

## The Legal Journal on Technology Public Comment: Case 2025-044-FB-UA, 2025-045-FB-UA, 2025-046-FB-UA, 2025-047-FB-UA (Content Moderation in Somaliland Cases)

**Summary:** In the "Reporting on Somaliland Current Affairs" case before Meta's Oversight Board, The Legal Journal on Technology (TJLT) argues that protecting media freedom in non-English speaking and politically sensitive regions requires systemic reform to prevent unjust censorship of journalistic content. As outlined in TJLT's response, the landscape in Somaliland is marked by the persistent suppression of independent media, despite constitutional promises of free expression. Journalists face intimidation, legal pressure, and the constant threat of digital repression, particularly on global platforms. TJLT contends that Meta's automated moderation tools, predominantly trained on Western datasets, routinely misclassify critical journalism in low-resource languages like Somali.

This not only silences vital watchdog voices but also perpetuates structural marginalization. TJLT identifies the need for procedural fairness, transparency, and culturally competent content moderation. The submission advocates for concrete safeguards such as including mandatory human review for public-interest news pages and a novel "Public-Interest Media Flag" to protect high-risk journalists in crisis zones. It recommends clear notification protocols, specialized appeal pathways, and transparency measures to restore trust when errors occur. Ultimately, TJLT asserts that rapid, tailored remedies and preventative mechanisms are essential to uphold journalistic freedom and due-process rights in digital spaces, ensuring media pluralism is preserved where it is needed most.

## ISSUE 1: MEDIA FREEDOM AND SAFETY OF JOURNALISTS IN SOMALILAND, THE ROLE OF SOCIAL MEDIA AND THE SITUATION FOR FREEDOM OF EXPRESSION.

TJLT Response: Somaliland is a self-declared republic that broke away from Somalia in 1991, which currently operates independently with de facto autonomy. Freedom of independent media and journalist safety continue to be long-awaited ideals and dreams that are yet to be realized. While Somaliland citizens have been formally promised freedom of expression and press independence in the Constitution, these protections are still only on paper.<sup>1</sup> On the ground, journalists face significant barriers in operating independently and disseminating information. Over the past decade, the Somaliland government has adopted an increasingly restrictive stance toward independent media, particularly in electronic platforms. Private radio stations have effectively been banned, based on the government's claim that open broadcasting would be disastrous in Somaliland's "highly argumentative society." This rationale, unfortunately, has been used to justify widespread suppression. One prominent example of this is Radio Horyaal, a privately-run station that was forced to relocate to Belgium to continue operating, following repeated harassment and denial of registration in Somaliland.<sup>2</sup>

Broadcasting via television does not fare well either. Although television broadcasters like Universal TV and Horn Cable TV are technically permitted to operate, they function under constant state scrutiny. They are routinely accused of bias and subjected to political pressure. The broader pattern is clear: journalists and media platforms that do not align with the narrative set by the state are met with interference,

<sup>&</sup>lt;sup>1</sup> Constitution of the Republic of Somaliland 2000, art 32.

<sup>&</sup>lt;sup>2</sup> Nicole Stremlau, 'Hostages of peace: the politics of radio liberalization in Somaliland' (*Somaliland Economic*, 13 May 2025) <a href="https://somalilandeconomic.com/hostages-of-peace-the-politics-of-radio-liberalization-in-somaliland/">https://somalilandeconomic.com/hostages-of-peace-the-politics-of-radio-liberalization-in-somaliland/</a> accessed 9 July 2025.

intimidation, or legal consequences intended to hold control over information. These restrictions are maintained through aggressive tactics. According to the Somali Journalists Syndicate, at least 18 journalists were arrested in 2019 alone, many in incidents involving excessive force.<sup>3</sup> There are documented cases of police firing on journalists, seizing equipment, and deleting footage of politically sensitive events. Such practices have created a climate of fear and self-censorship among journalists.

The judiciary further aggravates these challenges. Many judges frequently uphold convictions under vague or extremely broad charges like 'insulting public officials' or 'spreading false news.' Sentences in such cases range from several months to multiple years, further eroding confidence in legal protections for journalistic work. Social media platforms, especially Facebook, have emerged as vital tools for expression; however, they have simultaneously become dangerous spaces of repression as well. Journalists have faced imprisonment for critical posts addressing government policies, officials, and even public events.<sup>4</sup>

This digital repression goes beyond individual content. News websites have been blocked without a formal process, and digital reporters are often arrested for the material they publish online. In some cases, state authorities have relied on traditional clan structures to suppress dissent - by pressuring elders to serve as guarantors of a journalist's future conduct.<sup>5</sup> This practice by Somaliland exploits community norms to bypass legal accountability and silence critical voices.

For global platforms operating in regions like Somaliland, including Meta, it is essential to engage with local dynamics in a nuanced and informed manner. Wrongful content enforcement, especially against local journalism in politically sensitive and non-English-speaking contexts, risks reinforcing the very forms of repression it aims to avoid.

### ISSUE 2. CHALLENGES IN PREVENTING WRONGFUL ENFORCEMENT AGAINST JOURNALISTIC CONTENT IN NON-ENGLISH CONTEXTS

**TJLT Response:** The submission further argues that freedom of press and media in non-English speaking regions faces risks due to wrongful enforcement. Automated content moderation systems, which operate on Artificial Intelligence whose data sets are often western-centric, limiting their ability to understand context, language nuance and journalistic expressions from other languages and cultures.<sup>6</sup> This leads to frequent misclassification of legitimate reporting as harmful or illegal content, resulting in the takedown of journalists' posts and accounts.<sup>7</sup>

In regions where freedom of press is restricted, these moderation errors assume greater significance. In a place where a journalist is less likely to cover a politically sensitive topic due to authoritarian governments, the recourses are also limited against content moderation by Meta. Under-resourced content moderation teams, often unequipped for determining the validity of texts from other languages other than English,

<sup>6</sup> Kinga Pázmándi, 'Theory and Practice of Social Media's Content Moderation by Artificial Intelligence in Light of European Union's AI Act and Digital Services Act' (2025) *EJ Politics* <a href="https://www.ej-politics.org/index.php/politics/article/view/165">https://www.ej-politics.org/index.php/politics/article/view/165</a> accessed 15 July 2025.

<sup>&</sup>lt;sup>3</sup> Mohammed Dhaysane, 'Somalia: Attacks against journalists continued in 2019' Anadolu Agency (Ankara, 3 Jan 2020)

<sup>&</sup>lt;sup>4</sup> Reporters Without Borders, 'Three journalists held incommunicado and tortured in Somaliland' RSF (Paris, 11 January 2024) <a href="https://rsf.org/en/three-journalists-held-incommunicado-and-tortured-somaliland">https://rsf.org/en/three-journalists-held-incommunicado-and-tortured-somaliland</a> accessed 9 July 2025.

<sup>&</sup>lt;sup>5</sup> National Union of Somali Journalists, Report on Somaliland: Media Freedom Kept within Bounds (2009)

<sup>&</sup>lt;a href="https://rsf.org/sites/default/files/NUSOJ">https://rsf.org/sites/default/files/NUSOJ</a> Somaliland Report.pdf> accessed 9 July 2025.

<sup>&</sup>lt;sup>7</sup> Akshaya J and Nithya Sambasivan, 'Lost in Translation: How Content Moderation Fails Tamil Speakers Online' (2025) *arXiv preprint* <a href="https://arxiv.org/html/2501.13836v1#S1">https://arxiv.org/html/2501.13836v1#S1</a> accessed 15 July 2025.

frequently confuse between hate speech and regional journalism, silencing critical and marginal voices.<sup>8</sup> In these scenarios, journalists may often find themselves locked out of their audiences, with little recourse and lack of timely redressal even if the issue is recognised. This vulnerability is amplified by historical power imbalances and the continued marginalization of non-English languages within platform policies.

Effective solutions require investment in linguistic and cultural expertise, transparent appeal mechanisms, and multi-stakeholder collaboration to ensure moderators, both human and automated. Without corrective action, wrongful enforcement will continue to threaten media freedom and pluralism, particularly where journalism is already under siege.

# ISSUE 3. GOOD PRACTICES FOR ENSURING ACCESS TO ADEQUATE REMEDIES FOR JOURNALISTS AND MEDIA ORGANIZATIONS LOCKED OUT OF PAGES OR ACCOUNTS AS A RESULT OF WRONGFUL ENFORCEMENT.

**TJLT Response:** In the present case, Meta restored the content, acknowledging an error. What remains unresolved is the structural question: *How should platforms respond to wrongful enforcement, ensuring journalists can access effective remedies?* 

The posts reported on official events in Somalia, a politically sensitive and contested region. <sup>10</sup> In such a context, reporting on state activity often carries political or symbolic significance. However, as Meta later confirmed, contained no hate speech or incitement and were removed in error. The user stated clearly that their intention was to share public information, not to discriminate or provoke.

The implications of enforcement in such cases go beyond temporary account disruption. When a journalist loses access to their page and audience, it disrupts their ability to report on political developments, particularly in regions where traditional media is constrained. Remedy must therefore reflect not only a technical correction but the public role of journalistic speech and the disproportionate risks involved when removed without cause.

Meta's Community Standard on Hateful Conduct prohibits attacks based on protected characteristics, such as ethnicity and national origin. However, it also recognises that not all references to identity are inherently hateful. Where journalistic content engages with political or social themes, enforcement must consider the context and the intention of the speaker. This shows a disjuncture between the policy's intention and its application. Where enforcement results in wrongful unpublishing of journalistic pages, good practice should focus on procedural fairness, transparency, and institutional learning. To that end, several practices are worth adopting.

Meta's Oversight Board has identified *treating users fairly* and *refining automated enforcement* as strategic priorities.<sup>13</sup> In line with this, the first recommendation is to ensure mandatory human review of any

[3]

<sup>&</sup>lt;sup>8</sup> Amrutha Mohan and V Padmaja, 'Lost in Translation: How Content Moderation Fails Tamil Speakers Online' (*Tech Policy Press*, 10 January 2024) <a href="https://www.techpolicy.press/lost-in-translation-how-content-moderation-fails-tamil-speakers-online/">https://www.techpolicy.press/lost-in-translation-how-content-moderation-fails-tamil-speakers-online/</a> accessed 15 July 2025.

<sup>&</sup>lt;sup>9</sup> 'Reporting on Somaliland Current Affairs' (*Oversight Board*, 1 July 2025) <www.oversightboard.com/pc/reporting-on-somaliland-current-affairs/> accessed 15 July 2025.

<sup>&</sup>lt;sup>10</sup> Reporters Without Borders, 'Somalia' (RSF, 2024) <a href="https://rsf.org/en/country/somalia">https://rsf.org/en/country/somalia</a> accessed 15 July 2025.

<sup>&</sup>lt;sup>11</sup> 'Community Standards: Hateful Conduct' (*Meta Transparency Center*, May 2023) <a href="https://transparency.meta.com/engb/policies/community-standards/hateful-conduct/">https://transparency.meta.com/engb/policies/community-standards/hateful-conduct/</a> accessed 15 July 2025.

<sup>&</sup>lt;sup>12</sup> UN Human Rights Committee, 'General Comment No 34: Article 19: Freedoms of Opinion and Expression' (12 September 2011) UN Doc CCPR/C/GC/34, paras 42 and 46 <www.refworld.org/legal/general/hrc/2011/en/83764> accessed 15 July 2025. <sup>13</sup> 'Oversight Board Announces Seven Strategic Priorities' (*Oversight Board*, 13 March 2024)

<sup>&</sup>lt; www.oversightboard.com/news/543066014298093-oversight-board-announces-seven-strategic-priorities/> accessed 15 July

enforcement targeting media actors before removing pages that regularly post public-interest journalism. This can be triggered by posting patterns, regardless of verification status. Such pages should be routed into a specialised review path that includes contextual analysis by trained policy staff.

Second, the affected journalists should receive a policy-specific notice that clearly cites the Community Standard applied, identifies whether enforcement was triggered by user reports or automated systems, and outlines what kind of review (human or automated) will follow. Additionally, Meta should operationalise a "fast-track review lane" for accounts that

- (a) consistently report on political events, and
- (b) have no prior violations.

When such accounts are locked, they should automatically be routed to human moderators with training in public-interest content review.

Third, a log of reversed media-related enforcement decisions should be maintained internally and audited regularly to improve reviewer training and automated enforcement thresholds. In regions with heightened political sensitivity, Meta can utilize *local advisory input* from human rights groups and digital journalism networks to *improve classifications of political reporting*.

Finally, when a media page is unpublished and later reinstated, restorative remedy should not stop at access. Journalists should be offered the option to issue a reinstatement notice that explicitly explains the review process and affirms the content's compliance with platform standards. This feature wouldn't serve as an apology but rather as a transparency measure, helping journalists clarify to their followers that their reporting remains accurate and compliant, even after an inadvertent removal. Such clarity helps reduce reputational damage and restore trust.

Ensuring that journalists and media organizations affected by wrongful enforcement receive proper remedies, beyond simply restoring access, requires thoughtful, transparent, and respectful practices. Looking forward, strengthening remedies for journalists on social media is essential to uphold freedom of expression and press, ensuring digital spaces remain safe, credible forums for truth, accountability, and public voices.

#### RECOMMENDATIONS

Somaliland's media ecosystem is already under siege: at least forty-one Somali-language reporters have been arrested, intimidated, or assaulted since mid-March 2025, most for coverage of security or foreignpolicy dossiers. 14 When Meta mistakenly "unpublishes" a 90,000-follower news page, it extinguishes a watchdog voice for roughly 1 in 70 citizens of the six-million-strong polity. 15 The Oversight Board's own cross-check opinion underscores that current mistake-prevention privileges politicians and celebrities, while local newsrooms remain unprotected. Below are two lowest-friction, yet world-class safeguards that can be deployed inside one product cycle to curb such collateral damage.

1. Lightweight Public-Interest Media Flag (PIMF)

<sup>&</sup>lt;sup>14</sup> Committee to Protect Journalists, 'Alarming escalation: At least 41 journalists targeted since March in Somalia' (CPJ, 15 May 2025) <a href="https://cpj.org/2025/05/alarming-escalation-at-least-41-journalists-targeted-since-march-in-somalia/">https://cpj.org/2025/05/alarming-escalation-at-least-41-journalists-targeted-since-march-in-somalia/</a> accessed 15 July 2025.

<sup>&</sup>lt;sup>15</sup> M.A. Egge, 'Somaliland's population reaches 6.2 million' (*Horndiplomat*, 19 April 2024) <www.horndiplomat.com/2024/04/somalilands-population-reaches-6-2-million/> accessed 15 July 2025.

Meta already parses page metadata and runs language-ID scores for integrity workflows. A simple heuristic can elevate high-risk journalism without a line of new machine-learning code. A page is auto-labelled PIMF when it (a) self-identifies with "journalist/news/media" in its bio or domain e-mail, and (b) has  $\geq 50,000$  followers and  $\geq 60$  % of recent posts in a low-resource news language, such as Somali. Any proposed "unpublish" for a PIMF page is automatically paused until a *second* Somali-speaking reviewer confirms the violation. Because the rule is declarative, policy engineers can ship it via an Integrity rule-file update, avoiding the protracted model-training cycles that often stall Global-South fixes. The Santa Clara Principles demand upstream due-process guardrails; PIMF supplies them with almost no engineering lift. <sup>16</sup>

### 2. Twelve-Hour Soft-Stop with Shadow Appeal

When a PIMF post is suspected of breaching policy, Meta should apply a soft-stop rather than an immediate takedown. The content stays live behind a click-through interstitial and its reach is throttled by 75%, while a twelve-hour countdown notifies both editors and reviewers. Editors see a prominently placed *one-tap shadow-appeal* button that forwards the case to a priority queue without extending the deadline. If no definitive decision emerges within twelve hours, the post automatically regains full distribution and the strike is voided. This mirrors proven "hold-and-verify" patterns from financial fraud prevention, minimizing public-interest disruption while still containing potential harm, and squarely fulfils Access Now's call for accelerated human-review pathways in crisis regions.<sup>17</sup>

Together, PIMF functions as a proactive brake, while the soft-stop/appeal protocol operates as a reactive cushion. They translate the Santa Clara trilogy- *numbers, notice, appeals*- into safeguards tailored for low-resource languages, without the overhead of new neural architectures or extensive policy rewrites. Implementing these two layers would convert ad-hoc reversals into a rights-preserving workflow, strengthening Meta's compliance with international due process norms and materially improving the safety of frontline journalists in one of the world's most information-starved regions.

٤

<sup>&</sup>lt;sup>16</sup> 'Santa Clara Principles on Transparency and Accountability in Content Moderation' (*Santa Clara Principles*, 2022) <a href="https://santaclaraprinciples.org/">https://santaclaraprinciples.org/</a> accessed 15 July 2025.

<sup>&</sup>lt;sup>17</sup> Access Now, 'New: Content Governance in Crises Declaration' (*Access Now*, 20 February 2024) <a href="https://www.accessnow.org/publication/new-content-governance-in-crises-declaration/">www.accessnow.org/publication/new-content-governance-in-crises-declaration/</a> accessed 15 July 2025.