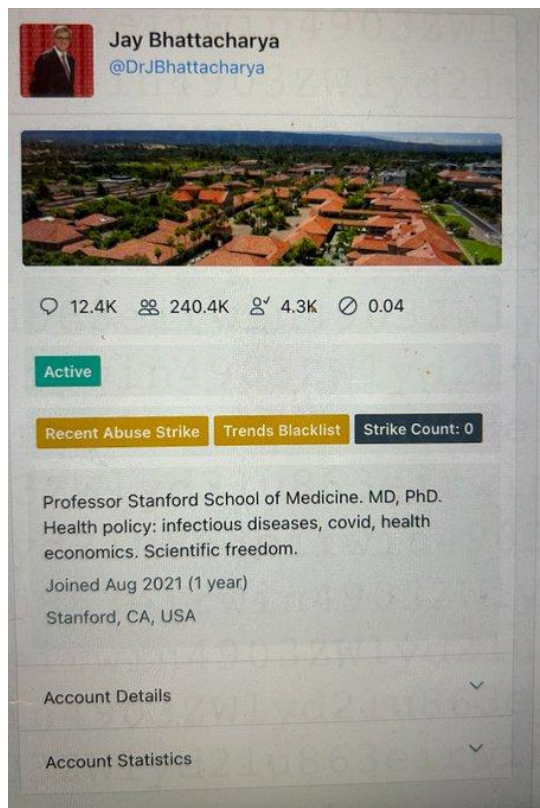


Bari Weiss

@bariweiss

THREAD: THE TWITTER FILES PART TWO. TWITTER'S SECRET BLACKLISTS.

1. A new [#TwitterFiles](#) investigation reveals that teams of Twitter employees build blacklists, prevent disfavored tweets from trending, and actively limit the visibility of entire accounts or even trending topics—all in secret, without informing users.
2. Twitter once had a mission “to give everyone the power to create and share ideas and information instantly, without barriers.” Along the way, barriers nevertheless were erected.
3. Take, for example, Stanford’s Dr. Jay Bhattacharya ([@DrJBhattacharya](#)) who argued that Covid lockdowns would harm children. Twitter secretly placed him on a “Trends Blacklist,” which prevented his tweets from trending.



4. Or consider the popular right-wing talk show host, Dan Bongino (@dbongino), who at one point was slapped with a "Search Blacklist."

The image shows a screenshot of Dan Bongino's Twitter profile page. At the top left is a profile picture of Dan Bongino. To its right, the name "Dan Bongino" is displayed with a blue verification checkmark, and the handle "@dbongino" is below it. A large banner image for "THE DAN BONGINO SHOW" is positioned below the header. Underneath the banner, statistics are shown: 33.1K replies, 3.7M retweets, 1.4K likes, and 0.08 followers. Below these are two status boxes: "Verified" (blue) and "Active" (green). A row of four labels follows: "Notifications Spike" (yellow), "Search Blacklist" (yellow), "Twitter Blue Verified" (dark grey), and "Multiple Accounts" (dark grey). Below that is another row: "Strike Count: 0" (dark grey), "NSFW View" (dark grey), and "SPMA" (dark grey). The bio section contains the text "Public Enemy #1", "Joined Jan 2011 (12 years)", "Planet 'Banned By YouTube'", and the website link "http://www.bongino.com". At the bottom left, the text "Account Details" is visible with a downward arrow icon to its right.

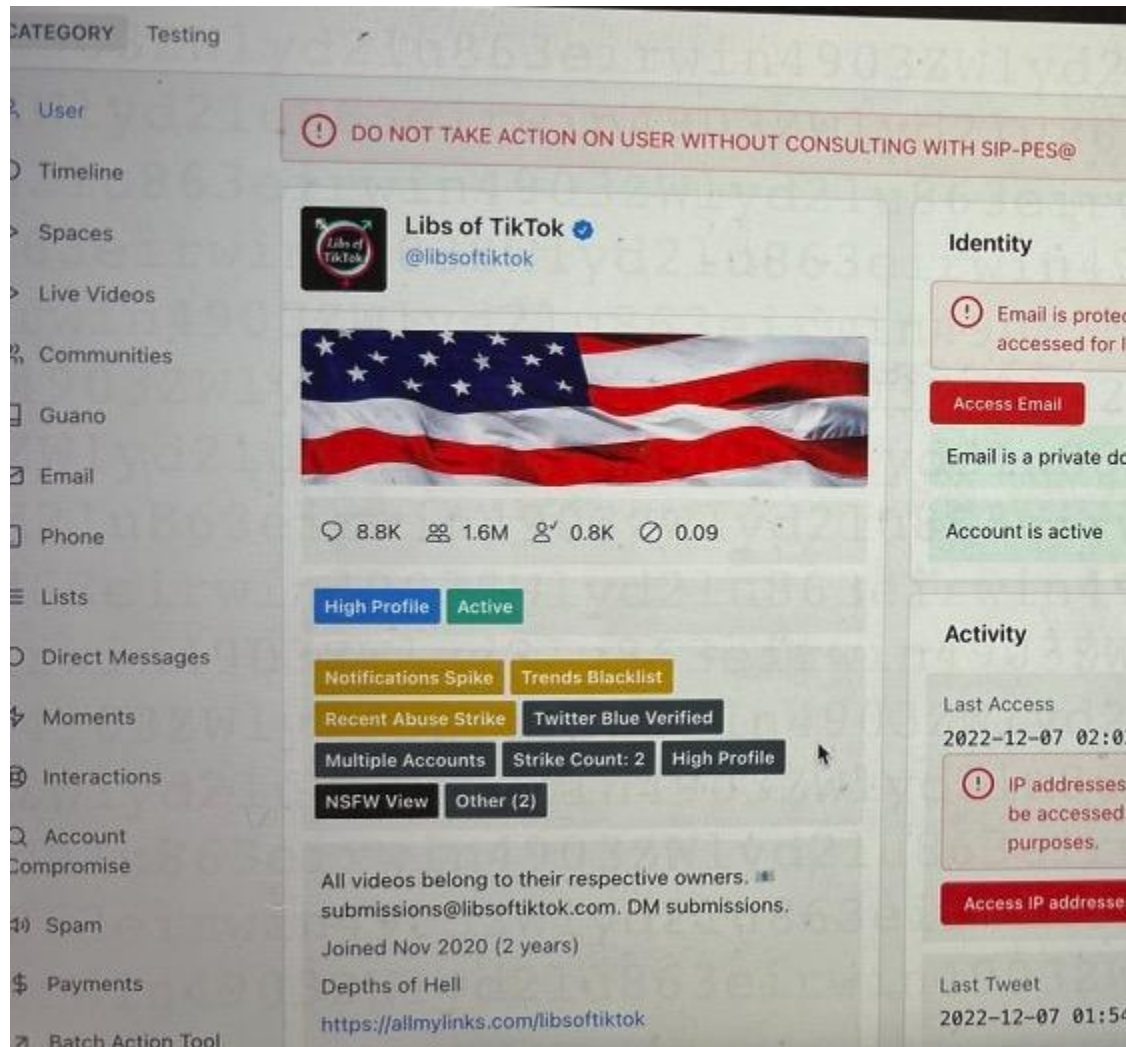
5. Twitter set the account of conservative activist Charlie Kirk (@charliekirk11) to “Do Not Amplify.”

The image shows a screenshot of Charlie Kirk's Twitter profile. At the top left is a profile picture of Charlie Kirk. To its right, the name "Charlie Kirk" is displayed with a blue verification checkmark, followed by the handle "@charliekirk11". Below the name is a large video thumbnail showing a man on a stage with a large audience and pyrotechnics. Underneath the video are statistics: 53K replies, 1.9M retweets, 192.8K likes, and 0.17 mentions. Below the statistics are two status labels: "Verified" (blue) and "Active" (green). A row of four yellow labels follows: "Recent Abuse Strike", "Notifications Spike", "Strike Count: 0", and "Do Not Amplify". Below these are two dark blue labels: "NSFW View" and "Other (2)". At the bottom of the profile, there is a bio: "Founder & President: @TPUSA • Host: The Charlie Kirk Show • Click the link below to subscribe 🇺🇸". Below the bio are the following details: "Joined May 2011 (12 years)", "Phoenix, AZ", and a link: "https://apple.co/2VCxGsh".

6. Twitter denied that it does such things. In 2018, Twitter's Vijaya Gadde (then Head of Legal Policy and Trust) and Kayvon Beykpour (Head of Product) said: “We do not shadow ban.” They added: “And we certainly don't shadow ban based on political viewpoints or ideology.”

7. What many people call “shadow banning,” Twitter executives and employees call “Visibility Filtering” or “VF.” Multiple high-level sources confirmed its meaning.
8. “Think about visibility filtering as being a way for us to suppress what people see to different levels. It’s a very powerful tool,” one senior Twitter employee told us.
9. “VF” refers to Twitter’s control over user visibility. It used VF to block searches of individual users; to limit the scope of a particular tweet’s discoverability; to block select users’ posts from ever appearing on the “trending” page; and from inclusion in hashtag searches.
10. All without users’ knowledge.
11. “We control visibility quite a bit. And we control the amplification of your content quite a bit. And normal people do not know how much we do,” one Twitter engineer told us. Two additional Twitter employees confirmed.
12. The group that decided whether to limit the reach of certain users was the Strategic Response Team - Global Escalation Team, or SRT-GET. It often handled up to 200 “cases” a day.
13. But there existed a level beyond official ticketing, beyond the rank-and-file moderators following the company’s policy on paper. That is the “Site Integrity Policy, Policy Escalation Support,” known as “SIP-PES.”
14. This secret group included Head of Legal, Policy, and Trust (Vijaya Gadde), the Global Head of Trust & Safety (Yoel Roth), subsequent CEOs Jack Dorsey and Parag Agrawal, and others.
15. This is where the biggest, most politically sensitive decisions got made. “Think high follower account, controversial,” another Twitter employee told us. For these “there would be no ticket or anything.”

16. One of the accounts that rose to this level of scrutiny was @libsoftiktok—an account that was on the “Trends Blacklist” and was designated as “Do Not Take Action on User Without Consulting With SIP-PES.”



17. The account—which Chaya Raichik began in November 2020 and now boasts over 1.4 million followers—was subjected to six suspensions in 2022 alone, Raichik says. Each time, Raichik was blocked from posting for as long as a week.

18. Twitter repeatedly informed Raichik that she had been suspended for violating Twitter’s policy against “hateful conduct.”

19. But in an internal SIP-PES memo from October 2022, after her seventh suspension, the committee acknowledged that “LTT has not directly engaged in behavior violative of the Hateful Conduct policy.” See here:

Site Policy Recommendation

Site Policy recommends placing @LibsOfTikTok ([LTT] 1.3M followers, not verified) in a 7-day timeout at the account level [meaning, not for a specific Tweet] based on the account’s continued pattern of indirectly violating Twitter’s Hateful Conduct Policy by tweeting content that either leads to or intends to incite harassment against individuals and institutions that support LGBTQ communities. At this time, Site Policy has not found explicitly violative Tweets, which would result in a permanent suspension of the account.

This type of enforcement action [repeated 7-day timeouts at the account-level] will not lead to permanent suspension, however: should LTT engage in any other direct Tweet-level violations of any of Site Policy’s policies, we will move forward with permanent suspension.

Assessment

Since its most recent timeout, while LTT has not directly engaged in behavior violative of the Hateful Conduct policy, the user has continued targeting individuals/allies/supporters of the LGBTQIA+ community for alleged misconduct. The targeting of at least one of these institutions

20. The committee justified her suspensions internally by claiming her posts encouraged online harassment of “hospitals and medical providers” by insinuating “that gender-affirming healthcare is equivalent to child abuse or grooming.”

21. Compare this to what happened when Raichik herself was doxxed on November 21, 2022. A photo of her home with her address was posted in a tweet that has garnered more than 10,000 likes.

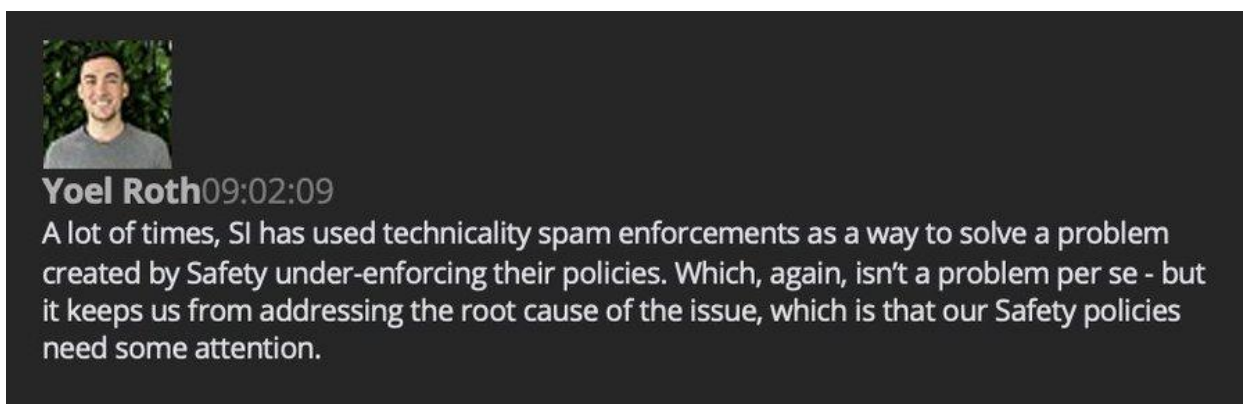
22. When Raichik told Twitter that her address had been disseminated she says Twitter Support responded with this message: “We reviewed the reported content, and didn’t find it to be in violation of the Twitter rules.” No action was taken. The doxxing tweet is still up.

Hello,

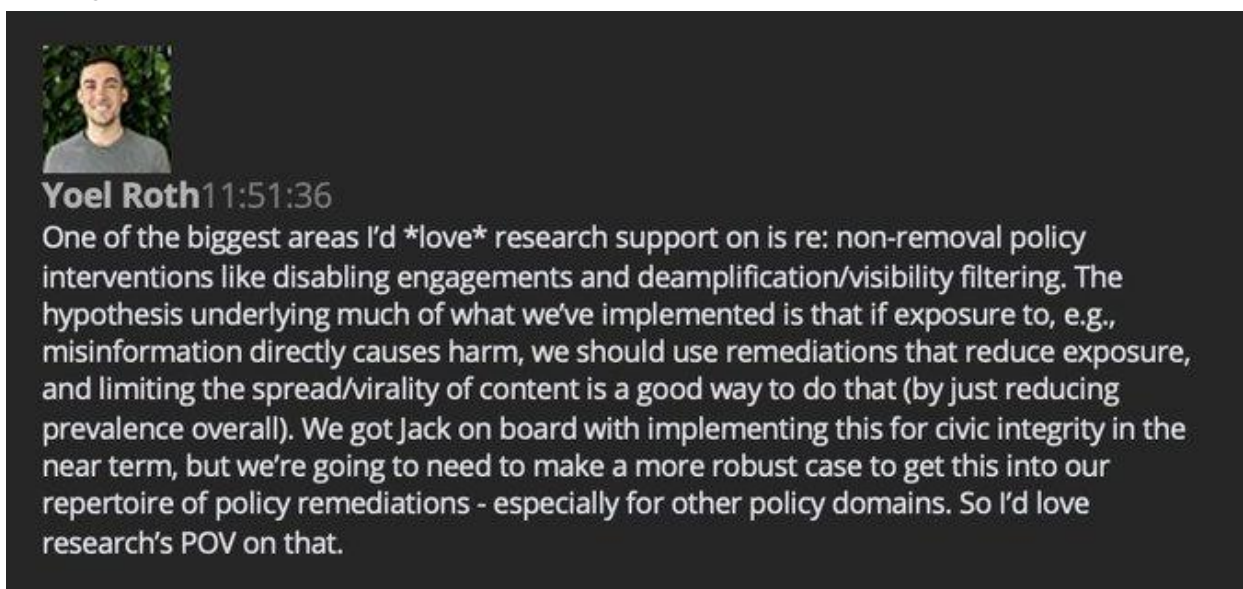
Thanks for reaching out. We reviewed the reported content, and didn’t find it to be in violation of the [Twitter rules](#). In this case, no action will be taken at this time.

If you have further concerns about intellectual property, your privacy, or your personal safety, the following guidelines can be of assistance:

23. In internal Slack messages, Twitter employees spoke of using technicalities to restrict the visibility of tweets and subjects. Here's Yoel Roth, Twitter's then Global Head of Trust & Safety, in a direct message to a colleague in early 2021:



24. Six days later, in a direct message with an employee on the Health, Misinformation, Privacy, and Identity research team, Roth requested more research to support expanding “non-removal policy interventions like disabling engagements and deamplification/visibility filtering.”



25. Roth wrote: “The hypothesis underlying much of what we've implemented is that if exposure to, e.g., misinformation directly causes harm, we should use remediations that reduce exposure, and limiting the spread/virality of content is a good way to do that.”

26. He added: “We got Jack on board with implementing this for civic integrity in the near term, but we’re going to need to make a more robust case to get this into our repertoire of policy remediations – especially for other policy domains.”

27. There is more to come on this story, which was reported by [@abigailshrier](#) [@shellenbergermd](#) [@nelliebowles](#) [@isaacgrafstein](#) and the team The Free Press [@thefp](#).

Keep up with this unfolding story here and at our brand new website: [thefp.com](#).