# CIS' Comments to the Christchurch Call

October 2019

By **Tanaya Rajwade, Elonnai Hickok, and Raouf Kundil Peedikayil**

Edited by **Gurshabad Grover and Amber Sinha**

**The Centre for Internet and Society, India**

# Introduction

In the wake of the Christchurch terror attacks, the Prime Minister of New Zealand, Jacinda Ardern, and the President of France, Emmanuel Macron co-chaired the Christchurch Call to Action in May 2018 to "bring together countries and tech companies in an attempt to bring to an end the ability to use social media to organise and promote terrorism and violent extremism."[1] Fifty one supporters, including India, and eight tech companies have jointly agreed to a set of non-binding commitments and ongoing collaboration to eliminate violent and extremist content online.[2] Facebook, Microsoft, Twitter, Google, and Amazon are all among the online service provider signatories that released a joint statement welcoming the call and committing to a nine-point action plan.

The Call has been hailed by many as a step in the right direction, as it represents the first collaboration between governments and the private sector companies to combat the problem of extremist content online at this scale. However, the vagueness of the commitments outlined in the Call and some of the proposed mechanisms have raised concerns about the potential abuse of human rights by both governments and tech companies.

This response is divided into two parts - Part One examines the call through the lens of human rights, and Part Two thinks through the ways in which India can adhere to the commitments in the Call, and compares the current legal framework in India with the commitments outlined in the Call.

# Part One: Response to Christchurch Call

## High level comments

### Enabling Civil Society Participation

The civil society response to the Call expressed concerns with the exclusion of civil society in the drafting and finalisation of the pledge.[3] As a stakeholder representing the public interest including the right to privacy and freedom of expression, any policy changes and attempts at defining terrorist and violent extremist content born out of the Call should be subject to consultations and independent review by civil society as far as possible.

### Ensuring Clarity of Purpose

---

[1] The Christchurch Call, <https://www.christchurchcall.com/christchurch-call.pdf>.

[2] ibid.

[3] 'Civil Society Positions on Christchurch Call Pledge' (EFF, 16 May 2019) <https://www.eff.org/files/2019/05/16/community_input_on_christchurch_call.pdf>.

It is critical that the Call clearly states its purpose and objectives to avoid function creep and politicisation of the issue. The call should also publish the framework and criteria followed for acceptance of countries for membership.

## Need for a shared understanding and metrics of success

The Call creates commitments for governments and companies, and draws upon the expertise of civil society. These stakeholders come from different backgrounds and contexts. While the Call aims to 'eliminate terrorist and extremist content online', a realistic shared understanding of success and parameters for its assessment should be evolved based on the roles played by these stakeholders in society, their respective constraints and obligations, as well as the differential levels of resources at their disposal.

## Placing the Call in the context of other initiatives

While the Call references other initiatives that exist, such as the EU Internet Forum, the G20 and G7, the Global Internet Forum to Counter Terrorism, the Global Counterterrorism Forum, Tech Against Terrorism, and the Aqaba Process - it is unclear how the Call complements, contrasts, or builds upon these different initiatives.

*There are also other international efforts that the Call can look to as guidance and best practice.* For example, the United Nations has set out detailed guidance on establishing CVE strategies.[4] Additionally, the Ankara Memorandum on Good Practices for a Multi-Sectoral Approach to CVE serves as guidance for GCTF members to devise domestic strategies in line with international human rights frameworks while paying due regard to their diverse histories and cultural contexts.[5] The Zurich-London Recommendations on Preventing and Countering Violent Extremism and Terrorism Online,[6] which have been created by the Global Counterterrorism Forum, provide a blueprint for combatting online terrorism and violent extremism while respecting rights.

## Need for transparency and accountability on actions taken by signatories

The commitments in the Call are voluntary, devoid of accountability mechanisms to measure fulfilment of the commitments by companies and governments.  With a view to encouraging

---

[4] 'Plan of Action to Prevent Violent Extremism' (United Nations) <https://www.un.org/ga/search/view_doc.asp?symbol=A/70/674>.

[5] 'Ankara Memorandum on Good Practices for a Multi-Sectoral Approach to Countering Violent Extremism' (Global Counterterrorism Forum, 1 September 2016) <https://www.thegctf.org/Portals/1/Documents/Framework%20Documents/A/GCTF-Ankara-Memorandum-ENG.pdf?ver=2016-09-01-114735-333>

[6] Countering Violent Extremism (CVE) Working Group, 'Zurich-London Recommendations on Preventing and Countering Violent Extremism and Terrorism Online' (Global Counterterrorism Forum, 15 September 2017) <https://www.thegctf.org/Portals/1/Documents/Framework%20Documents/A/GCTF%20-%20Zurich-London%20Recommendations%20ENG.pdf?ver=2017-09-15-210859-467>

debate and fostering accountability on actions taken under the Call, information like timeframes for actions, progress and action taken reports, allocation of funding, etc. could potentially be shared by signatories with the public.

## Need for focus on the differentiation of online service providers based on the manner of transmission

Extreme content can be created and spread via services offered by online service providers. It is important for the Call to recognise the myriad ways in which extremist content can be created and disseminated while emphasising the need for appropriate mechanisms tailored for different types of service providers. For example, ISPs should not be subject to the same requirements as social media companies. The Call is an opportunity to commit governments to developing regulation that is appropriate for the type of intermediary being targeted.

## Need for governments to define online extremist content in line with international human rights standards to prevent abuse

The high standards placed by the International Covenant on Civil and Political Rights[7] on any legislation curtailing freedom of expression makes it difficult to arrive at a sufficiently precise definition of extremism which satisfies legality.[8] As such, a broadly-accepted and internationally-agreed definition of extremism has not been reached, though there are multiple definitions developed at national and international levels. Most international approaches to defining violent extremism take note of the various types of extremist actions and the objective-driven nature of the action.

While it is not within the scope of the Call to arrive at a global definition of online extremism or extremist content, the Call should ask governments to arrive at a clear national definition for 'extremist content' which is in line with international human rights standards. Establishing a precise legal definition of the extremist content targeted under the Call is necessary to safeguard human rights from overreach and misuse by governments, and provides important clarity.

The guidance regarding the scope of the Christchurch Crisis Protocol may be useful in this regard; while it does not define terrorist or violent extremist content, it lays down parameters for the definition.[9] It creates explicit exceptions for political speech and non-violent protest,

---

[7] International Covenant on Civil and Political Rights, art 19.3

[Freedom of expression] will therefore be subject to certain restrictions, but these shall only be such as are provided by law and are necessary:

       (a) For respect of the rights or reputations of others;

       (b) For the protection of national security or of public order (ordre public), or of public health or morals.

[8] Shepherd, 'Extremism, Free Speech and the Rule of Law: Evaluating the Compliance of Legislation Restricting Extremist Expressions with Article 19 ICCPR' (2017) 33(85) Utrecht Journal of International and European Law 62.

[9] Defining a Terrorist Attack

and advocates for an action-based approach as opposed to an entity-based approach. This resolves the problems of counter-terrorism laws being weaponised against dissent[10] and the criminalisation of medical assistance.[11]

# Detailed Comments

**Need for greater focus on the role of the recommendation algorithms in extremist content dissemination**

The commitment in the Call for companies to 'review of algorithms or processes which drive users towards [harmful] content' is a welcome step. Content recommendation algorithms have been criticised for radicalising users[12] by basing advertising revenues on engagement, which is measured through metrics such as watch time.[13] These algorithms could potentially drive the

---

"There are a wide range of different definitions of terrorist and violent extremist content among participating Governments and Online Service Providers. These guidelines will not seek to create a common definition, but they will look for coherency in the interpretation of what is regarded as a terrorist attack for the purposes of these guidelines.

Individual participants should refer to a definition that includes: direct violence or the threat of harm to the civilian population with the intention of provoking terror or compelling a Government or International Organisation to do or to refrain from doing something, linked to a high risk of online platforms being instrumentalised in such an attack, including an intent to disseminate content as part of the attack.

     Political speech, protest, or other non-violent forms of expression are not valid reasons for action under these guidelines;

     Bearing in mind that terrorist attacks may be perpetrated by individuals or entities that are not formally listed by the United Nations Security Council Consolidated Sanctions List, and/or that the identity of perpetrators may not be immediately known, it is important to focus on assessing the actions of the perpetrators, and accomplices as well as those supporting such actions."

[10] United Nations Human Rights Council 'Human Rights Council discusses the protection of human rights while countering terrorism, and cultural rights' (1 March 2018) <http://www.ohchr.org/EN/HRBodies/HRC/Pages/NewsDetail.aspx?NewsID=22742&LangID=E>.

[11] Marine Buissonniere, Sarah Woznick, and Leonard Rubenstein, 'The Criminalisation of Healthcare' <https://www1.essex.ac.uk/hrc/documents/54198-criminalization-of-healthcare-web.pdf>

[12] Mathew Ingram, 'YouTube's secret life as an engine for right-wing radicalization' (Columbia Journalism Review, 19 September 2018) <https://www.cjr.org/the_media_today/YouTube-conspiracy-radicalization.php>; 'Opinion | YouTube, the Great Radicalizer - The New York Times' (New York Times, 10 March 2018) <https://www.nytimes.com/2018/03/10/opinion/sunday/YouTube-politics-radical.html>.

[13] Paul Covington, Jay Adams and Emre Sargin, 'Deep Neural Networks for YouTube Recommendations' <https://storage.googleapis.com/pub-tools-public-publication-data/pdf/45530.pdf>

users towards increasingly extremist content.[14] Recommendation algorithms may also have a role to play in building extremist networks on social media platforms.[15]

While attempts to address this have included updating community guidelines,[16] deployment of additional content reviewers,[17] and altering of recommendation algorithms,[18] the impact of such policy changes is yet to be assessed. Moreover, the implications of hostile content removal policies leading to relocation of extremists to secure, encrypted channels of communication, particularly on law enforcement, require further research.[19]

*The Call is an opportunity to commit companies to researching the role that an engagement driven business model can play in the dissemination of extremist content and develop specific measures, including alternative business models, to prevent the same. Companies can create tools that bring transparency to how algorithms work and give users granular control over the content that they view and the way content is disseminated on the platform. For example, giving individuals the choice to opt out of receiving content via a recommendation algorithm.[20]*

---

[14] Rodriguez, 'YouTube's recommendations drive 70% of what we watch' (Quartz, 13 January 2018) <https://qz.com/1178125/YouTubes-recommendations-drive-70-of-what-we-watch/>; Alastair Reed, Joe Whittaker, Fabio Votta and Seán Looney, 'Radical Filter Bubbles -Social Media Personalisation Algorithms and Extremist Content' (RUSI, 26 July 2019) <https://rusi.org/sites/default/files/20190726_grntt_paper_08_0.pdf>; Josephine B Schmitt, Diana Rieger, Olivia Rutkowski and Julian Ernst, 'Counter-messages as Prevention or Promotion of Extremism?! The Potential Role of YouTube: Recommendation Algorithms' (2018) 68(4) Journal of Communication 780 <https://doi.org/10.1093/joc/jqy029>; Derek O'Callaghan, Derek Greene, Maura Conway, Joe Carthy and Padraig Cunningham, 'Down the (White) Rabbit Hole: The Extreme Right and Online Recommender Systems' <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.822.7369&rep=rep1&type=pdf>. See Carole Cadwalladr, 'Google, democracy and the truth about internet search' (The Guardian, 4 December 2016) <https://www.theguardian.com/technology/2016/dec/04/google-democracy-truth-internet-search-facebook> for an understanding of how search engines optimisation is leveraged to push extremist content.

[15] Carole Cadwalladr, 'Google, democracy and the truth about internet search' (The Guardian, 4 December 2016) <https://www.theguardian.com/technology/2016/dec/04/google-democracy-truth-internet-search-facebook>.

[16] 'Our Ongoing Work to Tackle Hate', (Official YouTube Blog, 25 June 2019) <https://YouTube.googleblog.com/2019/06/our-ongoing-work-to-tackle-hate.html>; 'Limited features for certain videos' (YouTube Help) <https://support.google.com/YouTube/answer/7458465>.

[17] Rishika Chatterjee and Paresh Dave, 'YouTube set to hire more staff to review extremist video content' (The Independent, 5 December 2017) <https://www.independent.co.uk/news/business/news/YouTube-extremist-video-content-isis-neo-nazis-jihadi-white-supremacist-alphabet-a8092311.html>.

[18] Elizabeth Dwoskin, 'YouTube is changing its algorithms to stop recommending conspiracies' (The Washington Post) <https://www.washingtonpost.com/technology/2019/01/25/YouTube-is-changing-its-algorithms-stop-recommending-conspiracies/?noredirect=on&utm_term=.1262e4533055>.

[19] Mia Bloom, Hicham Tiflati, and John Horgan, 'Navigating ISIS's Preferred Platform: Telegram' (2019) 31(6) Terrorism and Political Violence 1242 <https://www.tandfonline.com/doi/abs/10.1080/09546553.2017.1339695?journalCode=ftpv20>.

[20] The use of an opt-out option was suggested by a former YouTube engineer. Kelly Weill, 'How YouTube Built a Radicalization Machine for the Far-Right' <https://www.thedailybeast.com/how-YouTube-pulled-these-men-down-a-vortex-of-far-right-hate>.

## Need for metrics and ethical frameworks for counter narratives and re-direct methods

The Call commits companies to *"[r]eview the operation of algorithms and other processes that may drive users towards and/or amplify terrorist and violent extremist content…this may include using algorithms and other processes to redirect users from such content or the promotion of credible, positive alternatives or counter-narratives."*

Though 'counter narratives' are one form of response that have been adopted by governments to discredit extremist narratives, experts have questioned the efficacy[21] and ethical implications[22] of such measures seeking to intentionally shape user behaviour.  As with counter-narratives, the efficacy of 're-direct campaigns', which direct users towards de-radicalising content,[23] has not been assessed.

In addition to the 'redirect' method,  companies can consider placing extremist or violent topics on a 'Do not amplify' list or algorithmically suppress extremist content using quality indicators such as legitimacy of content creators, the channels of dissemination, and the content matter.[24] However, since such methods  identify the types of speech to be suppressed, it is problematic for these to be created unilaterally. Creation of such standards should be a collaborative process involving the private companies, the government, the civil society and the broader public, and should be subjected to a periodic review process to ensure relevance. Open standards development processes could serve as a useful model on in this regard.

*As the Call encourages companies to research and employ these methods, it ought to emphasise the need for research into understanding the effectiveness of these efforts and establishment of a framework for evaluating ethical questions arising therefrom, alongside mechanisms for oversight, accountability, and redress.*


## Establish appropriate frameworks and transparency mechanisms prior to utilizing automated means to filter content

The Call commits companies to *"accelerate research into and development of technical solutions to prevent the upload of and to detect and immediately remove terrorist and violent extremist content online…"*

---

[21] Eric Rosand and Emily Winterbotham, 'Do counter narratives actually reduce violent extremism' (Brookings, 20 March 2019)  <https://www.brookings.edu/blog/order-from-chaos/2019/03/20/do-counter-narratives-actually-reduce-violent-extremism/>.

[22]  Elsevander Berg, 'Jigsaw's Redirect Method Brainwashing the Brainwashed' (Medium) <https://medium.com/@ElsevanderBerg/jigsaws-redirect-method-brainwashing-the-brainwashed-fe281733b9c3>; Anne Speckhard and Ardian Shajkovci, 'Breaking the ISIS brand: Counter narratives Part II- Ethical Considerations in Fighting ISIS Online' (Vox Pol, 7 March 2018) <https://www.voxpol.eu/breaking-the-isis-brand-counter-narratives-part-ii-ethical-considerations-in-fighting-isis-online/>.

[23] The Redirect Method <https://redirectmethod.org/>;  Andy Greenberg, 'Google's Clever Plan to Stop Aspiring ISIS Recruits' (WIRED, 9 July 2016) <https://www.wired.com/2016/09/googles-clever-plan-stop-aspiring-isis-recruits/>.

[24]  Renee Diresta, 'Up Next: A Better Recommendation System' (WIRED, 4 November 2018) <https://www.wired.com/story/creating-ethical-recommendation-engines/>.

Significant efforts are being undertaken by the industry to develop automated solutions to identifying and removing illegal or otherwise problematic content.[25] The Terrorist Content Regulation currently being worked on by the European Parliament mandated automated content filters in its earlier drafts but it was dropped after criticism from the civil society.[26] India is considering legally mandating companies to have these tools in place through provisions pertaining to intermediary liability.[27]

While automated content removal can be a powerful tool to tackle the challenge of scale of content on these platforms, there is a significant amount of work that needs to be done to ensure that algorithms are responsive to social and cultural context.[28] Inadequacies in training data used for developing algorithms, especially in the case of rarer types of content[29] can be acute in a country like India with  diversity in languages and cultures. Gaps in the data can also lead to blind spots in the algorithm where it encounters new content, as seen in the Christchurch attack, where filters failed to categorise the video as violent because of the  head-mounted view provided by the camera worn by the gunman.[30]

Without the correct context, videos that are violent but legitimate, such as those uploaded by activists and journalists for documenting bringing human rights abuses, can be removed by automated filters. It is unclear whether automated algorithms have the context necessary to separate such content from that disseminated by extremist actors and thus, could inadvertently prevent the data collection efforts to monitor and prosecute human rights abuses.[31] Deletion of non-extremist violent content erases parts of public memory and often serves to whitewash rampant human rights abuses that are not reported by the traditional

---

[25] 'Copyright Match Tool' (YouTube Help) <https://support.google.com/YouTube/answer/7648743?hl=en>; Evan Engstrom and Nick Feamster, 'The Limits of Filtering: A Look at the Functionality and Shortcomings of Content Detection Tools' (Engine, March 2017) <http://www.engine.is/the-limits-of-filtering/>; 'Microsoft Photo DNA' <https://news.microsoft.com/download/presskits/photodna/docs/photoDNAFS.pdf>.

[26] 'EU Parliament deletes the worst threats to freedom of expression proposed in the Terrorist Content Regulation' (EDRI, 17 April 2019) <https://edri.org/eu-parliament-deletes-worst-threats-to-freedom-of-expression-terrorist-content-regulation/>.

[27] The Information Technology [Intermediaries Guidelines (Amendment) Rules] 2018, s 3(9) <https://meity.gov.in/writereaddata/files/Draft_Intermediary_Amendment_24122018.pdf>.

[28] Sam Levin, Carrie Wong and Luke Harding, 'Facebook backs down from 'napalm girl' censorship and reinstates photo' (The Guardian, 9 September 2016) <https://www.theguardian.com/technology/2016/sep/09/facebook-reinstates-napalm-girl-photo>.

[29] Bird and Bird, 'Can AI content moderation keep us safe online?' (Lexology, 15 May 2019) <https://www.lexology.com/library/detail.aspx?g=e7f3ea70-2476-4f03-bc02-6bfbfadd4624>.

[30] 'Facebook Executive Testifies on AI Failure to Detect the Christchurch Mosque Shooting Video' (Fortune, 24 April 2019) <https://fortune.com/2019/04/24/facebook-new-zealand-terrorism-artificial-intelligence-ai/>

[31] 'Tech Advocacy' (Syrian Archive) <https://syrianarchive.org/en/tech-advocacy>. Over 200,000 videos documenting human rights abuses were made unavailable between 2011 and 2019. Human rights organisations have reported similar takedown of content in the context of conflicts in Yemen and Ukraine. See Electronic Frontier Foundation, 'Caught in the Net: The Impact Of "Extremist" Speech regulations On Human Rights Content' (EFF, May 2019) <https://www.eff.org/files/2019/06/03/extremist_speech_regulations_and_human_rights_content_-_eff_syrian_archive_witness.pdf>.

media.[32] This could also affect justice mechanisms which are increasingly relying on digital and social media-based evidence while adjudicating rights violations.[33]

Additionally, opportunities to address the root causes of extremism are extinguished when platforms automatically delete extremist content without reaching out to investigators.[34] We should also note that these algorithms may be no better than human moderators at distinguishing between counter-speech, satire and political activism pertaining to extremism and extremist content per se.[35]

*The Call should emphasise the need for creating regulatory frameworks that provides transparency into the parameters of algorithms as well as their deployment, which also contain mechanisms for oversight, accountability, and redress.*

Additionally, online service providers have undertaken to '*enforce those community standards or terms of service in a manner consistent with human rights and fundamental freedoms, including by prioritising moderation of terrorist and violent extremist content, however identified*'.

Apart from the concerns with the enforcement of these 'community standards' as well as the processes undertaken towards developing them, this commitment raises other questions relating to the extent of moderation and the moderators involved. Online service providers employ human content moderators to spot and regulate online nudity, violence, gore, etc.[36] In the recent past, reports have highlighted the inhumane working conditions of these moderators[37] with these companies failing to hold contractors accountable for not enforcing workplace safety standards.[38] Moderators have reported suffering from PTSD[39] and other forms of emotional and psychological trauma as a consequence of their jobs.

---

[32] EFF (n 31)

[33] The Prosecutor v Mahmoud Mustafa Busayf Al-Werfalli, Case No.ICC-01/11-01/17-2, Warrant of Arrest (15August 2017) <https://www.icc-cpi.int/Pages/record.aspx?docNo=ICC-01/11-01/17-2>.

[34] 'Civil Society Positions on Christchurch Call Pledge' (EFF, 16 May 2019) 6 < https://www.eff.org/files/2019/05/16/community_input_on_christchurch_call.pdf>.

[35] ibid.

[36] Andrew Arsht and Daniel Etcovitch, 'The Human Cost of Online Content Moderation' (Journal of Law and Technology, 2 March 2018) <https://jolt.law.harvard.edu/digest/the-human-cost-of-online-content-moderation>.
[37] Casey Newton, 'Bodies in Seats' (The Verge, 19 June 2019) <https://www.theverge.com/2019/6/19/18681845/facebook-moderator-interviews-video-trauma-ptsd-cognizant-tampa> Trigger Warning.

[38] Burns Charest LLP, 'California Lawsuit Claims Facebook Fails to Properly Protect Content Moderators' (CISION, 24 September 2018) <https://www.prnewswire.com/news-releases/california-lawsuit-claims-facebook-fails-to-properly-protect-content-moderators-300717600.html> Trigger Warning.
[39] 'Scola v Facebook' (RegMedia, 24 September 2018) <https://regmedia.co.uk/2018/09/24/scola_v_facebook.pdf> Trigger Warning.

*While we acknowledge that AI-based identification and takedown of extremist content is neither ideal nor viable, the Call is an opportunity to advocate for humane and fair working conditions for all content moderators employed by online service providers.*

## Need for robust transparency commitments

The Call commits service providers to *"provide greater transparency in the setting of community standards or terms of service including by: outlining and publishing the consequences of sharing terrorist and violent extremist content and describing policies and putting in place procedures for detecting and removing terrorist and violent extremist content"* and to *"implement regular and transparent public reporting in a way that is measurable and supported by clear methodology on the quantity and nature of terrorist and violent extremist content being detected and removed."*

Transparency is a key enabler for accountability and redress. It is particularly critical today as online content moderation is a rapidly evolving and contested process. Increased examination and input from multiple stakeholders will help to address emerging challenges and create an important check and balance to actions by the government and private companies. Currently, post facto aggregate transparency only allows for other stakeholders to call out statistics and retrospectively provide feedback to company and government action. Increased transparency by the government and the private sector could take the form of:

● *Government transparency of requests*: Governments should publish requests made to service providers to take down content or restrict services, as well as the rationale for the request as they are issued. This transparency will be important in ensuring the government does not misuse its power to block dissenting voices online.

● *Company transparency of government requests and user flags*: Companies should tag and allow for real time tracking of government and user flags received, complied with or rejected, and reasons for the decision. This would allow the public and the government to keep a check on what happens to requests, as well as how they are prioritised.

● *Company transparency of violations of community guidelines*: Given the current lack of transparency around content moderation,[40] companies should incorporate best practices around content moderation and transparency.[41]

● *Transparency of parameters and audit of algorithms*: Following research and impact assessment around automated technologies as recommended earlier, companies should

---

[40] Global Net Policy, 'The Santa Clara Principles on Transparency & Content Moderation' (Global Net Policy, May 2018) <http://globalnetpolicy.org/wp-content/uploads/2018/05/Santa-Clara-Principles_t.pdf>; Gennie Gebhart 'Who Has Your Back? Censorship Edition 2019' (EFF, 12 June 2019) <https://www.eff.org/wp/who-has-your-back-2019>

[41] The Santa Clara Principles on Transparency and Accountability in Content Moderation <https://santaclaraprinciples.org/>; 'Who Has Your Back' (n 40); Ryan Budish, Liz Woolery, and Kevin Bankston, 'The Transparency Reporting Toolkit' <https://www.newamerica.org/oti/policy-papers/the-transparency-reporting-toolkit/>.

share the parameters used for identifying extremist content in both manual or automated content moderation, as well as allow periodic public audits of algorithms used to filter content.

● *Transparency of policies and associated changes*: Companies should clearly explain their terms of service and community guidelines, and archive and present the changes in a way that users can clearly see how these evolve over time.

## Need for stronger and more explicit commitment to redress mechanisms for actions taken against content and users

The Call commits service providers to *"enforce those community standards or terms of service in a manner consistent with human rights and fundamental freedoms, including by: Providing an efficient complaints and appeals process for those wishing to contest the removal of their content or a decision to decline the upload of their content."*

Robust complaint and appeal mechanisms are critical in ensuring the protection of users' freedom of expression. Even as appeal mechanisms provided by the platforms have improved, a lack of meaningful notice for content takedowns undermines this progress by creating a gap whereby users lack information needed for appeal the decisions.[42]

The Call could go a step further and also outline cornerstones for these mechanisms including notice to the user when action is taken against their content, the specific guideline violated by their content, how they can contest the decision, and timeframes for companies to respond to these appeals. The Call should further commit governments to protect the right of the user to seek legal recourse against the violation of their right of expression, if the remedy mechanism provided by the online service provider proves unsatisfactory or inadequate.

## Need to commit governments to develop legislation and measures of enforcement in line with international human rights law

The Call commits governments to *"ensure effective enforcement of applicable laws that prohibit the production or dissemination of terrorist and violent extremist content, in a manner consistent with the rule of law and international human rights law, including freedom of expression."*

Prioritising efforts to remove or block content has been found to be insufficient in controlling extremist content, and comes with the risk of stifling free speech, as has been noted by civil society actors in response to the Call.[43] Laws like the German NetzDG have not demonstrated

---

[42] Who Has Your Back (n 41).

[43] Civil Society Positions on Christchurch Call Pledge (n 3).

a significant impact on elimination of extremist content,[44] with internet companies being aggressive in censoring non-extremist content.[45]

## Need to explore a range of technological and UI solutions to preventing the weaponisation of platforms and features

The Call commits companies to implement immediate, effective measures to mitigate the specific risk that terrorist and violent extremist content is disseminated through livestreaming, including identification of content for real-time review.

Companies have committed to identifying appropriate checks on livestreaming, enhanced vetting measures (such as streamer ratings or scores, account activity, or validation processes) and moderation of certain livestreaming events. Checks on livestreaming would be tailored for the specific livestreaming services, depending on the audience, the nature of the livestreaming service, and the likelihood of exploitation.[46]

Despite these measures, this commitment to immediate, effective measures remains unrealistic given the state of AI used for monitoring live streaming, especially in the case of terrorist events where community reporting remains the first line of defense relied on by live streaming platforms to trigger accelerated review.[47] However, given the unreliability of user reporting, mounting regulatory[48] and media[49] pressure against live streaming, it has been argued[50] that live streaming ought to be discontinued until AI technology is sufficiently developed.

As opposed to fully discontinuing live streaming, we recommend that friction be added to such services by restricting them to accounts with a certain minimum number of subscribers,[51]

---

[44] Echikson and Knodt, 'Germany's NetzDG: A Key Test for Combating Online Hate' (CEPS Policy Insight, 22 November 2018) <https://ssrn.com/abstract=3300636>.

[45] Mark Scott and Janosch Delcker, 'Free speech vs. censorship in Germany' (Politico, 1 April 2019) <https://www.politico.eu/article/germany-hate-speech-netzdg-facebook-YouTube-google-twitter-free-speech/>.

[46] Microsoft, 'Joint Statement in Support of Christchurch Call' (Microsoft Blog, 5 May 2019) <https://blogs.microsoft.com/wp-content/uploads/prod/sites/5/2019/05/Christchurch-Call-and-Nine-Steps.pdf>.

[47] Facebook, 'A Further Update on New Zealand Terrorist Attack' (Facebook Newsroom) <https://newsroom.fb.com/news/2019/03/technical-update-on-new-zealand/>.

[48] 'Australia targets tech firms with "abhorrent material" laws' (BBC News, 4 April 2019). <https://www.bbc.com/news/world-australia-47809504>.

[49] Newton, 'The world is turning against live streaming' (The Verge, 4 April 2019) <https://www.theverge.com/interface/2019/4/4/18294951/australia-live-streaming-law-facebook-twitter-periscope>.

[50] Doffman, 'Facebook Admits It Can't Control Facebook Live- Is This The End For Live Streaming?' (Forbes, 24 March 2019) <https://www.forbes.com/sites/zakdoffman/2019/03/24/could-this-really-be-the-beginning-of-the-end-for-facebook-live/>.

[51] 'Restrictions on Live Streaming' (YouTube Help) <https://support.google.com/YouTube/answer/2853834?hl=en>.

adding behavioural interventions to slow down dissemination of content by requiring users to pass captcha tests, and reminding users to verify content.

# Part Two: Indian Framework and the Call

This section analyses the commitments that India has signed on to under the Call and existing efforts in India. In doing so, it makes recommendations on potential steps that India could take as a signatory to the Call.

**Need to emphasise comprehensive, coordinated, and multi-faceted efforts via national strategies for countering violent extremism (CVE)**

The Call commits governments to *"[c]ounter the drivers of terrorism and violent extremism by strengthening the resilience and inclusiveness of our societies to enable them to resist terrorist and violent extremist ideologies, including through education, building media literacy to help counter distorted terrorist and violent extremist narratives, and the fight against inequality."*

It is commendable that the Call recognises that the roots of online extremism lie outside the internet and the need for more than online measures to truly address extremism. However, it is difficult to pinpoint specific single drivers, and correlations between structural development factors and terrorism are often counter-intuitive.[52] *The signatories should encourage a robust research ecosystem around drivers of extremism, which should inform national CVE strategies.*

Inspiration can be sought from the counter terrorist recruitment films developed by the Maharashtra Anti-Terrorism Squad and played before movies in the state, to create similar content in regional languages to reach audiences across India. The Pradhan Mantri Digital Saksharta Abhiyan (PMGDISHA) is India's primary digital literacy initiative, which aims to impart IT training to at least one member in every eligible household.[53] While it seeks to develop capacities to digitally access services and provide training for employability,[54] its scope can be expanded to educating users to consider information critically and perform basic fact-checking. The Call is an opportunity for the government to build a societal culture of giving legitimacy to counter narratives.

---

[52] JM Berger, 'Making CVE work - A Focussed Approach Based on Process Disruption, International Centre for Counter Terrorism - The Hague' (ICCT Netherlands, 1 May 2016) <https://www.icct.nl/wp-content/uploads/2016/05/J.-M.-Berger-Making-CVE-Work-A-Focused-Approach-Based-on-Process-Disruption-.pdf>; D Gambetta and H Steffen, 'Uncivil Engineers: The surprising link between education and jihad' (Foreign Affairs, 10 March 2016) <https://www.foreignaffairs.com/articles/2016-03-10/uncivil-engineers>.

[53] ANI, 'Cabinet approves PMGDISHA under Digital India Programme' (Business Standard, 8 February 2017) <https://www.business-standard.com/article/news-ani/cabinet-approves-pmgdisha-under-digital-india-programme-117020801404_1.html>.

[54] Factly Media and Research (Factly) and Internet and Mobile Association of India (IAMAI), 'Countering Misinformation in India- Solutions and Strategies' <https://2nafqn3o0l6kwfofi3ydj9li-wpengine.netdna-ssl.com/wp-content/uploads//2019/02/Countering-Misinformation-Fake-News-In-India.pdf>.

**Ensure effective enforcement of applicable laws that prohibit the production or dissemination of terrorist and violent extremist content, in a manner consistent with the rule of law and international human rights law, including freedom of expression.**

India has a gamut of laws pertaining to terrorism, ranging from the Unlawful Activities (Prevention) Act, the Prevention of Money Laundering Act[55] and related rules,[56] and legislation on hate speech.[57] Additionally, there is a legal regime to block, remove, and prevent access to networks- this includes s 144 of the Code of Criminal Procedure[58] and the Telegraph Act (with the rules thereunder)[59] which are used to order internet shutdowns; s 69A of the Information Technology Act which is used by the government to order the blocking of content,[60] and often criticised for a lack of transparency and judicial accountability;[61] and the framework under s 79 of the Information Technology Act which requires intermediaries to prohibit unlawful content through their terms of use and enables content takedowns, which would extend to terrorist content and creates a framework of intermediary liability and content removal.[62] The amended rules proposed under s 79 have been criticised for imposing excessive restraints on speech, in

[55] The Prevention of Money Laundering Act 2002 (PMLA 2002) <https://enforcementdirectorate.gov.in/PreventionOfMoneyLaunderingAct2002.pdf?p1=117211488412800032Preamble>; Nikesh Tarachand Shah v Union of India Writ Petition (Criminal) No. 67 OF 2017; Abhinav Sekhri, 'Not So Civil- The Money Laundering Act and Article 20' (2017) 4 NLUD Student Law Journal 72 <https://nludelhi.ac.in/download/publication/2017/NLUD%20-%20SLJ%20(Vol.%204).pdf>.

[56] Prevention of Money Laundering Rules, r 9; Justice K S Puttuswamy and Anr v Union of India and Ors Writ Petition (Civil) No. 494 of 2012 [436] <https://www.sci.gov.in/supremecourt/2012/35071/35071_2012_Judgement_26-Sep-2018.pdf>.

[57] For a comprehensive analysis, refer to Centre for Communication Governance, Hate Speech Laws in India (NLU Delhi Press 2018) <https://drive.google.com/file/d/1pDoIwlusnM3ys-1GAYbnTPmepU22b2Zr/view>.

[58] 'Legality of Internet shutdowns under Section 144 CrPC' (SFLC, 2 October 2016) <https://sflc.in/legality-internet-shutdowns-under-section-144-crpc>.

[59] Unified Telecom Licence, cl 38.1 <http://dot.gov.in/sites/default/files/Unified%20Licence_0.pdf>; Temporary Suspension of Telecom Services (Public Emergency or Public Safety) Rules 2017 <http://dot.gov.in/sites/default/files/Suspension%20Rules.pdf>.

[60] Shreya Singhal v Union of India AIR 2015 SC 1523 (Supreme Court of India) <https://meity.gov.in/writereaddata/files/Honorable-Supreme-Court-order-dated-24th-March%202015.pdf>; Jon Russell, 'India's Government Asks ISPs to Block GitHub, Vimeo and 30 Other Websites' (TechCrunch, 1 February 2014) <https://techcrunch.com/2014/12/31/indian-government-censorsht/>.

[61] Geeta Hariharan, 'What 66A Judgment Means For Free Speech Online' (CIS) <https://cis-india.org/internet-governance/blog/huffington-post-geetha-hariharan-march-26-2015-what-66-a-judgment-means-for-free-speech-online>; Jaddie Lannon, 'IT (Amendment) Act 2008, 69A Rules: Draft and Final Version Comparison' (CIS) <https://cis-india.org/internet-governance/blog/it-amendment-act-69-a-rules-draft-and-final-version-comparison> Nikhil Pahwa, 'Why Section 69A of the IT Act should have been changed by the Supreme Court' (Medianama, 25 March 2015) <https://www.medianama.com/2015/03/223-section-69-it-act-india>.

[62] Information Technology (Intermediaries Guidelines) Rules 2011 <https://www.wipo.int/edocs/lexdocs/laws/en/in/in099en.pdf>. Rule 3(2)(i) obligates intermediaries to prohibit users from sharing information that "threatens the unity, integrity, defence, security or sovereignty of India, friendly relations with foreign states, or public order or causes incitement to the commission of any cognisable offence or prevents investigation of any offence or is insulting any other nation".

addition to their failure to tailor regulation according to the functions of intermediaries and imposing obligations violating the principles laid down in *Shreya Singhal v Union of India*.[63]

*Any legislative efforts to adhere to tackle violent extremist and terrorist speech should be in consonance with the existing legal ecosystem. The content targeted should be clearly defined following research on the drivers of terrorism and multi-stakeholder dialogue, and should inform the holistic CVE programme suggested earlier. The commitments in the Call can be realised through regulation that is nuanced and appropriate for different types of intermediaries, is narrow in scope, adheres to the principles laid down in Shreya Singhal, and is in line with constitutional limitations on the freedom of expression.*

## Encourage creation of robust redress mechanisms for actions taken against content and users

The Rules framed under the Information Technology Act and the Telegraph Act constitute a Review Committee[64] which is responsible for reviewing whether blocking of access to information[65] and suspension of telecom services[66] has been carried out in compliance with the law. However, these procedures have been criticised for lacking judicial oversight[67] as well as transparency regarding blocked content.[68] Content blocking mechanisms should be amended to subject orders issued under s 69A of the IT Act and s 5 of the Telegraph Act to judicial review in order to boost accountability. Additionally, users should be informed about the manner in which content has been blocked, the reasons for such blocking, as well as the procedure for filing appeals in a transparent fashion.

The rules under S 79 of the IT Act do not include a notice mechanism informing users regarding their content which has been taken down by platforms, nor an appeals framework for those who wish to challenge these takedowns. To mitigate the danger of over-censorship, the user creating or posting content should be allowed to contest the intermediary's decision to take down content based on another user's complaint. The counter-notice model under the Digital

---

[63] Grover, Hickok and others, 'Response to the Draft of The Information Technology [Intermediary Guidelines (Amendment) Rules] 2018' (CIS, 31 January 2019) <https://cis-india.org/internet-governance/resources/Intermediary%20Liability%20Rules%202018.pdf>.

[64] Telegraph Rules 1951, r 419A (16).

[65] Information Technology (Procedure and Safeguards for Blocking Access of Information by Public) Rules 2009, r 1<https://www.meity.gov.in/writereaddata/files/Information%20Technology%20%28%20Procedure%20and%20safeguards%20for%20blocking%20for%20access%20of%20information%20by%20public%29%20Rules%2C%202009.pdf>.

[66] Temporary Suspension of Telecom Services (Public Emergency or Public Safety) Rules 2017 drafted under Telegraph Act 1885, s 7.

[67] Committee of Experts under the Chairmanship of Justice B.N. Srikrishna, A Free and Fair Digital Economy: Protecting Privacy, Empowering Indians <https://meity.gov.in/writereaddata/files/Data_Protection_Committee_Report-comp.pdf>

[68] Pahwa (n 61).

Millennium Communications Act could be taken as inspiration- it permits users to challenge takedowns and requires intermediaries to notify users prior to any takedowns.[69]

## Encourage media outlets to apply ethical standards when depicting terrorist events online, to avoid amplifying terrorist and violent extremist content

Media coverage, and reporting on social media in particular, may amplify the impact of terrorist attacks and violent extremist content if not done responsibly.[70] Therefore, it is important to adhere to standards of ethical reporting in order to guard against the dissemination of fear, distrust and extremist propaganda through journalistic reporting. The Press Council of India has published norms of journalistic conduct requiring journalists to refrain from displaying content that may create terror, incite violence, or glorify acts of violence or the ideologies behind them.[71]

The Press Council's norms are silent on the subject of social media and the norms to be followed while reporting via social media or using content shared over these online platforms. The RTDNA Guidelines For Using User-Generated Content can be studied to devise norms for responsible reporting in this regard.[72] Industry standards, such as the National Broadcasters' Association's  guidelines for responsible coverage,[73] can be adapted to suit the challenges of social media reporting, where users also contribute content. As enforcement of such norms is difficult, a self-regulatory approach would be the most effective, and norms prescribed by the RTDNA,[74] the UNESCO,[75] and the 'Voluntary Code of Ethics for the 2019 General Election'.[76]

---

[69] 17 USC 512(g)(1).

[70] Charlie Beckett, Fanning the Flames: Reporting Terror in a Networked World (Tow Centre for Digital Journalism 2016) <https://www.cjrarchive.org/img/posts/Reporting%20on%20Terror%20in%20a%20Networked%20World%20%28Beckett%29.pdf>.

[71] Press Council of India, Norms of Journalistic Conduct (2010) <http://presscouncil.nic.in/OldWebsite/NORMS-2010.pdf>.

[72] 'RTDNA Guidelines: User Generated Content' <https://www.rtdna.org/content/user_generated_content?_ga=2.29162684.1014258970.1520199730-1562739271.1519614513>.

[73] 'NBA to form emergency protocol for coverage of crisis situations' (Exchange4media, 8 August 2019) <https://www.exchange4media.com/advertising-news/nba-to-form-emergency-protocol-for-coverage-of-crisis-situations-33391.html>.

[74] Radio Television Digital News Association, 'RTDNA Guidelines: Mass Shootings' <https://rtdna.org/content/rtdna_guidelines_mass_shootings>; RTDNA Guidelines: User Generated Content' (n 73).

[75] United Nations educational, Scientific and Cultural Organisation, Terrorism and the Media <https://unesdoc.unesco.org/ark:/48223/pf0000247074/PDF/247074eng.pdf.multi>.

[76] Social Media Platforms present "Voluntary Code of Ethics for the 2019 General Election" to Election Commission of India' (ECI, 20 March 2019) <https://eci.gov.in/files/file/9467-social-media-platforms-present-voluntary-code-of-ethics-for-the-2019-general-election-to-election-commission-of-india/>.

**Adapt measures taken in furtherance of the Call to complement existing diplomatic and counter-terrorism initiatives**

India has historically been a part of multiple global initiatives to counter extremism. It is a founding member of the Global Counterterrorism Forum, an international forum to promote a strategic, long-term approach to counter terrorism and violent extremism. It was also a member of the G20 Osaka Leaders' Statement on Preventing Exploitation of the Internet for Terrorist and Violent Extremism Conducive to Terrorism in June 2019. This statement references the Christchurch shootings and echoes certain aspects of the Christchurch Call such as development of technology to automatically identify and remove content, and responsibility of companies to enforce their ToS on their platform.[77] Additionally, it has partnered with the US to exchange terrorist screening information[78] and to cooperate on bilateral counter terrorism efforts.[79] India has also participated in the Department of State's Antiterrorism Assistance program and received training in internet, dark web, mobile devices, and advanced digital forensics capabilities.[80] India is also an active participant in the Aqaba process.

*India's foreign policy and CVE strategy should include participating in global fora which brings together multiple stakeholders to address online extremism, while reconciling the commitments under these initiatives with those under the Call.*

# Summary

### A. Recommendations for the Call

**I.** The Call needs to clarify its interaction with existing CVE initiatives like the EU Internet Forum, the G20 and G7, the Global Internet Forum to Counter Terrorism, the Global Counterterrorism Forum, Tech Against Terrotism, and the Aqaba Process.

**II.** The purpose and objective of the call needs to be clearly defined

**III.** The basis for accepting companies and governments as members of the call needs to be clearly defined. This can mitigate the danger of weaponisation of the Call.

---

[77] G20 Osaka Leaders' Statement On Preventing Exploitation Of The Internet For Terrorism And Violent Extremism Conducive To Terrorism (Vect) <http://www.g20.utoronto.ca/2019/FINAL_G20_Statement_on_Preventing_Terrorist_and_VECT.pdf>.

[78] PTI, 'India-US sign key pact on sharing info on terror' (Economic Times, 12 July 2018) <https://economictimes.indiatimes.com/news/defence/india-us-sign-key-pact-on-sharing-info-on-terror/articleshow/52555612.cms?from=mdr>.

[79] For instance, see Ministry of External Affairs, 'Joint Statement on the Inaugural India-U.S 2+2 Ministerial Dialogue' <https://mea.gov.in/bilateral-documents.htm?dtl/30358/Joint_Statement_on_the_Inaugural_IndiaUS_2432_Ministerial_Dialogue>

[80] United States Department of State, 'Country Reports on Terrorism' <https://www.state.gov/country-reports-on-terrorism/>.

**IV.** The metrics of success of the Call need to be defined clearly based on the expectations from, and capabilities of different stakeholders.

**V.** The Call should encourage the classification of online service providers based on the nature of services provided and content shared, user base and methods of content moderation

**VI.** The Call should seek commitments from governments to provide clear national definitions for "extremist content" which are in conformity with international human rights standards.

**VII.** Companies should commit to develop tools that increase transparency related to algorithms and content moderation. Users should be permitted to opt-out of content recommendation algorithms.

**VIII.** The commitment to integrate re-direct mechanisms should be re-evaluated in light of the ethical challenges arising therefrom. This should be accompanied by testing of the actual efficacy of these methods.

**IX.** The mandatory deployment of automated solutions should be deferred until fundamental flaws therein are cured. Oversight, transparency and accountability mechanisms should be in place to prevent free speech violations before applying these solutions.

## B. Recommendations for India

**I.** Devise a holistic CVE strategy tackling the political and socio-economic causes of extremism.

**II.** Civic education programmes teaching individuals to evaluate information critically and engage in basic fact-checking for curbing fake news, right from the school level, can be implemented.

**III.** In line with our recommendations, the draft Intermediary Guidelines Rules 2018 ought to regulate intermediaries based on their function and nature of content created/shared, instead of adopting a one-size-fits-all approach.

**IV.** Create a clear national definition for "extremist content" which is in conformity with international human rights standards. Particular attention ought to be paid to the implications of this definition for free speech.

**V.** Include a robust notice and appeals mechanisms to permit users to challenge incorrect and unlawful content takedowns. This should be operationalised by amending s 79 of the IT Act as well as the Blocking Rules under s 69A to create an appeals mechanism.

**VI.** Reform the legal framework for blocking online content and access to communication systems to introduce judicial oversight and transparency in decisions. This would require amending the Telegraph Rules, the Information Technology (Procedure and Safeguards for Blocking Access of Information by Public) Rules, and the Temporary Suspension of Telecom Services (Public Emergency or Public Safety) Rules.

**VII.** Adopt a self-regulatory approach to online coverage of terrorist and violent extremist content, and encourage industry associations to develop norms for responsible reporting to prevent amplification of violent ideologies. Global best practices and pre-existing norms like those created by the RTDNA, the UNESCO as well as the 'Voluntary Code of Ethics for the 2019 General Election' can serve as inspiration in this regard.

**VIII.** India's foreign policy should include participating in global fora which bring together multiple stakeholders to address online extremism, while reconciling the commitments under existing initiatives with those under the Call.

### C. Areas of Collaboration between Governments and Online Service Providers

**I.** Development of educational modules empowering individuals to engage meaningfully on the internet.

**II.** Multi-stakeholder consultations on initiatives undertaken under the Christchurch Call can be arranged to ensure that diverse perspectives and differential resource capacities are respected.

**III.** Governments can collectively create soft law and/or oversight mechanisms to prevent the weaponisation of the Call and related domestic CVE legislations for restricting free speech, political dissent and human rights reporting.

**IV.** Standards for algorithmic suppression of content should be developed through collaboration between private companies, the government, the civil society and the public. Review mechanisms for these standards can also be developed to ensure relevance, redress and increase accountability. This includes the development of oversight mechanisms for automated takedowns that allow for redress and access to training data in a transparent fashion. Questionable content that is heavily context-dependent can be referred to human rights and counter-extremism experts prior to takedown.

### D. Future Research Areas

**I.** The method of operation and design of intermediaries would have to be researched to understand their potential for misuse as a tool for disseminating extremist content. This research can inform the classification of these intermediaries for the purpose of developing CVE frameworks suited to the challenges created by the nature of the platform. Such classification could also aid the development of technological solutions or development of standards for tackling terrorist content.

**II.** Studying the interrelation between advertising-driven business models and the dissemination of extremist content towards informing the development of business models that respect the freedom of expression, privacy and autonomy of users, and preventing the sharing of extremist content.

**III.** Research is needed to understand the efficacy of different re-direct mechanisms.

**IV.**    Development of frameworks for evaluating ethical questions surrounding re-direct methods, and de-radicalization methods,

**V.**    Further research into the development of technical solutions for removing content that are built on robust training data with the ability to precisely identify and target violent extremist content.

**VI.**    Further research is needed to develop effective oversight and accountability mechanisms for algorithmic content moderation takes into consideration the interests of industry, users, and governments.

**VII.**    Further research is needed to understand the impact of different forms of content – i.e. is live streaming more harmful than video than images?

# Conclusion

The Christchurch Call is an important step in fostering cooperation between governments, tech companies, and the civil society to address extremist speech. Some key ideas espoused by the Call such as greater transparency from technology companies and the commitment to fundamental human rights are commendable. However, in the face of mounting pressure to respond forcefully to horrific incidents such as the Christchurch attacks, governments have a tendency to resort to strategies such as overbroad censorship that yield immediately tangible results.[81] The Call, by nature of its vagueness, has language which lends itself to such impulsive policymaking, such as advocating for more automated content filtering and "appropriate action to prevent [dissemination of extreme content]". Therefore, it is imperative that any policies or regulations arising from the Call should refrain from adopting reactionary, knee-jerk solutions which are not evidence-based.

Although the scope of the Call is limited to violent extremist and terrorist content, it is important to recognise that such content does not exist in a vacuum. As highlighted above, it has been established that CVE efforts are most effective when designed to comprehensively address other malicious content such as hate speech and disinformation as well. As a signatory of the Call, India should use this bolstered political will as an opportunity to devise a comprehensive CVE strategy building on existing regulations and community-based efforts. This momentum can be channeled into conducting multi-stakeholder discussions exploring CVE strategies that respect and uphold the domestic as well as international human rights framework.

---

[81] York and Cope, 'In Wake of Charlie Hebdo Attack, Let's Not Sacrifice Even More Rights' (EFF, 8 January 2015) <https://www.eff.org/deeplinks/2015/01/wake-charlie-hebdo-attack-lets-not-sacrifice-even-more-rights>.