

SoK: Hate, Harassment, and the Changing Landscape of Online Abuse

Kurt Thomas, Devdatta Akhawe, Michael Bailey, Dan Boneh, Elie Burszstein, Sunny Consolvo, Nicola Dell, Zakir Durumeric, Patrick Gage Kelley, **Deepak Kumar**, Damon McCoy, Sarah Meiklejohn, Thomas Ristenpart, Gianluca Stringhini

Content warning: Potentially triggering language and difficult subject material ahead.

What does online hate and harassment look like?

A Timeline of Leslie Jones's Horrific Online Abuse

By Anna Silman



Leslie Jones Photo: Owen Kolasinski/BFA.com

Coordinated campaigns of **toxic comments** on social media that attempt to silence voices.

Falsely reporting targets to authorities or platforms to take action against their person or accounts.

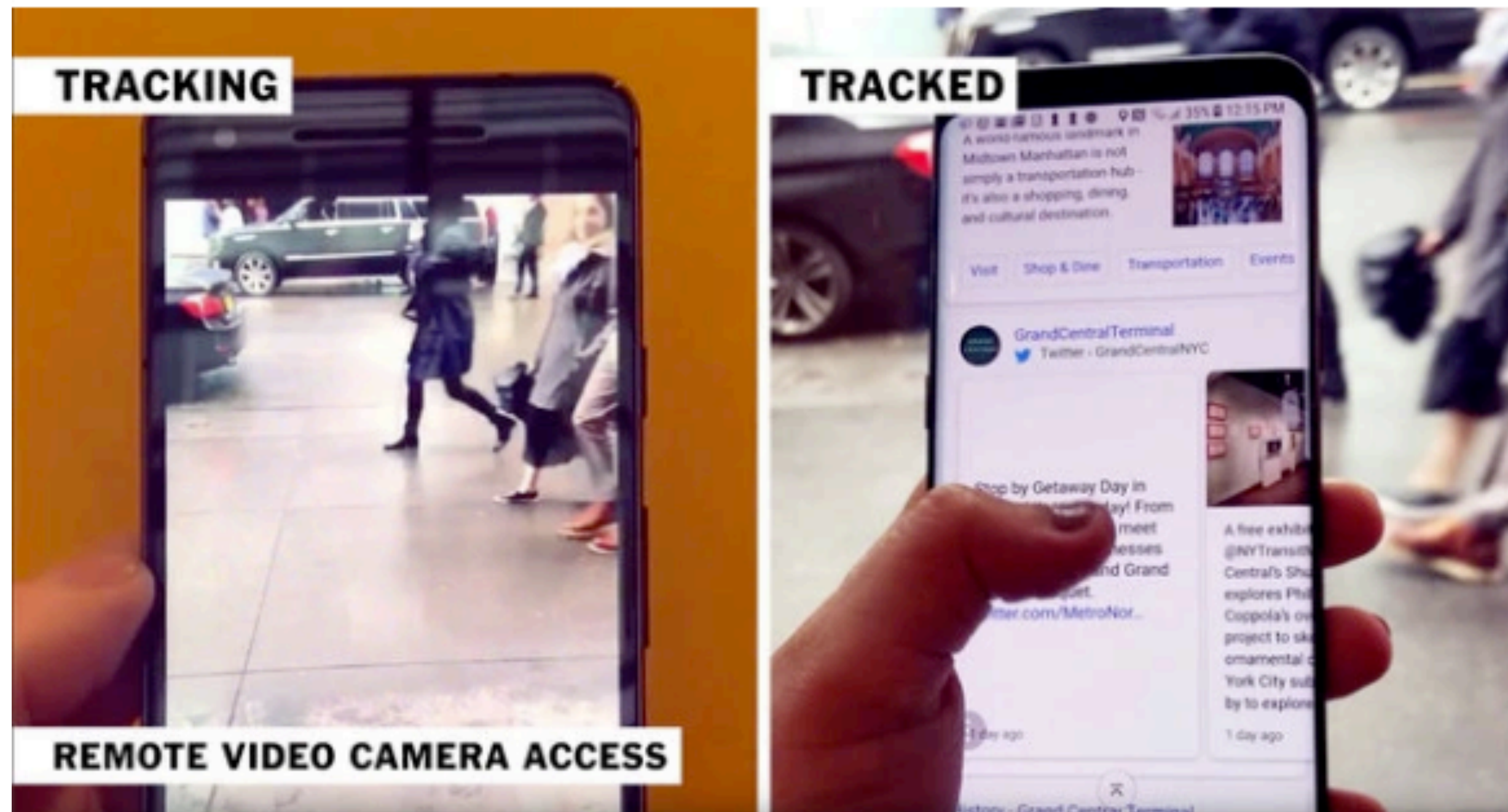
Twitch Streamer Nate Hill Swatted While Streaming Fortnite

A swatting incident is a terrifying event for all involved, which is why fans were concerned when streamer Nate Hill had to cut his stream suddenly.

BY MICHAEL LEE
PUBLISHED FEB 24, 2021



Hundreds of Apps Can Empower Stalkers to Track Their Victims



More than 200 apps and services offer would-be stalkers a variety of electronic capabilities, including basic location tracking, harvesting texts and secretly recording video. Drew Jordan/The New York Times

By Jennifer Valentino-DeVries

Spyware and tracking can aid in **surveilling** intimate partners through their devices and accounts.

Intent is to **inflict emotional harm,**
includes coercive control or instilling a
fear of sexual or physical violence.

Not just high profile targets



41% of people in US



40% of people globally



Source: PEW Research Center Online Harassment 2021, Microsoft Digital Civility Index

We should address online hate and harassment as a security problem.

Literature Review

- Examined the last five years of research and journalism on online hate and harassment
 - IEEE S&P, USENIX Security, CCS, CHI, CSCW, ICWSM, WWW, SOUPS, and IMC
 - Used related papers as a “seed set”, manually searched through related works, and expanded search to include findings from social sciences
 - Also included major news events (e.g., Gamergate) and related attacks and news coverage
- Reviewed over **150 news articles and research papers** in online hate and harassment

Threat Model: Targets and Attackers

Targets of harassment can be individuals or at-risk groups (e.g., LGBTQ+ people)

An attacker's main goal is to emotionally harm or coercively control the target.

Spouse,
family, peers

Anonymous
Internet user

Public figure,
media personality

Anonymous
mob



Types of Attackers

Differentiating Attacks

Research team synthesized criteria that differentiate attacks, falling into **three broad categories – Audience, Medium, Capabilities**

Category	Criteria
Audience	Intended to be seen by the target?
Audience	Intended to be seen by an audience?
Medium	Does attack use media, such as text or images?
Capabilities	Require deception of the audience?
Capabilities	Deception of a third-party authority?
Capabilities	Amplification?
Capabilities	Privileged access to information?

Differentiating Attacks – Audience

Category	Criteria	Examples
Audience	Intended to be seen by the target?	Bullying, Trolling
Audience	Intended to be seen by an audience?	Doxxing
Medium	Does attack use media, such as text or images?	Hate Speech
Capabilities	Require deception of the audience?	Impersonated profiles, Deepfakes
Capabilities	Deception of a third-party authority?	SWATing
Capabilities	Amplification?	Raiding, Dogpiling
Capabilities	Privileged access to information?	IPS, GPS monitoring

Differentiating Attacks – Medium

Category	Criteria	Examples
Audience	Intended to be seen by the target?	Bullying, Trolling
Audience	Intended to be seen by an audience?	Doxxing
Medium	Does attack use media, such as text or images?	Hate Speech
Capabilities	Require deception of the audience?	Impersonated profiles, Deepfakes
Capabilities	Deception of a third-party authority?	SWATing
Capabilities	Amplification?	Raiding, Dogpiling
Capabilities	Privileged access to information?	IPS, GPS monitoring

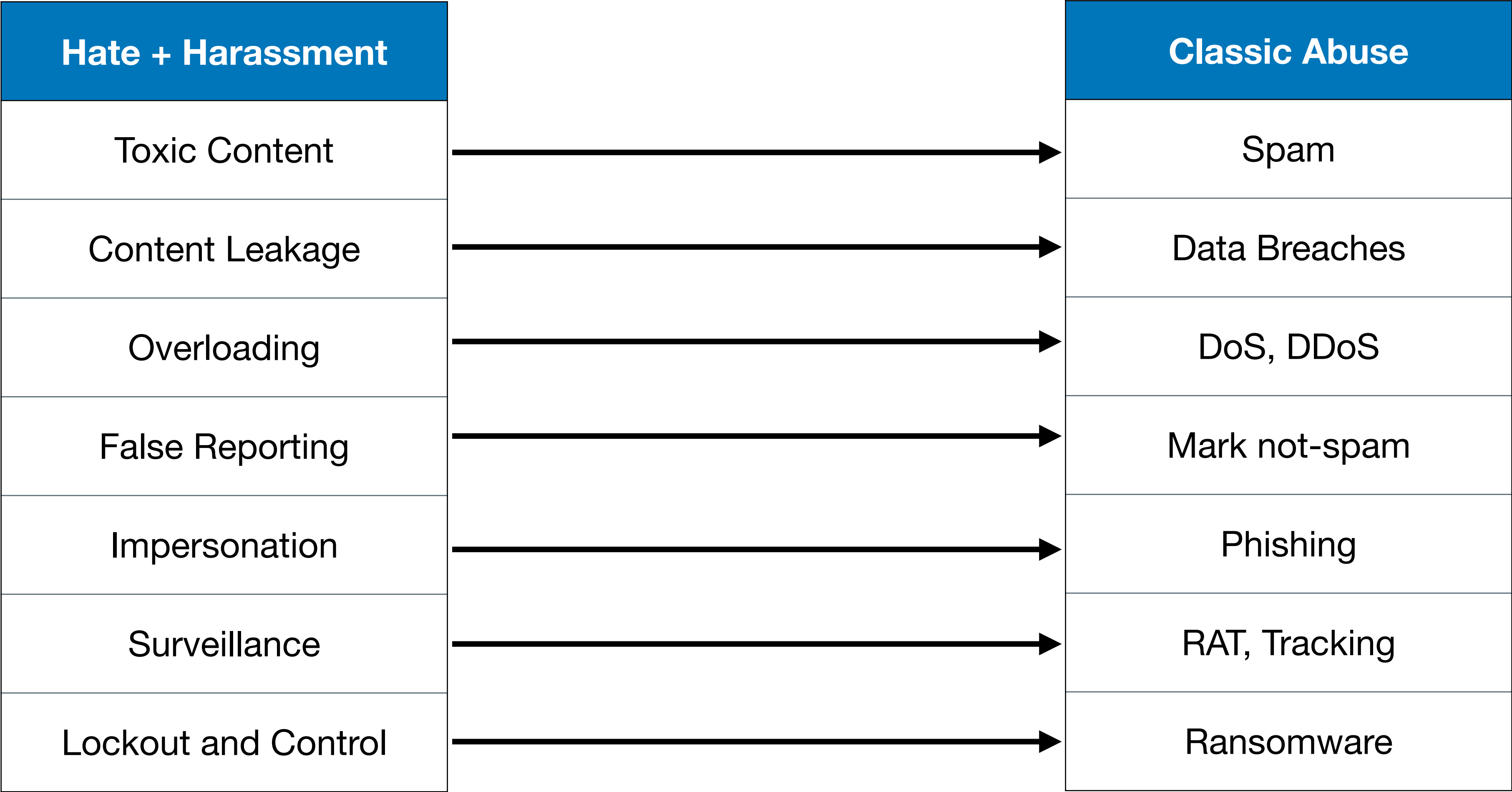
Differentiating Attacks – Capabilities

Category	Criteria	Examples
Audience	Intended to be seen by the target?	Bullying, Trolling
Audience	Intended to be seen by an audience?	Doxxing
Medium	Does attack use media, such as text or images?	Hate Speech
Capabilities	Require deception of the audience?	Impersonated profiles, Deepfakes
Capabilities	Deception of a third-party authority?	SWATing
Capabilities	Amplification?	Raiding, Dogpiling
Capabilities	Privileged access to information?	IPS, GPS monitoring

Seven Classes of Hate and Harassment

Attack Type	Security Principle
Toxic Content	Availability
Content Leakage	Confidentiality
Overloading	Availability
False Reporting	Integrity
Impersonation	Integrity
Surveillance	Confidentiality
Lockout and Control	Integrity, Availability

Parallels to Security Attacks



Parallels to Security Attacks

Hate + Harassment		Classic Abuse
Toxic Content	→	Spam
Content Leakage	→	Data Breaches
Overloading	→	DoS, DDoS
False Reporting	→	Mark not-spam
Impersonation	→	Phishing
Surveillance	→	RAT, Tracking
Lockout and Control	→	Ransomware

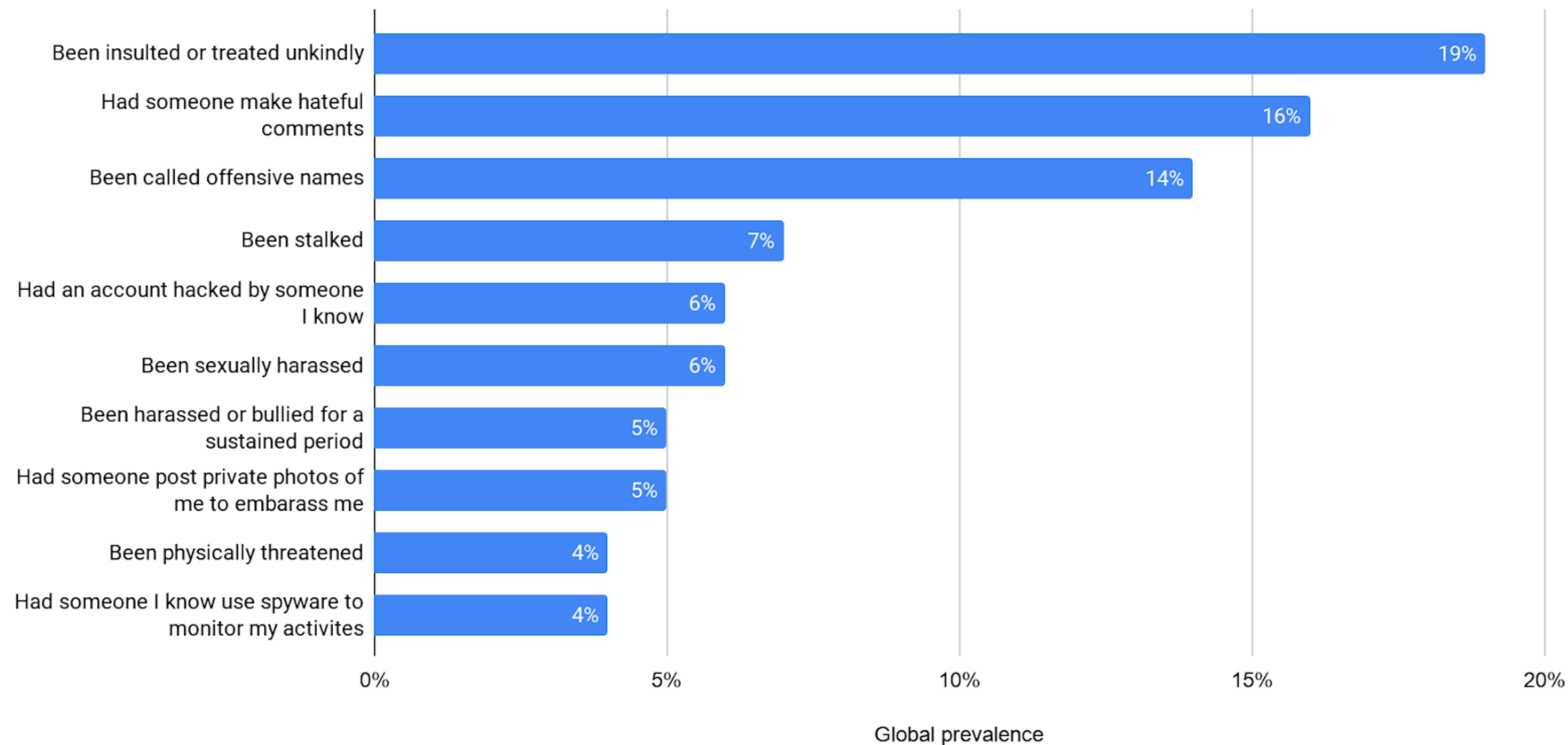
There is no single solution to address the diverse set of hate and harassment attacks.

**What are the lived experiences of
Internet users?**

Survey Instrument

- Surveyed ~1000 participants from 22 countries around the world for three years and asked about hate and harassment experiences
 - Survey was translated for countries that do not primarily speak English
 - Some countries do not appear for all three years to maximize unique countries
- Asked participants “Have you ever personally experienced any of the following online?”
 - Asked about hate and harassment experiences documented in prior work
 - Collected demographic data (e.g., gender, LGBTQ+ status, age, social media usage)

Breakdown of Harassment Experiences

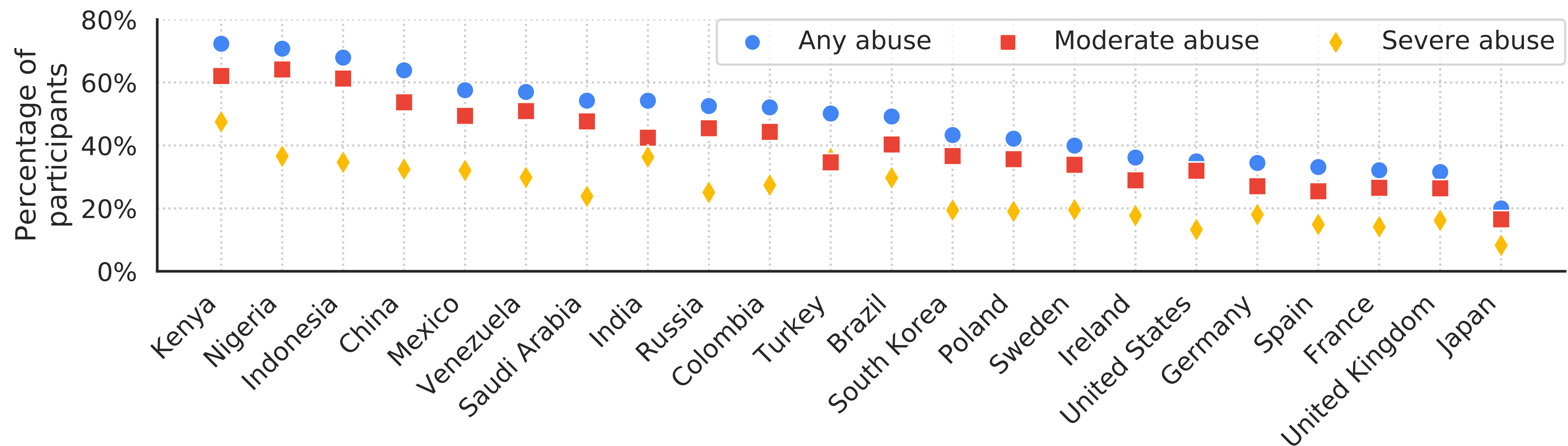


Breakdown of Harassment Experiences



Toxic content is one of the largest threats Internet users face.

Prevalence of Online Hate and Harassment



Measuring hate and harasssment outcomes

- Modeled experiencing any form of hate and harasssment as a binomial distribution
- Input variables are categorical demographic data

Demographic	Treatment	Reference	Odds
LGBTQ+	LGBTQ+	non-LGBTQ+	1.9x
Social Media Usage	Daily	Never	2.5x
	Weekly	Never	2.3x
Age	18 – 24	65 and up	4.0x
	25 – 34	65 and up	3.4x
Year	2017	2016	1.2x
	2018	2016	1.3x

Measuring hate and harassment outcomes

- Modeled experiencing any form of hate and harassment as a binomial distribution
 - Input variables are categorical demographic data
- Odds of experiencing online hate and harassment has increased over time!

Demographic	Treatment	Reference	Odds
LGBTQ+	LGBTQ+	non-LGBTQ+	1.9x
Social Media Usage	Daily	Never	2.5x
	Weekly	Never	2.3x
Age	18 – 24	65 and up	4.0x
	25 – 34	65 and up	3.4x
Year	2017	2016	1.2x
	2018	2016	1.3x

Measuring hate and harassment outcomes

- Modeled experiencing any form of hate and harassment as a binomial distribution
 - Input variables are categorical demographic data
- Odds of experiencing online hate and harassment has increased over time!
- Participants from *minority* groups experience more online hate and harassment

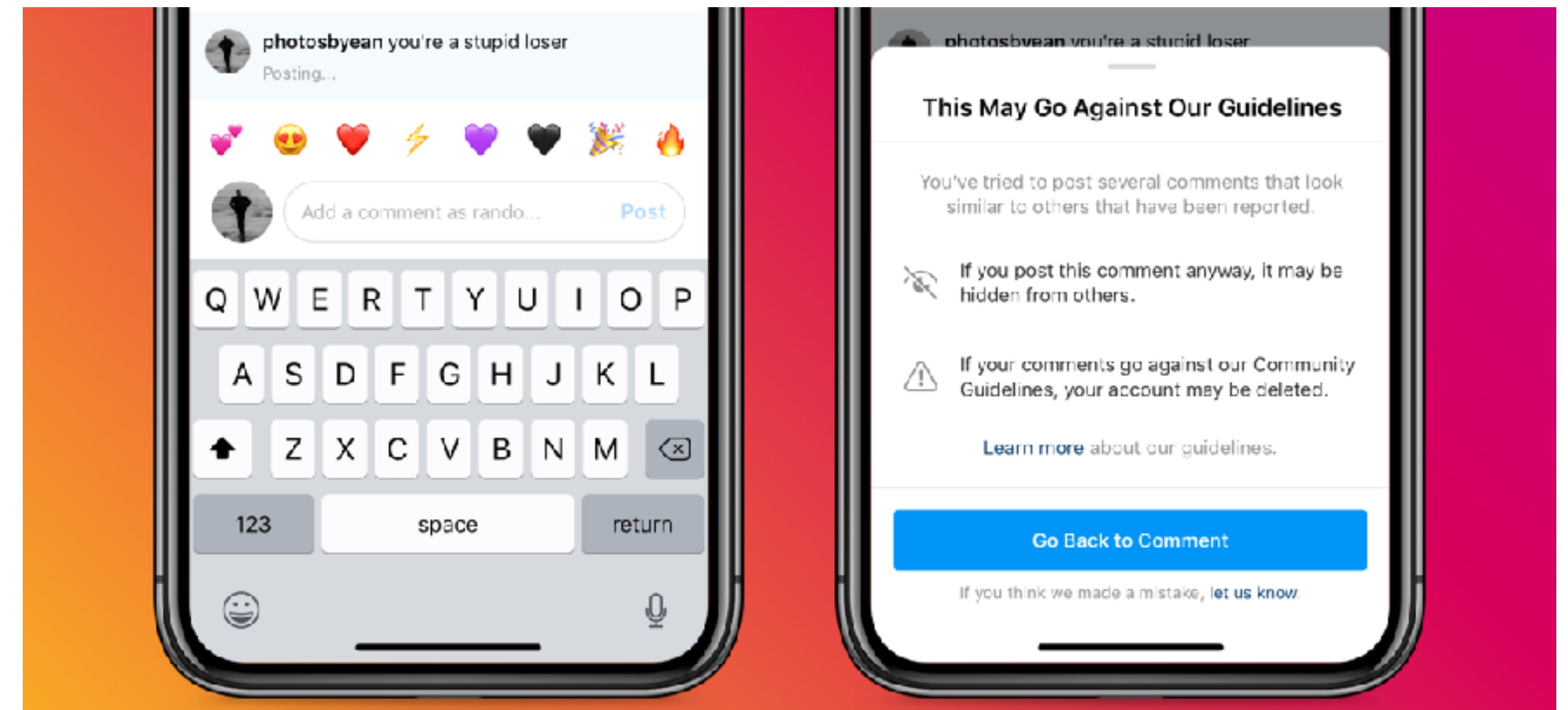
Demographic	Treatment	Reference	Odds
LGBTQ+	LGBTQ+	non-LGBTQ+	1.9x
Social Media Usage	Daily	Never	2.5x
	Weekly	Never	2.3x
Age	18 – 24	65 and up	4.0x
	25 – 34	65 and up	3.4x
Year	2017	2016	1.2x
	2018	2016	1.3x

Designing hate and harassment defenses
must take into account diverse online
experiences.

What can we do about it?

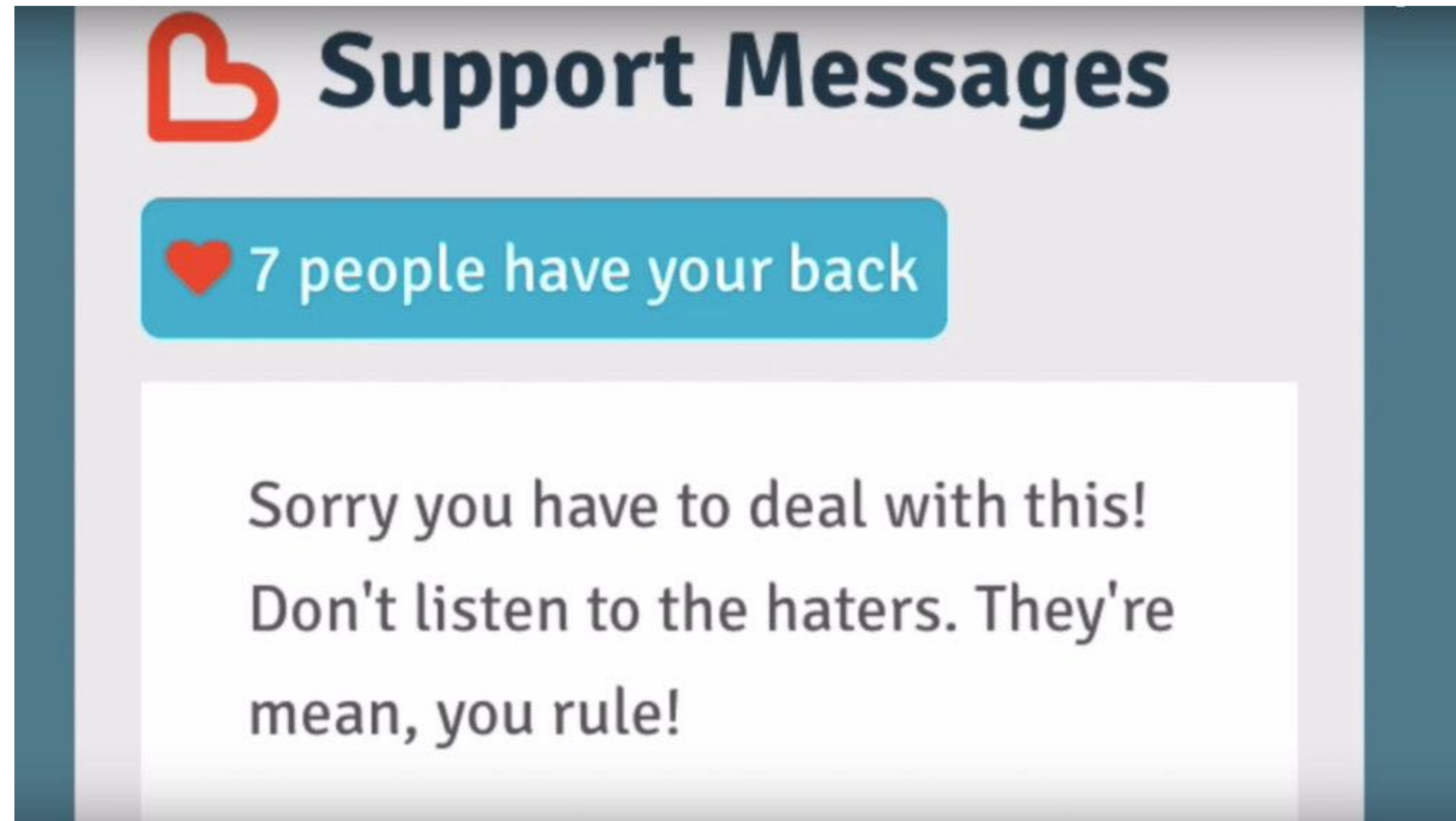
Towards Solutions and Interventions

- **Nudges, indicators, warnings**
- Human moderation, review, and delisting
- Automated detection
- Conscious design
- Policies, education, awareness



Towards Solutions and Interventions

- Nudges, indicators, warnings
- **Human moderation, review, and delisting**
- Automated detection
- Conscious design
- Policies, education, awareness



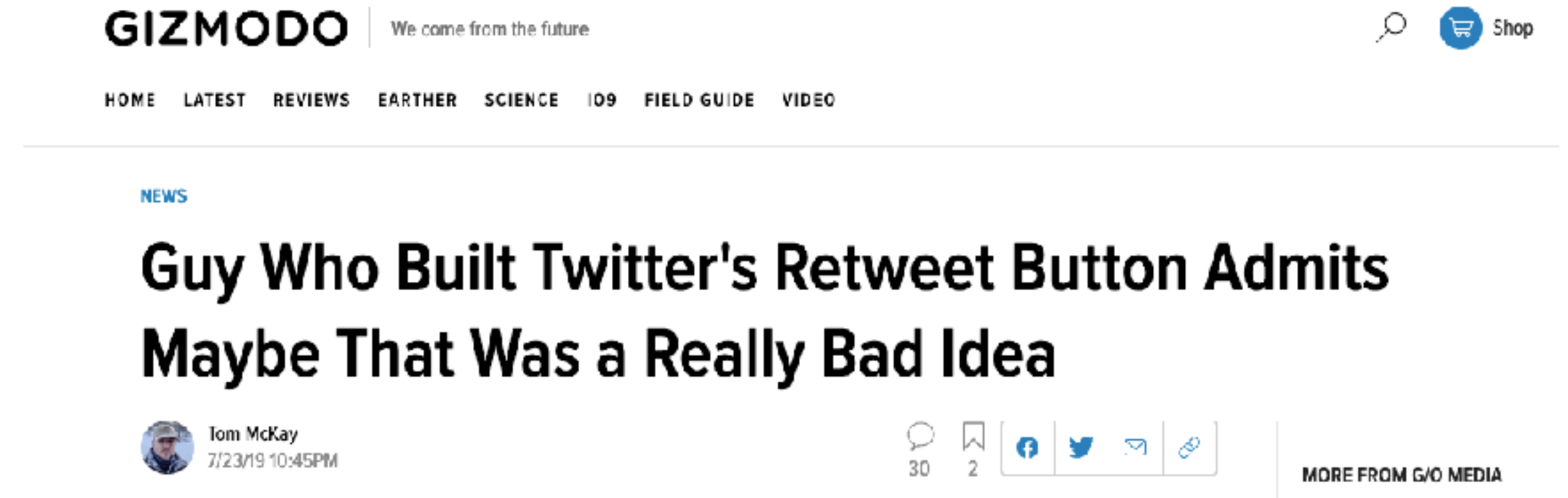
Towards Solutions and Interventions

- Nudges, indicators, warnings
- Human moderation, review, and delisting
- **Automated detection**
- Conscious design
- Policies, education, awareness



Towards Solutions and Interventions

- Nudges, indicators, warnings
- Human moderation, review, and delisting
- Automated detection
- **Conscious design**
- Policies, education, awareness



Towards Solutions and Interventions

- Nudges, indicators, warnings
- Human moderation, review, and delisting
- Automated detection
- Conscious design
- **Policies, education, awareness**

Here Are Twitter's Latest Rules for Fighting Hate and Abuse

Memo outlines steps Twitter plans to control hate and abuse on the service, including expanded definitions of nudity and more enforcement.

Twitter Has Made Changes To Its Policies Of Hateful Conduct In Order To Protect Its Users From Dehumanization

Twitch updates its hateful content and harassment policy after company called out for its own abuses

Facebook, in a reversal, will now ban Holocaust denial content under its hate-speech policy

Key Takeaways

- Online abuse is *changing*, the security community can and should work towards tackling the problem
- Online hate and harassment is growing over time and especially dangerous to Internet subcommunities
- Many techniques for defenses in research are already well studied in the security community, can draw on these for future research

Key Takeaways

- Online abuse is *changing*, the security community can and should work towards tackling the problem
- Online hate and harassment is growing over time and especially dangerous to Internet subcommunities
- Many techniques for defenses in research are already well studied in the security community, can draw on these for future research

Deepak Kumar

kumarde@cs.stanford.edu

@_kumarde