorithm to show examples of counter-speech to users that are searching for extremist content.²³ Each of these approaches would have a different effect on the rights of the speaker and allow for a more delicate balancing of the competing rights at play when regulating potentially harmful expression. These different tools have been categorized as “the concepts of *hard control* – a platform’s authority over what can be published online – and *soft control* – a platform’s authority over what we are likely to see, and is deprioritized in algorithms that govern a user’s view of posts on the network (the feed).”²⁴ Both types of control can and should be applied in the battle against hatred, exploitation, and misinformation and *soft control* comes with the benefit of not affecting the freedom of expression rights of users.

c. Content moderation: the hybrid approach

The scale of social media poses a unique challenge for implementing regulation. Facebook touts 2.89 billion monthly active users,²⁵ YouTube 2.3 billion,²⁶ and Twitter 463 million.²⁷ Facebook moderated 105 million pieces of content

²¹ See Kaleigh Rogers, “Facebook’s Fact-Checking Program only a Partial Solution to Disinformation, Report Says”, CBC (30 July 2019), online: <www.cbc.ca/news/science/facebook-fact-checking-full-fact-report-1.5230592>.

²² See Stefanie Ullmann & Marcus Tomalin, “Quarantining Online Hate Speech: Technical and Ethical Perspectives” (2019) 22 Ethics and Information Technology 69.

²³ See Daniel Kreiss & Matt Perault, “Four Ways to Fix Social Media’s Political Ads Problem-Without Banning them”, *The New York Times* (16 November 2019), online: <www.nytimes.com/2019/11/16/opinion/twitter-facebook-political-ads.html>.

²⁴ See Jillian C York & Ethan Zuckerman, “Moderating the Public Sphere” in Rikke Frank Jorgensen, *Human Rights in the Age of Platforms* (Cambridge Massachusetts: The MIT Press, 2019) 137 at 140.

²⁵ See Statistica, “Monthly Users”, *supra* note 7.

²⁶ See Statista Research Department, “Most Popular Social Networks Worldwide as of October 2021, ranked by number of active users”, *Statista* (16 November 2021), online: <www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/>.

²⁷ See *ibid.*

during the first quarter of 2020 alone and Instagram did the same for over 35 million posts.²⁸ There are 500 hours of content uploaded to YouTube every minute.²⁹ Regulating this volume of content requires a degree of systemization which complicates notions of speech as an individual's right. It is impossible for content moderation at this scale to be completely accurate and the costs associated with performing a discrete, legalistic, analysis in each case of moderation would be both economically and pragmatically impossible. The result is a system that utilizes machine learning tools to automate moderation at increasing rates. However, SMPs have also recognized that qualitative freedom of speech analysis is sometimes necessary, and they have developed their own court structures for dealing with this. The result is a blended system of moderation, using both algorithmic and human decisions, that is shaped by the principles of probability and proportionality.

i. Algorithmic moderation: Probability and proportionality

When dealing with SMPs with the scale of Facebook, Twitter, YouTube, and Instagram, Online Harms legislation with strict timelines for content removal are sure to increase the amount of content that is moderated using AI. As SMPs take a speed-over-accuracy approach to moderation to meet legislative deadlines. SMPs are decreasing the amount of time it takes to respond to flagged content using automated matching and predictive systems.³⁰ In doing so, the moderation opens itself up to a degree of error. Algorithmic moderation is not yet able to understand the context of speech well enough to determine the lines between essential political speech and hate speech, as in political activism by Palestinians, or the speech that on its face is harmful but is saved by the context, such as news reports showing terrorist footage.³¹ The result is a margin of error that the major SMPs are comfortable with.

²⁸ See Meta, News Release, "Community Standard Enforcement Report" (11 February 2021), online: <about.fb.com/news/2021/02/community-standards-enforcement-report-q4-2020/>.

²⁹ See L Ceci, "Hours of Video Uploaded to YouTube Every Minute as of February 2020", *Statista* (14 September 2021), online: <www.statista.com/statistics/259477/hours-of-video-uploaded-to-youtube-every-minute/>.

³⁰ See Douek, *supra* note 4 at 795.

³¹ See *ibid.*

SMPs employ two main types of AI tools to help with moderation: matching and classifying/predictive systems.³² Matching systems check content that is being uploaded against a backlog of content that is known to be illegal or against the sites guidelines. This is done through “hashing” of videos, which essentially gives each video a unique fingerprint that is then checked against the fingerprint of the content being uploaded.³³ The classification/predictive machine learning tool discerns whether a given post would fall into a category, say offensive/not offensive, hateful/not hateful, allowing moderators to keep a finger on the pulse of generally toxic content that is being hosted on the platform.³⁴ Both types of algorithmic moderation are prone to error; classification tools still struggle to understand the full context of speech³⁵ and matching tools can be duped by hackers that recreate the “fingerprint” of permissible content and impose it on illegal content.³⁶ In using these tools, SMPs accept that there will be both false positives and negatives in their moderation process, but the efficacy of these AI tools far outweigh the drawbacks.

Proportionality comes with the growing consensus that freedom of expression is typically not limited in isolation; there are countervailing rights which decision makers must balance when they determine whether a breach is justified or not. Hate speech might involve balancing the freedom of the expression of the speaker with the security of the person of the subject of that speech. Censoring false news about COVID-19 involves weighing societies interest in public health and safety against the speakers right to express themselves. The proportional approach posits that a person’s freedom of expression should only be limited to the extent proportional to the countervailing interest.³⁷ Using methods of *soft control* and new school speech regulation it is easier than ever before to implement a system of moderation which minimally impacts the speaker’s rights.

³² See Robert Gorwa, Reuben Binns & Christian Katzenbach, “Algorithmic content moderation: Technical and Political Challenges in the Automation of Platform Governance” (2020) 7:1 Big Data & Society 1 at 3–4.

³³ See *ibid.*

³⁴ See *ibid.*

³⁵ See *ibid.*

³⁶ See *ibid.*

³⁷ See Douek, *supra* note 4 at 781–82, 784–86.

ii. **Human moderation and the Facebook oversight board**

SMPs have come a long way over the past 20 years and their approach to speech regulation has had marked eras. As their scale increased, Facebook, Twitter, and YouTube hired American lawyers to craft their content policies.³⁸ It was easy to see the influence of the First Amendment education in the policies that each of the companies created, however, they realized quickly that there would be issues in transporting the American approach to freedom of expression abroad. In countries like Turkey, where it is illegal to make certain depictions of Atatürk³⁹, and Thailand, where it is illegal to make fun of the King,⁴⁰ the First Amendment approach would not work. These growing pains would pale in comparison to the effects that the social media algorithm would have during the 2016 and 2020 US elections, during which social media spurred a mob to descend on the capital and decry the legitimacy of the election of Joe Biden.⁴¹ The result was a softening of the First Amendment approach to speech. The major SMPs have realized that they got the equation wrong and that they will need to be more careful in balancing competing interests; the free market of ideas comes with its perils.⁴²

Today, SMPs employ an army of content moderators that review content.⁴³ This type of moderation falls into two categories: proactive and reactive. Proactive content moderation involves searching for target content and is mostly employed in the context of extremist or terrorist speech.⁴⁴ Conversely, reactive content moderation responds to content that has been flagged by other users. Until 2018 the process of how moderation happens remained a black box as no SMP had made their internal

³⁸ See Kate Klonick, "The New Governors: The People, Rules, and Processes Governing Online Speech" (2018) 131:6 Harv L Rev 1598 at 1619–21 [Klonick, "New Governors"].

³⁹ See *ibid* at 1624.

⁴⁰ See *ibid* at 1623.

⁴¹ See *supra* note 23.

⁴² See Klonick, "New Governors", *supra* note 38 at 1625–30.

⁴³ See Casey Newton, "The Trauma Floor: The secret lives of Facebook moderators in America", *The Verge* (25 February 2019), online: <www.theverge.com/2019/2/25/18229714/cognizant-facebook-content-moderator-interviews-trauma-working-conditions-arizona>.

⁴⁴ See Klonick, "New Governors", *supra* note 38 at 1638.

moderation guidelines public.⁴⁵ However, to increase transparency, Facebook has since made their community standards and content moderation policy public.⁴⁶ Using Facebook, or rather its parent company Meta as an example, their community standards are available online in their newly created “transparency” center.

The failure of a first-amendment approach to speech online not only incited a change to the moderation policies of the major SMPs but it also fueled a desire for more robust, transparent, and timely governance tools.⁴⁷ One such tool is the Facebook Oversight Board. Designed as “an external body that people can appeal to if they disagree with Meta’s content enforcement decisions.”⁴⁸ The Board is essentially a “private, independent arbitration system built by Facebook”⁴⁹ that enjoys some kind of deference from Facebook given that they “commit to the board’s independent oversight on content decisions and the implementation of those decisions.”⁵⁰ The Board is then appellate-court like in nature, having jurisdiction to select their cases and sitting in panels of five judges, and applying the Oversight Board charter in “a constitution-like”⁵¹ manner.

The Board was set up by an independent trust to insulate the board from accusations of being influenced by Facebook, however, the effectiveness of this has been called into question.⁵² Application for review to the Board may come from the poster of

⁴⁵ See Catherine Buni & Soraya Chemaly, “The Secret Rules of the Internet: The Murky History of Moderation, and How it’s Shaping the Future of Free Speech”, *The Verge*, online: www.theverge.com/2016/4/13/11387934/internet-moderator-history-youtube-facebook-reddit-censorship-free-speech.

⁴⁶ See Meta, News Release, “Publishing Our Internal Enforcement Guidelines and Expanding Our Appeals Process” (24 April 2018), online: about.fb.com/news/2018/04/comprehensive-community-standards/.

⁴⁷ See Kate Klonick, “The Facebook Oversight Board: Creating an Independent Institution to Adjudicate Online Free Expression” (2020) 129:8 Yale L J 2418 at 2448–51 [Klonick, “Oversight Board”].

⁴⁸ See Meta, “Oversight Board: The Purpose of the Board” (last visited 14 June 2022), online: transparency.fb.com/en-gb/oversight/.

⁴⁹ See Klonick, “Oversight Board”, *supra* note 47 at 2477.

⁵⁰ See Meta, “Bylaws”, (November 2021), online (pdf): about.fb.com/wp-content/uploads/2020/01/Bylaws_v6.pdf.

⁵¹ See Klonick, “Oversight Board”, *supra* note 47 at 2457.

⁵² See *ibid* at 2467.

the content, the person that flagged the content, or Facebook itself.⁵³ The subject matter of complaints under the Board's jurisdiction was left intentionally vague in the Board's charter, with clarifications coming through the by-laws and the Board's jurisprudence.⁵⁴ The decision-making process applied by the Board comes from Facebook's Values of "[v]oice, safety, privacy, authenticity, and dignity."⁵⁵ They are ostensibly embracing the new status quo: a proportional balancing approach to freedom of expression. Decisions need not be unanimous, and the Board is required to provide a specific determination, reasoning, and explanation for a given decision with room for concurring and dissenting opinions. The Board can also provide a policy-advisory statement to Facebook in their decisions.

d. Conclusion: Not your grandparents' marketplace of ideas

This section has laid out the substantial structural changes that should shape conversations about online harms governance. Its conclusions are as follows: firstly, the conception of free speech as a triangle should guide the architecture of governance options. Governments must find a way to regulate SMPs to secure the freedom of expression rights of the electorate and in turn safeguard the health of Canadian democracy. Governments must be wary of the point at which their regulations trigger the constitutional rights of end users, so as not to turn every post that gets taken down on Facebook into actionable *Charter* litigation. Secondly, speech regulation no longer needs to be dichotomous. Given the advent of new-school moderation techniques in the online sphere, SMPs are better placed to balance freedom of expression interests with dignity, safety, and privacy, than governments have ever been. Thirdly, the scale of the online marketplace necessitates automation and the acceptance of a degree of error. Fourthly, there are market incentives for private corporations to develop complex moderation tools, both algorithmic and juridical. Legislators should be looking to create synergies with SMPs to improve existing systems and to avoid taking on further responsibilities at the expense of taxpayers.

⁵³ See *ibid* at 2463.

⁵⁴ See *ibid* at 2462.

⁵⁵ See *ibid* at 2463.

3. Canada's new hate speech regime

a. Online Harms, criminal, and human rights law

Over the summer of 2021, the Trudeau liberals entered public consultations on their Online Harms legislation.⁵⁶ The legislation would create a new set of rules for SMPs and a government body to ensure compliance. The rules would place an obligation to remove five categories of harmful content on SMPs (hate speech, child sexual exploitation content, non-consensual sharing of intimate images, incitement to violence content, and terrorist content) within 24 hours of being flagged. It would also impose requirements of transparency and procedural fairness, and mandate that internet service providers (ISPs) block access to SMPs should they persistently fail to take down child sexual exploitation and terrorist content. The oversight framework would create an appeal mechanism for content moderation decisions made by SMPs and would order the removal of harmful content when SMPs fail to do so through their own content moderation frameworks. The legislation would also put an obligation on SMPs to preserve records of potentially illegal content and content of national security which could later be lawfully obtained for investigation.

The oversight body would be called the Digital Safety Commission of Canada and would be comprised of the Digital Safety Commissioner of Canada, the Digital Recourse Council of Canada, and the Advisory Board. The tripartite body would oversee the enforcement of the rules, research online safety, provide independent recourse for content removal, and provide avenues for experts, equity-deserving, and Indigenous interests to shape the regulation.

The proposed Online Harms legislation also claims to bolster the ability of law enforcement to deal with illegal activity in the online sphere. Canada has a criminal hate speech prohibition in the form of sections 318, 319(1), and 319(2), of the *Criminal Code of Canada* which prohibit advocating or promoting genocide, inciting hatred against any identifiable

⁵⁶ See Government of Canada, News Release, "Consultation closed: The Government's proposed Approach to address Harmful Content Online" (7 October 2021), online: <www.canada.ca/en/canadian-heritage/campaigns/harmful-online-content.html>.

group, and willfully promoting hatred against any identifiable group. Consent of the Attorney General is required to before charges are laid using these provisions to prevent the intimidation of those whose speech might be controversial. Uttering threats and harassment may also take place in the online sphere and hate is considered in aggravating factor in these offences. These provisions are largely considered to be ineffective at combatting hate given that they often go unused and police services have limited experience using them in the context of online harms.⁵⁷ There are also issues of clarity on when the Attorney General should or should not allow charges to be laid under the provisions.⁵⁸

Bill C-36, which was introduced in June 2021,⁵⁹ is seen as complementary to the Online Harms legislation. The Bill would reintroduce a provision to the *Human Rights Act* making it discriminatory to communicate hate speech by means of the internet or other means of telecommunication in a context in which the hate speech is likely to foment detestation or vilification of an individual or a group of individuals based on a prohibited ground of discrimination.

The legislation also introduces a definition of *hate speech*: *hate speech* means the content of a communication that expresses detestation or vilification of an individual or group of individuals. Though the legislation clarifies that speech which expresses dislike or disdain, or discredits, humiliates, hurts, or offends does not come within this definition. The government claims that this definition is consistent with how the Supreme Court of Canada has defined hate speech⁶⁰ through its jurisprudence, but the Canadian Civil Liberties Association has been critical of this, claiming that the definition remains vague and will result in a chilling effect on

⁵⁷ See Canada, House of Commons, Standing Committee on Justice and Human Rights, *Taking Action to End Online Hate* (June 2019) (Chair: Anthony Housefather) at 14–17, online (pdf): www.ourcommons.ca/Content/Committee/421/JUST/Reports/RP10581008/justrp29/justrp29-e.pdf.

⁵⁸ See *ibid.*

⁵⁹ See Bill C-36, *An Act to amend the Criminal Code and the Canadian Human Rights Act and to make related amendments to another Act (hate propaganda, hate crimes and hate speech)*, 2nd Sess, 43rd Parl, 2021.

⁶⁰ See Government of Canada, News Release, “Combating Hate Speech and Hate Crimes: Proposed legislative Changes to the Canadian Human Rights Act and the Criminal Code” (1 September 2021), online: www.justice.gc.ca/eng/csj-sjc/pl/chshc-lcdch/index.html.

free speech.⁶¹ The government has been explicit in saying that Bill C-36 as meant as a complementary pathway to the Online Harms legislation, giving victims of hate speech the choice of complaining against individuals or websites under the CHRA or against SMPs under online harms. However, it is unclear how reintroducing the section would eliminate the issues which called for it to be repealed.

The hate speech provision of the *Human Right Act* has a storied history. University of Windsor law Professor Richard Moon was commissioned by the Canadian Human Rights Commission to compile a report on section 13 in 2008.⁶² In the report, Moon advocates for more frequent use of the Criminal hate speech provisions,⁶³ the creation of provincial “Hate Crime Teams” made of police and Crown prosecutors⁶⁴, and for the repeal or amendment of section 13.⁶⁵ Today, Moon’s position remains largely the same given that “this is substantially the same provision ... it relies upon private citizens, organizations to initiate a complaint. And one of the problems with the old Section 13 was the incredible burden it places upon individuals or groups, both to do the basic investigation or inquiry, but also to see the complaint through the process.”⁶⁶ However, the debate was not one sided, the Canadian Bar Association opposed the repeal of the provision in 2012.⁶⁷ They argued that section 13 was a

⁶¹ See Karadeglija, *supra* note 1.

⁶² See Richard Moon, “Report to the Canadian Human Rights Commission Concerning Section 13 of the Canadian Human Rights Act and the Regulation of Hate Speech on the Internet” (October 2008), online (pdf): *Canadian Human Rights Commission* <deliverypdf.ssrn.com/delivery.php?ID=57202011103108407702808209311400510800604509106506300002809609806410212508902008410304501112003310612005309811503100410409811204704002701302309207706412710411210806802908311008808002902511511709008500411309206608109908>.

⁶³ See *ibid* at 2.

⁶⁴ See *ibid* at 2, 32–33.

⁶⁵ See *ibid* at 2.

⁶⁶ See Karadeglija, *supra* note 1.

⁶⁷ See Constitutional and Human Rights Law Section, “Bill C-304 *Canadian Human Rights Act* Amendments (Hate Messages)” (April 2012), online (pdf): *Canadian Bar Association* <www.cba.org/CMSPages/GetFile.aspx?guid=0aa9dfdd-92d4-429f-8946-8b2e08c28a10>.

reasonable limit under section 1 of the *Charter* and that it was a useful as a tool for groups targeted by hate speech.⁶⁸

I had the opportunity to learn about Canada's Online Harms legislation early in the consultation process during and internship with the British Columbia Civil Liberties Association. I was charged with developing a refined position on freedom of expression in the 21st century given the Trudeau's government's intention to change the legislative landscape with Bill C-36 and Online Harms. This paper is a continuation of the policy research I conducted over the summer of 2021.

b. Online harms: undertheorized and underinclusive?

Defining online harms can be difficult and questions as to how the Canadian government came to only five categories of harmful content is not obvious. While the legislation is pointed in its design to deal with issues of violence, exploitation, discrimination, and terrorism, which bleed from the offline sphere to the online, it is silent on several areas of harm which have been plaguing online spaces and SMPs. There is no mention of misinformation,⁶⁹ fake news,⁷⁰ conspiracy theories,⁷¹ surveillance,⁷² copyright infringement,⁷³ foreign interference in domestic politics,⁷⁴ or expression which is legal but nonetheless

⁶⁸ See *ibid.*

⁶⁹ See Evidence for Democracy, *supra* note 10 at 9–15.

⁷⁰ See Elizabeth Thompson, "Poll finds 90% of Canadians have Fallen for Fake News", CBC (11 June 2019), online: <www.cbc.ca/news/politics/fake-news-facebook-twitter-poll-1.5169916>.

⁷¹ See Belle Riley Thompson, "COVID and Conspiracy", *Open Canada* (29 September 2021), online: <opencanada.org/covid-and-conspiracy/>.

⁷² See Carole Cadwalladr & Emma Graham-Harrison, "Revealed: 50 million Facebook profiles Harvested for Cambridge Analytica in Major Data Breach", *The Guardian* (17 March 2018), online: <www.theguardian.com/news/2018/mar/17/cambridge-analytica-facebook-influence-us-election>.

⁷³ See Fara Tabatabai, "A Tale of Two Countries: Canada's Response to the Peer-to-Peer Crisis and What it Means for the United States" (2005) 73:5 *Fordham L Rev* 2321.

⁷⁴ See Catharine Tunney, "Canadian Voters are likely to face Foreign Cyber Interference in the Next Elections, says Cyber Spies", CBC (16 July 2021), online: <www.cbc.ca/news/politics/election-security-threats-cse-1.6104745/>. See also Canadian Centre for Cyber Security, News Release, "Cyber threats to Canada's democratic process: July 2021 update" (July 2021), online:

harmful. Not to mention harms that come distinctly from SMPs and flow to their users such as the misuse of user information,⁷⁵ the role of the social media algorithm in creating social media addiction,⁷⁶ and amplifying polarization.⁷⁷

It becomes necessary at this juncture to understand the relation between social media algorithms, democracy, and acts of hate in the real world. The social media algorithm prioritizes popularity and user-interest; giving greater visibility to posts which are popular and intrigue the user, causing them to click.⁷⁸ While this function is much simpler than, say an algorithm which performs short trades based on market variables, it is nonetheless prone to the same pitfalls. Chiefly, “algorithms (a) construct identity and reputation through (b) classification and risk assessment, creating the opportunity for (c) discrimination, normalization, and manipulation, without (d) adequate transparency, accountability, monitoring, or due process.”⁷⁹ Around the world, social media algorithms have exacerbated genocidal acts,⁸⁰ terrorist

cyber.gc.ca/en/cyber-threats-canadas-democratic-process-july-2021-update/>.

⁷⁵ See Moon, *supra* note 62.

⁷⁶ See Mike Wright, “One in five underage Facebook users are ‘addicted’ to social media app, whistleblower says”, *National Post* (26 October 2021), online: <nationalpost.com/news/world/one-in-five-underage-facebook-users-are-addicted-to-social-media-app-whistleblower-says>.

⁷⁷ See Jeff Horwitz & Deepa Seetharaman, “Facebook Executives Shut Down Efforts to Make the Site Less Divisive”, *The Wall Street Journal* (16 May 2020), online: <www.wsj.com/articles/facebook-knows-it-encourages-division-top-executives-nixed-solutions-11590507499>.

⁷⁸ See Oscar Alvarado & Annika Waern, “Towards Algorithmic Experience: Initial Efforts for Social Media Contexts” (Honourable mention delivered at the 2018 CHI Conference on Human Factors in Computing Systems, 4 April 2018) at 2–3.

⁷⁹ See Jack M Balkin, “2016 Sidley Austin Distinguished Lecture on Big Data Law and Policy: The Three Laws of Robotics in the Age of Big Data” (2017) 78:5 Ohio St LJ 1217 at 1239.

⁸⁰ See Paul Mozur, “A Genocide Incited on Facebook, With Posts From Myanmar’s Military”, *The New York Times* (15 October 2018), online: <www.nytimes.com/2018/10/15/technology/myanmar-facebook-genocide.html>.

recruitment⁸¹ and the destabilization of democracies.⁸² Or, as put by professors Elkin-Koren and Perel “a user who has watched the NBA championship is more likely to be offered additional sporting events. A user who has searched for information on extreme Islamic ideology might be offered videos on ISIS.”⁸³ A legislative scheme aimed at online harms is undermined by failing to account for the intersection between individual posts containing harmful content and the systems which allow for that content to be disseminated worldwide and fed directly to like-minded individuals, at the click of a button. It may be unfair to critique legislation beyond the breadth of its goals; effective governance of the online space is sure to be a multifaceted endeavor requiring different standards for ISPs, SMPs, and citizens. Expecting all this to come from the same piece of legislation is not feasible. Nonetheless, the literature supports a legislative scheme that incorporates systemic governance of algorithms with moderation of single pieces of content. These options will be explored in the final section of the paper.

4. Lessons from the European approach to online harms

As the problem of online harms became more prevalent over the course of the 21st century, not all governments have been as restrained as the United States. Notably, the EU and Germany have taken aggressive stances against the proliferation of harm

⁸¹ See Stewart Bell, “Extremist groups ‘actively recruiting’ military and police, Canadian intelligence report warns”, *Global News* (23 August 2021), online: [<globalnews.ca/news/8128463/extremist-groups-military-recruitment-report/>](https://globalnews.ca/news/8128463/extremist-groups-military-recruitment-report/).

⁸² See Kevin Roose, “Facebook reverses Postelection Algorithm changes that Boosted News from Authoritative Sources”, *The New York Times* (16 December 2020), online: www.nytimes.com/2020/12/16/technology/facebook-reverses-postelection-algorithm-changes-that-boosted-news-from-authoritative-sources.html>. See also, Catherine Kim, “Poll: 70 percent of Republicans don’t think the Election was Free and Fair”, *Politico* (11 September 2020), online: www.politico.com/news/2020/11/09/republicans-free-fair-elections-435488/>.

⁸³ See Niva Elkin-Koren and Maayan Perel, “Separation of Functions for AI: Restraining Speech Regulation by Online Platforms” (2020) 24:3 *Lewis & Clark L Rev* 857 at 888.

online. In 2016 the EU entered an agreement on a Code of Conduct with major SMPs.⁸⁴ The Code, an instrument of soft-law which requires companies to make good faith efforts, stipulates that SMPs should create rules and community standards prohibiting hate speech, create bureaucracies to review content, review flagged content within 24 hours, and promote transparency.⁸⁵ However, German lawmakers found that SMPs were not mobilizing quickly enough to achieve their obligations under the EU Code⁸⁶ and so decided to create further legal incentive through the Network Enforcement Act (Netzwerkdurchsetzungsgesetz or NetzDG) which came into effect in 2018. The NetzDG is an intermediary liability takedown regime, it requires that SMPs “(1) block access within Germany to ‘manifestly unlawful’ content within twenty-four hours of receiving notice of any such content, and (2) block access to other unlawful content (that is unlawful but not ‘manifestly’ unlawful) within seven days of receiving such notice.”⁸⁷ Should SMPs fail to comply with these obligations they can be fined up to 50 million euros.⁸⁸ The model was seemingly developed for political reasons with legal experts decrying the bill⁸⁹ and academics questioning its constitutional validity.⁹⁰

The legislation proposed in Canada seems to mirror the approach of Germany; a take-down regime of intermediary liability that requires SMPs to remove user content that falls within the legal definition of hate speech within 24 hours or be subject

⁸⁴ See Dawn Carla Nunziato, “The Marketplace of Ideas Online” *Notre Dame L Rev* 1519 at 1532.

⁸⁵ See Council of the European Union, News Release, “Assessment of the Code of Conduct on Hate on line State of Play” (27 September 2019), online (pdf): <ec.europa.eu/info/sites/default/files/aid_development_cooperation_fundamental_rights/assessment_of_the_code_of_conduct_on_hate_speech_on_line_-_state_of_play_0.pdf>.

⁸⁶ See Heidi J S Tworek, “Fighting Hate with Speech Law: Media and German Visions of Democracy” (2021) 25:2 *J of Holocaust Research* 106 at 110–11.

⁸⁷ See Nunziato, *supra* note 84 at 1533.

⁸⁸ See Tworek, *supra* note 86 at 112.

⁸⁹ See Wolfgang Schulz, “Regulating Intermediaries to Protect Privacy Online – The Case of the German NetzDG” (2018) in Marion Albers & Ingo Sarlet, *Personality and Data Protection Rights on the Internet* at 5 (forthcoming).

⁹⁰ See Victor Claussen, “Fighting Hate Speech and Fake News. The Network Enforcement Act (NetzDG) in Germany in the Context of European Legislation” (2019) 3 *Fake news, pluralismo informativo e responsabilità in rete* 110 at 119–24.

to government fines. Though there are a few key differences. The Canadian regime seems to take oversight a step further by allowing decisions to be appealed to the Digital Recourse Council of Canada and by omitting the distinction between manifestly unlawful and just unlawful content. This section will apply Balkin's free speech as a triangle approach to the proposed Canadian legislation and identify how the regime will promote censorship and create prior restraint and likely inciting constitutional challenges.

a. Collateral censorship and prior restraint

The system of incentives that is created by Online Harms legislation which holds SMPs liable for failing to moderate content in a timely manner is quite straightforward. When facing penalties that attack their bottom line, SMPs will be incentivized to censor content more liberally to avoid fines.⁹¹ Rather than incentivizing a system which encourages SMPs to develop accurate and proportionate tools for moderation, a takedown regime is a blunt instrument that does not take advantage of new school moderation tools widely and creates excess moderation to prevent liability.⁹² While in the early period of social media platforms it was difficult to discern the effects of collateral censorship on freedom of expression, new transparency standards have allowed for research into the effect of collateral censorship. Based on the limited quantitative data, it appears that intermediary regimes encourage over censorship and produce a chilling effect on speech.⁹³

Prior restraint has been defined as “[c]ensorship before publication. It commonly refers to attempts to restrain publication of material prior to a court adjudicating on whether the publication would give rise to any civil action such as defamation

⁹¹ See Balkin, “Triangle”, *supra* note 3 at 2016–17.

⁹² See *ibid.*

⁹³ See Daphne Keller, “Empirical evidence of over-removal by internet companies under intermediary liability laws: an updated list” (8 February 2021), online (blog): cyberlaw.stanford.edu/blog/2021/02/empirical-evidence-over-removal-internet-companies-under-intermediary-liability-laws.

or breach of confidence”⁹⁴ or “action that prohibits speech or their expression before the speech happens.”⁹⁵ However, digital prior restraint differs in a couple of ways. Firstly, the restraint in the case of the Online Harms legislation happens without any juridical authorization, in some cases, a human moderator may not have even examined the restrained post, as is the case with algorithmic moderation. There is no obligation for the lawfulness, or permissibility of the post to be considered in a timely manner, infringing a fundamental democratic right for an indeterminate amount of time with little or no due process. Secondly, digital prior restraint is done by private companies that have been deputized by nation-states and thirdly, the restraint may occur after the post has been online for a short time.

5. Alternative regulatory regimes

As the legal world has seen the growing appetite for governance in the online sphere and misguided attempts such as the German NetzDG, several middle ways have been proposed. Intermediary immunity, information fiduciaries, and the duty to act responsibly are three legal regimes which would better protect the Canadian public from online harms while maintaining freedom of expression guarantees and working in concert with existing systems of private governance.

a. Digital-information fiduciaries

A solution proposed by Balkin, is the imposition of fiduciary duties on SMPs. They compare the relationship between SMPs and their users to the relationship between lawyers and their clients, or doctors and their patients,⁹⁶ and it is true that Facebook has probably collected enough data to know most users better

⁹⁴ See Jeffery Berryman, *Canadian Online Legal Dictionary*, 2nd ed, (Toronto: Irwin Law, 2013) sub verbo “Prior Restraint”, accessed online (14 June 2022): Irwin Law <irwinlaw.com/cold/prior-restraint/>.

⁹⁵ See Legal Information Institute, “Prior Restraint” (last accessed 22 August 2022), online: Cornell Law <www.law.cornell.edu/category/keywords/freedom_of_the_press>.

⁹⁶ See Balkin, “Triangle”, *supra* note 3 at 2048.

than their family doctor. Balkin distinguishes SMPs on the basis that they offer their services for free in exchange for the right to collect user data and send targeted ads.⁹⁷ This creates a conflict of interest between users and SMPs, as the companies are incentivized to use user data to increase profits while users expect their data not to be used nefariously.⁹⁸ This conflict can be remedied through regulation.⁹⁹ Balkin therefore proposes that SMPs should have a duty of care, confidentiality, and loyalty towards their users to curtail the misuse of user information.¹⁰⁰ The duties of care and confidentiality would require SMPs to keep user's data secure and confidential. The duty of loyalty would require SMPs not to run awry of user's expectations of what is being done with their data. Preventing companies from selling data to companies that aim to politically manipulate end users, such as the Cambridge Analytica scandal, or to use user data to encourage social media addiction. In simple terms, "the most general obligation of digital-information fiduciaries is that they may not act like con artists."¹⁰¹

b. Intermediary immunity

Regimes of intermediary immunity give SMPs immunity from liability for content that is posted on their site, so long as they meet certain criteria regarding due process, transparency, and content moderation.¹⁰² This differs from intermediary liability in that it does not create a regime whereby the government will be constantly looking over the shoulder of SMPs, ready to fine them at the first sign of trouble, and therefore encouraging collateral censorship. Instead, immunity regimes use the incentive of protection to encourage private actors to improve their existing system while increasing transparency. This approach also minimally impacts the rights of corporations and users. In the United States, intermediary immunity is already a legislated

⁹⁷ See *ibid* at 2049

⁹⁸ See *ibid*.

⁹⁹ See *ibid*.

¹⁰⁰ See *ibid* at 2051

¹⁰¹ See *ibid* at 2053.

¹⁰² See *ibid* at 2047.

norm,¹⁰³ this concept would simply require the State to attach conditions to that immunity.

A drawback to intermediary immunity is that it is unclear under what circumstances the immunity may be pierced. Derivations from the conditions of immunity may be difficult to substantiate and it may be unclear to courts and companies what standards apply. This difficulty could be remedied by a hybrid approach which uses breach of the fiduciary duties as a roadmap for when to pierce an SMPs intermediary immunity.

b. Duty to act responsibly

The Canadian Commission on Democratic Expression expounded on possible policy approaches to addressing online harms at the end of 2020. The report of the Commission, established in 2020 by the Public Policy Forum, to “understand, anticipate, and respond to the effects of new digital technologies on public life and Canadian democracy”¹⁰⁴, was authored by both American free speech experts and Canadian legal behemoths.¹⁰⁵

The report takes a broad definition of online harms, recognizing traditional criminal acts which occur online (Luring a child via a computer, non-consensual distribution or transmission of an intimate image(s) with intent to harm or humiliate, extortion, criminal harassment, indecent/harassing communications, uttering threats, fraud, identity fraud, and making or distributing child pornography), in addition to broader societal harms that come from SMPs.¹⁰⁶ The report recognizes that “[w]hen a piece of harmful speech or disinformation goes viral, it is because it was

¹⁰³ See *ibid* at 2046. Referring to section 230 of the 1996 *Telecommunications Act* in the U.S.

¹⁰⁴ See Public Policy Forum, News Release, “Canadian Commission on Democratic Expression” (2021), online: <ppforum.ca/project/canadian-commission-on-democratic-expression/>.

¹⁰⁵ Notably, Adam Dodek (Dean of University of Ottawa Common Law Section), Jameel Jaffer (Executive Director of Knight First Amendment Institute at Columbia University), The Right Honourable Beverly McLachlin (Former Chief Justice of the Supreme Court of Canada).

¹⁰⁶ See Canadian Commission on Democratic Expression, “Final Report 2020–2021” (January 2021) at 17–22, online (pdf): ppforum.ca/wp-content/uploads/2021/01/CanadianCommissionOnDemocraticExpression-PPF-JAN2021-EN.pdf.

amplified by an algorithm.”¹⁰⁷ While the report ultimately recommends regulation of speech online, it is clear that it prefers “regulating the system rather than the content.”¹⁰⁸ The report also warns against exactly the approach considered by the Canadian government, cautioning that “an attempt to tick off an exhaustive list of harms, deal with them individually and move on would be fanciful, partial and temporary.”¹⁰⁹

The Committees approach includes six steps: (1) a statutory duty to act responsibly for SMPs, large messaging groups, search engines and ISPs; (2) a regulatory body to represent the public interest and oversee a Code of Conduct (3) creation of a Social Media Council for dialogue on policies and practices (4) create transparency mechanisms (5) create an avenue for appeal, favoring an e-tribunal; and (6) developing a quick response system to ensure rapid removal of content that creates a reasonable apprehension of an imminent threat to the health and safety of a targeted person or group.¹¹⁰ Reading the framework in isolation, it is unclear how exactly the proposed regime would function. The report clarifies that the duty would “require platforms to show that reasonable measures are being taken”¹¹¹ implying that failing to moderate an individual piece of harmful content would not result in a breach so long as “reasonable measures” were in place. The Code of Conduct is meant to define the content of the Duty to Act Responsibly, however the content of both of these key elements was not canvassed by the Committee. Leaving the definition of such a key element of the approach up to the legislator gave reason for the first amendment expert on the panel to demur from the recommendations of the report.¹¹² They found “it difficult to endorse the proposed Duty to Act Responsibly when the content of the duty is left almost entirely to Parliament and the new regulator to decide.”¹¹³ Breaches of the new duty could result in the regulator imposing fines or administrative penalties based on severity, frequency, repetition, and the size and scale of the

¹⁰⁷ See *ibid* at 18.

¹⁰⁸ See *ibid* at 22.

¹⁰⁹ See *ibid*.

¹¹⁰ See *ibid* at 27–28.

¹¹¹ See *ibid* at 29.

¹¹² See *ibid* at 48.

¹¹³ See *ibid*.

offending party. The NetzDG can impose fines of 40 million dollars to SMPs that fail to meet its strict timelines and the UK government approach would allow fines to equal 5 to 10% of global revenues.¹¹⁴ It is unclear what level of fines the committee recommends.

d. The benefits of soft-Law approaches to online harms

The regimes proposed by Balkin and the Commission on Democratic Expression largely avoid the problems of collateral censorship and prior restraint which the proposed Canadian legislation is sure to bring. The regimes also better fit the contours of the internet governance landscape by acknowledging points of friction within the expression as a triangle theory, engaging with the probabilistic nature of online content moderation, creating incentives to make systems more accurate instead of airing on the side of caution, and encouraging public-private cooperation to prevent the creation of yet another overburdened administrative law regime.

The crux of expression as a triangle is that states have a democratic imperative to regulate private companies while at the same time having a constitutional imperative not to unjustifiably impede the rights of said companies and their customers. Take-down regimes of intermediary liability get the equilibrium wrong in balancing these two imperatives. The ideal approach must focus on the system employed by the companies and not the speech of the individual users, should states hope to respect the public/private divide and shield their legislation from claims of unconstitutionality. Such an approach might still intervene in situations involving an individual pieces of content, but only to correct the system which allowed that content to be proliferated, such as the case of the January 6th insurrection.¹¹⁵ The approach is in essence substance neutral.

Soft-law approaches are less likely to create problems relating to the public/private divide because there is limited

¹¹⁴ See *ibid* at 33.

¹¹⁵ See Alan Suderman & Josua Goodman, "Amid the Capitol Riot, Facebook faced its own Insurrection", AP (23 October 2021), online: apnews.com/article/donald-trump-technology-business-social-media-media-07124025bdbeba98a7c7b181562c3c1a.

government involvement; the government acts only on high level governance issues and in extreme cases. While the duty to act responsibly and duties of information fiduciaries remain underdefined, they nonetheless seem to strike the appropriate balance by imposing flexible obligations of a systematic nature, which can be corrected through litigation by users or government in exceptional circumstances.

General obligations embrace the probabilistic and proportional nature of content moderation explored in the section 3. In employing a governance system which does not impose fines for every instance of harmful content that slips through the hybrid moderation system employed by major SMPs, the government encourages private actors to continue to hone the accuracy of those systems while reaping the pragmatic benefits of automation. The duty of responsibility regime also provides an important opportunity for public-private cooperation. The interplay between organisms like the Facebook Oversight Board and a Canadian regulator would likely result in valuable dialogue between the two decision-makers.

Impactful online harms legislation could likely incorporate elements of all three solutions. Intermediary immunity would be the rule with breaches of fiduciary duties, or the duty to act responsibly, allowing for an exception. The duty to act responsibly seems to have considerable overlap with Balkin's fiduciary duties; an SMP that is breaching their duty of care or loyalty is likely behaving irresponsibly. However, where the duty to act responsibly becomes murky is in the extent to which it would place substantive duties on SMPs. Fiduciary duties of care, confidentiality, and loyalty are, conversely, more procedural. They place obligations on fiduciaries to act in certain ways not necessarily to do certain things. At a minimum one could base a statutory fiduciary duty off of those that already exist in Canadian business law, requiring SMPs to:

- (a) Act honestly and in good faith with a view to the best interests of a democratic society; and
- (b) Exercise the care diligence and skill that a reasonably prudent person would exercise in comparable circumstances.¹¹⁶

¹¹⁶ This language comes directly from section 122(1) of the *Canada Business Corporations Act* (RSC 1985, c C-44), which articulates the directors/officers statutory duty of care towards the company.

The best interests of a democratic society should be defined in this legislation and should include considering, at the very least, the interests of: government, citizens, minorities, businesses, and the SMPs self-interest. Such a framework would allow for sufficient discretion for SMPs to make decisions about the type of moderation they would like to do, the type of product they would like to offer the public and give adequate deference to their expertise. Such a duty would only place a procedural obligation on them to consider how their business model would affect the democratic function. Under this definition, individual citizens could bring actions against SMPs as well as the public regulator, however it is unclear whether this would be a standard of negligence or strict liability. It may be overly difficult to show causation and damages for breaches of this duty, but conversely a strict liability standard may invite constitutional difficulties. This standard seems likely to be able to hold SMPs liable in situations where they act irresponsibly, covering egregious cases, while leaving individual content moderation decisions to the algorithm and individual moderators. Given this limited ambit, it follows that penalties should be on the higher side of the spectrum contemplated by the Commission for on Democratic Expression, in the ballpark of 5 to 10% of global revenues, to provide an undeniable incentive for SMPs to abide by their duties.

Conclusion

The marketplace of ideas is dead. The challenge of regulating online harms is the challenge of regulating the metaverse of ideas. Regulators must understand the structural differences of regulating speech in this new arena and must be attentive to new technologies that continue to transform the landscape. The age of free speech absolutism is becoming a thing of the past, despite the yearnings of the ideologies remaining adherents,¹¹⁷ and the attractive normative purity of the theory. The

¹¹⁷ See generally, Clyde Wayne Crews Jr, "The Case against Social Media Content Regulation: Reaffirming Congress' Duty to Protect Online Bias, 'Harmful Content', and Dissident Speech from the Administrative State" (28 June 2020), online (pdf): *Competitive Enterprise Institute* <deliverypdf.ssrn.com/delivery.php?ID=199020086001027023127014115127127099104015006077091033071072022106070081067003076103097114000125006036111066126071066125065126062015046052031093>

democratic imperative for regulating SMPs is simply too high. Not to mention the evolution of the right to freedom of expression over the past thirty years. When thirty years ago one would have to yell from a soapbox on a street corner to proliferate their ideas, they may now do so anonymously, at the click of a button. Citizens have an endless option of mediums through which to express themselves; a limitation on one of these mediums does not necessarily diminish the citizens ability to proliferate it through another medium. This new breadth of the content of the right of expression should help free speech absolutists to swallow the hard pill of content moderation.

In this vein, the regulation of online harms is complicated by jurisdiction. In what jurisdiction does speech occurring in the metaverse fall under? That question is beyond the scope of this paper; however, it is worth noting that the relationship between domestic and international responses to these harms is anything but clear. The diffuse nature of hosting online data has led to the premise that online harms require an international response, yet domestic legislation has also been seen to impact the norms that SMPs use to regulate globally.¹¹⁸ The result is a sometimes complementary sometimes counterintuitive relationship as seen by the move of dominant moderation ideology from an American First Amendment standpoint (domestic influence) to one that embraces proportionality (international influence).¹¹⁹

This paper makes three arguments: First, that the landscape of freedom of expression has been changed by the triangular nature of free speech, new school regulation, and hybrid moderation. Second, the Canadian government's approach fails to account for these structural variables instead opting for an approach that is underinclusive and creates unnecessary collateral censorship and prior restraint. And, third, that alternative approaches such as intermediary liability, information fiduciary duties, and the duty to act responsibly are better suited to the new freedom of expression landscape and

[09812201601506712011009503305309209211306508707000409708310307600608607602](https://doi.org/10.1009/503305309209211306508707000409708310307600608607602)>. See also, John Samples, "Why the Government Should Not Regulate Content Moderation of Social Media" (9 April 2019), online: *Cato Institute* <www.cato.org/policy-analysis/why-government-should-not-regulate-content-moderation-social-media>.

¹¹⁸ See Tworek, *supra* note 86 at 121.

¹¹⁹ See Klonick, *supra* note 38.

would better balance the democratic deficit caused by the social media phenomena.

Bibliography

LEGISLATION

Bill C-36, *An Act to amend the Criminal Code and the Canadian Human Rights Act and to make related amendments to another Act (hate propaganda, hate crimes and hate speech)*, 2nd Sess, 43rd Parl

SECONDARY SOURCES: JOURNAL ARTICLES

- Alvarado, Oscar & Annika Waern, "Towards Algorithmic Experience: Initial Efforts for Social Media Contexts" (Honourable mention delivered at the 2018 CHI Conference on Human Factors in Computing Systems, 4 April 2018)
- Balkin, Jack M., "Free Speech is a Triangle" (2018) 118:7 Colum L Rev 2011.
- Balkin, Jack M., "2016 Sidley Austin Distinguished Lecture on Big Data Law and Policy: The Three Laws of Robotics in the Age of Big Data" (2017) 78:5 Ohio St LJ 1217.
- Belli, Luca, Pedro A Francisco and Nicolo Zingales, "Law of the Land or Law of the Platform? Beware of the Privatisation of Regulation and Police" in Luca Belli & Olga Cavalli, *Internet Governance and Regulations in Latin America*, 1st edition (Rio de Janeiro: FGV Direito Rio, 2019) 423.
- Crews Jr., Clyde Wayne, "The Case against Social Media Content Regulation: Reaffirming Congress' Duty to Protect Online Bias, 'Harmful Content', and Dissident Speech from the Administrative State" (2020) 4 Competitive Enterprise Institute: Issue Analysis.
- Claussen, Victor, "Fighting hate speech and fake news. The Network Enforcement Act (NetzDG) in Germany in the context of European Legislation" (2019) 3 110.
- Douek, Evelyn, "Governing Online Speech" (2021) 121:3 Colum L Rev 759.
- Elkin-Koren, Niva & Maayan Perel, "Separation of Functions for AI: Restraining Speech Regulation by Online Platforms" (2020) 24:3 Lewis & Clark L Rev 857.
- Gorwa, Robert, Reuben Binns & Christian Katzenbach, "Algorithmic content moderation: Technical and Political challenges in the

- automation of platform governance" (2020) 7:1 Big Data & Society 1.
- Klonick, Kate, "The New Governors: The People, Rules, and Processes Governing Online Speech" (2018) 131:6 Harv L Rev 1598.
- Klonick, Kate, "The Facebook Oversight Board: Creating an Independent Institution to Adjudicate Online Free Expression" (2020) 129:8 Yale L J 2418 at 2448–2451.
- Lamo, Madeline, & Ryan Calo, "Regulating Bot Speech" (2019) 66:4 UCLA L Rev 988.
- Luberda, Rachel, "The Fourth Branch of Government: Evaluating the Media's Role in Overseeing the Independent Judiciary" (2014) 22:2 Notre Dame J of L, Ethics & Pub Policy 507.
- Moon, Richard, "Report to the Canadian Human Rights Commission Concerning Section 13 of the Canadian Human Rights Act and the Regulation of Hate Speech on the Internet (2008)
- Nunziato, Dawn Carla, "The Marketplace of Ideas Online" Notre Dame L Rev 1519
- Pasquale, Frank, "Toward a Fourth Law of Robotics: Preserving Attribution, Responsibility, and Explainability, in an Algorithmic Society" (2017) 78:5 Ohio St LJ 1243.
- Schulz, Wolfgang, "Regulating Intermediaries to Protect Privacy Online – The Case of the German NetzDG" (2018) in Marion Albers and Ingo Sarlet, *Personality and Data Protection Rights on the Internet*, Forthcoming.
- Samples, John, "Why the Government Should Not Regulate Content Moderation of Social Media" (2019) 865 Cato Institute: Policy Analysis.
- Tabatabai, Fara, "A Tale of Two Countries: Canada's Response to the Peer-to-Peer Crisis and What it Means for the United States" (2005) 73:5 Fordham L Rev 2321.
- Tworek, Heidi J. S. "Fighting Hate with Speech Law: Media and German Visions of Democracy" (2021) 25:2 The J of Holocaust Research 106.
- Ullmann, Stefanie & Marcus Tomalin, "Quarantining online hate speech: technical and ethical perspectives" (2019) 22 Ethics and Information Technology 69.

York, Jillian C. and Ethan Zuckerman, "Moderating the Public Sphere" in Rikke Frank Jorgensen, *Human Rights in the Age of Platforms*, (Cambridge Massachusetts: The MIT Press, 2019) 137.

SECONDARY SOURCES: NEWSPAPERS

Bell, Stewart, "Extremist groups 'actively recruiting' military and police, Canadian intelligence report warns", *Global News* (23 August 2021), online: <globalnews.ca/news/8128463/extremist-groups-military-recruitment-report/>.

Buni, Catherine, & Soraya Chemaly, "The Secret Rules of the Internet: The Murky History of Moderation, and How it's Shaping the Future of Free Speech", *The Verge*, online: <www.theverge.com/2016/4/13/11387934/internet-moderator-history-youtube-facebook-reddit-censorship-free-speech>.

Cadwalladr, Carole & Emma Graham-Harrison, "Revealed: 50 million Facebook profiles harvested for Cambridge Analytica in major data breach", *The Guardian* (17 March 2018), online: www.theguardian.com/news/2018/mar/17/cambridge-analytica-facebook-influence-us-election.

Ghosh, Dipayan, "How the Free Market Incentivized Facebook's Harmful Monopoly", *Centre for International Governance Innovation* (6 January 2021), online: <www.cigionline.org/articles/how-free-market-incentivized-facebooks-harmful-monopoly/>.

Horwitz, Jeff, & Deepa Seetharaman, "Facebook Executives Shut Down Efforts to Make the Site Less Divisive", *The Wall Street Journal* (16 May 2020), online: <www.wsj.com/articles/facebook-knows-it-encourages-division-top-executives-nixed-solutions-11590507499>.

Karadeglija, Anja, "New hate law could have chilling effect, free speech advocates say", *National Post* (4 June 2021), online: <nationalpost.com/news/politics/new-hate-law-could-have-chilling-effect-free-speech-advocates-say>.

Kim, Catherine, "Poll: 70 percent of Republicans don't think the election was free and fair", *Politico* (11 September 2020), online: <www.politico.com/news/2020/11/09/republicans-free-fair-elections-435488>.

- Kreiss, Daniel & Matt Perault, "Four Ways to Fix Social Media's Political Ads Problem-Without Banning them", *The New York Times* (16 November 2019), Online: www.nytimes.com/2019/11/16/opinion/twitter-facebook-political-ads.html.
- Keller, Daphne, "Empirical evidence of over-removal by internet companies under intermediary liability laws: an updated list" (8 February 2021), online (blog): cyberlaw.stanford.edu/blog/2021/02/empirical-evidence-over-removal-internet-companies-under-intermediary-liability-laws.
- Mozur, Paul, "A Genocide Incited on Facebook, With Posts From Myanmar's Military", *The New York Times* (15 October 2018), online: www.nytimes.com/2018/10/15/technology/myanmar-facebook-genocide.html.
- Newton, Casey, "The Trama Floor: The secret lives of Facebook moderators in America", *The Verge* (25 February 2019), online: www.theverge.com/2019/2/25/18229714/cognizant-facebook-content-moderator-interviews-trauma-working-conditions-arizona.
- Rogers, Kaleigh, "Facebook's fact-checking program only a partial solution to disinformation, report says", *CBC* (30 July 2019), online: www.cbc.ca/news/science/facebook-fact-checking-full-fact-report-1.5230592.
- Roose, Kevin, "Facebook reverses postelection algorithm changes that boosted news from authoritative sources", *The New York Times* (16 December 2020), online: www.nytimes.com/2020/12/16/technology/facebook-reverses-postelection-algorithm-changes-that-boosted-news-from-authoritative-sources.html.
- Suderman, Alan & Joshua Goodman, "Amid the Capitol riot, Facebook faced its own insurrection", *AP* (23 October 2021), online: apnews.com/article/donald-trump-technology-business-social-media-media-07124025bdbeba98a7c7b181562c3c1a.
- Tunney, Catharine, "Canadian voters are likely to face foreign cyber interference in the next elections, says cyber spies", *CBC* (16 July 2021), online: www.cbc.ca/news/politics/election-security-threats-cse-1.6104745.
- Thompson, Elizabeth, "Poll finds 90% of Canadians have fallen for fake news", *CBC* (11 June 2019), online:

<www.cbc.ca/news/politics/fake-news-facebook-twitter-poll-1.5169916>.

Thompson, Belle Riley, "COVID and conspiracy", Open Canada (29 September 2021), online: <opencanada.org/covid-and-conspiracy/>.

Wright, Mike, "One in five underage Facebook users are 'addicted' to social media app, whistleblower says", *National Post* (26 October 2021), online: <nationalpost.com/news/world/one-in-five-underage-facebook-users-are-addicted-to-social-media-app-whistleblower-says>.

SECONDARY SOURCES: NEWS RELEASES

Canadian Centre for Cyber Security, News Release, "Cyber threats to Canada's democratic process: July 2021 update" (July 2021), online: <cyber.gc.ca/en/cyber-threats-canadas-democratic-process-july-2021-update/>.

Council of the European Union, News Release, "Assessment of the Code of Conduct on Hate on line State of Play" (27 September 2019), online (pdf): <ec.europa.eu/info/sites/default/files/aid_development_cooperation_fundamental_rights/assessment_of_the_code_of_conduct_on_hate_speech_on_line_-_state_of_play_0.pdf>.

Government of Canada, News Release, "Consultation closed: The Government's proposed approach to address harmful content online" (7 October 2021), online: <www.canada.ca/en/canadian-heritage/campaigns/harmful-online-content.html>.

Government of Canada, News Release, "Combatting hate speech and hate crimes: Proposed legislative changes to the *Canadian Human Rights Act* and the *Criminal Code*" (1 September 2021), online: <www.justice.gc.ca/eng/csj-sjc/pl/chshc-lcdch/index.html>.

Public Policy Forum, News Release, "Canadian Commission on Democratic Expression" (2021), online: <ppforum.ca/project/canadian-commission-on-democratic-expression/>.

SECONDARY SOURCES: MISCELLANEOUS

- Canada, House of Commons, Standing Committee on Justice and Human Rights, *Taking Action to End Online Hate* (June 2019) (Chair: Anthony Housefather), online: <www.ourcommons.ca/Content/Committee/421/JUST/Reports/RP10581008/justrp29/justrp29-e.pdf>.
- Canadian Online Legal Dictionary, “Prior Restraint” *Irwin Law*, online: <irwinlaw.com/cold/prior-restraint/>.
- Canadian Commission on Democratic Expression, “Final Report 2020-2021” *Public Policy Forum* (January 2021), online: <ppforum.ca/wp-content/uploads/2021/01/CanadianCommissionOnDemocraticExpression-PPF-JAN2021-EN.pdf>
- Evidence for Democracy, “Misinformation in Canada: Research and Policy Options” (2021), online: <evidencefordemocracy.ca/sites/default/files/reports/misinformation-in-canada-evidence-for-democracy-report_.pdf>.
- Legal Information Institute, “Prior Restraint” (last visited 22 August 2022), online: Cornell Law <www.law.cornell.edu/wex/prior_restraint#:~:text=In%20First%20Amendment%20law%2C%20prior,expression%20before%20the%20speech%20happens.%20>.
- L. Ceci, “Hours of video uploaded to YouTube every minute as of February 2020”, *Statista* (14 September 2021), online: <www.statista.com/statistics/259477/hours-of-video-uploaded-to-youtube-every-minute/>.
- Meta, News Release, “Community Standard Enforcement Report” (11 February 2021), online: <about.fb.com/news/2021/02/community-standards-enforcement-report-q4-2020/>.
- Meta, News Release, “Publishing Our Internal Enforcement Guidelines and Expanding Our Appeals Process” (24 April 2018), online: <about.fb.com/news/2018/04/comprehensive-community-standards/>.
- Meta, “Oversight board: The purpose of the board” (last visited 22 August 2022), online: <transparency.fb.com/en-gb/oversight/>.
- Meta, “Bylaws”, (November 2021), online: <about.fb.com/wp-content/uploads/2020/01/Bylaws_v6.pdf>.

Statista Research Department, "Worldwide digital population as of January 2021", *Statista* (10 September 2021), online: [<www.statista.com/statistics/617136/digital-population-worldwide/>](https://www.statista.com/statistics/617136/digital-population-worldwide/).

Statista Research Department, "Facebook: number of monthly active users worldwide 2008-2021", *Statista* (1 November 2021), online: [<www.statista.com/statistics/264810/number-of-monthly-active-facebook-users-worldwide/>](https://www.statista.com/statistics/264810/number-of-monthly-active-facebook-users-worldwide/).

Statista Research Department, "Most popular social networks worldwide as of October 2021, ranked by number of active users", *Statista* (16 November 2021), online: [<www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/>](https://www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/).

Statista Research Department, "Distribution of Instagram users worldwide as of October 2021, by age group", *Statista* (23 November 2021), online: [<www.statista.com/statistics/325587/instagram-global-age-group/>](https://www.statista.com/statistics/325587/instagram-global-age-group/).

Annex A: Technical Discussion Paper on Online Harms Legislation



Technical Discussion Paper

Online Harms Legislation

Minister of Canadian Heritage

Minister of Public Safety and Emergency
Preparedness

Minister of Justice and the Attorney General



Purpose

Share our vision for a safer, more inclusive online environment

Garner feedback on a technical, complex proposal

- Engage Canadians on a novel approach to regulating social media platforms
- Present a specific, detailed legislative direction
- Seek Canadians' views on the elements of the proposal

Use feedback to help design legislation to introduce in the Fall

Context

94% of Canadian adults have an account on at least one social media platform

Harmful content online – a serious and growing problem



1 in 5 Canadians have experienced some form of online hate

58%

of **women** in Canada have been **victims of violence online**

3x

Racialized Canadians are almost **three times** more likely to have experienced harmful behaviour online

1,106%

Increase in online child sexual exploitation reports received by the RCMP National Child Exploitation Crime Centre between 2014 to 2019

Context

Canadians want something to be done

60%

of Canadians think
there should be **more
regulation of online
hate speech**

80%

of Canadians support
requirements **to remove
racist or hateful content
within 24 hours**

There is a clear role for Government

- Efforts by social media platforms are inconsistent and not enough
- Like-minded countries have developed their own approaches

The Government has committed to act and has
developed a proposal

Vision



Serve the public interest online



Support safe and inclusive digital expression



Provide additional tools to confront online harms

Proposal

Module 1

A new legislative and regulatory framework for social media



Module 2

Modifying Canada's existing legal framework



Set new rules and define scope of new legislation



Create new regulatory bodies



Explore how to engage law enforcement & CSIS



Update Mandatory Reporting Act



Explore options to update CSIS Act

Module 1: A new legislative and regulatory framework for social media

Set new rules and define scope of new legislation



Set new rules for social media platforms

- Obligation to remove 5 categories of harmful content
- Harmful content to be removed within 24 hours of being flagged
- Transparency, reporting and preservation requirements
- Procedural fairness for users, victims, and advocacy groups
- Direct internet service providers (ISPs) to block access in Canada as a last resort with a court order, for platforms that persistently do not comply with orders to take down child sexual exploitation and terrorist content

Provide checks on platform decisions

- Provide an appeal mechanism for content moderation decisions by platforms
- Order the removal of harmful content when platforms get it wrong

Module 1: A new legislative and regulatory framework for social media

Set new rules and define scope of new legislation



Target **five** categories of harmful content, drawing on *Criminal Code*:

- 1 Hate speech
- 2 Child sexual exploitation content
- 3 Non-consensual sharing of intimate images
- 4 Incitement to violence content
- 5 Terrorist content

Module 1: A new legislative and regulatory framework for social media

Set new rules and define scope of new legislation



Legislation would apply to 'Online Communication Service Providers (OCSPs)

OCSPs:



Exemptions for private communications and telecommunications

Excluded:



Legislation would not apply to products and services that are not OCSPs

**Not
OCSPs:**



Module 1: A new legislative and regulatory framework for social media

Create new regulatory bodies

A new **Digital Safety Commission** of Canada supporting three new bodies:

1 Digital Safety Commissioner of Canada

- Oversee and enforce new rules
- Set norms and build a base of research for online safety

2 Digital Recourse Council of Canada

- Provide independent recourse through a digital tribunal system
- Make binding decisions on content removal

3 Advisory Board

- Provide expert advice and guidance to the Commissioner and the Recourse Council
- Bring expert, equity-deserving, and Indigenous interests to social media regulation

Module 1: A new legislative and regulatory framework for social media

Engaging law enforcement and CSIS



Set new preservation requirements:

- Require platforms to preserve potentially illegal content and content of national security concern falling within the five categories of harmful content
- Prevent platforms from deleting content and important identifying information that could be lawfully obtained (i.e. judicial authorizations) for use in future investigations

Explore 2 options to alert law enforcement and CSIS of certain forms of harmful content under the five categories

<i>OPTION</i>	<i>Scope of Content</i>	<i>Information sent by platforms</i>
Notify law enforcement where the content suggests an imminent risk of serious harm	Where there are reasonable grounds to suspect that there is an imminent risk of serious harm to any person or to property	The content itself plus any additional public-facing information as prescribed by the GiC regulations to law enforcement
Report prescribed content of criminal concern to law enforcement and content of national security concern to CSIS	Certain types of potentially criminal content and content of national security concern – thresholds and specific offences to be set through Governor-in-Council regulations	The content itself plus any additional public-facing information as prescribed by the GiC regulations to law enforcement and/or CSIS

Module 2: Modifying Canada's existing legal framework

Update and modernize the Mandatory Reporting Act

- Centralize and clarify the legal requirements for the mandatory reporting of child pornography offences
- Clarify the application and scope of the law, including requiring information to assist in promoting compliance with the Act
- Extend the legally required preservation period for information related to child pornography offences

Explore 2 options to require ISPs to report certain information in their mandatory reporting when a child pornography offence is clearly evident

OPTION	What's included
Require provision of transmission data in mandatory reports	IP address, date, time, type, origin, and destination associated with the material in question
Require provision of Basic Subscriber Information in mandatory reports	Transmission data + Customer name, address, phone number, billing information associated with the IP address



Module 2: Modifying Canada's existing legal framework

Explore amending the CSIS Act:

- Provide CSIS with a new judicial authorization for obtaining Basic Subscriber Information, akin to a law enforcement *Criminal Code* production order
- Enable CSIS to quickly identify perpetrators behind threats to national security in a rapidly-evolving online environment
- Could be used for investigating national security threats beyond terrorist content, including foreign interference and espionage
- Would be subject to checks and balances, including Ministerial oversight and possible review by the National Security and Intelligence Review Agency

Hate speech: Linkages with Bill C-36

New legislation is designed to synchronize with Bill C-36



Bill C-36: Amending the *Canadian Human Rights Act* and *Criminal Code*

- Provides for re-enactment of section 13 of the *Canadian Human Rights Act* (CHRA), making it a discriminatory practice to communicate hate speech online
- Section 13 would not apply to social media platforms regulated under online harms legislation

Online harms legislation would target hate speech on social media platforms

- Definition of hate speech to be aligned with definition in Bill C-36

Separate but complimentary tracks for addressing hate speech online

- Complaints against *social media platforms* → Digital Safety Commissioner
- Complaints against *individuals and websites* → CHRA

Seeking Input

The Government is publishing:

- Narrative description of proposal
- Technical discussion paper containing elements of a legislative proposal

Comments now open on these documents

- Comments open until September 25, 2021
- Send input to: pch.icn-dci.pch@canada.ca

What comes next: Fall 2021

- Use feedback to help design legislation