

# 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



# Online Hate and Harassment is Ubiquitous



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

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

**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, Web, 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

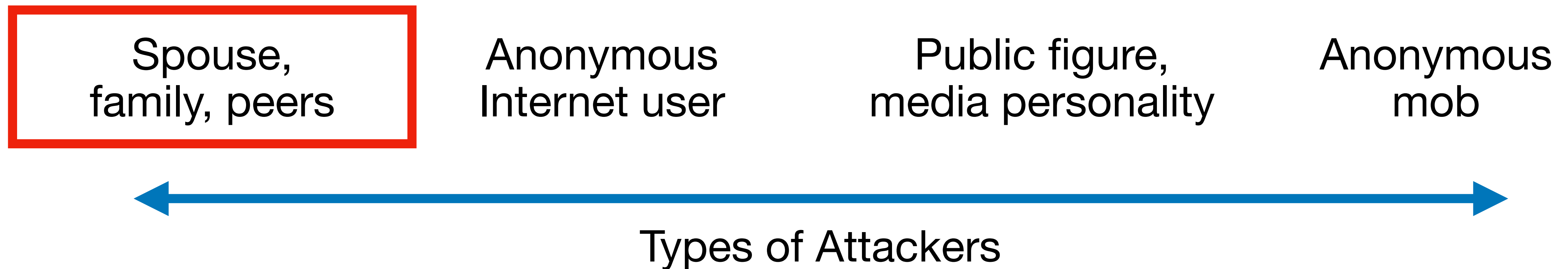


Types of Attackers

# Threat Model: Targets and Attackers

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

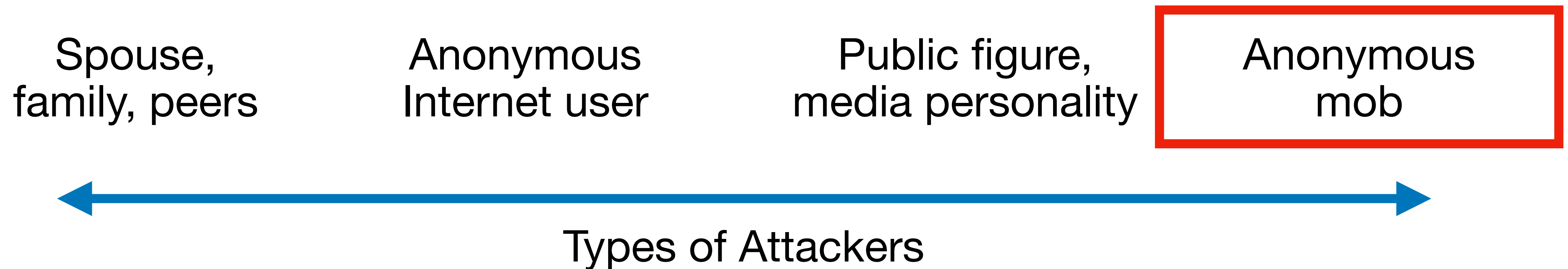
*An attacker's main goal is to emotionally harm or coercively control the target.*



# 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.*



# 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
<b>Medium</b>	<b>Does attack use media, such as text or images?</b>	<b>Hate Speech</b>
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 Online 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

# Seven Classes of Online Hate and Harassment

Attack Type	Security Principle		Classic Abuse
Toxic Content	Availability	→	Spam
Content Leakage	Confidentiality	→	Data Breaches
Overloading	Availability	→	DoS, DDoS
False Reporting	Integrity	→	Mark not-spam
Impersonation	Integrity	→	Phishing
Surveillance	Confidentiality	→	RAT, Tracking
Lockout and Control	Integrity, Availability	→	Ransomware



# Seven Classes of Online Hate and Harassment

Attack Type	Security Principle		Classic Abuse
Toxic Content	Availability	→	Spam
Content Leakage	Confidentiality	→	Data Breaches
Overloading	Availability	→	DoS, DDoS
False Reporting	Integrity	→	Mark not-spam
Impersonation	Integrity	→	Phishing
Surveillance	Confidentiality	→	RAT, Tracking
Lockout and Control	Integrity, Availability	→	Ransomware

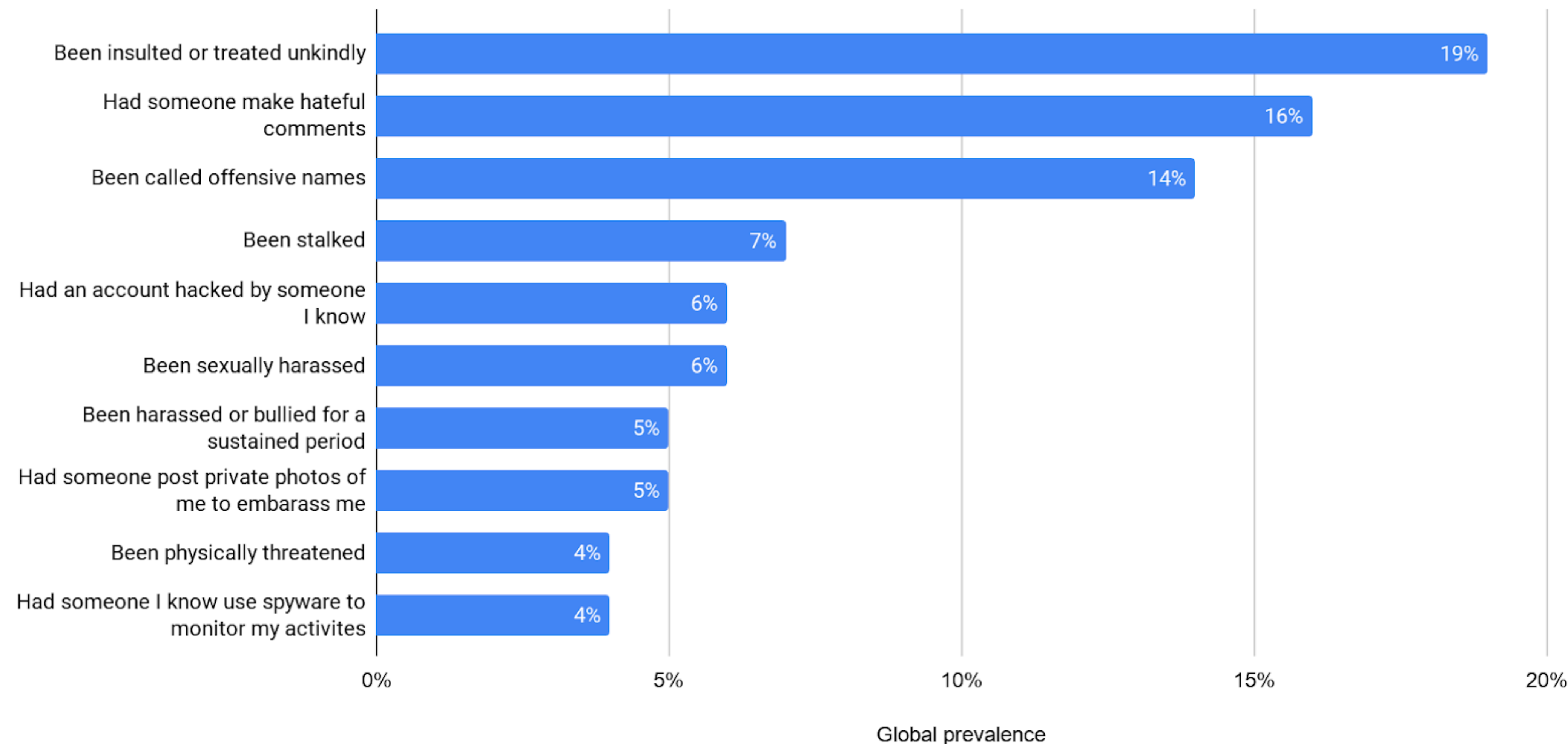
**There is no single solution to address the diverse set of hate and harassment attacks.**

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

# Survey Instrument

- Surveyed ~1000 participants from 22 countries each 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



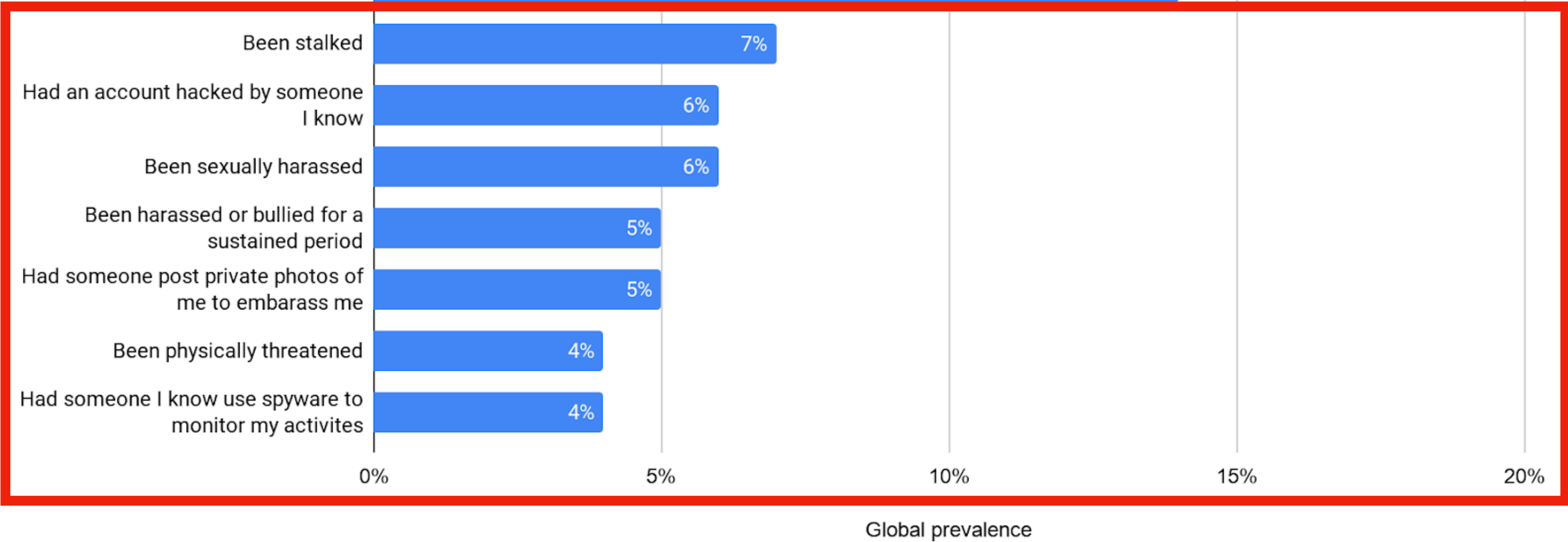


# Breakdown of Harassment Experiences

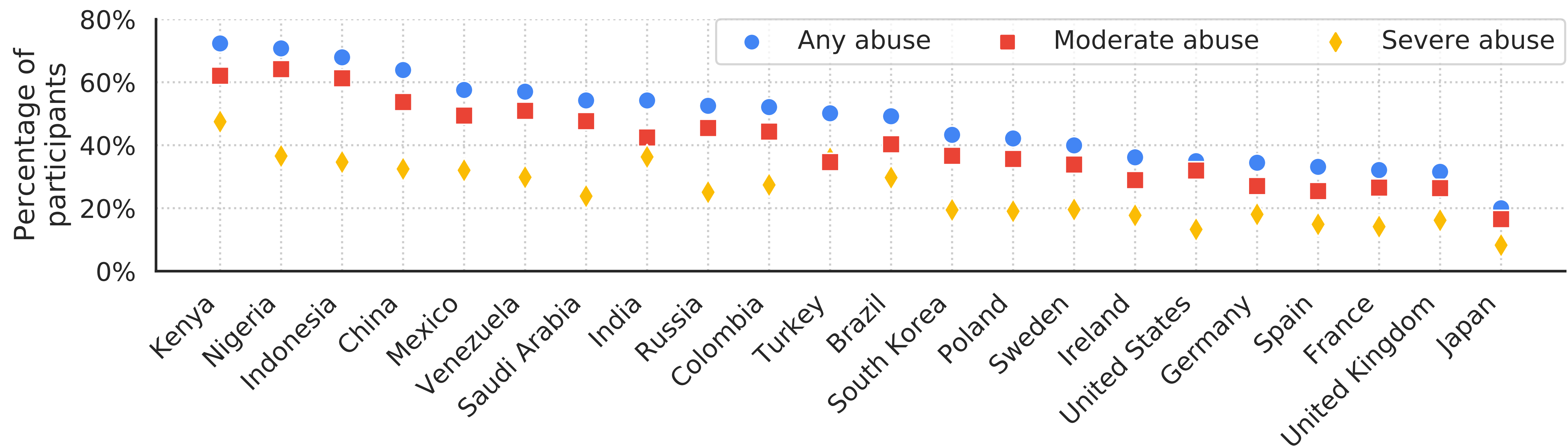


**Toxic content is one of the largest threats Internet users face.**

# Breakdown of Harassment Experiences



# Prevalence of Online Hate and Harassment



# Measuring hate and harassment outcomes

- Modeled experiencing any form of hate and harassment 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
- 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



# Measuring hate and harasssment outcomes

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

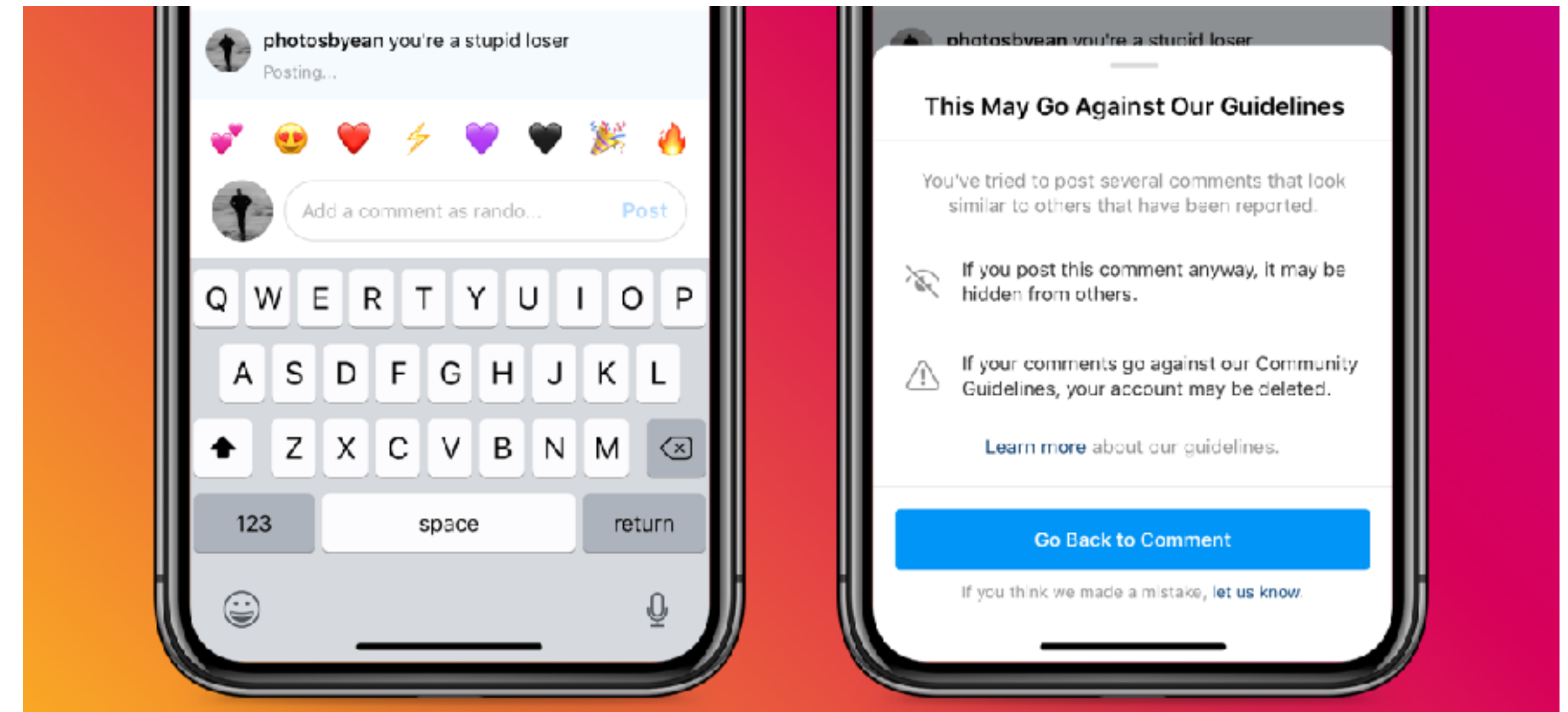
Demographic	Treatment	Reference	Odds
LGBTQ+	LGBTQ+	non-LGBTQ+	1.9x
Social Media Usage	Daily	Never	2.5x
	Weekly	Never	2.3x
Age	18 – 24	65 and up	4.0x
	25 – 34	65 and up	3.4x
Year	2017	2016	1.2x
	2018	2016	1.3x

**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



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

# Tensions and Challenges

- How do we empower targets of abuse instead of burdening them with choice?
- How do you balance moderation with filter bubbles and free speech?
- How do we enable both privacy and accountability?

## THE TRAUMA FLOOR

*The secret lives of Facebook moderators in America*

**TikTok Admits It Suppressed Videos  
by Disabled, Queer, and Fat  
Creators**

# 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 some Internet users
- Many techniques and defenses 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 some Internet users
- Many techniques and defenses are already well studied in the security community, can draw on these for future research

Deepak Kumar

[kumarde@cs.stanford.edu](mailto:kumarde@cs.stanford.edu)

@\_kumarde