# UK: TECHNICAL EXPLAINER ON X'S RECOMMENDER SYSTEM AND THE 2024 RACIST RIOTS

**Amnesty International is a movement of 10 million people which mobilizes the humanity in everyone and campaigns for change so we can all enjoy our human rights. Our vision is of a world where those in power keep their promises, respect international law and are held to account. We are independent of any government, political ideology, economic interest or religion and are funded mainly by our membership and individual donations. We believe that acting in solidarity and compassion with people everywhere can change our societies for the better.**

**amnesty.org**

*Cover photo:* People attend a counter demonstration against an anti-immigration protest called by far-right activists, close to an Immigration Solicitors' office in Westcliff, near Southend-on-Sea, eastern England on August 7, 2024. (Photo by STF/AFP via Getty Images)

**AMNESTY INTERNATIONAL**

# CONTENTS

# 1. EXECUTIVE SUMMARY

## 1.1 OVERVIEW

- After a triple murder in Southport, false claims regarding the perpetrator's identity spread rapidly on X (formerly Twitter) contributing to violent racist riots which erupted in several locations across the UK.
- Amnesty International's technical assessment of X's publicly available algorithm reveals a platform that systematically promotes and amplifies content that provokes strong reactions and controversy, often with insufficient safeguards to mitigate resulting harms.
- Amnesty International's analysis finds that X's design and policy choices created fertile ground for inflammatory racist narratives to thrive in the wake of the Southport attack – with significant human rights consequences for migrants, Muslims, and racialised communities.

## 1.2 EXECUTIVE SUMMARY

In the immediate aftermath of a tragic triple murder in Southport on 29 July 2024, social media platform X (formerly Twitter) became a hotspot for racist, Islamophobic and xenophobic rhetoric. False claims alleging the perpetrator of the attack was a Muslim immigrant or asylum-seeker gained significant traction online. As hateful narratives spread on X, offline violence erupted, with mobs targeting mosques, refugee shelters, and Asian, Black and Muslim communities amid a wave of violent racist riots which swept across multiple UK towns and cities. As outlined in the UN Guiding Principles on Business and Human Rights (UN Guiding Principles), companies like X have a responsibility to respect human rights. This includes taking steps to avoid causing or contributing to human rights harms through their design and operations, and to address those impacts when they occur.

One year after the racist riots[1] took place, Amnesty International's technical assessment of X's publicly available recommender algorithm reveals that, through its algorithmic design, the social media platform played a central role in the spread of false narratives and content which fuelled and contributed to violence against racialised communities. Analysis also showed that the platform's policy and design choices meant it was more likely to provoke heated exchanges, reactions, and engagement, often without adequate safeguards to prevent or mitigate harm.

In the critical window after the Southport attack, X's engagement-driven system meant that inflammatory posts, even if entirely false, went viral, outpacing efforts to correct the record or de-amplify harmful content - some of which amounted to advocacy of hatred that constitutes incitement to discrimination or violence. These engagement-first design choices contributed to heightened risks amid a wave of anti-Muslim and anti-

---

[1] Amnesty International, UK *Government Must Address Root Cause of Racism That Plagues Our Society*, 2 August 2024, https://www.amnesty.org/en/latest/news/2024/08/uk-government-must-address-root-cause-of-racism-that-plagues-our-society/

migrant violence observed in several locations across the UK at the time, and which continues to present a serious human rights risk today.

This technical explainer analyses the role of X's design and policy choices and their contribution to adverse human rights impacts in the context of the 2024 racist riots. It provides background on the architecture of X's recommendation system, highlights specific design choices made by the company, and explains why X's engagement-first approach to content ranking results in significant human rights risks.

The platform-fuelled racist violence which took place in the UK reinforces longstanding concerns about Big Tech and the way their platforms operate. Amnesty International's previous research on 'Toxic Twitter' and investigations into Meta's role in contributing to violence in Myanmar and Ethiopia has consistently shown a pattern: when platforms prioritise engagement and profit over human rights, marginalized groups often pay the price.[2]

[2] Amnesty International, *#Toxictwitter: Violence and abuse against women online*, 21 March 2018, https://www.amnesty.org/en/documents/act30/8070/2018/en/; Amnesty International, *The Social Atrocity: Meta and the right to remedy for the Rohingya,* 28 September 2022, https://www.amnesty.org/en/documents/asa16/5933/2022/en/; Amnesty International, *A Death Sentence for My Father: Meta's contribution to human rights abuses in northern Ethiopia*, 11 October 2023, https://www.amnesty.org/en/documents/afr25/7292/2023/en/

**X'S RECOMMENDER SYSTEM TECHNICAL EXPLAINER**
X'S DESIGN CHOICES AMPLIFIED FALSE AND HARMFUL CONTENT TARGETING RACIALISED PEOPLE, INCLUDING MIGRANTS, REFUGEES AND MUSLIMS, DURING THE 2024 RACIST RIOTS

Amnesty International                                                                                                          5

# 2. INTRODUCTION

On 29 July 2024, three young girls - Alice da Silva Aguiar, Bebe King and Elsie Dot Stancombe -were tragically murdered in Southport, England, by 17-year-old Axel Rudakubana.[3] Before official accounts were shared, false statements and Islamophobic narratives about the incident began circulating on social media. Accounts posing as news outlets and influential far-right figures alleged online, without evidence, that the crime was an "Islamist" attack or part of a wider "anti-Western" plot. Almost immediately after the incident, misinformation and falsehoods about the perpetrator's identity, religion, and immigration status flooded social media platforms, incorrectly naming the perpetrator as "Ali Al-Shakati" and claiming he was an asylum-seeker who had arrived by boat. These false narratives spread across all social media platforms, and were prominent on X, a platform that under Elon Musk's ownership since late 2022, has significantly loosened content moderation measures, verification requirements, cut trust and safety teams and reinstated previously banned accounts.[4] These changes created an environment in which racist and discriminatory content could thrive on X with minimal restraint.

Despite local police promptly debunking the speculative claims about the attacker's identity and clarifying the suspect was a UK born 17-year-old British male who is not an asylum-seeker - the damage was done.[5] The false claims, anger and paranoia – fuelled not only by days of online vitriol but also by years of entrenched structural racism in the UK, including anti-Muslim racism - had already spilled into the streets.[6] Within hours and over the ensuing days following the attack, online rhetoric was accompanied by offline violence as racist riots erupted in Southport and across several UK towns and cities, with the worst violence taking place in England and Northern Ireland. Mobs, many mobilised by rumours on social media, targeted Muslims or those perceived as Muslim, and migrant communities: a local mosque in Southport was attacked and vandalised, bricks were hurled at police, and businesses associated with immigrants were destroyed. Rioters set up improvised "checkpoints" on roads, and attempted arson against facilities housing asylum seekers.[7]

The unrest continued through late July and into August 2024, marking some of the worst racist violence the UK had seen in years.[8] The riots demonstrated the pernicious role of social media algorithms in fuelling online falsehoods and offline harms.

As authorities tried to contain the unrest, the role of social media, and X in particular, came under scrutiny. Prime Minister Keir Starmer warned that violent disorder "whipped up online" was criminal and happening

---

[3] Liverpool Crown Court, *R v. Axel Rudakubana: Sentencing Remarks,* 27 January 2025, https://www.judiciary.uk/wp-content/uploads/2025/01/R-v-Axel-Rudakubana.pdf

[4] eSafety Commissioner, "Report reveals the extent of deep cuts to safety staff and gaps in Twitter/X's measures to tackle online hate", 11 January 2024, https://www.esafety.gov.au/newsroom/media-releases/report-reveals-the-extent-of-deep-cuts-to-safety-staff-and-gaps-in-twitter/xs-measures-to-tackle-online-hate.

[5] Telegraph, "Merseyside police say suspect named online in Southport stabbings is wrong man", 30 July 2024, https://www.telegraph.co.uk/news/2024/07/30/merseyside-police-southport-suspect-name-online-wrong/

[6] Amnesty International UK, "Failure to Tackle Institutional Racism Is Root Cause of Racist Violence Unfolding on Our Streets", 12 August 2024, https://www.amnesty.org.uk/press-releases/uk-failure-tackle-institutional-racism-root-cause-racist-violence-unfolding-streets

[7] BBC News, "Footage shows panic inside Rotherham asylum hotel during riot", 16 August 2024, https://www.bbc.co.uk/news/videos/c0e8lx2dvvyo; BBC News, "Violent Southport protests reveal organising tactics of the far-right", 2 August 2024, https://www.bbc.co.uk/news/articles/cl4y0453nv5o

[8] BBC News, "Did social media fan the flames of riot in Southport?", 1 August 2024, https://www.bbc.co.uk/news/articles/cd1e8d7llg9o

"on [social media companies'] premises," stating "the law must be upheld everywhere".[9] In the aftermath of the racist riots British authorities responded with a series of arrests and pursued criminal charges against individuals who used X and other platforms to incite violence or spread malicious falsehoods. Some perpetrators received prison sentences for their social media posts.[10]

UK government officials and regulators publicly expressed concern over the role of social media platforms — including X — in fuelling the July-August 2024 riots.[11] Following the riots, the UK Communications Regulator (Ofcom) undertook a fact-finding exercise. Ofcom's Chief Executive, Dame Melanie Dawes, sent an open letter to the government in October 2024 outlining the regulator's observations.[12] Ofcom found that illegal content and malicious falsehoods "spread widely and quickly online" after the Southport murders, and that some posts had "malicious intent, seeking to influence public reaction." Crucially, Ofcom concluded that social media posts appeared "to have contributed to the significant violent disorder" that ensued. The letter highlighted specific failures: accounts with over 100,000 followers had falsely identified the attacker as a Muslim asylum-seeker and spread unverified claims about his background.[13] In response, officials stressed that under the new Online Safety Act social media firms must swiftly remove illegal incitement and disinformation, adding, that "they [social media companies] shouldn't be waiting for the Online Safety Act for that" and that companies will be held to account if they fail to curb such content.[14] In November 2024, the House of Commons Science, Innovation and Technology Select Committee launched an inquiry into the role of social media algorithms in the riots.[15] The committee published its findings on 11 July 2025.[16]

[9] Reuters, "Keir Starmer warns social media firms after Southport misinformation fuels riots", 2 August 2024 https://www.reuters.com/world/uk/pm-starmer-warns-social-media-firms-after-southport-misinformation-fuels-uk-2024-08-01
[10] Crown Prosecution Service, "Man jailed just two days after posting online during public disorder", 9 August 2024, https://www.cps.gov.uk/cps/news/man-jailed-just-two-days-after-posting-online-during-public-disorder: BBC News, "'Keyboard warrior' jailed for part in UK disorder", 16 August 2024, https://www.bbc.co.uk/news/articles/c5y3gre3y9yo
[11] Politico, "Starmer slams social media (but not Farage) as far right roils UK", 1 August 2024, https://www.politico.eu/article/keir-starmer-nigel-farage-social-media-far-right-roils-uk/; Telegraph, "Sir Keir Starmer to review social media laws in wake of riots", 9 August 2024, https://telegraph.co.uk/politics/2024/08/09/sir-keir-starmer-to-review-social-media-laws-wake-of-riots/; Politico, Race riots put Britain on collision course with Elon Musk, 5 August 2024, available at: https://www.politico.eu/article/united-kingdom-riots-elon-musk-x-twitter-fake-news-disinformation
[12] Ofcom (Dame Melanie Dawes), "Letter to the Secretary of State", 22 October 2024, https://www.ofcom.org.uk/siteassets/resources/documents/about-ofcom/public-correspondence/2024/letter-from-dame-melanie-dawes-to-the-secretary-of-state-22-october-2024.pdf
[13] Ofcom (Dame Melanie Dawes), "Letter to the Secretary of State" (previously cited)
[14] BBC News, "Clear link between online posts and violent disorder", 22 October 2024, https://www.bbc.co.uk/news/articles/c70w0ne4zexo; Guardian, "No 10 criticises Elon Musk for 'civil war is inevitable' post on England riots", 5 August 2024, https://www.theguardian.com/uk-news/article/2024/aug/05/no-10-criticises-elon-musk-post-x-riots
[15] UK Parliament, Science, Innovation and Technology Committee "What are the links between social media algorithms, generative AI, and the spread of harmful content online?", 20 November 2024, https://committees.parliament.uk/work/8641/social-media-misinformation-and-harmful-algorithms/news/203870/what-are-the-links-between-social-media-algorithms-generative-ai-and-the-spread-of-harmful-content-online/
[16] UK Parliament, Science, Innovation and Technology Committee, *Social Media, Misinformation and Harmful Algorithms: Second Report of Session 2024–25,* 11 July 2025, https://committees.parliament.uk/publications/48745/documents/255558/default/

# 3. REMOVING THE GUARDRAILS: CHANGES AT X

The Southport tragedy and ensuing racist riots unfolded against a backdrop of major policy and personnel changes at X. Since Elon Musk's takeover in late 2022, X has dismantled or weakened many of its safety guardrails aimed at curbing harmful content and disinformation, from mass layoffs of content moderation staff to the reinstatement of banned accounts.[17] Elon Musk disbanded Twitter's Trust and Safety advisory council, fired trust and safety engineers, restored numerous accounts banned for hate or harassment (including that of Stephen Yaxley-Lennon, also known as Tommy Robinson), and publicly declared a more permissive approach to speech on the platform.[18] These actions coincided with a documented rise in "hate speech" on the platform. A study conducted by researchers at UC Berkeley, University of California, Los Angeles (UCLA) and University of Southern California (USC) found that the "increase in hate speech just before Musk bought X persisted until at least May of 2023, with the weekly rate of hate speech being approximately 50% higher than the months preceding his purchase, although this increase cannot be directly attributed to any policy at X. The increase was reported across multiple dimensions of hate, including racism, homophobia, and transphobia. Moreover, there is a doubling of hate post "likes," indicating increased engagement with hate posts."[19]

There is to date, no evidence X conducted any human rights due diligence around these major policy shifts, despite the company having a responsibility to identify and address risks in line with international business and human rights standards. Amnesty International put this allegation to X in a letter dated 18 July 2025. X did not respond. X should have evaluated how restoring figures like Tommy Robinson, who had been previously banned for violating hate speech rules, might result in human rights harms.
Elon Musk publicly espoused a "free speech absolutist" stance, indicating that virtually any speech would be allowed on the platform.[20]  This approach was reflected in X's new policy - 'Freedom of Speech, Not Reach'-

---

17 ABC News, "Potential mass layoffs at Twitter could cripple content moderation, some experts say", 23 October 2022, https://abcnews.go.com/Business/potential-mass-layoffs-twitter-cripple-content-moderation-experts/story?id=91856973
18 Associated Press, "Elon Musk's Twitter dissolves trust and safety council", 13 December 2022, https://apnews.com/article/elon-musk-twitter-inc-technology-business-a9b795e8050de12319b82b5dd7118cd7; eSafety Commissioner, "Report reveals the extent of deep cuts to safety staff and gaps in twitter/x's measures to tackle online hate", 11 January 2024, https://www.esafety.gov.au/newsroom/media-releases/report-reveals-the-extent-of-deep-cuts-to-safety-staff-and-gaps-in-twitter/xs-measures-to-tackle-online-hate; Washington Post, "'Opening the gates of hell': Musk says he will revive banned accounts", 24 November 2022, https://www.washingtonpost.com/technology/2022/11/24/twitter-musk-reverses-suspensions
19 Berkeley News, "Study finds persistent spike in hate speech on X", 13 February 2025, https://news.berkeley.edu/2025/02/13/study-finds-persistent-spike-in-hate-speech-on-x/; Daniel Hickey, Daniel M. T. Fessler, Kristina Lerman, Keith Burghardt, *X under Musk's leadership: Substantial hate and no reduction in inauthentic activity*, PLoS One , 12 February 2025, Volume (18 (3)), Available at: https://pubmed.ncbi.nlm.nih.gov/39937728/
20 Elon Musk (@elonmusk), X post: "Starlink has been told by some governments (not Ukraine) to block Russian news sources", 5 March 2022, https://x.com/elonmusk/status/1499976967105433600; Independent, "Elon Musk fires back at Twitter censorship critic: 'You're such a numbskull'", 29 May 2023, https://www.independent.co.uk/tech/elon-musk-twitter-turkey-censorship-b2347509.html; Vanity Fair, "Elon Musk, self-proclaimed 'free speech absolutist,' says Twitter will treat 'cisgender' as a slur", 21 June 2023, https://www.vanityfair.com/news/2023/06/elon-musk-free-speech-twitter-will-treat-cisgender-slur; NPR, "Elon Musk calls himself a 'free speech absolutist.' what could Twitter look like under his control?", 8 October 2022, https://www.npr.org/2022/10/08/1127689351/elon-

X'S RECOMMENDER SYSTEM TECHNICAL EXPLAINER
X'S DESIGN CHOICES AMPLIFIED FALSE AND HARMFUL CONTENT TARGETING RACIALISED PEOPLE, INCLUDING MIGRANTS, REFUGEES AND MUSLIMS, DURING THE 2024 RACIST RIOTS
Amnesty International

8

introduced in April 2023.[21] The Southport case starkly demonstrated the failures of this approach: far from being de-amplified, false and harmful narratives gained viral reach on X, with little evidence to suggest any interventions from the company succeeded at slowing them down. It is unclear if X pre-emptively down-ranked the topic "Ali al-Shakati" or if it added friction measures to its sharing, even after officials refuted the false alleged identity of the perpetrator.[22] Amnesty International put this question to X in a letter dated 18 July 2025. X did not respond.

One of the starkest examples of the impact of these policy change choices and removal of safeguards was the case of Lucy Connolly, a woman who posted on X in the midst of the riots calling for *"Mass deportation now, set fire to all the f\*\*ing hotels full of the bastards for all I care"*, which was not initially taken down by X. This was because X's moderators (or automated systems) decided that her post did not violate the platform's rules on violent threats, therefore leaving it online.[23] That post, explicitly inciting arson against asylum-seeker accommodations, remained public and was viewed 310,000 times, despite her account having less than 9,000 followers at the time. The post was seen far beyond the reach of her immediate network, suggesting algorithmic distribution likely played at least some role in aiding this. Lucy Connolly was later convicted under the Public Order Act for intentionally stirring up racial hatred with that message, receiving a 31-month prison sentence. Yet X's failure to remove it promptly is telling, as it highlights how the platform's policy choices combined with the algorithm's prioritisation of engagement, meant dangerous content faced few barriers to virality.[24]

One of the most concerning changes introduced following Elon Musk's acquisition of Twitter (now X) was the rollout of X Premium (formerly Twitter Blue), a paid subscription model launched in late 2023. The service introduced tiered access Basic, Premium, Premium+ that includes enhanced platform features, most notably, reply "boosts", which increase the visibility of a subscriber's content in conversations and ranked feeds. The higher the subscription tier, the greater the amplification a user's replies and posts receive.[25]

By mid-2024, the new model no longer required identity verification. At the time of the racist riots in July–August 2024, users could access verification through X Premium by providing only a display name, a profile picture, a confirmed phone number, and recent activity within 30 days, regardless of identity.[26]

Whilst the company maintained that basic integrity requirements remained, the reduced verification threshold and the monetisation of visibility enabled a pay-to-amplify model. In practice, this meant that users, including those spreading disinformation or engaging in inflammatory rhetoric, could secure increased reach for their content simply by subscribing. X's transition from identity-based verification to a paid "blue tick" model enabled the algorithmic boosting of content posted by paying users, irrespective of the type of content.

musk-calls-himself-a-free-speech-absolutist-what-could-twitter-look-like-un; X Corp., About offensive content, https://help.x.com/en/safety-and-security/offensive-posts-and-content (accessed 17 July 2025)

[21] X Corp., "Freedom of speech, not reach: An update on our enforcement philosophy", 17 April 2023, https://blog.x.com/en_us/topics/product/2023/freedom-of-speech-not-reach-an-update-on-our-enforcement-philosophy

[22] Sky News, "Southport attack misinformation fuels far-right discourse on social media", 31 July 2024, https://news.sky.com/story/southport-attack-misinformation-fuels-far-right-discourse-on-social-media-13188274; Independent, "How a few Twitter posts on Elon Musk's X helped fan the flames of unrest and rioting across the UK", 6 August 2024, https://www.independent.co.uk/tech/elon-musk-uk-riots-southport-twitter-x-b2591725.html; Merseyside Police, "Update on Major Incident in Southport", 30 July 2024, https://www.merseyside.police.uk/news/merseyside/news/2024/july/update-on-major-incident-in-southport.

[23] BBC News, "Racial hatred post did not break X rules", 3 September 2024, https://www.bbc.co.uk/news/articles/cn8ljjjmpg5o.

[24] Royal Courts of Justice, *Lucy Connolly v. The King: Approved Judgment*, 15 May 2025, https://www.judiciary.uk/wp-content/uploads/2025/05/Lucy-Connolly-v-The-King.pdf

[25] X Corp., "Verified on X", https://verified.x.com/en (accessed on 17 July 2025)

[26] X Corp., "About Verified accounts", https://help.x.com/en/managing-your-account/about-x-verified-accounts (accessed on 17 July 2025)

# 4. ISLAMOPHOBIC, RACIST AND XENOPHOBIC RHETORIC ON X AFTER THE SOUTHPORT ATTACK

In the 48 hours following the Southport murders, incendiary posts by far-right influencers went viral on X. Within hours of the stabbing on 29 July, #Southport, #Stabbing and #EnoughisEnough among other related hashtags trended on X, as users raced to offer their own narratives with unverified claims, based on racist and xenophobic stereotypes and prejudices, that the attacker was "a Muslim" who "came to UK by boat". An account on X called "Europe Invasion", known to publish anti-immigrant and Islamophobic content[27] posted shortly after news of the attack emerged that the suspect was "alleged to be a Muslim immigrant".[28] That post garnered over four million views.[29] Within 24 hours, all X posts speculating that the perpetrator was Muslim, a refugee, a foreign national or arrived by boat were tracked to have an estimated 27 million impressions[30] in total.[31] This claim gained further traction through X's own features when the fake name 'Ali al-Shakati' appeared amongst X's top trending topics[32] and was surfaced in user search results. One later analysis found that posts sharing the name Ali al-Shakati were seen over 420,000 times and had an enormous potential reach across the social media platform.[33]

As false claims that the Southport attacker was a Muslim and/or an asylum-seeker trended on X, prominent far-right accounts rushed to amplify the narrative and spread anti-Muslim racism on the platform. Stephen Yaxley-Lennon a far-right activist, also known as Tommy Robinson, told his 840,000 X followers at the time there was "more evidence to suggest Islam is a mental health issue rather than a religion of peace". [34]

---

[27] Radio Free Europe, "The many faces of Europe invasion, the anonymous x account stoking anti-immigrant narratives", 17 December 2024, https://www.rferl.org/a/europe-invasion-x-account-disinformation-xenophobia-immigration-x-account/33239067.html; Marc Owen Jones (Associate Professor at Northwestern University in Qatar), "Written evidence submitted by Marc Owen Jones", 12 March 2025, https://committees.parliament.uk/writtenevidence/138332/html/

[28] Shayan Sardarizadeh, Senior journalist at BBC Verify, X post: "A reminder that far-right account 'Europe Invasion' was among the most influential accounts that circulated misinformation about the identity of the Southport stabbing suspect", 25 August 2024, https://x.com/Shayan86/status/1827694952413114394

[29] BBC Bitesize, "Timeline of how online misinformation fuelled UK riots", August 2024, https://www.bbc.co.uk/bitesize/articles/zshjs82

[30] A Twitter impression refers to the number of times people view a particular tweet organically on the X platform. Some authors describe impressions as the number of times a tweet appears on users' screens.

[31] Guardian, "How false online claims about Southport knife attack spread so rapidly", 31 July 2024, https://www.theguardian.com/uk-news/article/2024/jul/31/how-false-online-claims-about-southport-knife-attack-spread-so-rapidly

[32] Institute for Strategic Dialogue, "From rumours to riots: How online misinformation fuelled violence in the aftermath of the Southport attack", 31 July 2025, https://www.isdglobal.org/digital_dispatches/from-rumours-to-riots-how-online-misinformation-fuelled-violence-in-the-aftermath-of-the-southport-attack/

[33] Independent, "How fake claims over Southport suspect spread like wildfire with false name seen more than 420,000 times", 17 August 2024, https://www.independent.co.uk/news/uk/crime/southport-riots-uk-false-identity-misinformation-suspect-b2594042.html

[34] Tommy Robinson, X post: "There's more evidence to suggest islam is a mental health issue rather than a religion of peace " 30 July 2024, https://x.com/TRobinsonNewEra/status/1818328857587413359

Andrew Tate, a notorious online influencer, posted a video falsely asserting the attacker was an "undocumented migrant" who "arrived on a boat" and lamented that "when the invaders slaughter your daughters, you do absolutely f****** nothing".[35] These incendiary comments racked up millions of views on X and poured fuel on an already volatile situation. Notably, both Tommy Robinson and Andrew Tate had previously been banned from Twitter for hate speech and harmful content, but their accounts were reinstated in late 2023 under Elon Musk's "amnesty" for suspended users.[36] By 2024, they once again had amassed a large following on X. X's owner Elon Musk – whose personal X account has over 220 million followers[37] – became a notable amplifier of the false narratives being exchanged regarding the racist riots. On 5 August 2024, as riots spread, he replied under a video of the disorder originally posted by Ashley St Clair with the statement: "civil war is inevitable".[38] In another post on 5 August 2024, when Prime Minister Keir Starmer, urged protection of Muslim communities,[39] Elon Musk publicly retorted: "Shouldn't you be concerned about attacks on *all* communities?". He further replied to Tommy Robinson's posts about the violence with enthusiastic punctuation "!!".[40] Elon Musk's engagement with these posts is consequential not only because of their content, but rather because it likely triggered algorithmic cascades due to both the size of his audience and the engagement-based design of the platform. By interacting with those spreading inflammatory narratives, X's leadership not only lent visibility to their claims but likely contributed to their further algorithmic amplification.

On 6 August 2024 Elon Musk reposted a video and tweet by Andy Ngo claiming that ""Armed "Muslim patrol" members surrounded and attacked a pub after marching around the area looking for white right-wingers to attack. (The rumoured "far-right" protest never materialized.)"[41]– adding "Why aren't all communities protected in Britain? @Keir_Starmer" echoing a slogan alleging preferential treatment for Muslims.[42] X's public view counter states the post had 68 million views as of 31 July 2025.[43] He then further commented "#TwoTierKeir" , a post that's received 7.2 million impressions, with the thread subsequently trending under "#TwoTierKeir".[44]  Because Elon Musk's account was and remains the platform's most followed and carries a Premium status, any reply, quote tweet or "like" from him triggers a high weight engagement signal within X's ranking model.[45]

On X, when accounts with large follower numbers or verified Premium accounts, such as Elon Musk's, engage with content, for example by replying to Tommy Robinson or commenting on riot videos, the platform's ranking system tends to push that content to millions who would likely not otherwise see it. Amnesty International's data on impression counts supports this: the reach of Tommy Robinson's and others' posts far exceeded their follower numbers and typical expected engagement metrics (likes, views, reposts, etc.), suggesting that algorithmic distribution played a role in distributing this content beyond the followers of these accounts. Amnesty International researchers analysed engagement metrics on 292 posts (including reposts) from Tommy Robinson's X account from 3 Aug 2024 to 8 Aug 2024. 15% of these posts exceeded 1 million views by August 8th, all gaining a wider audience than his 900,000 followers as of August 2024.[46] The net result was a staggering amplification of hate and anti-migrant sentiment. Amnesty International's analysis shows that in the two weeks following the Southport attack, Tommy Robinson's posts

[35] Al Jazeera, "Southport stabbing: What led to the spread of disinformation?", 2 August 2024, https://www.aljazeera.com/news/2024/8/2/southport-stabbing-what-led-to-the-spread-of-disinformation
[36] BBC News, "The haters and conspiracy theorists back on Twitter", 8 March 2023, https://www.bbc.co.uk/news/technology-64554381
[37] Elon Musk (@elonmusk), https://x.com/elonmusk (accessed on 1 August 2025)
[38] Elon Musk (@elonmusk), X comment :"Civil War is inevitable.", 4 August 2024, https://x.com/elonmusk/status/1819933223536742771
[39] Elon Musk (@elonmusk), X comment: "Shouldn't you be concerned about attacks on *all* communities?", 5 August 2024, https://x.com/elonmusk/status/1820502825308504499
[40] Elon Musk (@elonmusk), X comment:"!!", 1 August 2024, https://x.com/elonmusk/status/1819062565395439767
[41] Andy Ngo (@MrAndyNgo), X Post: "Armed 'Muslim patrol' members surrounded and attacked a pub after marching around the area looking for white right-wingers to attack.", 5 August 2024, https://x.com/MrAndyNgo/status/1820564628163670159
[42] Elon Musk (@elonmusk), X Post: "Why aren't all communities protected in Britain? @Keir_Starmer", 6 August 2024, https://x.com/elonmusk/status/1820790297233592361
[43] Elon Musk (@elonmusk), X Post: "Why aren't all communities protected in Britain? @Keir_Starmer", 6 August 2024, https://x.com/elonmusk/status/1820790297233592361
[44] Elon Musk (@elonmusk), X comment:"#TwoTierKeir,", 6 August 2024, https://x.com/elonmusk/status/1820805621534400786
[45] X's ranking model refers to one stage of the recommendation pipeline, namely where a list of potential 'candidate' tweets are scored according to their relevance to the user in question. This process determines the content that each user receives on their 'For You' timeline; Twitter, "Twitter's recommendation algorithm", 31 March 2023, https://blog.x.com/engineering/en_us/topics/open-source/2023/twitter-recommendation-algorithm
[46] Telegraph, "The podcast that rocketed Tommy Robinson into the mainstream", 7 August 2024, https://www.telegraph.co.uk/news/2024/08/07/how-tommy-robinson-became-mainstream/

on X received over 580 million views[47] - an unprecedented reach for a figure banned on most mainstream platforms for breaching hate speech rules.[48]

Our analysis of the recommender system and the removal of safeguards at X demonstrate that the recommender system is set up to amplify this type of harmful content and that there are few safeguards to prevent it.  It is important to acknowledge that heightened user activity during periods of public crisis, and user-to-user resharing of content, could both also contribute to increased views, and therefore, algorithmic distribution may not be the sole factor driving this extensive reach.

The combination of these factors led to the spread of anti-migrant and anti-Muslim rhetoric online, which entailed significant adverse human rights impacts for British Muslim and migrant communities, including discrimination and threats to personal security.

---

[47] Evidence of X posts on file with Amnesty International.
[48] Guardian, "Tommy Robinson banned from Facebook and Instagram", 26 February 2019, https://www.theguardian.com/uk-news/2019/feb/26/tommy-robinson-banned-from-facebook-and-instagram.

# 5. AN ALGORITHM ENGINEERED FOR ENGAGEMENT

X's recommender system, the algorithm behind the "For You" timeline, is an engagement-driven ranking system that decides which posts each user sees. Rather than showing tweets chronologically or only from followed accounts, the "For You" feed uses machine-learning models and heuristics to maximize user engagement. X's platform constantly asks: "What content will this user likely interact with?" and then boosts posts that fit that prediction.

The longer a user remains active on the platform, the more behavioural and engagement data can be collected about them, which in turn enables increasingly precise targeting for both content delivery and advertising. This design is not unique to X but rather reflects a broader industry trend of social media algorithms optimised for engagement and profit.

Amnesty International has previously found that this model incentivises platforms to optimise for content that provokes strong engagement – including material that is inflammatory, discriminatory, or harmful– as this tends to sustain user attention and generate the most interaction.[49] On X, this dynamic is shaped by personalised recommendation systems that determine which posts, hashtags, and accounts are promoted to each user. Notably, the system does not limit itself to content from accounts a user follows.

Based on the publicly available information as of 2023, every time a user opens X, the system pulls in a large pool of candidate tweets (up to roughly 1500) from various sources. This includes In-Network content - tweets by accounts the user follows - and Out-of-Network content - popular or relevant tweets by others. [50] Once the candidate tweets are gathered, X's "Home Mixer" service ranks them using a machine-learning model, often called the "heavy ranker". This ranking model is a neural network (reportedly with around 48 million parameters) that evaluates thousands of features about the user, the tweet, and their interaction history.[51] It essentially tries to predict the probability of various engagement actions for each tweet – such as the likelihood the user will like, reply, retweet, unlike, click on, or spend time reading that tweet. X then records every interaction and reply as discrete explicit signals or feedback that feed into both sourcing and ranking models. Actions such as bookmarking a tweet and sharing it via the share button also inform the

---

[49] Amnesty International, *Surveillance Giants: How the Business Model of Google and Facebook Threatens Human Rights*, 21 November 2019, https://www.amnesty.org/en/documents/pol30/1404/2019/en/; Amnesty International, "Big Tech platforms are playing an active role in fuelling racist violence", 6 August 2024, https://www.amnesty.org/en/latest/news/2024/08/uk-big-tech-platforms-play-an-active-role-in-fuelling-racist-violence/; Amnesty International, "Written Evidence submitted to the Science, Innovation and Technology Committee: Social Media, Misinformation, and Harmful Algorithms inquiry", 14 May 2025, https://committees.parliament.uk/writtenevidence/141492/default/
[50] Twitter, "Twitter's recommendation algorithm" (previously cited)
[51] A neural network is a type of machine learning model inspired by the structure of the human brain. It consists of layers of interconnected nodes (or "neurons") that learn patterns in large datasets. In X's algorithm, a neural network is used to predict which tweets a user is most likely to engage with, based on complex relationships between user behaviour, tweet content, and social signals.

algorithm about which content a user values most. In addition to the explicit feedback, X tracks what is referred to as "tweet clicks". These include actions such as opening a tweet detail view, watching videos (measured by duration or percentage viewed), and profile visits - all collected through the User Signal Service. These signals help the system infer interest even when a user does not actively like or retweet. The model outputs a relevance score for each candidate tweet based on these predicted probabilities. However, not all user actions are treated equally by X's algorithm. In general, active, and conversational interactions are heavily rewarded, whereas passive approval such as retweets and likes or negative feedback carry little weight (or even incur penalties).

By design, X's algorithm reaches far beyond the accounts followed by users. The algorithm actively seeks out tweets that are trending or engaging across the platform, so it can recommend them to new audiences. For example, if many people in your extended network liked or replied to a particular tweet, or if a topic similar to your interests is suddenly trending, the system will consider those tweets for your feed even if you do not follow the account which posted the original content. This approach, powered by X's Social Graph Traversal and Clustering Algorithms, means that viral content can jump from fringe communities into mainstream feeds if it generates enough engagement.[52]

◎ ↓FIGURE 1: WEIGHTINGS FOR EACH PREDICTED PROBABILITY[53]

| FEATURE | WEIGHT | DESCRIPTION |
|---|---|---|
| FAVOURITE (LIKE) | 0.5 | Predicted probability of the user "favouriting" (liking) a tweet: very low influence on the final ranking score. |
| RETWEET | 1.0 | Predicted probability of the user retweeting: a light signal, only marginally more than a like. |
| REPLY | 13.5 | Predicted probability of the user replying: strongly boosts tweets that spark direct conversation. |
| GOOD PROFILE CLICK | 12.0 | Probability the user clicks into the author's profile and then likes/replies: valued nearly as much as a reply. |
| VIDEO PLAYBACK >50% | 0.005 | Probability the user watches more than 50% of a video: effectively zero impact on ranking. |
| REPLY ENGAGED BY AUTHOR | 75.0 | Probability the user replies, and the author subsequently engages: highest reward for sustained back-and-forth. |
| GOOD CLICK (CONVERSATION OPEN) | 11.0 | Probability the user opens the conversation view and then likes/replies: signals deep conversational interest. |
| GOOD CLICK V2 (2-MIN CONVERSATION VIEW) | 10.0 | Probability the user stays more than two minutes in the conversation view: strong indicator of engagement depth. |
| NEGATIVE FEEDBACK | −74.0 | Probability of negative feedback (for example, "show less," block or mute): heavily penalizes disliked or unwanted content. |
| REPORT | −369.0 | Probability the user reports the tweet: significantly demotes content deemed offensive or problematic. |

Amnesty International's analysis of X's open-source recommender algorithm uncovered systemic design

---

[52] Social Graph Traversal is a method used by platforms to explore relationships between users, based on their interactions and connections. For example, if many of your friends engage with a post, the algorithm may promote it to your feed even if you don't follow the original poster because the system "traverses" the social graph to identify relevant content. Clustering Algorithms are techniques used to group similar data points together in this case, grouping users or tweets into topical communities based on shared interests or behaviours. X uses clustering to identify which tweets might interest you, even from accounts you don't follow, by recommending content popular in communities you're algorithmically associated with.
[53] Github, "Heavy Ranker", https://raw.githubusercontent.com/twitter/the-algorithm-ml/main/projects/home/recap/README.md#:~:text%3Dcontributes%20a%20near%2Cyou%20can%20run%20the%20model

choices that favour contentious engagement over safety.[54] Under X's "heavy ranker" model – the machine-learning system that decides which posts get promoted – certain interaction signals carry high weightings. The system's top priority is to drive "conversation", regardless of the nature of the content. For example, a user reply to a tweet is valued at 26 times more than a like, and if the original author replies back, it receives a 75x boost[55]. By contrast, passive signals of approval are minimal, assigned just 1x and 0.5x respectively. In practical terms, X's code values a single reply thread as much as dozens of reshares, propelling tweets that spark argumentative back-and-forth to the top of recommendation lists.

Amnesty International could find no evidence that the ranking algorithm assesses a tweet's substance or potential harm (e.g. hate speech, misinformation) before boosting a post, aside from excluding known policy-violating material via separate filtering rules applied after scoring. A tweet that generates outrage and replies can be accelerated across the platform long before user reports or "negative feedback" can take effect to slow it down. Whilst X's algorithm does down-rank content that receives many reports or blocks, assigning a –369 weight to "report" signals, those safeguards depend on users to report harmful content, and only typically come into play after harmful posts have already gone viral. Amnesty International put this allegation to X in a letter dated 18 July 2025 but received no response.

In practice, this means if a post does not trigger an immediate automated takedown for company policy violations, the recommender system treats its virality as a proxy for value. As long as a tweet drives engagement, the algorithm appears to have no mechanism for assessing the potential for causing harm - at least not until enough users themselves report it. This apparent absence of a precautionary check is critical. Content that is discriminatory or perpetuates hateful narratives can achieve maximum visibility in the crucial early hours, well before any manual content moderation or user reporting can be effective. In the Southport case, the outrage-fuelled claims about a "Muslim asylum seeker" attacker were already trending across X and racking up millions of impressions long before officials or community leaders could effectively intervene.

Further Amnesty International analysis of X's recommender algorithm code also uncovered built-in amplification biases favouring certain users. Notably, X's code includes a feature that boosts content from "Premium" (formerly Blue) verified subscribers, which are paid accounts.[56] Tweets by paying users receive a multiplier on their ranking score roughly 4x for posts shown to one's followers and 2x greater for posts shown algorithmically to others.[57] This product decision by X's leadership to incentivise subscriptions means that Premium subscribers enjoy artificial amplification over ordinary users. In practice, many of the platform's most controversial and influential figures – including owner Elon Musk and several high-profile accounts, such as Tommy Robinson and Andrew Tate, reinstated after his takeover in 2022 – carry Blue verification badges, amplifying their reach by design.

These algorithmic weights, disclosed in X's own source code published in March 2023, create an environment which favours content that is likely to provoke a response and delivers it at scale. When divisive content drives replies, and replies drive ranking, then falsehoods, irrespective of their harm, may be prioritised and surface more quickly in timelines than verified information. These design features provided fertile ground for inflammatory racist narratives to thrive on X in the wake of the Southport attack. Posts that trigger frenzied replies to threads or come from Blue-verified accounts are systematically elevated, regardless of accuracy or harm. X's recommender system, which has an outsized influence on the user experience, therefore risks amplifying harmful content during highly polarising events.

In the days in which the riots took place, Amnesty International monitored trending metrics and engagement on the platform. Based on that, it is not clear if X pre-emptively downranked the trending false "Ali al-Shakati" topic or added friction measures to its sharing. X should have had an adequate crisis protocol for such a situation, similar to how some platforms claim to deploy "break the glass measures" during emergencies.[58] If such a protocol or measures existed at X, it either was not triggered, was triggered too late, or was ineffective. Amnesty International wrote to X on 18 July 2025 and asked about whether such measures were deployed, but received no response.

---

[54] Github, "Twitter/The Algorithm", https://github.com/twitter/the-algorithm; Twitter, "Twitter's recommendation algorithm" (previously cited)
[55] Github, "Heavy Ranker" (previously cited)
[56] X Corp., "About X Premium", https://help.x.com/en/using-x/x-premium (accessed 1 August 2025)
[57] This rule is found in a feature switch for "isBlueVerified". It is a product decision to incentivize subscription, and it's one of the areas where we could clearly see how business logic mixes with ranking. Github, "Twitter/The Algorithm" (previously cited); Twitter, "Twitter's recommendation algorithm" (previously cited)
[58] Tech Policy Press, "We know a little about Meta's 'break glass' measures. we should know more", 1 October 2024, https://www.techpolicy.press/we-know-a-little-about-metas-break-glass-measures-we-should-know-more/; Facebook, "Our comprehensive approach to protecting the US 2020 elections and inauguration day", 22 October 2021, https://about.fb.com/news/2021/10/protecting-us-2020-elections-inauguration-day

Amnesty International's analysis of X's technical architecture raises serious concerns about how the platform's recommender system functions and how it can fuel human rights abuses - especially during moments of crisis. The way the system weights, ranks, and boosts content, particularly posts that generate heated replies or are shared or created by "Blue" or "Premium" accounts, can result in harmful material being surfaced to large audiences. Designed in this way, X's recommendation engine can operate as a powerful accelerant of harmful content once it begins to circulate. Without meaningful safeguards to assess risk or harmfulness before distribution, the likelihood increases that inflammatory or hostile posts will gain traction, particularly during periods of heightened social tension. Where such content targets racial or religious groups, portrays marginalised communities as threatening or violent, and circulates in high volumes, X's failure to mitigate these foreseeable risks constitutes a failure to respect human rights, especially in contexts where the content may amount to advocacy of hatred that incites discrimination, hostility or violence. This failure reflects broader concerns on whether X's recommender system is designed in a manner consistent with the company's human rights responsibilities.

# 6. TRANSPARENCY BARRIERS AND LIMITATIONS

It is important to acknowledge the limitations in analysing X's algorithm and the transparency gaps that persist. Whilst this open-source release provided valuable insight into the recommender system's design at the time of its release in March 2023, X has not since published any updates to its underlying models or ranking weights. As such, Amnesty's analysis is based on the last verifiable configuration, and it remains unclear whether subsequent changes, if any, have addressed the system's vulnerability to amplifying harmful content. Amnesty International posed this question to X in a letter dated 18 July 2025, but did not receive any response.

Amnesty International's assessment of the recommender system relies exclusively on components of X's recommender system that have been made public such as the open-source code on X's GitHub and official engineering blog posts. It therefore represents a snapshot and does not encompass the platform's entire live production system. In fact, significant parts of the recommendation pipeline remain proprietary or absent from public documentation, for example, the actual learned weights of the machine-learning models, the data used to train them, various internal threshold values and configuration settings, the trust and safety enforcement algorithms that detect or downrank policy-violating content, and the modules governing advertisement recommendation.[59]

Because these components have not been disclosed, it is not fully possible to reconstruct or simulate how content was ranked or filtered during specific incidents such as the racist riots in real-time using only the available open-source materials. Accordingly, Amnesty International's findings are limited to the design choices and ranking signals that are observable in the released code or described in public sources. These transparency gaps themselves pose a serious barrier to accountability. Regulators, researchers and civil society seeking to analyse the full picture of the extent to which X's design features impact human rights or whether X's algorithms are currently operating in line with the company's human rights responsibilities to prevent the amplification of harmful content - in a way that respects and upholds human rights and user safety – cannot fully do so when much remains opaque. Notably, since Elon Musk's takeover, the platform has also erected barriers to scrutiny by severely restricting its API access, thereby limiting the ability of researchers and civil society to independently monitor the spread of harmful content.[60] Whilst X has published systemic risk assessments under the EU's Digital Services Act, it has not provided adequate detail or disclosures on the specific effects of its recommender systems or safeguard practices particularly in high-risk contexts.[61] The absence of meaningful transparency means we often only see the consequences after the fact, without a clear window into the algorithmic processes that led there. By curtailing transparency, X

---

[59] Twitter, "Twitter's recommendation algorithm" (previously cited)
[60] Guardian, "Techscape: Why Elon Musk's Twitter API changes could kill off useful bots – and researchers' access too", 7 February 2023, https://www.theguardian.com/technology/2023/feb/07/techscape-elon-musk-twitter-api
[61] X Corp., "Twitter International Unlimited Company: DSA Systemic Risk Assessment Report 2024", August 2024, https://transparency.x.com/content/dam/transparency-twitter/dsa/dsa-sra/dsa-sra-2024/TIUC-DSA-SRA-Report-2024.pdf

has effectively insulated itself from the very forms of scrutiny that make human rights due diligence enforceable and meaningful.

Despite these transparency gaps, there is enough evidence that suggests X's risk controls are insufficient. Key mitigation mechanisms, notably X's proprietary trust and safety filters that flag or downrank "abusive" and "toxic" content, remain a black box outside public scrutiny. However, the observable outcomes strongly indicate these measures are failing to curb harmful content. The open-source portions of X's 2023 algorithm showed an engagement-driven ranking system with only limited safeguards e.g. mild downranking after negative feedback, and X has not published any updates to suggest it has enhanced those safeguards since.

**X'S RECOMMENDER SYSTEM TECHNICAL EXPLAINER**
X'S DESIGN CHOICES AMPLIFIED FALSE AND HARMFUL CONTENT TARGETING RACIALISED PEOPLE, INCLUDING MIGRANTS, REFUGEES AND MUSLIMS, DURING THE 2024 RACIST RIOTS

Amnesty International                                                                                                          18

# 7. TIME FOR ACCOUNTABILITY: RECOMMENDATIONS

The racist riots of 2024 stand as a stark reminder of how social media design choices can contribute to real-world harms. Whilst primary responsibility for violent acts lies with the perpetrators, X's engagement-driven algorithm and weakened safeguards created a foreseeable risk of serious human rights impacts – a risk the company did not adequately address. One year on, the factors that fuelled the spread of Islamophobic and anti-migrant narratives on X remain a serious concern. Without significant reforms, X's platform could again serve to fuel the spread of hatred in future crises. Amnesty's analysis concludes that X's recommender system contributed to human rights harms during the 2024 racist riots, within the meaning of Principle 13 of the UN Guiding Principles. This contribution triggers responsibilities to cease such involvement and enable effective remedy, in line with Principle 22.

Whilst regulatory frameworks such as the UK's Online Safety Act (OSA) and the EU's Digital Services Act (DSA) now establish legal obligations for platforms to assess and mitigate some systemic risks, these obligations must be robustly enforced to have any effect. X's design choices and opaque practices continue to pose human rights risks that demand greater accountability, not just scrutiny. Regulatory authorities must now take action to ensure that X fulfils its responsibility to respect human rights. X must radically improve transparency in relation to its algorithm given that this algorithm is having measurable human rights impacts. Amnesty International is calling for independent accountability measures to address the systemic harms created by X's design choices.

Without systemic reform to X's design and transparency practices, the platform risks contributing to similar human rights harms in future crises. The company's algorithmic architecture and transparency practices must be reformed in order for the company to meet its human rights responsibilities and to avoid becoming a vector for future outbreaks of violence.
In line with international human rights standards, X has a responsibility to respect human rights and prevent its platform from facilitating incitement to discrimination or violence. The UN Guiding Principles state that companies should identify, prevent, mitigate, and account for how they address their impacts on human rights. This includes remediating harms to which they have contributed to.

Accordingly, Amnesty International makes the following recommendations:

To X:
-   X should institute human rights due diligence across all aspects of its operations including recommender systems and content moderation policies, in line with the UN Guiding Principles. This process should span the design, deployment, and continuous evolution of these systems, ensuring that potential human rights risks are identified, assessed, and addressed proactively as an integral part of X's operations rather than only being addressed after harms occur. X should also be transparent about how risks and impacts are identified and addressed.

- X should redesign its ranking system to ensure compliance with its human rights responsibilities. Concretely, the algorithm should include checks that de-amplify content that poses human rights risks without requiring user reports to do so.
- X should reduce human rights risks by adapting how content is ranked, promoted, or distributed on the platform. X should establish mechanisms for increasing its ability to prevent harm, through changes in design and platform rules, especially when serious risks are identified.
- X should proactively share data and insights about its algorithms and content flows, particularly around major incidents. This means restoring access for researchers by ensuring access to the latest API and/or creating a safe data-sharing mechanism so that external experts can monitor trends like spikes in discriminatory content in real time.
- X should publish information such as the prevalence of harmful content in "For You" feeds, the effectiveness (or failure rate) of its content moderation in crises, and the outcomes of any alterations to its algorithms.
- Independent third-party audits of X's recommender system and human rights due diligence processes should be facilitated on a regular basis, with results made public.
- In cases where X's algorithm is found to have amplified content that contributed to human rights harms, X should provide effective remedies for affected communities.

To Regulators:
- Ofcom should formally request comprehensive information from X regarding the spread of illegal and harmful content during the 2024 racist riots, details of algorithmic design and configurations at that time, and details of the specific policy and risk mitigation measures introduced subsequently.
- Ofcom should use this information to rigorously assess whether X now meets its statutory duties under the Online Safety Act, including minimizing illegal content such as incitement to racial hatred. Should current or future breaches be identified, Ofcom should impose appropriate sanctions on X to reinforce compliance expectations.
- Regulators should use their powers to compel greater transparency from X concerning its algorithmic practices, content moderation, and policy adjustments since August 2024. Even though Ofcom cannot sanction or investigate specific actions from 2024, it can legitimately ask for this information to inform its ongoing regulatory oversight and ensure that appropriate safeguards are in place.

To the UK Government:
- The Government should address major gaps in the current online safety regime to hold platforms like X accountable for the broader harms caused by their algorithms. This should include introducing mandatory human rights due diligence legislation which requires social media companies to proactively identify and address risks their services pose to public safety and human rights, especially in high-risk contexts. Where platforms fail to assess and act on foreseeable harms, they should be held accountable.
- The Government should ensure that those harmed by online-amplified violence or hatred have access to justice and support. This includes exploring avenues for remedy and restitution when platform failures contribute to offline harm.

# AMNESTY INTERNATIONAL IS A GLOBAL MOVEMENT FOR HUMAN RIGHTS. WHEN INJUSTICE HAPPENS TO ONE PERSON, IT MATTERS TO US ALL.