

## Covid-19 emergency demands platforms to

- **Publicly commit to halt exclusively AI driven content moderation after the sanitary crisis.** Algorithms can be an aid in content moderation, but must not make any decisions on the removal of content as they are not able to assess compliance with standards on freedom of expression and the context of content, and are therefore prone to misidentify legal content.
- **Establish mechanisms to notify illegal contents, and increase the visibility of such mechanisms.** Notification mechanisms must be transparent, user friendly and easily understandable.
- **Strengthen mechanisms for appeal against content removal decisions.** These mechanisms must as well be transparent, user friendly and easily understandable
- **Reporting mechanisms and appeal mechanisms should not be lengthened to discourage users from using them.**
- **Publish a post-covid-19 transparency report.** This report should include data on moderation operations carried out at the request of governments, users or on their own initiative.

## RSF recommendations to platforms

- **Platforms must comply strictly with their duty of care,** and the law should strengthen their obligations in this regard, in order to make sure they do evaluate how their activities and services affect the rights of their users, and take actions to mitigate this impact. The results of this assessment must be made public.
- **Content moderation operations,** whether carried out by technological and/or human means, **must comply with international human rights standards.** They must not be allowed to restrict freedom of expression in a way that is excessive as regards permissible limitations provided for by article 19 of the ICCPR.

## Moderation and freedom of information

- **Ensure a balance between protecting users from hateful content and respecting their freedom of expression.** Platforms must be under an obligation to do their best efforts to remove illegal content, taking into consideration their users' right to freedom of expression.
- **Commit to ensuring that in each moderation operation of content notified as illegal, a human is involved in the moderation process** in order to assess the context of the content. To effectively assess the context of contents, platforms should involve members of communities most affected by hate speech, such as journalists.
- **Journalists should also be invited to contribute to reflexion and studies over hate speech** on platforms and develop the appropriate solutions.

- **Mechanisms to protect legitimate contents against notifications in bad faith**, and to sanction such notifications in bad faith must be put in place. Platforms must be careful that their rules are not misused to silence journalists.
- **Journalistic contents must benefit from a special protection**, to ensure they cannot be removed by digital services providers in application of their terms of use or to respond to a notification. Journalists and media should have the ability to seize the judge for an urgent provisional decision on the legitimacy of the removal of a content and interim measures
- **Put in place visible and easily actionable mechanisms for reporting illegal content.**
- **Put in place visible and easily actionable mechanisms to appeal content removal decisions.** The removal of a content must be appealable before the platform, and the decision of the platform over this appeal must be open to a recourse before a court of law - or before an independent body (such as a public regulator) under the control of a judge

### Transparency obligations

- **Be transparent about their content moderation rules** and specify the details of the application of these moderation rules.
- **Be transparent about moderation operations** carried out at the request of governments, users or on their own initiative. These results should be published periodically and include the percentage of requests that are acted upon, the reasons why the platform decided to act upon these requests or not, and the moderation operations related to the reporting of hateful content.
- **Increase the transparency of their actions against online harassment.**

### Combating cyber violence

- **Collaborate actively with judicial authorities in the investigation of cyber-violence against journalists** - provided the requests by such authorities comply with international standards on free speech and due process, in particular by:
  - responding to requests from the judicial authorities, in particular with regard to enabling the identification and prosecution of those responsible for illegal content;
  - removing illegal content at the request of the judicial authorities.
- **Strengthen the fight against coordinated online harassment campaigns**, including those perpetrated by bots. Bots' accounts should be clearly marked as such and platforms must provide effective mechanisms to report suspicious malicious bots.
- **Develop communication and awareness-raising campaigns about online violence** specifically targeting journalists, especially women.
- **All digital services providers must have a legal representative in all the countries they are operating**, in order for individuals to be able to sue them in their country of residence for the personal harm they may have suffered because of platforms activities.