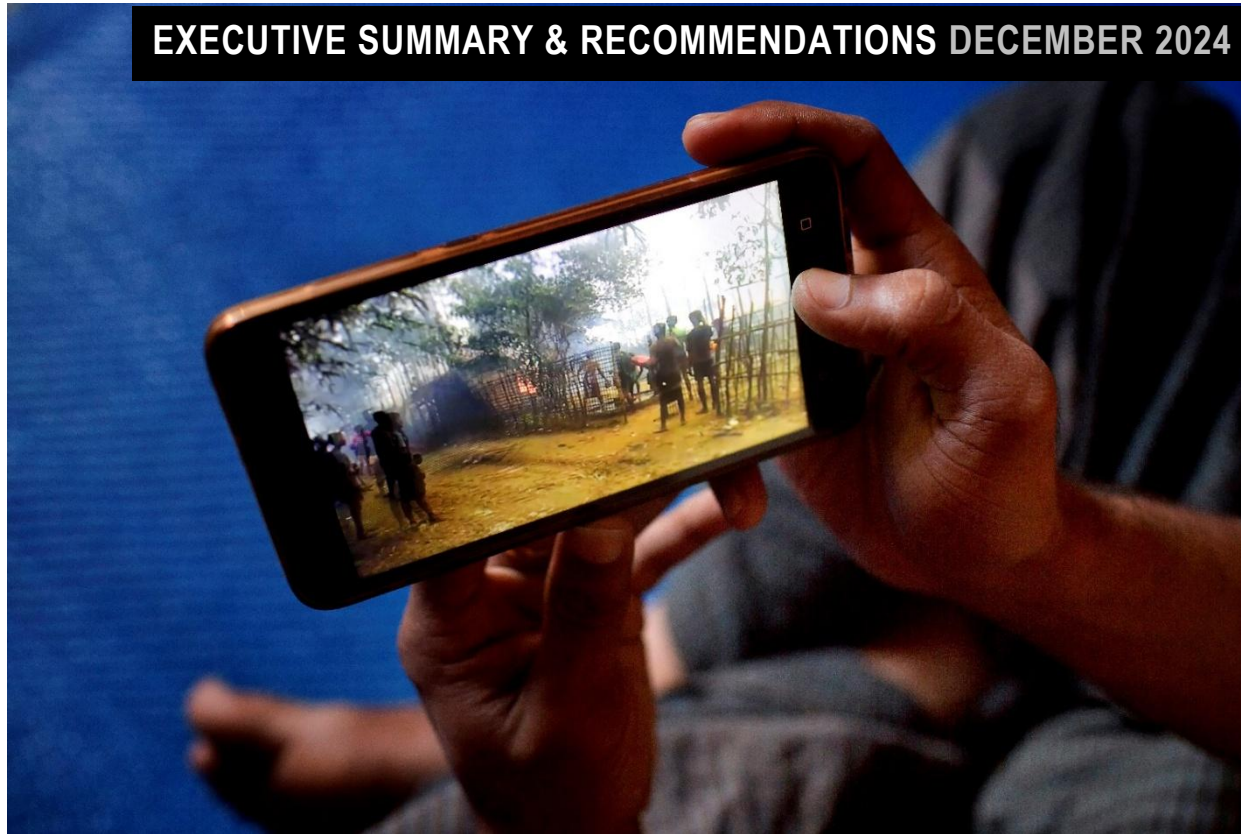# LEVERAGING SOCIAL MEDIA FOR GENOCIDE AND MASS ATROCITY PREVENTION

## Understanding the Digital Toolbox



**EXECUTIVE SUMMARY & RECOMMENDATIONS** DECEMBER 2024

UNITED STATES HOLOCAUST MEMORIAL MUSEUM

SIMON-SKJODT CENTER
FOR THE PREVENTION OF GENOCIDE

# EXECUTIVE SUMMARY

Much of the existing literature discussing social media focuses on how it might fuel or incite mass atrocities, drawing from experiences in contexts such as Sri Lanka and Burma. But there is significantly less awareness of how tools deployed or developed by social media companies might *reduce* the risk of mass violence and contribute constructively to atrocity prevention efforts.

This report aims to address this gap by focusing on how social media tools can support two core atrocity prevention strategies:
    (1) protecting vulnerable civilian populations at risk of mass atrocities, and
    (2) degrading potential perpetrators' capacity to commit mass atrocities.

It provides a landscape assessment of the suite of social media product, policy, and operational interventions that may offer potential to support these strategies and articulates some of the associated limitations, risks, and important considerations when these tools are deployed.

This report is primarily aimed at those inside social media companies with authority to develop or deploy tools in moments of heightened atrocity risk (which may include trust and safety professionals, human rights or crisis response teams, and senior leadership), as well as atrocity prevention experts and policy makers who may be able to encourage or incentivize the use of digital tools to support atrocity prevention. Select tools may also be of interest or use for humanitarian and civil society advocacy organizations that operate in atrocity risk settings.

The objective of this report is to fill a gap by expanding the understanding of both policy makers and social media platform representatives about the available tools in the digital realm to support atrocity prevention efforts, to stimulate future research in this space, and to broaden our collective imaginations in designing modern atrocity prevention policy strategies that leverage digital tools and opportunities.

This report is based on a series of semi-structured expert consultations, held under the Chatham House Rule of non-attribution, with more than 30 current and former representatives of social media companies, academics and practitioners specialized in technology and atrocity prevention, and members of at-risk communities who lent their experiences and insights to support this project.

The report concludes that expanding the atrocity prevention toolbox to include digital tools and interventions offers an opportunity to develop more modern atrocity prevention strategies to meet the challenges of the moment.

It identifies the following categories of interventions as offering potential to support civilian protection:

- **Protecting online privacy**: tools or interventions aimed at restricting the visibility of digital content that may put civilians at risk in atrocity risk settings
- **Securing social media accounts**: interventions aimed at protecting social media users against hacking, impersonation, and account takeover efforts
- **Surfacing crisis resources and credible information**: interventions aimed at connecting social media users to crisis resources and/or amplifying credible information

- **Disseminating early-warning information**: interventions that make use of social media to communicate warnings about atrocity risks
- **Enhancing communication and coordination capabilities**: interventions that enhance civilians' ability to communicate and coordinate in atrocity risk settings

This report also identifies the following categories of interventions as offering potential to degrade the capacity of atrocity perpetrators:

- **Preventing perpetrators from gaining a foothold of platforms at scale**: interventions aimed at preventing perpetrators from setting up a large presence on social media platforms
- **Disrupting perpetrators from organizing and coordinating**: interventions aimed at disrupting perpetrators from using social media to coordinate and organize the commission of violence
- **Limiting the presence or visibility of dangerous content in atrocity risk settings**: interventions aimed at reducing the presence or visibility of potentially inflammatory digital content during periods of heightened atrocity risk
- **Contextualizing perpetrator content**: interventions aimed at providing additional information or context around inflammatory digital content
- **Preventing perpetrators from mobilizing bystanders**: interventions aimed at reducing the incentives for bystanders or third-party enablers to inadvertently contribute to narratives and ideologies being advanced by perpetrators
- **Implementing last resort or "break glass" measures**: interventions that temporarily and intentionally degrade or disable social media features in moments of heightened atrocity risks

For each of the preceding categories, this report sets out specific considerations and preliminary recommendations on how they might be developed and implemented. It also sets out the following as general recommendations to platforms seeking to constructively contribute to atrocity prevention efforts:

- Platforms should invest in building internal atrocity prevention capacity and expertise. They should ensure they have a dedicated crisis response function that can define and categorize potential atrocity risk situations according to a principled risk assessment process and should develop clear protocols on when various interventions and policies will be deployed.
- Platforms should invest in research and development on social media tools that hold potential to help prevent mass atrocities. The inventory of tools in this report offers a starting point for both deepening understanding of when and how different tools can address mass atrocity risks and expanding the range of available tools.
- Platforms should invest heavily in local partnerships that can support awareness of atrocity risk dynamics. These relationships should be established well in advance of moments of crisis, and platforms should explore providing training on relevant product and policy interventions so they can be rolled out more effectively in at-risk communities.
- Platforms should build their awareness on how their products are being used in atrocity risk settings to create a baseline for further assessment of risks and opportunities.
- Platforms should localize all resources to ensure accessibility and ease of use for affected communities. Any tools or interventions developed for use by individuals in at-risk communities must be made available in the relevant local languages of affected populations.
- Platforms should hold tabletop or scenario-based simulations to prepare for atrocity risk settings.
- Platforms should preserve digital evidence of mass atrocities and, where appropriate, share information to assist in the investigation and prosecution of atrocity crimes. They should also clarify their policies on data preservation in atrocity risk and conflict settings, and consult with

civil society organizations (and, where feasible, affected communities) to identify content relevant to international justice and accountability efforts.

Finally, this report sets out recommendations to policy makers, urging them to assess both risks and opportunities to leverage the digital environment to address the risks of mass violence and to explore opportunities to incorporate social media tools and interventions into atrocity prevention strategies.

# SUMMARY OF TOOLS AND RECOMMENDATIONS

## A. Recommendations for Platforms

As discussed, this report sets out the landscape of social media tools and interventions that may be able to support either (a) protecting vulnerable civilian populations or (b) degrading perpetrator capacity. Because many of these interventions are within the control of platforms, most of the resulting recommendations are directed at social media companies.

### 1. Preliminary Recommendations: Interventions to Support Civilian Protection

First, preliminary recommendations on specific tools and interventions that may be able to contribute to the protection of vulnerable civilian populations are as follows:

**Atrocity Prevention Strategy: Protect Vulnerable Civilian Populations**

**Key Assumptions**

Social media can enable vulnerable civilian populations to access critical information and coordinate actions to protect themselves in moments of crisis.

At the same time, information available on social media can place civilians at greater risk of physical attack.

Social media can enable communication between members of affected communities about emerging atrocity risks, and from affected communities to policy makers.

**Mechanisms**

Safeguarding sensitive information about vulnerable civilian populations (for example, by *locking profiles or increasing account security measures to prevent hacking*)

Coordinating and facilitating self- or external-protection efforts (for example, by *users communicating warnings on unsafe locations or circulating information on humanitarian aid access points*)

Supporting access to essential information that could be used for protection

## TOOLS & INTERVENTIONS TO PROTECT VULNERABLE CIVILIAN POPULATIONS

| **Protect Online Privacy**<br><br>*Tools or interventions aimed at restricting the visibility of digital content that may put civilians at risk in atrocity risk settings* | **THEORY OF CHANGE:**<br><br>If digital content could be used to target civilians, restricting the visibility of that content can contribute to civilian protection. | **EXAMPLES:**<br><br>• Facebook's locked profile feature, which limits the ability to view various elements of a person's social media account, or similar interventions to limit the ability to view a user's affiliations or friends lists<br>• Obscuring users' previously shared location information<br>• Reviewing features to which users may be added without their consent that could make them more readily visible to perpetrators<br>• Creating channels for users' social media accounts to be secured or locked down in case of detention or arrest<br>• Proactively sharing instructions on the deletion or deactivation of social media accounts | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br><br>• Platforms should explore interventions to proactively restrict the visibility of digital information that could be used to target civilians in atrocity risk settings, such as their affiliations or location history.<br>• Privacy interventions aimed at protecting civilians should be carefully balanced against their potential interests in sharing information in atrocity risk settings. Wherever feasible, civilians should be afforded agency over their digital presence.<br>• Platforms should carefully review features through which civilians' digital information may be visible without their consent, or where they may not realize they gave prior consent.<br>• Platforms should ensure that vulnerable civilian populations can readily understand how to temporarily deactivate or delete their social media accounts should they deem it necessary for their protection.<br>• Platforms should communicate available privacy tools to vulnerable populations in advance of crises, and should clearly articulate relevant limitations to avoid overpromising to people who are at risk. |

| **Secure Social Media Accounts**<br><br>*Interventions aimed at protecting social media users against hacking, impersonation, and account takeover efforts* | **THEORY OF CHANGE:**<br>Civilian protection includes ensuring that civilians' digital information cannot be obtained and used against them through hacking and impersonation campaigns. This can in turn protect others who may be misled by hacked and impersonated accounts. | **EXAMPLES:**<br>• Account security push notifications, deployed in Ukraine<br>• End-to-end encryption channels | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br>• Platforms should ensure they put in place and stringently enforce policies prohibiting impersonation in atrocity risk settings.<br>• Platforms should explore opportunities, such as through push notifications or prompts, to proactively communicate information to civilians about how to best secure their online accounts.<br>• Platforms should clearly communicate their choices around the use of encrypted or unencrypted features, so that users readily understand the security of the tools they use in atrocity risk settings. |
|---|---|---|---|
| **Surface Crisis Resources and Credible Information**<br><br>*Tools or interventions aimed at connecting social media users to crisis resources, amplifying credible information, or both* | **THEORY OF CHANGE:**<br>Ensuring that civilians can access information about crisis resources can contribute to protection by helping them avoid or withstand attacks.<br>OR<br>Ensuring that civilians can access reliable information about evolving developments can prevent misinformation and disinformation from inciting violence. | **EXAMPLES:**<br>• Creating centralized landing pages or information hubs that compile authoritative information in atrocity risk settings<br>• Modifying approaches to ranking and amplification of information to align with needs in atrocity risk settings<br>• Amplifying content from credible accounts, such as reliable media or civil society organizations<br>• Providing ad credits to credible local organizations | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br>• In atrocity risk settings, the way information is ranked and prioritized takes on heightened importance. Platforms should review approaches to the ranking and amplification of information to align with needs in atrocity risk settings.<br>• Platforms should develop principled approaches for how they will identify credible and useful information for civilian protection, and under what circumstances specific information will be amplified.<br>• Platforms should consider affording users choice in how content is prioritized in user feeds so they can quickly identify resources and information important to their protection. |

| | | • Using push or pop-up notifications, or "nudges," to direct people to important resources or news items | • Platforms should, in partnership with relevant organizations and humanitarian agencies, explore opportunities to direct users to credible information or resources that could support their protection.<br>• Platforms should explore opportunities to afford vulnerable communities greater control and agency in efforts to surface crisis resources on social media and mitigate risks associated with information sharing. |
|---|---|---|---|
| **Disseminate Early-Warning Information**<br><br>*Interventions that make use of social media to communicate warnings about atrocity risks* | **THEORY OF CHANGE:**<br>Social media may be used to communicate warnings (either to civilians at risk or to policy makers), with a view to influencing outcomes on civilian protection. | **EXAMPLES:**<br>• Publishing emergency air raid alerts on social media<br>• Using social media posts to warn people about safe/unsafe locations in Libya, or potential targets for air strikes in Syria | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br>• In light of the use of social media for early warning, platforms operating in atrocity risk settings should review the way that content moderation policies on graphic media or violent content are applied and enforced, with consideration to the needs of affected populations to understand emerging events and risks of violence. Platforms should also explore the use of technical interventions to mitigate psychosocial harm associated with the viewing of graphic content, such as the use of interstitials, grayscale, or image blurring.<br>• Platforms should explore opportunities to support early-warning initiatives by trusted third-party entities, but they also should implement safeguards to carefully assess information credibility and timeliness. |

| **Enhance Communication and Coordination Capabilities**<br><br>*Interventions that expand or enhance civilians' ability to communicate and coordinate* | **THEORY OF CHANGE:**<br><br>Supporting open communication and coordination between civilians will enable them to better avoid or withstand atrocities. | **EXAMPLES:**<br><br>• "Groups" or "Communities" features on social media<br>• Social media features that enable group messaging<br>• Features that help users connect to social media platforms via proxy servers, bypassing restrictions on internet access | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br><br>• Platforms should be mindful of the value of features that enable group discussion for coordination and communication between civilians in atrocity risk settings. When considering modifying or updating these features, platforms should take particular care in assessing the needs of those in atrocity risk settings who may be using them for this purpose.<br>• Platforms should pay particular attention to the moderation of group discussion forums in settings where there is a heightened risk of mass atrocities. This may include ensuring that the administrators of discussion forums have the tools and resources they need to support moderation, such as the ability to approve or remove members from digital spaces.<br>• Platforms should ensure that the visibility or privacy of group discussion forums is clearly communicated to all participants.<br>• In atrocity risk settings where restrictions on communication are heightened, such as through internet blackouts or platform shutdowns, platforms should consider exploring opportunities to expand access to social media, particularly for vulnerable or isolated communities. |
|---|---|---|---|

## 2. Preliminary Recommendations: Interventions to Degrade Perpetrator Capacity

Preliminary recommendations on specific tools and interventions that may be able to contribute to degrading the capacity of atrocity perpetrators are as follows:

**Atrocity Prevention Strategy: Degrade Potential Perpetrators' Capacity to Commit Atrocities**

**Key Assumptions**

- Mass atrocities depend on perpetrators having certain material and operational capacities. In many countries at risk of mass atrocities today, perpetrators may use social media as a resource for facilitating systematic attacks.
- Social media can enable potential perpetrators to communicate rapidly and persuasively with large audiences in ways that may contribute to atrocity risk, by inciting violence, spreading exclusionary ideologies, or disseminating disinformation or misinformation about a particular group.
- Social media can also play a role in the planning and organization of attacks, such as by providing forums for recruitment or weapons sales.
- Tools that make social media platforms less effective or efficient means of advancing perpetrators' goals can therefore contribute to degrading their overall capacity to commit atrocities.

**Mechanisms**

- Decreasing the speed and audience-reach efficiency of social media features for potential perpetrators (for example, via *content moderation policies on crisis misinformation, rate limits, or nudges suggesting users think twice before re-sharing certain content*)
- Decreasing the persuasiveness of inciting, misleading, or otherwise dangerous content (for example, via contextualizing content or labeling the source of posts, such as state-affiliated media)
- Disrupting digital spaces in which perpetrators are organizing or planning the commission of atrocities (*for example, weapons sales and recruitment*)
- Denying potential perpetrators access to social media platforms entirely or to specific social media features or platforms (for example, via *detection and removal of coordinated networks of accounts of potential perpetrators, deplatforming violent organizations, or disabling social media features in moments of heightened atrocity risk*)

| TOOLS & INTERVENTIONS TO DEGRADE POTENTIAL PERPETRATORS' CAPACITY TO COMMIT ATROCITIES | | | |
|---|---|---|---|
| **Prevent Perpetrators Gaining Foothold on Platforms at Scale**<br><br>*Interventions aimed at preventing perpetrators from setting up a large presence on social media platforms* | **THEORY OF CHANGE:**<br><br>Preventing perpetrators from establishing or maintaining extensive networks of accounts will make them less able to weaponize social media in furtherance of atrocities (such as to incite or coordinate violence). | **EXAMPLES:**<br><br>• Preventing perpetrators from registering social media accounts<br>• Expanding detection of coordinated networks of accounts of potential perpetrators<br>• Designating and banning perpetrators under policies governing violent individuals and organizations<br>• Deplatforming perpetrators, or subjecting them to heightened monitoring against content moderation policies | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br><br>• Platforms should explore interventions that can prevent atrocity perpetrators from setting up networks of inauthentic accounts. This should include both efforts to prevent the registration of new accounts, and review of older accounts that may exhibit suspicious behavior.<br>• Platforms should enhance their in-house investigative capacities to detect and remove coordinated networks of inauthentic accounts that may be used by atrocity perpetrators.<br>• Platforms should, in atrocity risk settings, proactively review potential perpetrators against criteria for designation under violent organizations policies.<br>• Platforms should explore heightened monitoring of accounts of atrocity perpetrators and more stringent enforcement of content moderation policies given these individuals' offline behavior. |

| **Disrupt Perpetrators from Coordinating and Organizing on Social Media**<br><br>*Interventions aimed at disrupting perpetrators from using social media to coordinate and organize the commission of violence* | **THEORY OF CHANGE:**<br><br>To the extent that digital spaces are being used to coordinate and organize violence, disrupting perpetrators' ability to use social media to advance the planning and organization of violence will degrade their overall capacity to commit atrocities. | **EXAMPLES:**<br><br>• Enforcement of content moderation policies prohibiting weapons sales or to promote criminal activities<br>• Heightened monitoring of online spaces where perpetrators may be organizing violent activities, such as groups or pages<br>• Policies prohibiting the use of social media for surveillance | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br><br>• Platforms should ensure that they have policies in place that prohibit the abuse of their platforms for the coordination and organization of mass violence, including but not limited to the purchase and sale of weapons and recruitment to violent organizations.<br>• Platforms should also ensure that the enforcement of these policies is sufficiently resourced in atrocity risk settings, particularly in online spaces where perpetrators may be gathering.<br>• Platforms should explore interventions to prevent perpetrators from readily collecting information on social media about potential targets, and ensure policies are in place prohibiting the use of social media data for surveillance. |

| Limit the Presence or Visibility of Dangerous Content in Atrocity Risk Settings | THEORY OF CHANGE: | EXAMPLES: | PRELIMINARY GUIDANCE FOR PLATFORMS: |
|---|---|---|---|
| *Interventions aimed at reducing the presence or visibility of potentially inflammatory digital content during periods of heightened atrocity risk* | Reducing the presence, audience reach, or visibility of potentially inflammatory digital content, limits potential perpetrators' ability to use social media to incite atrocities or further societal divisions. | • Having policies governing how platforms manage dangerous misinformation in crisis settings, such as limiting it from appearing on users' home feed or timeline or limiting its ability to be re-shared<br>• Deamplifying content that could create a serious risk of harm, such as potentially dehumanizing language or exclusionary ideologies<br>• Implementing rate limits or forwarding limits that reduce the number of people a user can forward content to at scale | • Platforms should ensure they have policies in place to manage dangerous misinformation in atrocity risk settings, perhaps by limiting users' ability to share, recommend, or amplify unverified and potentially harmful information.<br>• Platforms should also ensure that they have policies in place prohibiting the incitement of violence and that these policies are rigorously enforced in atrocity risk settings. These policies should also be developed and enforced with an understanding of behaviors and patterns around the commission of mass violence, such as the use of dehumanization, hate speech, and coded language or "dog-whistling" to incite violence.<br>• Where dangerous misinformation remains online in atrocity risk settings, platforms should explore the use of "soft interventions" to reduce the risk of misinformation contributing to violence, such as placing warning labels over the content.<br>• Platforms should engage in further research on the benefits, risks, and unintended consequences of deamplifying dangerous content (such as dehumanizing language or derogatory terms) in atrocity risk settings, but they should be transparent about their approach.<br>• In atrocity risk settings, platforms should explore reasonable rate limits or |

| | | | requirements that users accumulate some indicia of trustworthiness before they are permitted broad reach and engagement on the platform, to prevent perpetrators from reaching other users en masse.<br>• Platforms should explore opportunities to link indicia of trustworthiness to the ability to use features like ads in atrocity risk settings, or to prohibit the use of ads outright in certain contexts. To the extent ads are permitted, they should be rigorously scrutinized against policies prohibiting hate speech and incitement to violence. |
|---|---|---|---|
| **Contextualize Perpetrator Content**<br><br>*Interventions aimed at providing additional information or context around inflammatory digital content posted on social media by potential perpetrators, where it is not removed outright* | **THEORY OF CHANGE:**<br><br>Situating inflammatory digital content in the context of credible, factual information can reduce perpetrators' ability to spread and persuade people of dangerous rumors or incite violence. | **EXAMPLES:**<br>• Placing warning labels or interstitials over potentially inflammatory digital content, sharing further context about what is depicted or asserted<br>• Verifying and labeling accounts belonging to certain types of users, such as government officials, electoral candidates, or state-affiliated media<br>• Providing further context on or labeling the provenance of misleading media<br>• "Prebunking" or inoculating users against dangerous misinformation | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br>• In atrocity risk settings, platforms should explore labeling and verifying certain categories of accounts, such as those belonging to government officials or electoral candidates, to prevent users from being persuaded by impersonation attempts.<br>• Platforms should explore the use of interstitials, paired with deamplification, for a small subset of high-risk, high-visibility content in atrocity risk contexts. They should also explore options to communicate the provenance of misleading media, so users better understand the source of content they encounter.<br>• In partnership with local organizations, platforms should explore the use of prebunking to reduce the potency of mis/disinformation in atrocity risk contexts, |

| | | | |
|---|---|---|---|
| | | | and support independent research on the efficacy of these efforts.<br>• Where local partnerships are absent, platforms should explore the possibility of user-led or community-to-community interventions that would enable users to flag misinformation themselves. |
| **Prevent Perpetrators from Mobilizing Bystanders**<br><br>*Interventions aimed at reducing the incentives for bystanders or third-party enablers to inadvertently contribute to narratives and ideologies being advanced by perpetrators* | **THEORY OF CHANGE:**<br><br>By reducing the likelihood that third-party enablers contribute to the dissemination of dangerous narratives and ideologies advanced by perpetrators, reduce perpetrators' ability to weaponize social media to incite or fuel atrocities. | **EXAMPLES:**<br>• "Nudges" suggesting users think twice before re-sharing certain content on social media<br>• Prompts warning users if they are about to share a potentially harmful or hurtful reply or comment<br>• Interventions to interrupt the user interface to make it more difficult to rapidly re-share content that may contribute to violence | **PRELIMINARY GUIDANCE FOR PLATFORMS:**<br>• Platforms should, in atrocity risk settings, explore the use of "nudges" to encourage critical thinking, and they should make it more difficult for bystanders to rapidly re-share information that could contribute to violence.<br>• In settings where atrocities have already begun, platforms may want to consider suspending interventions that add friction to users' ability to rapidly share content that may be necessary for their protection. |

| Last Resort or "Break Glass" Measures

Interventions that temporarily and intentionally disable or degrade social media features in moments of heightened atrocity risk | THEORY OF CHANGE:

Where social media features are at risk of being abused to contribute to atrocities, disabling features reduces the tools available to perpetrators. | EXAMPLES:

• Intentionally disabling features that allow users to share hashtags, to avoid inciting violence in Ethiopia
• Intentionally slowing down or degrading the functionality of certain features (i.e., adding friction) to prevent content from rapidly circulating on social media | PRELIMINARY GUIDANCE FOR PLATFORMS:

• In light of the gravity and irremediability of mass atrocities, platforms should keep on the table interventions that would temporarily degrade or disable platform features at risk of severe abuse by atrocity perpetrators.
• At the same time, because of the dual-use nature of most social media features, these measures should typically be used as a last resort or "break glass" measure, deployed only after assessing relevant limitations and trade-offs. |
|---|---|---|---|

[**The full version of this report** contains an additional category of recommendations, "General Recommendations to Platforms," that has been omitted from this document, as it is already outlined in the executive summary. For the complete Recommendations section, please see pp.46-62 in the full report.]

## B. Recommendations for Policy Makers

Most of the recommendations set out in this report are aimed at platforms as the primary actors in conceptualizing, developing, and deploying the types of features and interventions described herein. This report, however, is not solely aimed at platforms, but also at atrocity prevention policy makers responsible for developing strategies that could better integrate digital tools and interventions. First and foremost, policy makers should ensure that atrocity prevention strategies include an assessment of both risks and opportunities in the digital environment, taking into account how both at-risk communities and perpetrators are using social media. Further, policy makers should consider taking the following actions:

• Partner with social media platforms to research the benefits and risks of specific interventions in atrocity risk settings;
• Establish dedicated channels for communication between the atrocity prevention community and social media companies;
• Engage in greater information sharing with social media companies on settings where there is a heightened risk of mass atrocities, with the aim of raising awareness of the need for digital interventions;
• Explore opportunities to share atrocity prevention expertise with platforms, to support the development and deployment of interventions focused on prevention; and
• Explore opportunities to incorporate social media tools and interventions into atrocity prevention strategies.

A nonpartisan federal, educational institution, the **UNITED STATES HOLOCAUST MEMORIAL MUSEUM** is America's national memorial to the victims of the Holocaust, dedicated to ensuring the permanence of Holocaust memory, understanding, and relevance. Through the power of Holocaust history, the Museum challenges leaders and individuals worldwide to think critically about their role in society and to confront antisemitism and other forms of hate, prevent genocide, and promote human dignity.

ushmm.org/connect

**UNITED STATES HOLOCAUST MEMORIAL MUSEUM**

**SIMON-SKJODT CENTER FOR THE PREVENTION OF GENOCIDE**

**100 Raoul Wallenberg Place, SW  Washington, DC 20024-2126**  ushmm.org