# Society Against Violence and Abuse

A set of proposed ideas for making Instagram a safer platform

# Critical Research

Orlikowski and Baroudi classify research as critical where a critical stance is taken toward taken-for-granted assumptions about organizations and information systems, and where the aim is to critique to status quo "through the exposure of what are believed to be deep-seated, structural contradictions within social systems"

Klien, Hienz, and Michael Myers. A SET OF PRINCIPLES FOR CONDUCTING CRITICAL RESEARCH IN INFORMATION SYSTEMS. MIS Quarterly

# Globally 1 out of 3 women have experienced gender based violence or abuse in their lifetime.

This form of gender based hate crime has a very long history, and has taken numerous forms with varying intensity over the years.
Yet even today, less than 40% women seek help of any sort, less than 10% report to the police and choose to disregard or deal with the issue themselves.
We have come to a point where there are multiple organisations and more help around us than ever before, but VAWG ( Violence and Abuse against women and girls ) is still rampant.
Hence it becomes necessary to understand how and where we are failing women and girls, and what are the challenges that abide.

*We conducted critical research on the following two subjects, and understood what works and what doesn't in both approaches towards the issue*

# MAVA

Men Against Violence and Abuse ( MAVA ), is an organisation that engages men and boys instead of women and girls, in conversations about gender inequality and male dominance in our society, and also how it affects them at a personal level. They view men as equal partners and stakeholders, not just as perpetrators in this issue of VAWG.

# TOXIC TWITTER

Amnesty International is a group of human rights activists, that took a different turn on twitter's toxicity. They used data science and crowd sourcing to conduct research and study on abusive tweets that women recieve every 30 seconds. They found patterns and fatual data from the masses of tweets, that twitter needed to address and hasn't. They held the platform responsible for turning it into a battlefeild for women.

Next →

# MAVA

Some of our insights and critique on MAVA are :

## Archaeology and Geneology of Knowledge

Traditional efforts to tackle gender-based violence against women have concentrated on empowering women to assert themselves.

## Cultural Capital

Effects of mainstream media (films, ads and internet) and also stories of gods in religion like hinduism, glorifies male dominance, violence and gender norms. This has a direct impact to the way men are brought up

## Discipline

Using the existing power structures as an advantage to enable power sharing among the oppressed

MAVA primarily focuses on adolescent men and not men in power that hold a bigger percentage for social reproduction

## Habitus

Socially conditioned by the ideals of masculinity also makes men significantly blind towards the contributions of women

MAVA wants oppressors to be a part of conversations and movements where the oppressed are usually the only ones fighting at the forefront, the majority often does not take onus of protecting the minority.

## Economic Capital

Women's vulnerability to male dominance has been long associated and impacted by their economic disempowerment.

This has a direct impact on women being unable to accept positive changes in men

# TOXIC TWITTER

Some of our insights and critique on Toxic Twitter campaign are :

## Archaeology and Geneology of Knowledge

V and A against women on Twitter is not new, its simply any extension of existing and systematic discrimination against women that has found its way into the digital sphere.

## Cultural Capital

Instead of women using their voices to 'impact the world', many are instead being pushed backwards to a culture of silence

India's linguistic diversity, local slang makes it difficult for Twitter's reporting system to analyse and block reported content

## Discipline

Amnesty, as a Human Rights group calls out powerful groups by investigating and exposing facts. It uses media as a tool to unite people from everywhere and seek for justice.

## Habitus

Lack of real world consequences tend to bring out the worst in people

Amnesty focused on crowdsourcing information, primary and secondary research and data analysis compared to individual emancipation/strategic action.

## Economic Capital

Twitter also has workplace effects for many women. Due to constant abuse, their work is also unfairly hindered thus affecting their economic capital

For female journalists, being active on Twitter or any other social media platform is quite necessary, yet they're forced to prepare themselves for any kind of backlash not just to do their job but to fight for their space as women on on the internet.

# Objective

Our main objective with the project was to make the internet a safe space for all genders by attempting to curb online abuse. But we narrowed it down to :

**"Making social media platforms safer safer by attempting to curb online hate and abuse"**

# Social Media Reporting Systems

## Facts and Figures :

From the 40 responses we recieved on our survey,
- 92.5% have reported gender based abuse online
- 85% have reported on Instagram, as compared to other social media platforms
- 55% recieved an uncertain response to their report, and 22% recieved no action at all
- 50% people didn't choose to report because they didn't believe it would make an impact

Link to the form : https://forms.gle/d1byqw5i23eR55nW6

**Instagram was chosen as the working space for the project**

# Ideation

These are some of our ideas from our SCAMPER brainstorming session

Substitute the action of blocking the offender with a learning process or a tedious game that involves a learning process to get back into instagram

**S**

Combine online reporting and police reporting systems to ensure real world consequences

**C**

Verify each account for Instagram at the time of creating one. Ask users to update and verify their account periodically. (like tinder/bumble)

**A**

Minify the engagement that abuser gets on their account when people encourage others to report them.

**M**

Using reporting system in emergency situations (SOS) that could perhaps help save someone in a real world scenario

**P**

Eliminate the process of a user having to choose to view sensitive content without enough context

**E**

Using reverse pyschology on offenders to reconsider their comment

**R**

View all our research and ideation on this miro board :
https://miro.com/welcomeonboard/Qtf9rBIDVDNPsjJOhP2Nk222q3qMnBY2Q20QaT4uErAsvT94LRkeG8PuaxJlza48

Final Outcome →

# Idea 1

## Description

*Modifying the account creation process, which would make users take the policies and code of conduct seriously*

Verifying account history and linking accounts to a person's biometrics helps Instagram background check their users and debars them from creating multiple fake accounts. Customization during account creation process is a new add-on which allows users to choose from a given set of hashtags. Selecting them under the customization section, Instagram would make sure to not show any posts related to those hashtags. Thus, providing users a comfortable user experience. A clear explanation of the policies and code of conduct at the beginning (during the account creation process) is necessary. We understood from our research that regional slurs and abuse is not detectable by Instagram. Which is why mentioning the region you're in, would modify the experience accordingly and make the reporting procedure customized to the region.

## Instagram benefits

1. Conduct a background check on its users to understand their history of locations to get a better sense of their communicative attributes.

2. Opens the door to communicating advertisements in regional language to its users.

3. Better filtration system is equally proportional to easier recognition of the user's attributes.

4. Provide similar sponsorship based on the interests of the user, on all their accounts. Easier to recognize its the same person on multiple accounts.

5. Removing the person from all their accounts after being reported on one. Makes sure Instagram is a safe space.
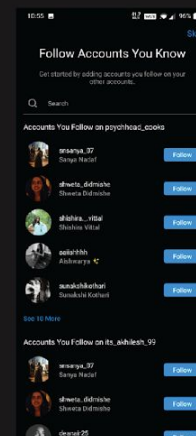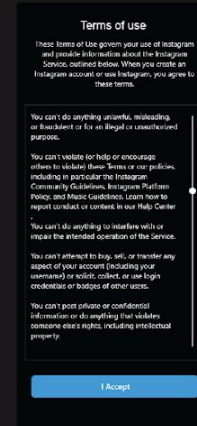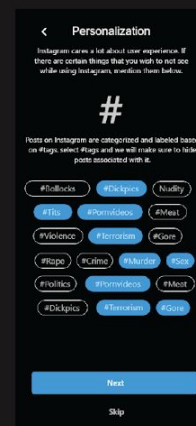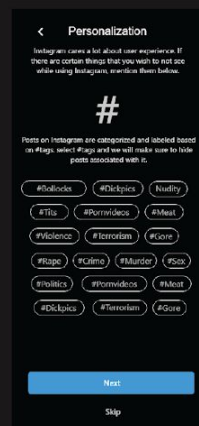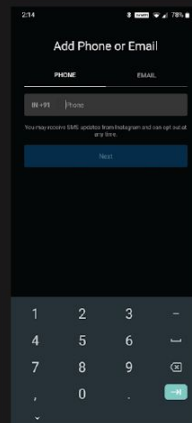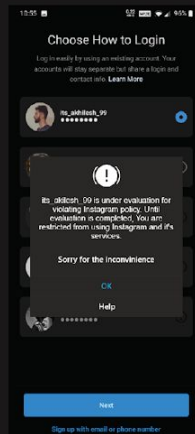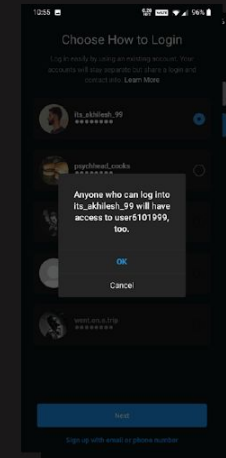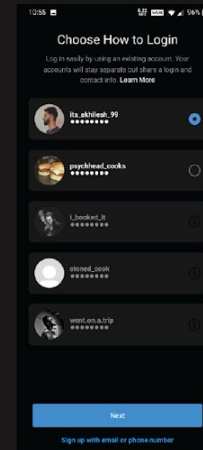
## CR concepts

Panopticon : Self Surveillance

Discourse

# Instagram

## Account creation process update

# Idea 2

## Description

*Sensitive content warning with descriptors and keywords, to give more context to the same and refining feed by giving users more control*

This idea looks at posts labeled as 'sensitive content' on IG, where posts that can trigger users with explicit imagery are blurred. However, this becomes counter-productive since the blurriness of the image only make users more curious to tap on it and end up getting scandalized with the graphic images/video, sometimes to not be able to report it too. In this case, having more context through keywords about what the 'sensitive content' is really about, gives users more information, making them aware and allowing them to make an informed decision. Taking a step further, we're proposing a "Hide" feature under each post, where users can opt to not see posts of a certain kind. We're also extending one of the ideas from the account creation process here, where users decide what kind of content they don't wish to see by customizing hashtags.
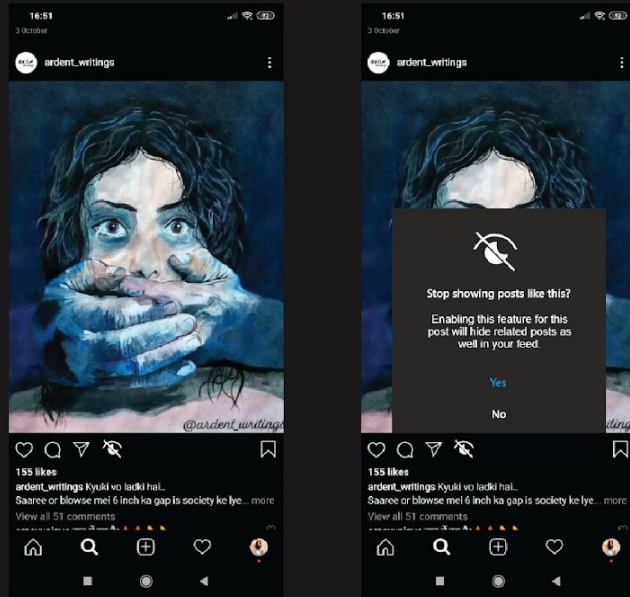
## Instagram benefits

1. Keywords can be used to improve image recognition, for it to describe the content without the need of a person.

2. Descriptors can gain traction over those specific issues with relevant audiences

3. Better filtration system increases instagram's archive and narrows down on the range of ads that will grab the user
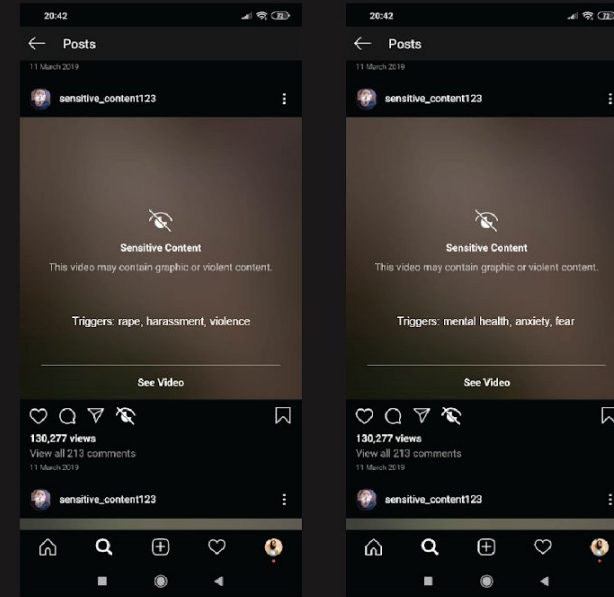
## CR concepts
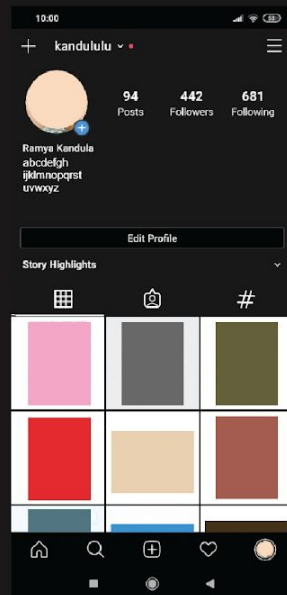
Panopticon :
Self Surveillance

Cognitive Interests

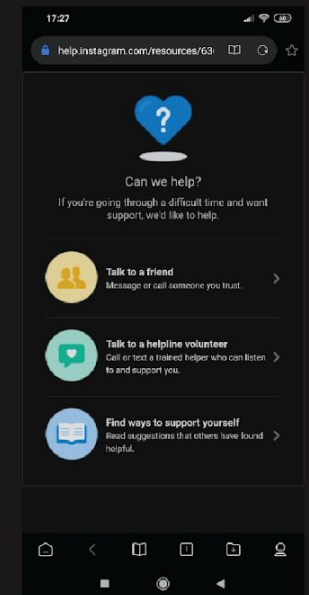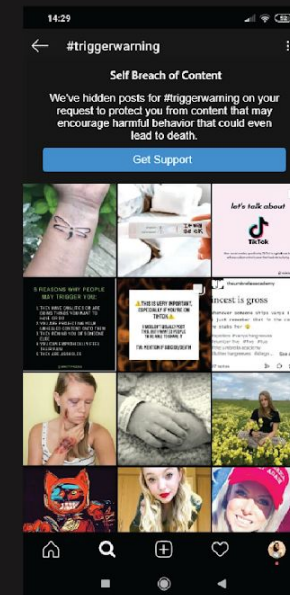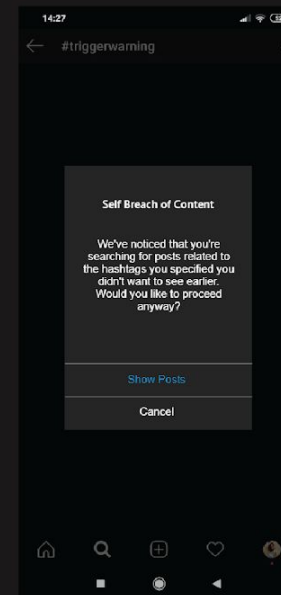## Option for users to choose the kind of posts they wish to see
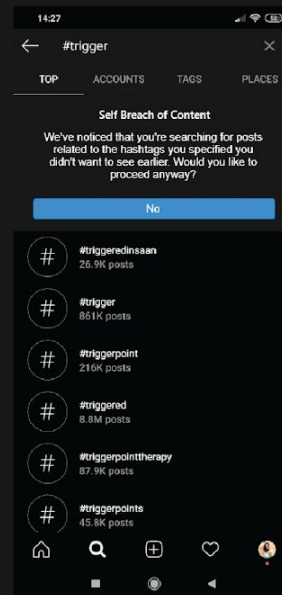


## Keywords on sensitive content to give user context about the post



## Customize content through hashtags



## Cautioning user when self breach of content occurs

# Idea 3

## Description

*Reverse reporting system by putting abuser through similar harrowing steps that usually abused person needs to to go through*

The number of steps to post an abusive comment is three while the number to report one is six. In adopting survivor-centric preventative measures in response to GBV, it is unfair that the onus to seek justice for the abused individuals/communities lies only with themselves. It is possible that survivors of online abuse have already been traumatized and have had this affect their mental health negatively. This proposed new feature for posting abusive comments aims at eliminating and minimalizing the need to report by attempting to nip the abuse in the bud while also educating the abuser. By providing the same sense of uncertainty, hopelessness and discouragement that is otherwise felt by the abused person even after reporting the abuser, would eventually lead to self-surveillance in the abuser.

## Instagram benefits

1. Paid tie ups with multiple news websites, organizations and ad revenue, when promoted during the process
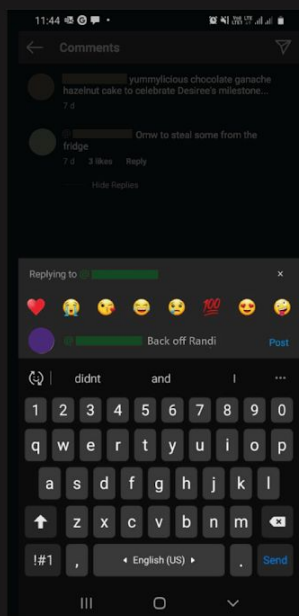
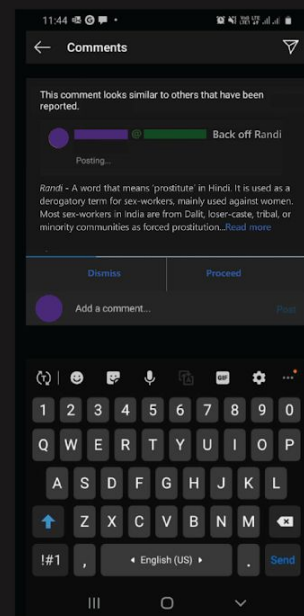## CR concepts

Panopticon :
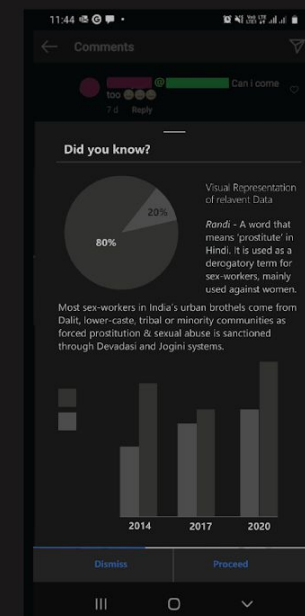Self Surveillance

Communicative
Action

Habitus

Archaeology of
Knowledge

Abuser types in comment

AI detects abusive words/phrases/ tones. User is asked if they want to still post comment

Suggests alternative use of words and shows info related to the abusive language

If they wish to go ahead, they're presented with additional info tabs about the effect of the particular kind of abuse

Users get an option to retract and end process anytime which when clicked will dismiss the comment

If user goes through all steps/tabs they receive a an uncertain response saying only that their comment will be reviewed and only then if it doesn't go against CG will it be posted

# Idea 4

## Description

*Substituting reporting options with examples of abuse for the user to understand which category to choose*

The purpose of this idea is to simplify and humanize the reporting process for the user by replacing overwhelming terminology with more accessible language as well as examples of abusive content to help the user categorize the post better. The list of types of abuse or harassment with examples would take into account the 'blocked tags' from the user's account customization and t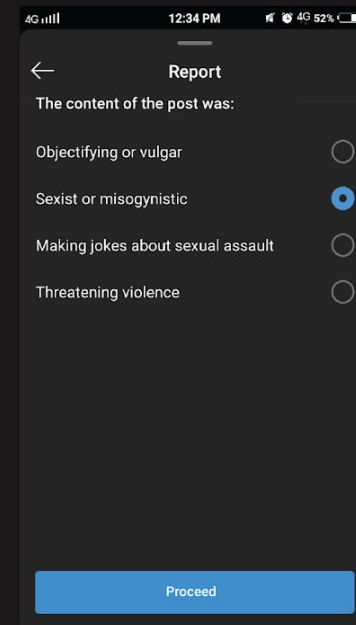he type of content that they have chosen not to see or have reported in the past. The options could also be region specific and display categories and examples related to the user's geographic location and local language. At the end of the report, along with a message that notifies the user of their report being sent and an option to block the reported user, there would be a link to a feature that allows them to keep track of the progress of consequences on the accounts that they reported.

## Instagram benefits

1. More engagement on the reporting system by making it more accessible

2. Lesser content that users dislike and therefore spend more time on the platform

3. Trust in the system with a better tone of voice

## CR concepts

Social and Cultural Capital

An additional category to report beyond 'inappropriate content', to humanize the process.

**Screen 1 — Explore / Report**

politicalindians • Follow

Report

Why are you reporting this post?

It's spam

It violates community guidelines

It is triggering or upsetting to me

**Screen 2 — Report**

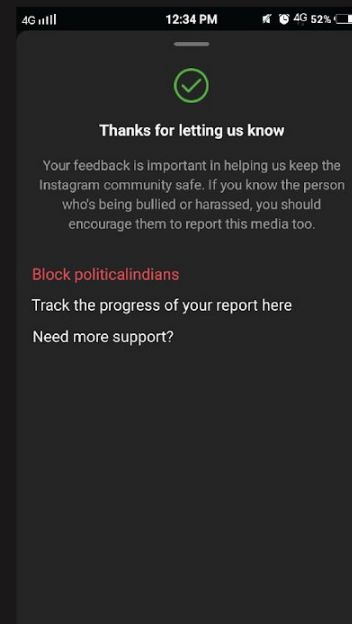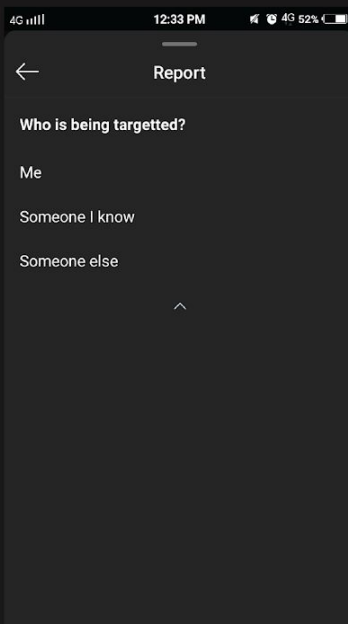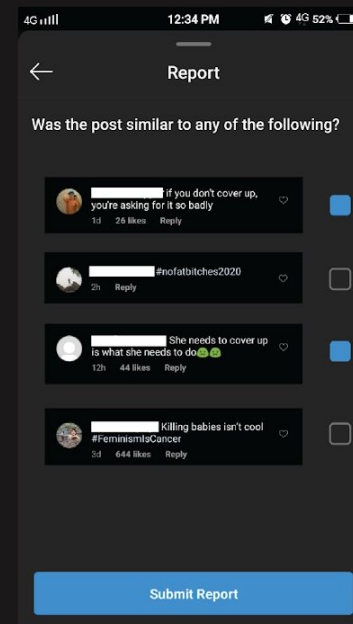Why was the post triggering or upsetting?

It included images or words related to:
#selfharm #homophobia #casteism
#blood #rape #violence #transphobia
#death (tags blocked by user)

Bullying or hateful content

Discrimination against someone's identity

Sexual harassment

**Screen 3 — Report**

The post discriminated against someone's:

Race or ethnicity

Gender or sexual orientation ●

Religion

Caste

Disability

Proceed

**Screen 4 — Report**

The content of the post was:

Objectifying or vulgar

Sexist or misogynistic ●

Making jokes about sexual assault

Threatening violence

Proceed

**Screen 5 — Report**

Was the post similar to any of the following?

So  it's just  a dude then ?
1d   996 likes   Reply

You can't switch genders it's impossible
1d   1149 likes   Reply

LGBT and Demokkkrats 2020
14h   19 likes   Reply

If he can't even figure out what gender he is then how is he supposed to make political decisions
1d   152 likes   Reply

Submit Report

Examples if selected 'Discrimination based on gender identity'

**Screen 6 — Report**

Was the post similar to any of the following?

if you don't cover up, you're asking for it so badly
1d   26 likes   Reply

#nofatbitches2020
2h   Reply

She needs to cover up is what she needs to do
12h   44 likes   Reply

Killing babies isn't cool #FeminismIsCancer
3d   644 likes   Reply

Submit Report

Examples if selected 'Sexist or misogynist' content

**Screen 7 — Report**

Who is being targetted?

Me

Someone I know

Someone else

**Screen 8 — Final**

Thanks for letting us know

Your feedback is important in helping us keep the Instagram community safe. If you know the person who's being bullied or harassed, you should encourage them to report this media too.

Block politicalindians

Track the progress of your report here

Need more support?

Final screen.
Track progress will take the user to the reports history/activity feature on their profile.

# Idea 5

## Description

*History of a user's abusive or offensive behavior can be transferred to other platforms or employers.*

This feature came into being, to create an archive of reports. This allows the users to keep track of their reports and hold Instagram accountable in case the reported content doesn't check mark their guidelines, and have a constant scope for improvement. The other part of this archive is the history of the user's own reports. This space was created for an ethical background check for employees of a workspace. While this creates a set of real world consequences to their online behavior, it also tends to be a reality check. This archive is a space for self surveillance of the user's behavior online, to make them more conscious about their words and actions. Rewarding positive changes and conscious behavior online is equally important. Which is why curated highlights on user's profile is a small but meaningful addition - like being "report free". Instagram can create a joint database with other platforms as well, create guidelines for reporting content and have the power to remove hate online.
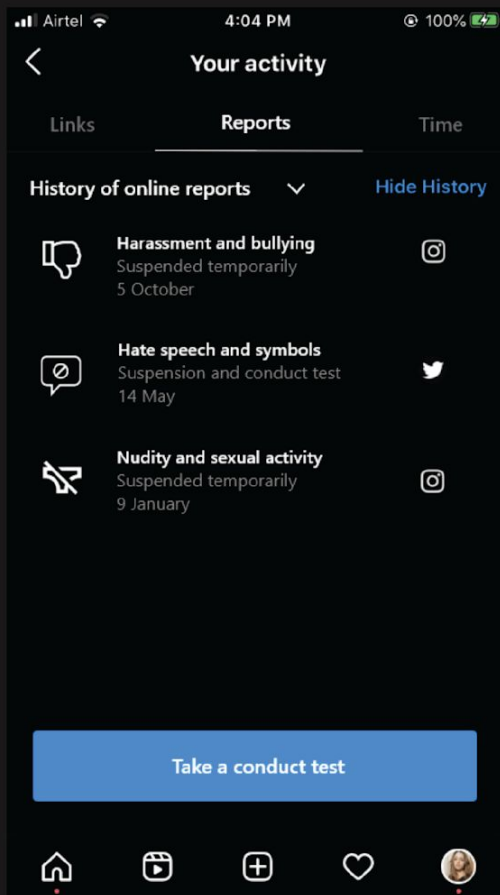
## Instagram benefits

1. Joint database of reporting guidelines from all platforms can make the Ai better in recognizing hateful content

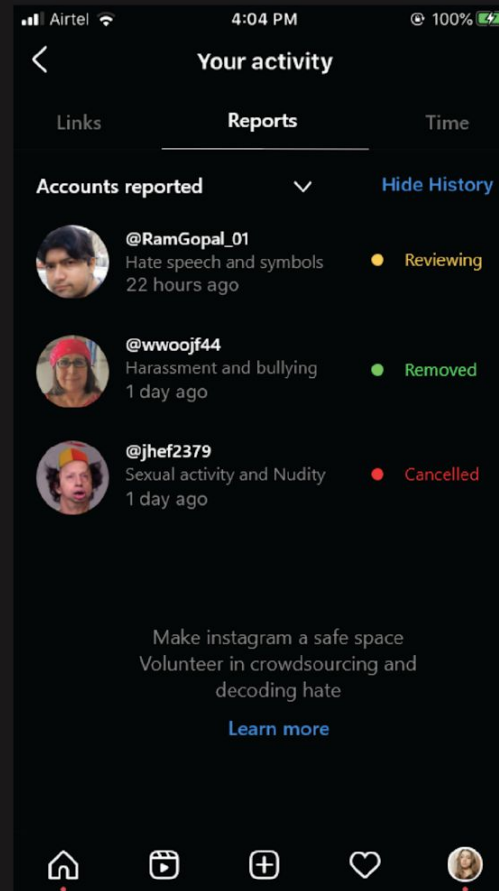2. Increased screen time opportunity by adding a tracking status of the report
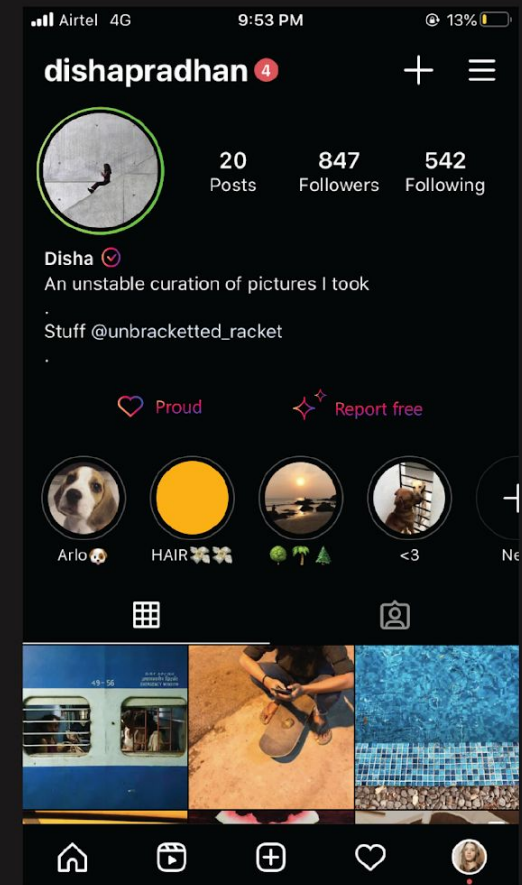
## CR concepts

Cultural Capital

Panopticon :
Self Surveillance

History of user's online reports
+ Option to redeem by taking
conduct test to qualify for it

Tracking content reported

Verify : Bumble approach to verify the
user is the same person they are
impersonating
+ Profile highlights

Adobe XD link :
https://xd.adobe.com/view/b681509e-4404-40e0-af77-0744ab9b7987-e3d2/?fullscreen

# Thank You.