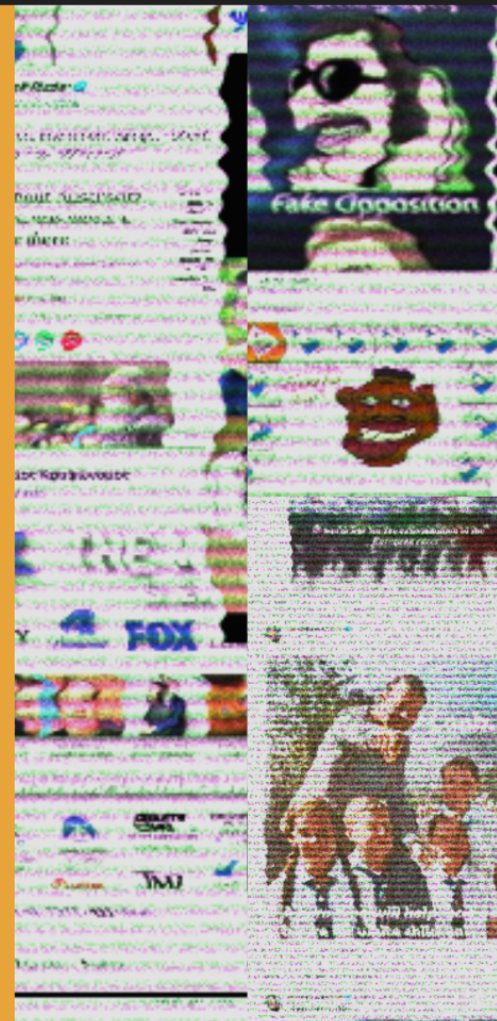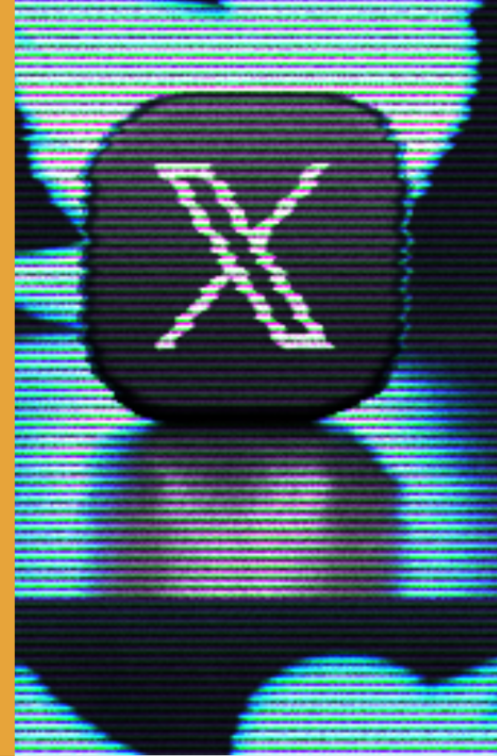## CCDH Center for Countering Digital Hate

# X CONTENT MODERATION FAILURE

## How Twitter/X continues to host posts reported for extreme hate speech

The Center for Countering Digital Hate works to stop the spread of online hate and disinformation through innovative research, public campaigns and policy advocacy.

Our mission is to protect human rights and civil liberties online.

Social media platforms have changed the way we communicate, build and maintain relationships, set social standards, and negotiate and assert our society's values. In the process, they have become safe spaces for the spread of hate, conspiracy theories and disinformation.

Social media companies erode basic human rights and civil liberties by enabling the spread of online hate and disinformation.

At CCDH, we have developed a deep understanding of the online harm landscape, showing how easily hate actors and disinformation spreaders exploit the digital platforms and search engines that promote and profit from their content.

We are fighting for better online spaces that promote truth, democracy, and are safe for all. Our goal is to increase the economic and reputational costs for the platforms that facilitate the spread of hate and disinformation.

**If you appreciate this report, you can donate to CCDH at counterhate.com/donate. In the United States, Center for Countering Digital Hate Inc is a 501(c)(3) charity. In the United Kingdom, Center for Countering Digital Hate Ltd is a non-profit company limited by guarantee.**

**Please Read**

This report includes content on the following themes which may be distressing to readers:

- Violence
- Antisemitism
- Holocaust denial
- Racism

Reader discretion is advised.

**Contents**

**X continues to host 86% of 300 posts reported for extreme hate speech**

Research has found that X, formerly Twitter, continues to host 86% of 300 posts reported for hate speech, despite a commitment to combat content "motivated by hatred, prejudice or intolerance".

Researchers collected a total of 300 posts promoting hate from 100 X accounts, amounting to three posts per account. These posts were categorized as promoting antisemitism, anti-black hatred, other racism, neo-Nazism and white supremacism. Together the accounts identified have a sum total of 1,060,106 followers.

Posts were independently assessed by two researchers to determine whether they represent a violation of X's policies on hate speech, with only those that both agreed were violations included in the final set of 300 posts. Posts were reported to X using its user reporting tools on August 30 and 31, and reviewed for action taken against them one week later on September 7.

CCDH has provided a full database of the posts it selected, and has taken screenshots. They will be provided, on request, to any journalist or researcher enquiring.

One week after reporting, researchers found that X had continued to host 86.33% (259) of the 300 posts and 90.00% (90) of the 100 accounts remained active. X continued to host posts that were:

- Promoting racist caricatures of black people and Jewish people
- Denying the Holocaust or mocking victims of the Holocaust
- Labeling Hitler as "A hero who will help secure a future for white children!"
- Claiming "Blacks don't need provoking before becoming violent. It's in their nature."
- Encouraging users to "Stop Race mixing" and "break up with your non-white gf today"
- Promoting conspiracies that Jews promote mass migration and "control the blacks"
- Claiming Jews exert "victimhood and cult-like behavior" and suffer from "schizophrenia"
- Memes accusing Black people of being harmful to "A quiet functioning society"

X continued to host these posts even after they were reported, despite them clearly violating the platform's policies against hateful content, which prohibit racist slurs, dehumanization, and hateful imagery such as the Nazi swastika.

Researchers received notifications from X that three accounts reported as part of the study had been "locked", stating that "they can't post, repost, or Like content, and we'll ask them to remove the reported content if they want to regain full access to their account." As of September 7, this has resulted in posts that promote Holocaust denial, including one that mocks Anne Frank's status as a victim of the Holocaust, remaining publicly visible on the platform.

Musk recently posted: "I'm pro free speech, but against anti-Semitism of any kind", whilst X's new CEO, Linda Yaccarino, claimed that X is "a much healthier and safer platform than it was a year ago". However, previous research from the CCDH found that the volume of posts containing hate speech surged after Elon Musk's acquisition of Twitter.

**X continues to host antisemitic content including Holocaust denial**

Researchers identified a total of 140 posts that promoted antisemitism, including racist caricatures of Jewish people and claims that Jews control the world. X continued hosting this content in 85.00% (119) of cases.

X continued to host 16 posts that researchers identified as containing Holocaust denial, including posts that mock victims of Holocaust and present death camps as benign.

This comes as X's owner Elon Musk declared his intention to sue the Anti-Defamation League (ADL), America's oldest organization working to counter antisemitism, over their criticism of X for allowing a boom in anti-Jewish hate on the platform.

**Researchers identify 38 ads placed next to hateful posts**

CCDH researchers also identified 38 ads from companies, such as Apple and Disney, that were placed adjacent to hateful posts. The ads appeared either on the "For You" feed or on the profiles of accounts promoting hate speech.

On the point of brand safety on X, Yaccarino said: "since acquisition, we [X] have built brand safety and content moderation tools that have never existed before at this company."

In hopes of winning back advertisers, X recently signed a new brand safety deal with digital ad-tech firm Integral Ad Science and introduced an industry-standard blocklist that claims to "protect advertisers from appearing adjacent to unsafe keywords in the Home Timeline".

**Methodology**

- The full set of posts used in this analysis is available below. **Note that content includes extreme antisemitism, anti-black hatred, white supremacism and neo-Nazism.**

  *X Content Moderation dataset, CCDH, 7 September 2023,*
  *https://docs.google.com/spreadsheets/d/1Uu4yh9YeeYN2uWEBsqwohf33D9Y2Jikfnr0*
  *3_InrRYc/edit?usp=sharing*

- All 300 posts, X accounts and ads included in this analysis can be found in the link above. Posts were reported to X on August 30th and 31st, and reviewed to see what action X had taken at least one week later on September 7, examining both notifications from X on action against reports, and relevant posts or accounts themselves to examine whether they had been removed or otherwise acted upon.

- For each post, researchers recorded whether the following actions had been taken:
  - No action
  - Actioned, noting the action taken
    - Post removed
    - Account suspended
    - Account locked (see below for further explanation)

- Where it was unclear whether a post or account had been removed by X or the user controlling the account, this was counted as "actioned" to ensure that all possible actions taken by X were accounted for.

- Researchers identified hateful accounts and content by examining the content being shared by known hateful accounts and account recommendations presented by X on hateful profiles under the "Who to follow" and "You might like" panels.

- Each post was categorized according to the types of hate it promotes, and then checked by a second researcher to ensure it represented a clear violation of X's policies on hate speech (see below). Researchers stopped collecting posts when they amassed 300 posts from 100 unique accounts that met these criteria.

- X continued to host posts identified by researchers that clearly violate the platform's policies against hateful conduct, which prohibits:
  - "[attacks on] other people on the basis of race, ethnicity, national origin"
  - "dehumanization" on the basis of race, religion, nationality and other qualities

- ○ Hateful imagery such as the Nazi swastika

- The policy also contains a commitment to "combating abuse motivated by hatred, prejudice or intolerance, particularly abuse that seeks to silence the voices of those who have been historically marginalized."

    *Hateful Conduct, X, April 2023,*
    *https://help.twitter.com/en/rules-and-policies/hateful-conduct-policy*

**Background**

- The new CEO of X, Linda Yaccarino, recently claimed in an interview that X is "a much healthier and safer platform than it was a year ago".

    *"X Corp. now a much healthier and safe platform than a year ago says Lina Yaccarino", CNBC, 10 August 2023, 01:11,* [*https://www.cnbc.com/video/2023/08/10/x-corp-now-a-much-healthier-and-safer-platform-than-a-year-ago-says-linda-yaccarino.html*](https://www.cnbc.com/video/2023/08/10/x-corp-now-a-much-healthier-and-safer-platform-than-a-year-ago-says-linda-yaccarino.html)

- However, previous research from the CCDH found that the volume of posts containing hate speech surged after Elon Musk's acquisition of Twitter.

    *The Musk Bump: Quantifying the rise in hate speech under Elon Musk, CCDH, 6 December 2022,* [*https://counterhate.com/blog/the-musk-bump-quantifying-the-rise-in-hate-speech-under-elon-musk/*](https://counterhate.com/blog/the-musk-bump-quantifying-the-rise-in-hate-speech-under-elon-musk/)

- To brand safety on X, Yaccarino said that "since acquisition, we [X] have built brand safety and content moderation tools that have never existed before at this company."
    *"X Corp. now a much healthier and safe platform than a year ago says Lina Yaccarino", CNBC, 10 August 2023, 01:18,* [*https://www.cnbc.com/video/2023/08/10/x-corp-now-a-much-healthier-and-safer-platform-than-a-year-ago-says-linda-yaccarino.html*](https://www.cnbc.com/video/2023/08/10/x-corp-now-a-much-healthier-and-safer-platform-than-a-year-ago-says-linda-yaccarino.html)

- In hopes of winning back advertisers, X recently signed a new brand safety deal with digital ad-tech firm Integral Ad Science and introduced an industry-standard blocklist that alleges to "protect advertisers from appearing adjacent to unsafe keywords in the Home Timeline".

    *A safer X is a better X, X, accessed 1 September 2023,* [*https://business.twitter.com/en/help/ads-policies/brand-safety.html*](https://business.twitter.com/en/help/ads-policies/brand-safety.html)
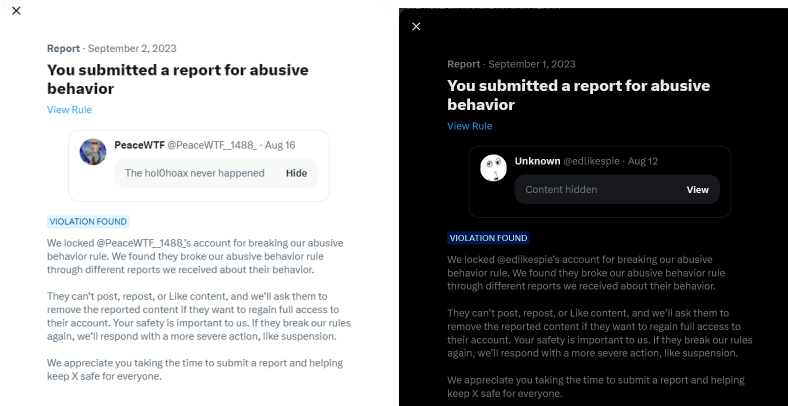
- X's owner Elon Musk has declared his intention to sue the Anti-Defamation League (ADL), America's oldest organization working to counter antisemitism, over their criticism of X for allowing a boom in anti-Jewish hate on the platform.

    *"Elon Musk blames the ADL for 60% ad sales decline at X, threatens to sue", CNN, 5 September 2023,* [*https://www.cnn.com/2023/09/05/tech/elon-musk-adl-lawsuit/index.html*](https://www.cnn.com/2023/09/05/tech/elon-musk-adl-lawsuit/index.html)

**X's decision to "lock" accounts in response to reports**

- Researchers received notifications from X that three accounts reported as part of the study had been "locked", stating that "they can't post, repost, or Like content, and we'll ask them to remove the reported content if they want to regain full access to their account." Note that the date in screenshots is the date the notification was received.



- As of September 7, this has resulted in posts that promote Holocaust denial and mock victims of the Holocaust remaining publicly visible on the platform.



*@P33Hole, X, 14 August 2023,*
*https://drive.google.com/file/d/1EuksNb39yr–J3Mbr6UcENZ9tvsmnOeib/view*
*@oldhead_scott, X, 2 August 2023,*
*https://drive.google.com/open?id=1g7pJ27dkmogS9CiQXGVdo1CFmxempaWj*

**Examples of hateful posts still hosted by X as of September 7**

- Promoting racist caricatures of black people and Jewish people.



*@wayotworld, X, 21 July 2023,*
*https://drive.google.com/open?id=14i6BYE9Vs4Sbp3O2jkciQclHDzeYCCA7*



*IiiLothrop, X, 25 August 2023,*
*https://drive.google.com/file/d/1oysZhmdjFGNYJjjF-HEtOecJqVpUoROk/view*

*Nordicfrost, X, 27 July 2023,*
*https://drive.google.com/file/d/1dFrDfbBG9271CQbDdhqeObgSOwUiQv64/view*

- Promoting conspiracies that Jews promote mass migration and "control the blacks"



*Databased_01, X, 16 August 2023,*
*https://drive.google.com/file/d/1b9EtaNGUdhDxsnZ2AlR4XZo-Tx1XOcPw/view*



*NorthernGirl248, X, 23 May 2023,*
*https://drive.google.com/file/d/1oKNgv4hJLkVccbiYHcDFfmvrvt3dXC4T/view*

1:03 PM · Aug 19, 2023 · **149** Views

- Claiming Jews exert "victimhood and cult-like behavior" and suffer from "Schizophrenia"



*IiiLothrop, X, 30 July 2023,*
*https://drive.google.com/file/d/1KRsu1–5WEgQKog7LEc3YRmpOLKO5YAQ4/view*

- Denying the Holocaust or mocking victims of the Holocaust



*P33Hole, X, 14 August 2023,*
*https://drive.google.com/file/d/1EuksNb39yr–J3Mbr6UcENZ9tvsmnOeib/view*

- Labeling Hitler as "A hero who will help secure a future for white children!"



*WilliamA_33, X, 18 August 2023,*
*https://drive.google.com/file/d/1rHifb5djnJP–JHtg84yHwaNUVStB85so/view*

- Claiming "Blacks don't need provoking before becoming violent. It's in their nature."



*NorthernGirl248, X, 2 March 2023,*
*https://drive.google.com/file/d/1OgPoSyaClPXnY4LYVt3FsdDh_fKgNPBo/view*

- Encouraging users to "Stop Race mixing" and "break up with your non–white gf today"



*Makeeuropasnow, X, 8 August 2023,*
*https://drive.google.com/file/d/12jdR9XVDDVilx5Kk87wBleLpekj57zRC/view*

- Memes accusing Black people of being harmful to "A quiet functioning society"

**Examples of ads placed next to hateful posts**

- The full set of ads used in this analysis is available below.
  *X Content Moderation dataset, CCDH, 7 September 2023,*
  *https://docs.google.com/spreadsheets/d/1Uu4yh9YeeYN2uWEBsqwohf33D9Y2Jikfnr0*
  *3_InrRYc/edit#gid=726502028*

- An ad from Walt Disney World placed below a post mocking Black Americans



*AntiWhiteWatch1, X, 13 August 2023,*
*https://twitter.com/AntiWhiteWatch1/status/1690535047349350400*
*WaltDisneyWorld, X, 21 July 2023,*
*https://twitter.com/WaltDisneyWorld/status/1682492538018750464*

- An ad from Apple placed above a post alluding to Holocaust Denial



*Wayotworld, X, 13 August 2023,*
*https://twitter.com/Apple/status/1641602026856669186*
*Apple, X, 31 March 2023, https://twitter.com/Apple/status/1641602026856669186*

● An ad from Supermicro placed between two posts praising Nazis

21

- This ad for UFC was placed next to a post praising Hitler.



*@I_Loooove_Pizza, X, 16 August 2023,*
*https://twitter.com/I_Loooove_Pizza/status/1691995395302314043*
*@UFC, X, 16 August 2023, https://twitter.com/ufc/status/1691917749914877975*

CCDH

Center for
Countering
Digital Hate