# Moderating the Moderators:
# A Look Inside Bluesky's Content Moderation

**November 2025**

**Abhisek Dash**

# Caution : Some slides may contain explicit content.

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

# 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**



FREDAG 9. september 2016   Uke 36 · Nr. 253 · 157. årgang · Løssalg kr 40   (Levert hjem fra kr 12. Bestill på ap.no/abo)

Facebook krever at Aftenposten fjerner dette historiske bildet fra vår Facebook-side. Her er Aftenpostens svar:

## Dear Mark Zuckerberg

Jeg skriver til deg for å fortelle hvorfor Aftenposten ikke vil etterkomme Facebooks krav om å fjerne eller redigere dette viktige dokumentarbildet.

**Espen Egil Hansen, sjefredaktør**

NYHETER · DEL 1 · SIDE 2-5 · KULTUR · DEL 2 · SIDE 2-3

# 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* **strike a difficult balance between enabling expression and protecting our community and ended up making a mistake***… "*

# Even today...

**Bluesky Safety** ✓
@safety.bsky.app

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

⊕ Follow

*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.

# 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

## Community Standards

Our Community Standards apply to everyone, all around the world, and to all types of content, including AI-generated content.

Each section of our Community Standards starts with a "Policy Rationale" that sets out the aims of the policy followed by specific policy lines that outline:

🚫 Content that's not allowed; and

⚠️ Content that requires additional information or context to enforce on, content that is allowed with a warning screen or content that is allowed but can only be viewed by adults aged 18 and older.

Coordinating Harm and Promoting Crime

Dangerous Organizations and Individuals

Fraud, Scams, and Deceptive Practices

Restricted Goods and Services

Violence and Incitement

Adult Sexual Exploitation

Bullying and Harassment

Child Sexual Exploitation, Abuse, and Nudity

Human Exploitation

Suicide, Self-Injury, and Eating Disorders

Adult Nudity and Sexual Activity

Adult Sexual Solicitation and Sexually Explicit Language

Hateful Conduct

Privacy Violations

Violent and Graphic Content

Account Integrity

Authentic Identity Representation

Cybersecurity

Inauthentic Behavior

Memorialization

Misinformation

Spam

Third-Party Intellectual Property Infringement

Using Meta Intellectual Property and Licenses

Additional Protection of Minors

Locally Illegal Content, Products, or Services

# Terms and conditions : Community standards

## Community Standards

Our Community Standards apply to everyone, all around the world, and to all types of content, including AI-generated content.

Each section of our Community Standards starts with a "Policy Rationale" that sets out the aims of the policy followed by specific policy lines that outline:

🚫 Content that's not allowed; and

⚠️ Content that requires additional information or context to enforce on, content that is allowed with a warning screen or content that is allowed but can only be viewed by adults aged 18 and older.

Coordinating Harm and Promoting Crime
Dangerous Organizations and Individuals
Fraud, Scams, and Deceptive Practices
Restricted Goods and Services
Violence and Incitement
Adult Sexual Exploitation
Bullying and Harassment
Child Sexual Exploitation, Abuse, and Nudity
Human Exploitation
Suicide, Self-Injury, and Eating Disorders
Adult Nudity and Sexual Activity
Adult Sexual Solicitation and Sexually Explicit Language
Hateful Conduct
Privacy Violations
Violent and Graphic Content

### ✖ Help Center

Using X Managing your account Safety and security Rules and policies

participate in the public conversation freely and safely.

## Safety

**Violent Content:** You may share graphic media if it is properly labeled, not prominently displayed and is not excessively gory or depicting sexual violence, but explicitly threatening, inciting, glorifying, or expressing desire for violence is not allowed. Learn more.

**Violent & Hateful Entities:** You can't affiliate with or promote the activities of violent and hateful entities. Learn more.

**Child Safety:** We have zero tolerance for any forms of child sexual exploitation and remove certain media depicting physical child abuse to prevent the normalization of violence against children. Learn more.

**Abuse/Harassment:** You may not share abusive content, engage in the targeted harassment of someone, or incite other people to do so. Learn more.

**Hateful conduct:** You may not attack other people on the basis of race, ethnicity, national origin, caste, sexual orientation, gender, gender identity, religious affiliation, age, disability, or serious disease. Learn more.

**Perpetrators of Violent Attacks:** We will remove any accounts maintained by individual perpetrators of terrorist, violent extremist, or mass violent attacks, and may also remove posts disseminating manifestos or other content produced by perpetrators. Learn more.

**Suicide:** You may not promote or encourage suicide or self-harm. Learn more.

**Adult Content:** You may share consensually produced and distributed adult nudity or sexual behavior, provided it's properly labeled and not prominently displayed. Learn more.

**Illegal or Certain Regulated Goods or Services:** You may not use our service for any unlawful purpose or in furtherance of illegal activities. This includes selling, buying, or facilitating transactions in illegal goods or services, as well as certain types of regulated goods or services. Learn more.

### ♪ TikTok

Overview
Community Principles
Youth Safety and Well-Being
Safety and Civility
Mental and Behavioral Health
Sensitive and Mature Themes
Integrity and Authenticity
Regulated Goods and Commercial Activities
Privacy and Security
For You feed Eligibility Standards
Accounts and Features
Enforcement

**Content Moderation**

Keeping our platform safe, trustworthy, and vibrant requires balancing creative expression and preventing harm. We use a combination of safety approaches to strike the right balance:

🚫 **Remove** content that we do not allow

Everyone who joins TikTok has the ability to freely share content on the platform. However, we remove content--whether posted publicly or privately--when we find that it violates our rules.

⛔ **Restrict** content that is not suitable for youth

We allow a range of content on our platform, but also recognize that not all of it may be suitable for younger audiences. We restrict content that may not be suitable so that it is only viewed by adults (18 years and older). A summary of restricted content categories can be found here.

⚠️ **Make ineligible** for the FYF content that does not meet our recommendation standards

The FYF is an opportunity to discover new content and reach new audiences, but it is not guaranteed that all content will be recommended. Content that does not meet our standards will be ineligible for the FYF. A summary of these standards can be found here.

🔲 **Empower our community with information, tools, and resources**

...rience on TikTok. We may add labels, "opt-in" ...r out content with specific hashtags or ...nd in-app features with safety resources.

### YouTube Help Describe your issue

**Sensitive content**

We hope to protect viewers, creators, and especially minors. That's why we've got rules around keeping children safe, sex & nudity, and self harm. Learn what's allowed on YouTube and what to do if you see content that doesn't follow these policies.

- Nudity & sexual content policies
- Thumbnails policy
- Child safety policy
- Suicide, self-harm, and eating disorders policy
- Vulgar language policy

**Violent or dangerous content**

Hate speech, predatory behavior, graphic violence, malicious attacks, and content that promotes harmful or dangerous behavior isn't allowed on YouTube.

- Harmful or dangerous content policies
- Violent or graphic content policies
- Violent criminal organizations policy
- Hate speech policy
- Harassment & cyberbullying policies

**Regulated goods**

Certain goods can't be sold on YouTube. Find out what's allowed and what isn't.

- Sale of illegal or regulated goods or services policies
- Firearms policy

**Misinformation**

Certain types of misleading or deceptive content with serious risk of egregious harm are not allowed on YouTube. This includes certain types of misinformation that can cause real-world harm, like promoting harmful remedies or treatments, certain types of technically manipulated content, or content interfering with democratic processes.

- Misinformation policies
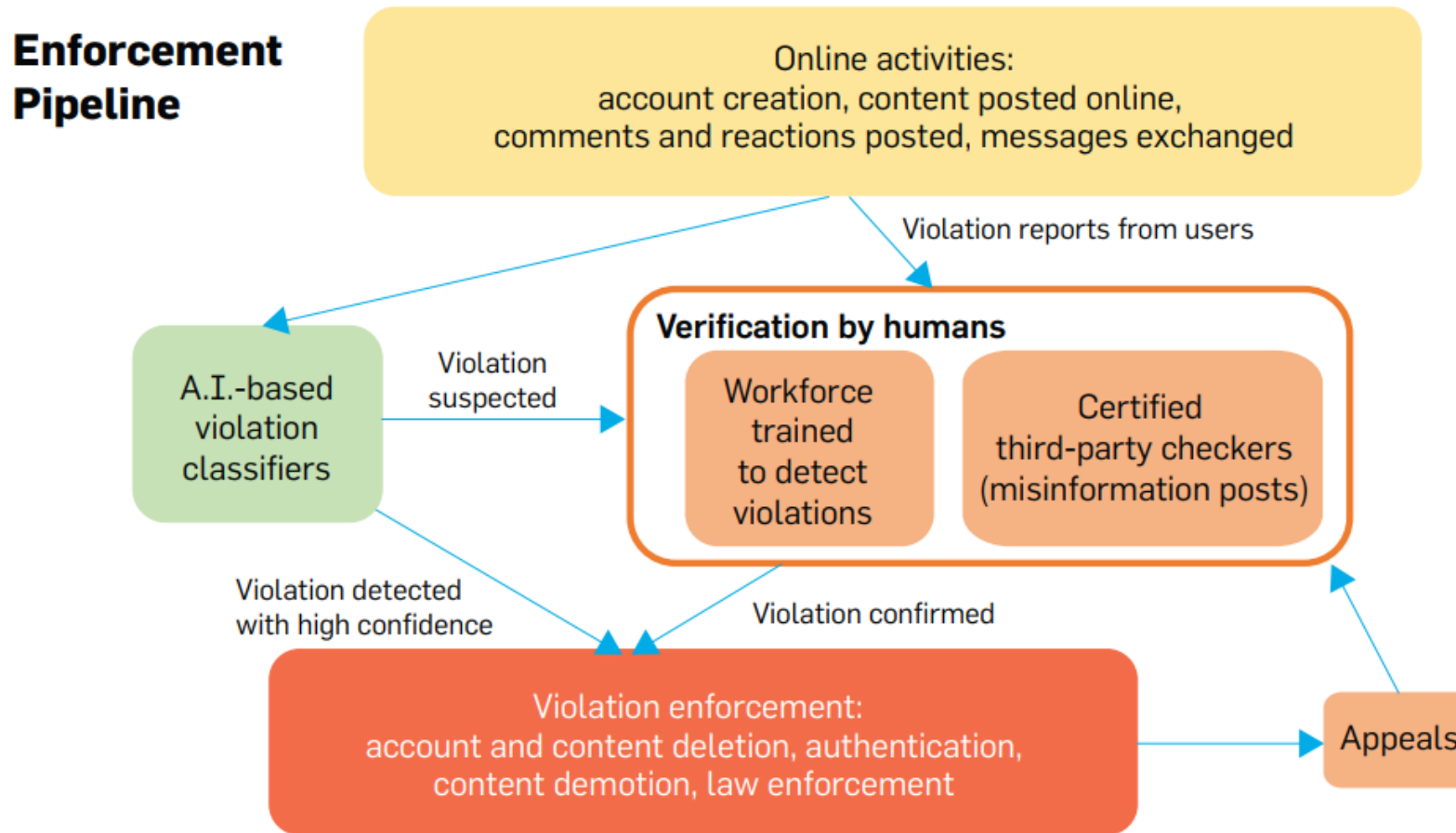- Elections misinformation policies
- Medical misinformation policies

# 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
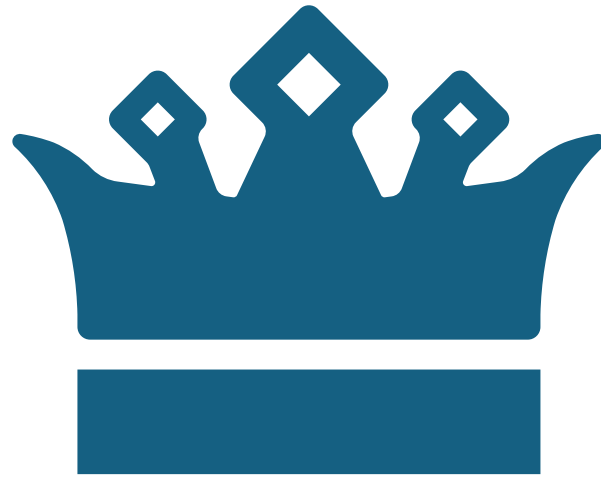- Misinformation

… … …

# A typical content moderation pipeline



**Enforcement Pipeline**

Online activities:
account creation, content posted online,
comments and reactions posted, messages exchanged

Violation reports from users

A.I.-based violation classifiers

Violation suspected

**Verification by humans**

Workforce trained to detect violations

Certified third-party checkers (misinformation posts)

Violation detected with high confidence

Violation confirmed

Violation enforcement:
account and content deletion, authentication, content demotion, law enforcement

Appeals

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

# Although thoughtful, some systemic flaws exist

**Opaque by design**

**Centralized power**

**Zero user agency**

**Arbitrary / inconsistent enforcement**

# XChecks : Above the community standards

Source: https://www.wsj.com/articles/the-facebook-files-11631713039#_=_

# Takeaways

Centralized moderation is an incredibly complex, expensive, human-in-the-loop socio-technical system, built with thoughtful policies.
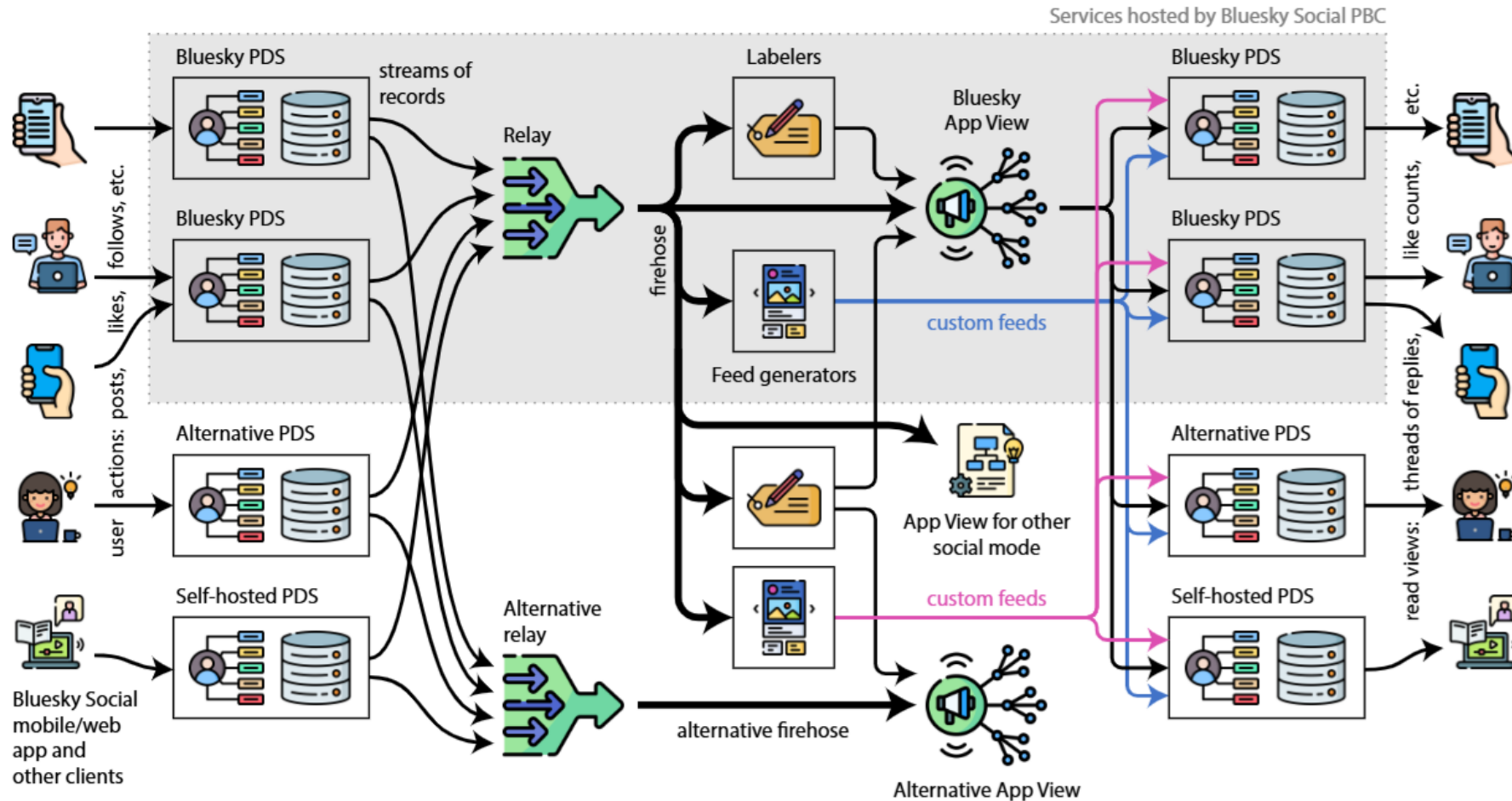
However, its opaque and centralized nature leads to arbitrary enforcement, denying users agency and eroding trust.

# A decentralized alternative: Bluesky

# Bluesky

18

# 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.
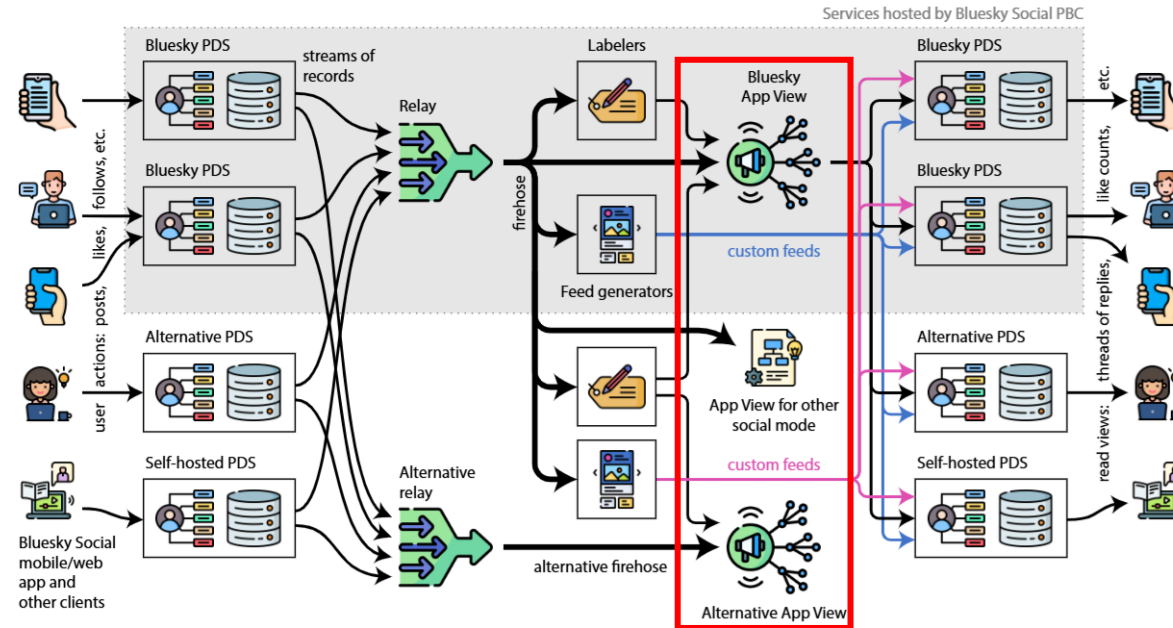
# Indexing I : Relay



❑ 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.

# Composable content moderation on Bluesky

# Labelers on Bluesky



Types of labelers

Default labeler

Country specific labelers

Community labelers

# 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?

❑ Who apply these labels : Algorithms or Humans?

❑ How do the labelers operationalize abstract principles?

❑ How consistent are these operationalizations?

# Labeler outcomes are publicly accessible



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

❑ 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).

# Timeliness of application of labels



Bluesky Labels: Volume vs. Reaction Time

# Automated vs. Human-in-the-loop labels



Bluesky Labels: Volume vs. Reaction Time

# Automated vs. Human-in-the-loop labels



Bluesky Labels: Volume vs. Reaction Time

# 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

# 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

# Understanding operationalization from label stream

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

❑ Many of the clusters correspond to gender, religion, race, health status, politics etc.

# How does Bluesky apply these labels?

We can see *what* was labeled, but not *why*. The implementation policies remain <span style="color:red">opaque</span>, while the operationalization outcome is <span style="color:#1f3a5f">transparent</span>.

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

# How consistent are these operationalizations?

❑ Let's look at one of the clusters : "Only two genders exist" cluster

# "Only two genders exist" cluster

Should the posts on the right be labelled?

| | |
|---|---|
| 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 labels

| 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

# LLM based consistency evaluation

**Prompt (for Llama 3 70B Instruct) :**

You are a human annotator. The posts given are labeled as *intolerant* in a semantic subcategory of "{cluster_name}":

{examples_str; randomly chosen at most 50 examples}


Following only these labeled examples as implicit rules to ascertain what gets labeled as *intolerant*, would you label the following post as 'intolerant' or 'not intolerant' and give reasoning based on already labeled content:


Post Text: {target_text}


Your label ('intolerant' or 'not intolerant'):

# An anecdote

- LLM Label : Intolerant



bluenomatterwho69.bsky.social
@bluenomatterwho69.bsky.social                + Follow

There are only 2 genders. Male or female. You can be confused all you want but you are one or the other.

20 March 2025 at 16:30    Everybody can reply

- Reason: The post asserts that there are only two genders, male and female, and implies that anyone who identifies otherwise is "confused." This language is similar to other posts that have been labeled as intolerant, such as "Non-binary does not compute, biologically; it is ILLOGICAL. There are ONLY TWO GENDERS, despite the inane claims of the lgbqt alphabet Gestapo." The tone is also dismissive and condescending, implying that the author's view is the only correct one.

# Consistency in intolerance clusters

| Cluster | # LLM Intolerant | # Has Intolerant | Consistency (%) |
|---|---|---|---|
| cluster_0 *mental illness insults* | 67 | 0 | 0.00 |
| cluster_1 *anti-LGBTQ sentiment* | 34 | 2 | 5.88 |
| cluster_2 *anti-semitic hate speech* | 65 | 4 | 6.15 |
| cluster_3 *anti-Russian sentiment* | 73 | 0 | 0.00 |
| cluster_4 *abusive language and slurs* | 37 | 0 | 0.00 |
| cluster_5 *racial slurs and hate speech* | 73 | 0 | 0.00 |
| cluster_6 *only two genders exist* | 29 | 3 | 10.34 |
| cluster_7 *anti-Trump insults and profanity* | 93 | 2 | 2.15 |
| cluster_8 *derogatory slurs* | 7 | 0 | 0.00 |
| cluster_9 *Elon Musk is a retard* | 49 | 1 | 2.04 |
| cluster_10 *autism slurs and insults* | 16 | 0 | 0.00 |
| cluster_11 *anti-Catholic sentiment* | 26 | 1 | 3.85 |
| cluster_12 *anti-Islamic sentiment* | 49 | 2 | 4.08 |
| cluster_13 *anti-LGBTQ+ slurs* | 87 | 2 | 2.30 |
| cluster_14 *biological sex matters in sports* | 20 | 2 | 10.00 |
| cluster_15 *trans women are not women* | 18 | 1 | 5.56 |
| cluster_16 *name-calling and Nazi comparisons* | 95 | 0 | 0.00 |
| cluster_17 *derogatory insults* | 28 | 1 | 3.57 |
| cluster_18 *derogatory Trump supporters* | 89 | 1 | 1.12 |
| cluster_19 *abusive name-calling* | 91 | 12 | 13.19 |

# 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.
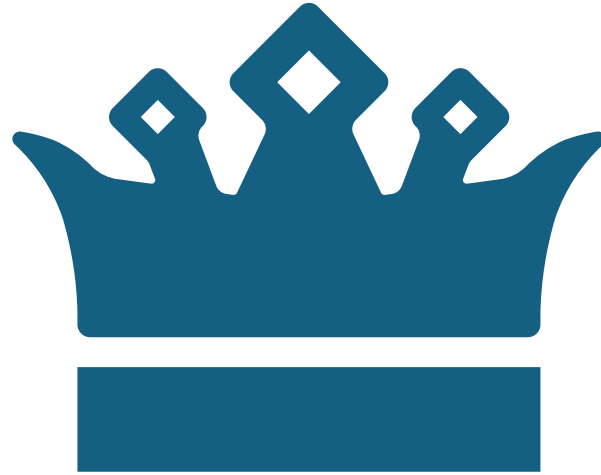
# Systemic flaws still exist

**Opaque by design**

↓

**Opacity has decreased**

**Centralized power**

↓

**Implementation is centralized**

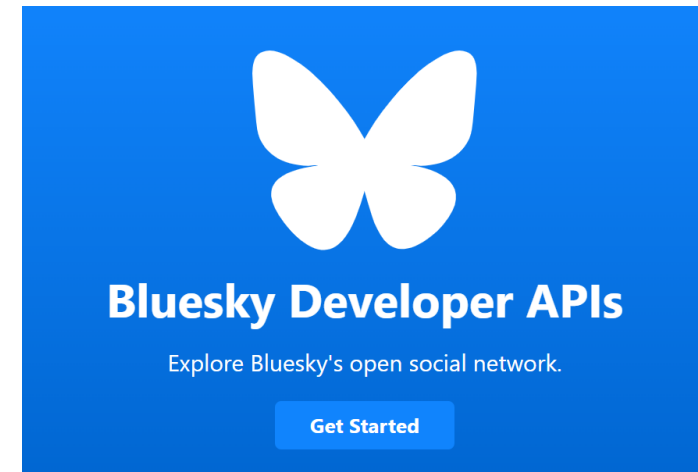**Zero user agency**

↓

**Better user agency**

**Inconsistent enforcement**

↓

**Now auditable**

43

# 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**

**Colab notebook to explore**