# Who Moderates the Moderators?
## A Look Inside Bluesky

**SICSS-Saarbrücken**

**September 2025**

**Abhisek Dash, Pushpdeep Singh**

# Caution : Some slides may contain explicit content.

However, I believe this will make us all appreciate the complexity.

# Plan for today

Content moderation

The centralized black-box

A decentralized alternative: Bluesky

Content moderation on Bluesky
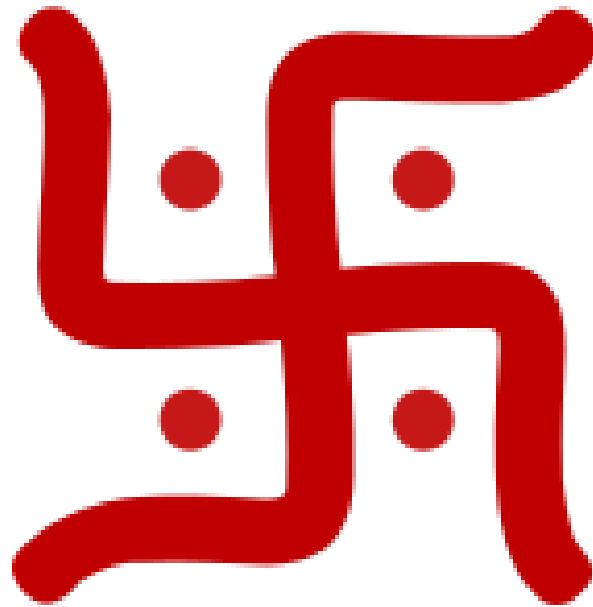
Our role as CSS researchers

Hands on session (in the afternoon)

# Content moderation

# Before we proceed ...

**If someone posted the following image on your favorite social media platform. <span style="color:red">Should the platform take down the post?</span>**

# Should the platform take down this image?

# Different meaning in different cultures

In Hinduism, the right-facing symbol (clockwise) (卐) is called swastika, symbolizing sun, prosperity and good luck.

# Should the platform take down this image?

# A bit of a history ...

❑ **Napalm Girl** : 1972 Pulitzer Prize winning photo by Nick Ut (AP).

❑ **2016** : Norwegian journalist Tom Egeland included it in an article reflecting on photos that changed the history of warfare.

❑ The picture contains graphic suffering and underage nudity.

❑ **Consequence** : Facebook moderators deleted the post.

❑ The shown snapshot is the outrage of the editor in chief of Aftenposten on the front page of the newspaper.



**Photo taken from Custodians of The Internet by Tarleton Gillespie**

# Statement from Facebook Vice President

*" These decisions aren't easy. In many cases, there's no clear line between an image of nudity or violence that carries global and historic significance and one that doesn't. Some images may be offensive in one part of the world and acceptable in another, and even with a clear standard, it's hard to screen millions of posts on a case-by-case basis every week…  In this case, we tried to* <span style="color:red">***strike a difficult balance between enabling expression and protecting our community and ended up making a mistake***</span>*… "*

# Even today...

**Bluesky Safety** ✓
@safety.bsky.app

+ Follow

Glorifying violence or harm violates Bluesky's Community Guidelines. We review reports and take action on content that celebrates harm against anyone. Violence has no place in healthy public discourse, and we're committed to fostering healthy, open conversations

September 11, 2025 at 12:26 AM   ⊘ Replies disabled

**495** reposts   **1** quote   **2.4K** likes   **36** saves

### The New York Times

**Protests in Nepal** | Updates | What to Know | Class Tensions | Censorship Playbook Fails | Longstanding Problems

## Nepal Bans 26 Social Media Platforms, Including Facebook and YouTube

Critics worry a new law could curb freedom of expression, affect tourism and cut communication with the many Nepalis who work abroad.

11

# Definition from the Digital Services Act (DSA)

Article 3(t) defines content moderation as:

"*content moderation means the activities undertaken by platforms aimed at* ***detecting, identifying and addressing*** *illegal content or information incompatible with their terms and conditions …*"

Source: https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32022R2065

# Takeaways

Content moderation is often striking a balance between enabling expression of the community and protecting the community.

Thus, it is notoriously difficult.

# The centralized black box

# Terms and conditions : Community standards

# Partial list of violations

- Abusive accounts
- Privacy violation
- Spam
- Financial scams
- Regulated goods
- Adult sexual exploitation
- Adult sexual solicitation
- Child nudity and sexual exploitation
- Human trafficking

- Self harm
- Violence and graphical content
- Societal harm
- Personal harm
- Bullying
- Credible threats of violence
- Hate speech
- Misinformation

… … …

# A typical content moderation pipeline



**Figure courtesy :** Halevy, Alon, et al. "Preserving integrity in online social networks." *Communications of the ACM* 65.2 (2022)

# A broader view of enforcement

**Remove clearly violating content**

**Remove / downrank offensive content**

**Personalize removal or downranking decisions**

**Minimize bad experiences on your newsfeed**

# Although thoughtful, some systemic flaws exist

**Opaque by design**

**Centralized power**

**Zero user agency**

**Arbitrary / inconsistent enforcement**

# XChecks : Above the community standards

# Takeaways

Centralized moderation is an incredibly complex, expensive, human-in-the-loop socio-technical system, built with thoughtful policies.

However, its opaque and centralized nature leads to arbitrary enforcement, denying users agency and eroding trust.

# A decentralized alternative: Bluesky

# Recent shift to decentralized alternatives



✓ Built on open protocols

✓ Power is decentralized

✓ More user agency and autonomy

✓ Foster innovation through open standards

# Recent shift to decentralized alternatives



✓Built on open protocols

✓Power is decentralized

✓More user agency and autonomy

✓Foster innovation through open standards

# Bluesky

# User data storage : Personal Data Servers



❑ A Personal Data Server (PDS) stores the user's repositories (i.e., all the actions they ever performed).

❑ Currently, there are 2,555 valid PDSs on Bluesky and 70 of them are run by Bluesky Social PBC.

# Indexing I : Relay



❑ Relay crawls the user repositories on all known PDSes and consumes the streams of updates that they produce

❑ Relay creates the firehose: an aggregated stream of updates.

# Content dissemination : Labeler and Feed generators



❑ Labelers consume the firehose and moderate the content.

❑ Feed generators consume the firehose and reorder the content.

# Indexing II & User interface : App view



❑ The App View is a service that consumes the firehose, and processes the records that are relevant to the Bluesky social app

❑ The App View is also responsible for enforcing moderation controls.

# Takeaways

Bluesky is a decentralized alternative where the power is diffused across different entities.

Moderation is no longer a top-down function of the platform, but a service in an open marketplace.

# Break (5 mins)

Content moderation

The centralized black-box

A decentralized alternative: Bluesky

Content moderation on Bluesky

Our role as CSS researchers

Hands on session (in the afternoon)

# Composable content moderation on Bluesky

# Labelers on Bluesky



**Bluesky Moderation Service**
@moderation.bsky.app

Official Bluesky moderation service. bsky.social/about/suppor...

Labels     Posts     Replies

Labels are annotations on users and content. They can be used to hide, warn, and categorize the network.

**Adult Content**
Explicit sexual images.    Hide
ⓘ Configured in *moderation settings.*

**Sexually Suggestive**
Does not include nudity.    Warn
ⓘ Configured in *moderation settings.*

**Non-sexual Nudity**
E.g. artistic nudes.    Off
ⓘ Configured in *moderation settings.*

**Sexually Suggestive (Cartoon)**
Art with explicit or suggestive sexual themes, including provocative imagery or partial nudity.
Off / **Warn** / Hide

**Graphic Media**
Explicit or potentially disturbing media.    Warn
ⓘ Configured in *moderation settings.*

**Self-Harm**
Promotes self-harm, including graphic images, glorifying discussions, or triggering stories.
Off / **Warn** / Hide

**Sensitive**
May be upsetting, covering topics like substance abuse or mental health issues, cautioning sensitive viewers.
**Off** / Warn / Hide

**Extremist**
Radical views advocating violence, hate, or discrimination against individuals or groups.
**Off** / Warn / Hide

**Intolerance**
Discrimination against protected groups.
Off / **Warn** / Hide

**Bluesky Moderation Service (Germany)**
@moderation-de.bsky.app

Offizieller Bluesky-Moderationsdienst. bsky.social/about/suppor...

Labels     Posts     Replies

Labels are annotations on users and content. They can be used to hide, warn, and categorize the network.

This labeler hasn't declared what labels it publishes, and may not be active.

## Types of labelers

**Default labeler**

**Country specific labelers**

**Community labelers**

**Blacksky Moderation**
@blackskyweb.xyz

Subscribe to Labeler

Creating a safe space for Black community building. We built the Blacksky feeds, this mod service, an atproto implementation called "rsky" and more.
www.blackskyweb.xyz

♡ Liked by 6,237 users

Labels     Lists     Posts     Replies

Labels are annotations on users and content. They can be used to hide, warn, and categorize the network.

Subscribe to @blackskyweb.xyz to use these labels:

**Synthetic Media**
Content which has been generated or manipulated to appear as though based on reality, when it is in fact artificial. Also referred to as manipulated media. Synthetic media may sometimes (but not always) be generated through algorithmic processes (such as artificial intelligence or machine learning). A deepfake is a form of synthetic media where an image or recording is altered to misrepresent someone doing or saying something that was not done or said.

**Violence**
Content that expresses violence via images or statements that target Black people. This may include, but is not limited to, threats of physical violence and sexual violence.

**White Supremacy or Antiblack Rhetoric**
Content that expresses statements that are rooted in white supremacy and anti-black rhetoric that does not fall under targeted anti-black harassment. This may include internalized anti-blackness that is harmful in nature.

**Doxxing**
The act of disclosing someone's personal, non-public information — such as a real name, home address, phone number or any other data that could be used to identify the individual — in an online forum or other public place without the person's consent.

**Non-Consensual Intimate Imagery**
Non-consensual image sharing, or non-consensual intimate image sharing (also called 'non-consensual explicit imagery' (NCEI) or colloquially called 'revenge porn'), refers to the act or threat of creating, publishing or sharing an intimate image or video without the consent of the individuals visible in it.

**Misogynoir**
Content that expresses hatred, bias, or prejudice against Black women, specifically where racism, sexism, and/or transphobia intersect. This includes sexual harassment, objectification, and targeted attacks on Black women's identity or appearance.

**Antiblack Harassment**
Content that targets individuals or groups based on their Black identity with derogatory, hateful, or dehumanizing language or imagery. This label is applied to content that perpetuates harmful stereotypes, slurs, or direct harassment aimed at Black individuals or communities.

# Users have more autonomy and agency

❑Bluesky and country specific labeler apply by default.

❑On the other hand, this architecture empowers users to

   ❑choose what should be added to their feed.

   ❑choose how labeled content should appear.

   ❑choose which other lebeler should affect their feed.

# Some questions that may arise

❑What kind of labels do these labelers apply?

❑What are the interpretations of these labels?

❑How do they apply these labels? Are the algorithms transparent?

# Labeler outcomes are publicly accessible



Top Labels in posts (for which data is available in firehose) json file

1521785 porn
500396 sexual
213873 !takedown
191673 spam
69654 nudity
23141 rude
19337 graphic-media
19107 sexual-figurative
9978 intolerant
3862 self-harm
3114 !hide
2662 threat
805 !warn
16 extremist
8 misinformation
6 misleading
4 gore
2 rumor
2 unsafe-link
1 !no-unauthenticated
1 inauthentic
1 needs-review
1 sensitive
1 impersonation

❑ Total distinct posts labelled : 2,528,786 (19 March '25 - 1 June '25 )

❑ The data reveals a strong focus on platform hygiene (spam, porn) and user safety issues (intolerance, rude).

# What do these labels even mean?

| Labels | Description |
| --- | --- |
| Intolerance | Discrimination against protected groups. |
| Threats | Promotes violence or harm towards others, including threats, incitement, or advocacy of harm. |
| Rude | Rude or impolite, including crude language and disrespectful comments, without constructive purpose. |

Source: https://bsky.app/profile/moderation.bsky.app

# How does a labeler apply these labels?

# How does Bluesky apply these labels?

The key takeaway: We can see *what* was labeled, but not *why*. The implementation policies remain <span style="color:red">opaque</span>.

Source: https://bsky.social/about/blog/01-17-2025-moderation-2024

# Intolerance: Discrimination against protected groups

Q: *How does Bluesky labeler operationalize 'protected group'?*

Information about someone's protected characteristics—*such as sexual orientation, gender identity, disability, caste, or immigration status*—for the purpose of targeting them or discriminating against them.

Source: https://bsky.social/about/support/community-guidelines

# Access to label stream gives some hope!

❑Here is a cluster of all the posts labelled as Intolerance.

❑Such analyses on the labelled text could give us a more nuanced view of what is happening behind the scenes.



Semantic Subclusters within 'intolerant' Posts

# Takeaways

Bluesky's composable moderation paradigm provides more autonomy and agency to users.

While procedural transparency is still lacking on Bluesky labeler, it is auditable thanks to the accessible label stream.

# Our role as Computational Social Scientists

# Systemic flaws still exist

**Opaque by design**

↓

**Opacity has decreased**

**Centralized power**

↓

**Power is decentralized**

**Zero user agency**

↓

**Better user agency**

**Inconsistent enforcement**

↓

**Now auditable**

44

# Opportunity 1: Data Access

❏Reliability of content moderation is now measurable.

❏We will see some examples of this in the hands-on session.

  ❏What kind of labels are applied?

  ❏Are they applied timely?

  ❏Are they applied consistently?



**Bluesky Developer APIs**

Explore Bluesky's open social network.

Get Started

# Opportunity 2: Open marketplace

❑For ages, researchers have been able to audit platform practices.

❑Now, we have the opportunity to rectify the drawbacks.

# Opportunity 3 : Stay ahead of the curve

Proactively think about regulations for better accountability, data

protection while not hindering innovation.

# Takeaways

Centralized moderation is opaque and unaccountable.

Decentralized moderation is more transparent and composable, giving users more autonomy. Yet it lacks procedural transparency.

This paradigm shift empowers us not only to measure the moderation service, but also to set the standards.

# Thank You!

✉ **adash@mpi-sws.org**

# Plan for today

Content moderation

The centralized black-box

A decentralized alternative: Bluesky

Content moderation on Bluesky

Our role as CSS researchers

Hands on session (in the afternoon)

# Hands-on session

# DID and PLC Directory

- Each user on Bluesky is identified via a unique DID (decentralised identifier) e.g., *did:plc:ljxfsne42aud2cd2iowkkqiz*.

- Each DID points to its associated **DID Document**, a document that stores service information about the user.

**did:plc**                                           Lookup   API   Specification   Code

## Resolve a did:plc Identifier

🔍  did:plc:ar7c4by46qjdydhdevvrndac

https://plc.directory/did:plc:ar7c4by46qjdydhdevvrndac

# DID and PLC Directory

**DID Document JSON**

```json
{
  "@context": [
    "https://www.w3.org/ns/did/v1",
    "https://w3id.org/security/multikey/v1",
    "https://w3id.org/security/suites/secp256k1-2019/v1"
  ],
  "alsoKnownAs": [
    "at://moderation.bsky.app"
  ],
  "id": "did:plc:ar7c4by46qjdydhdevvrndac",
  "service": [
    {
      "id": "#atproto_pds",
      "serviceEndpoint": "https://inkcap.us-east.host.bsky.network",
      "type": "AtprotoPersonalDataServer"
    },
    {
      "id": "#atproto_labeler",
      "serviceEndpoint": "https://mod.bsky.app",
      "type": "AtprotoLabeler"
    }
  ],
  "verificationMethod": [
    {
      "controller": "did:plc:ar7c4by46qjdydhdevvrndac",
      "id": "did:plc:ar7c4by46qjdydhdevvrndac#atproto",
      "publicKeyMultibase": "zQ3shoG4QW9B3zvKSSiRwwc1De7MFQLNBT9A71gr12GKwMgHu",
      "type": "Multikey"
    },
    {
      "controller": "did:plc:ar7c4by46qjdydhdevvrndac",
      "id": "did:plc:ar7c4by46qjdydhdevvrndac#atproto_label",
      "publicKeyMultibase": "zQ3shmV1BNcX17coaDbfen6zArEad6SCLT3jVWCbC6Y9iinTa",
      "type": "Multikey"
    }
  ]
}
```

Additional entry
for a labeler

53

# Getting all labelers

- Using the **sync.listRepos** call offered by the Bluesky Relay, we can obtain a list of all active Bluesky users and their DIDs.

"cursor": "10",
"repos": [
{ "active": true,
**"did": "did:plc:hlm3aibaluzuzuqqnzxg2urq"**,
"head": "bafyreigizoa62issah3ajgtmwgpi3wn4bsz3uchv5dihbchm2uubuelx2i",
"rev": "3lqmtafrzyv2r"
},
{ "active": true,
**"did": "did:plc:6qgqcyg5gfd6sxmghdclxmpi"**,
"head": "bafyreidnlloy7pc634f36a2a2q3a7mwhzdslzfx2vfyk2ugel4sjpxpvsq",
"rev": "3lqn6vruciq2n"
},
{ "active": true,
**"did": "did:plc:kbf77syjgjcjciabbsdm37qq"**,
...

Using https://relay1.us-west.bsky.network/xrpc/com.atproto.sync.listRepos endpoint

- Check plc directory to get DID document for each of these DID and see if they are a labeler.

# Getting labeler information

Suppose we know a particular user (identified by a DID) is a labeler, from DID document we have its service endpoint (we will explore where it is used, later in this session)

But, what about :

- their descriptions?

- which labels are issued? their definitions?

- default settings?

**API : app.bsky.labeler.getServices**

Colab notebook

# https://shorturl.at/ZicnE

# TASK 1 : Getting labeler info

TASK 2 : Getting label-stream from a labeler

# TASK 3 : Measuring labeling time

# TASK 4 : Analyzing labeled content & consistency

# "Only two genders exist" cluster

Should the posts on the right be labelled?

| Examples labeled as Intolerant by Bluesky | From Top-k posts from firehose | |
|---|---|---|
| Only two genders exist: ○ Male ○ Female ~Anything other than that is pure confusion. | There's only 2 genders | ✅ |
| There's only two genders | There's only two genders | ✅ |
| There are only 2 genders | Есть только два гендера | ✅ |
| Il existe deux genres, mâle et femelle. Le reste c'est de la psychiatrie. | There are exactly two genders.  Just two. | ✅ |
| There are 3 genders: Male Female Mentally ill! | Gibt trotzdem nur 2 Geschlechter | ✅ |

Examples labeled as Intolerant by Bluesky

From Top-k posts from firehose

# "Only two genders exist" cluster

Bluesky label intolerant?

| Examples labeled as Intolerant by Bluesky | From Top-k posts from firehose | Bluesky label intolerant? |
|---|---|---|
| Only two genders exist: ○ Male ○ Female ~Anything other than that is pure confusion. | There's only 2 genders | ✅ |
| There's only two genders | There's only two genders | ❌ |
| There are only 2 genders | Есть только два гендера | ❌ |
| Il existe deux genres, mâle et femelle. Le reste c'est de la psychiatrie. | There are exactly two genders. Just two. | ❌ |
| There are 3 genders: Male Female Mentally ill! | Gibt trotzdem nur 2 Geschlechter | ✅ |

Semantic Subclusters within 'intolerant' Posts