# Analyzing Norm Violations in Real-Time Live-Streaming Chat

Jihyung Moon*, **Dong-Ho Lee***, Hyundong Cho, Woojeong Jin, Chan Young Park, Minwoo Kim, Jonathan May, Jay Pujara, Sungjoon Park
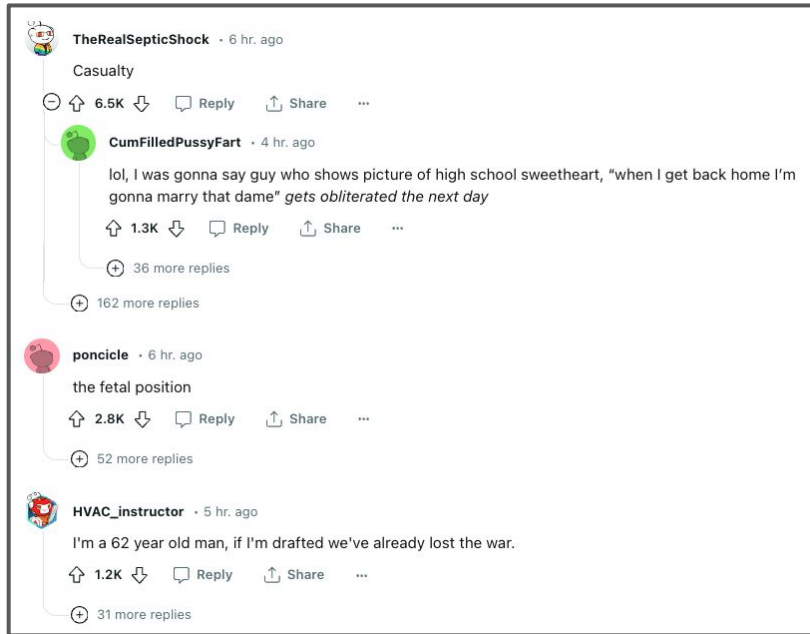
# Reddit vs. Twitch conversations

## Reddit



## Twitch



2

# **Asynchronous** vs. **Synchronous** conversations

## Reddit (Asynchronous)



## Twitch (Synchronous)

# **Asynchronous** vs. **Synchronous** conversations

## Reddit (Asynchronous)



## Twitch (Synchronous)



- Smaller message length.
- Smaller temporal gap between chats.
- Relationships between chats are less clearly defined.

4

# Existing Toxicity Detection Models

Asynchronous

A: □ □ □ □ □ □ □ □ □ □ □   *02:23*
  └ B: □ □ □ □ □ □ □ □ □ □   *11:07*
  └ C: □ □ □ □ □ □ □ □ □   *18:23*

Train

Perspective
OpenAI
Delphi

Existing models primarily trained on toxicity-labeled data
in the **asynchronous** domain.

# What if we use such models on **Synchronous** domain?



Asynchronous — Synchronous

A: □ □ □ □ □ □ □ □ □ □ □ *02:23*
└ B: □ □ □ □ □ □ □ □ □ *11:07*
└ C: □ □ □ □ □ □ □ □ □ *18:23*

A : □ □ □ *02:23:10*
B: □ □ □ □ *02:23:22*
C: □ □ *02:23:57*
B: □ □ □ *02:24:03*

Train

**Perspective**
**OpenAI**  Delphi

?

Test

Can these models perform well in
detecting violations also in the **synchronous** domain?

# Data Construction



**Twitch Top 200 Streamers**

Gather
chat rules

CHAT RULES
be nic
I will not be disrespectful to fellow chatters.
I will be a good memer in chat.
I will not rudely tell myth what to do.
I will have a good time. :)

**Okay, Got It!**

Define Norms

Map rules

| Norm | Rule Examples |
|---|---|
| Incivility | Be Nice |
| Discrimination | No Racism |
| … | … |

# Violation norms in **Twitch**

| Coarse | Fine-grained | Target | Rule Examples |
|---|---|---|---|
| Discrimination | Discrimination | - | No racism, sexism or homophobia. |
| HIB (Harassment, Intimidation, Bullying) | HIB | Broadcaster OIB OOB | No HIB towards broadcaster<br>No HIB towards viewers, moderators, etc.<br>No HIB towards other broadcasters, politicians, etc. |
| Privacy | Doxing | - | Please no personal questions about me.<br>Do not share any personal information about yourself or others. |
| Inappropriate Contents | NSFW<br>Self-destructive<br>Illegal<br>Spoiler | -<br>-<br>-<br>- | No NSFW content (e.g., Inappropriate ASCII arts).<br>No talk of suicide.<br>No drug discussion of any kind.<br>Do not give game spoilers. |
| Off Topic | Controversial Topic<br>Begging | -<br>- | No drama, politics or religion.<br>No begging for subscriptions or money. |
| Spam | Excessive & Repetitive<br>Advertisements | -<br>- | No walls of text.<br>No self promotion unless authorized. |
| Meta-Rules (Live streaming specific) | Backseating & Tall order<br>Mentioning other broadcasters<br>Specific language only | -<br>-<br>- | Don't tell me what to do. Don't ask for mod.<br>Don't talk down on other streamers.<br>English only. |
| Incivility (Miscellaneous) | Incivility | - | Be nice, Be civil |

# Violation norms in Twitch vs Reddit

| Coarse | Fine-grained |
|---|---|
| Discrimination | Discrimination |
| HIB (Harassment, Intimidation, Bullying) | HIB |
| Privacy | Doxing |
| Inappropriate Contents | NSFW Self-destructive Illegal Spoiler |
| Off Topic | Controversial Topic Begging |
| Spam | Excessive & Repetitive Advertisements |
| Meta-Rules (Live streaming specific) | Backseating & Tall order Mentioning other broadcasters Specific language only |
| Incivility (Miscellaneous) | Incivility |

## Reddit (NormVio, Park et al., 2021)

| | |
|---|---|
| **Incivility**: {Personality} | *"Be civil"* |
| **Harassment**: {Harassment, Doxxing} | *"Don't harass others"* |
| **Spam**:{Spam, Reposting, Copyright} | *"No excessive posting"* |
| **Format**: {Format, Images, Links} | *"Use the correct tags"* |
| **Content**:{Low-quality Content, NSFW, Spoilers} | *"No low-quality posts"* |
| **Off-topic** :{Off-topic, Politics} | *"Only relevant posts"* |
| **Hate speech**:{Hatespeech} | *"No racism, sexism"* |
| **Trolling**:{Trolling, Personal Army} | *"No trolls or bots"* |
| **Meta-rules**:{Voting, Moderation Enforcement, Reddiquette} | *"No Downvoting"* |

Detecting Community Sensitive Norm Violations in Online Conversations., **Park** et al., EMNLP 2021 Findings
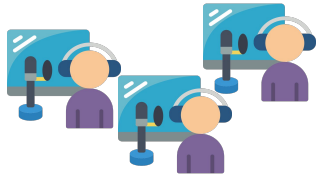
# Violation norms in Twitch vs Reddit

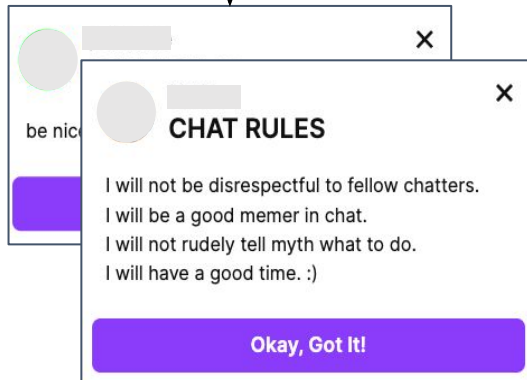| Trolling: {Trolling, Personal Army} | *"No trolls or bots"* |
|---|---|
| Format: {Format, Images, Links} | *"Use the correct tags"* |

| Coarse | Fine-grained |
|---|---|
| Discrimination | Discrimination |
| HIB (Harassment, Intimidation, Bullying) | HIB |
| Privacy | Doxing |
| Inappropriate Contents | NSFW Self-destructive Illegal Spoiler |
| Off Topic | Controversial Topic Begging |
| Spam | Excessive & Repetitive Advertisements |
| Meta-Rules (Live streaming specific) | Backseating & Tall order Mentioning other broadcasters Specific language only |
| Incivility (Miscellaneous) | Incivility |

| | |
|---|---|
| **Hate speech**:{Hatespeech} | *"No racism, sexism"* |
| **Harassment**: {Harassment, Doxxing} | *"Don't harass others"* |
| **Content**:{Low-quality Content, NSFW, Spoilers} | *"No low-quality posts"* |
| **Off-topic** :{Off-topic, Politics} | *"Only relevant posts"* |
| **Spam**:{Spam, Reposting, Copyright} | *"No excessive posting"* |
| **Meta-rules**:{Voting, Moderation Enforcement, Reddiquette} | *"No Downvoting"* |
| **Incivility**: {Personality} | *"Be civil"* |

Detecting Community Sensitive Norm Violations in Online Conversations., **Park** et al., EMNLP 2021 Findings

# Data Construction

Twitch API & IRC

Monitor Moderation Events
(Ban, Timeout, Delete)

**Annotators**

**Twitch Top 200 Streamers**

Gather
chat rules

CHAT RULES

be nic

I will not be disrespectful to fellow chatters.
I will be a good memer in chat.
I will not rudely tell myth what to do.
I will have a good time. :)

Okay, Got It!

Define Norms

Map rules

| Norm | Rule Examples |
|---|---|
| Incivility | Be Nice |
| Discrimination | No Racism |
| … | … |

# Data Construction



Twitch API & IRC

Monitor Moderation Events
(Ban, Timeout, Delete)

**Twitch Top 200 Streamers**

**Annotators**

Gather
chat rules

| Stage 1 | Moderated Message |
| Stage 2 | Moderated Message + Chat Context |
| Stage 3 | Moderated Message + Chat / Video Context |

CHAT RULES

be nic

I will not be disrespectful to fellow chatters.
I will be a good memer in chat.
I will not rudely tell myth what to do.
I will have a good time. :)

**Okay, Got It!**

Define Norms

Map rules

| **Norm** | **Rule Examples** |
|---|---|
| Incivility | Be Nice |
| Discrimination | No Racism |
| … | … |

12

**Annotation Interface –** Stage 1: Moderated Message

**Video Context** ⌃

**Bad Utterance** ⌃

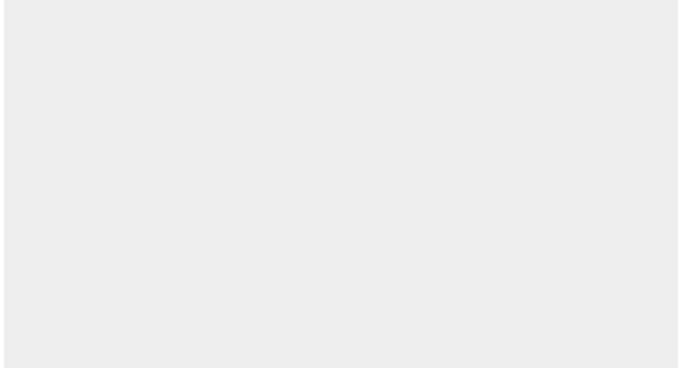██████ (██████) [2022-06-14 10:55:20]: **check this out chatr, playing league without hand - r3deyejedi_za**

**Channel Setting** ⌃

██████ (██████)

CAN THIS MAN BE SO GOOD WHILE BEING 80 YEARS OLD? COME AND SEE, SOLOQ, CLIMBING BACK TO CHALLENGER !YT !discord

League of Legends  English

**Single user context** ⌃

**Multi user context** ⌃

13

**Annotation Interface –** Stage 2: Moderated Message + Chat Context

**Video Context** ⌃

**Bad Utterance** ⌃

█████ (████████) [2022-06-14 10:55:20]: **check this out chatr, playing league without hand - r3deyejedi_za**

**Channel Setting** ⌃

████ (████)

CAN THIS MAN BE SO GOOD WHILE BEING 80 YEARS OLD? COME AND SEE, SOLOQ, CLIMBING BACK TO CHALLENGER !YT !discord
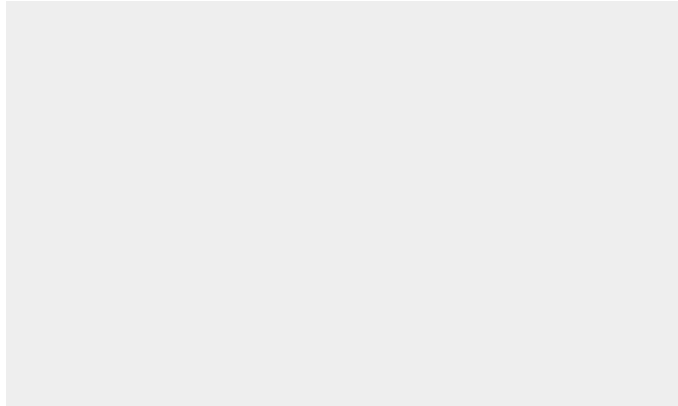
League of Legends `English`

**Single user context** ⌃

█████ (████████) [2022-06-14 10:54:42]: check this out chatr, playing league without hand - r3deyejedi_za {ID: bb4a8330-5b65-4f71-b08a-a166f8122a41}

█████ (████████) [2022-06-14 10:55:20]: **check this out chatr, playing league without hand - r3deyejedi_za** {ID: c37d60df-88ac-4198-9b83-838956f04e86}

████ ) [2022-06-14 10:55:37]: [ban]

**Multi user context** ⌃

RowdyJinx (rowdyjinx) [2022-06-14 10:54:36]: D:

████████████ [2022-06-14 10:54:36]: D:

█████████ [2022-06-14 10:54:37]: true

████████ [2022-06-14 10:54:37]: me too

███████████ [2022-06-14 10:54:37]: KEKW

██████ [2022-06-14 10:54:37]: D:

███████ [2022-06-14 10:54:38]: yes

█████████ [2022-06-14 10:54:40]: D:

███████ [2022-06-14 10:54:40]: D:

██████████ [2022-06-14 10:54:40]: @Jankos scene

███████████ [2022-06-14 10:54:42]: check this out chatr, playing league without hand - r3deyejedi_za

█████████ [2022-06-14 10:54:42]: D:

**Annotation Interface – Stage 3: Moderated Message + Chat/Video Context**



**Video Context**

**Bad Utterance**

[REDACTED] ([REDACTED]) [2022-06-14 10:55:20]: **check this out chatr, playing league without hand - r3deyejedi_za**

**Channel Setting**

[REDACTED] ([REDACTED])

CAN THIS MAN BE SO GOOD WHILE BEING 80 YEARS OLD? COME AND SEE, SOLOQ, CLIMBING BACK TO CHALLENGER !YT !discord

League of Legends    English

**Single user context**

[REDACTED] ([REDACTED]) [2022-06-14 10:54:42]: check this out chatr, playing league without hand - r3deyejedi_za {ID: bb4a8330-5b65-4f71-b08a-a166f8122a41}

[REDACTED] ([REDACTED]) [2022-06-14 10:55:20]: **check this out chatr, playing league without hand - r3deyejedi_za** {ID: c37d60df-88ac-4198-9b83-838956f04e86}

[REDACTED] ) [2022-06-14 10:55:37]: [ban]

**Multi user context**

RowdyJinx (rowdyjinx) [2022-06-14 10:54:36]: D:

[REDACTED] [2022-06-14 10:54:36]: D:

[REDACTED] [2022-06-14 10:54:37]: true

[REDACTED] [2022-06-14 10:54:37]: me too

[REDACTED] [2022-06-14 10:54:37]: KEKW

[REDACTED] [2022-06-14 10:54:37]: D:

[REDACTED] [2022-06-14 10:54:38]: yes

[REDACTED] [2022-06-14 10:54:40]: D:

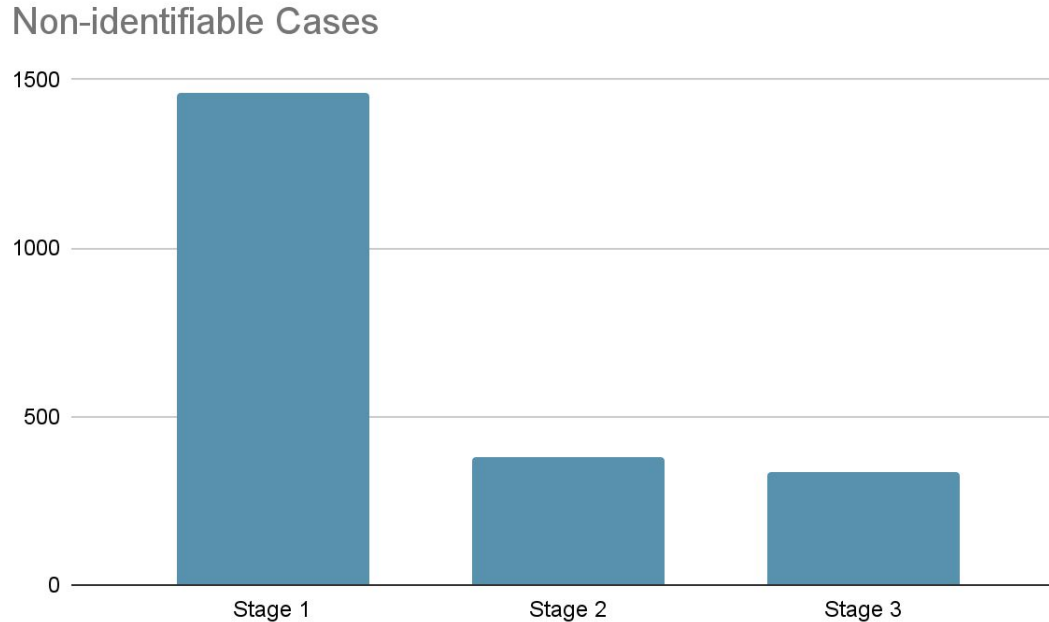[REDACTED] [2022-06-14 10:54:40]: D:

[REDACTED] [2022-06-14 10:54:40]: @Jankos scene

[REDACTED] [2022-06-14 10:54:42]: check this out chatr, playing league without hand - r3deyejedi_za

[REDACTED] [2022-06-14 10:54:42]: D:
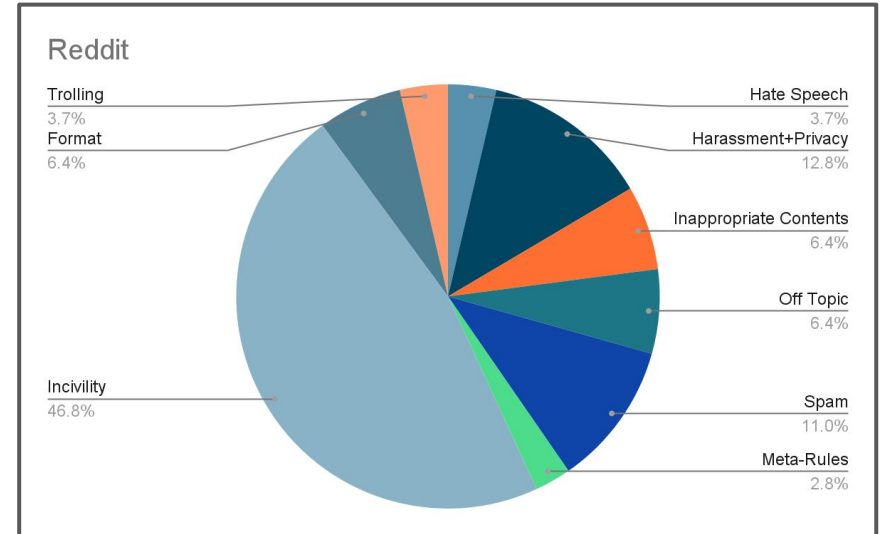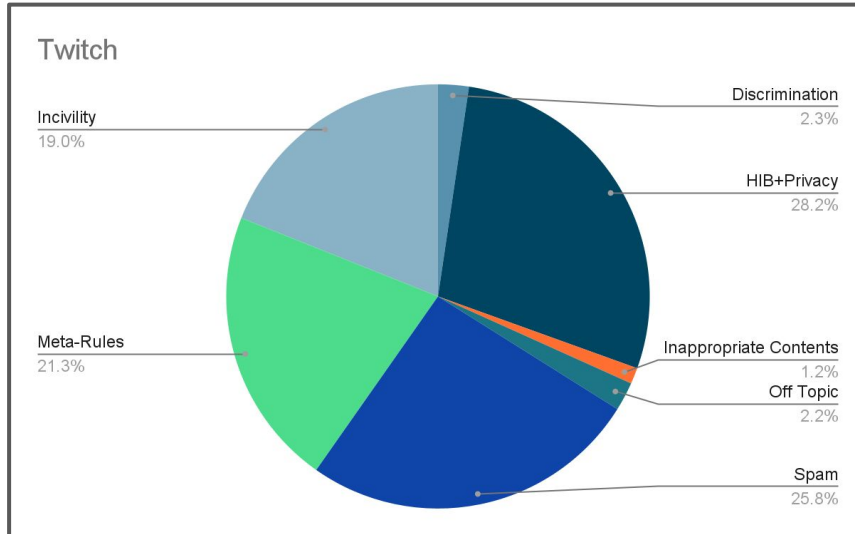
# Data Statistics

Number of non-identifiable cases



Non-identifiable Cases

**Providing context** will <u>lessen the ambiguity of the violation</u>.

# Data Statistics

Twitch vs. Reddit

Reddit (NormVio, Park et al., 2021)



**Twitch**

- Discrimination 2.3%
- HIB+Privacy 28.2%
- Incivility 19.0%
- Meta-Rules 21.3%
- Inappropriate Contents 1.2%
- Off Topic 2.2%
- Spam 25.8%



**Reddit**

- Trolling 3.7%
- Format 6.4%
- Hate Speech 3.7%
- Harassment+Privacy 12.8%
- Inappropriate Contents 6.4%
- Off Topic 6.4%
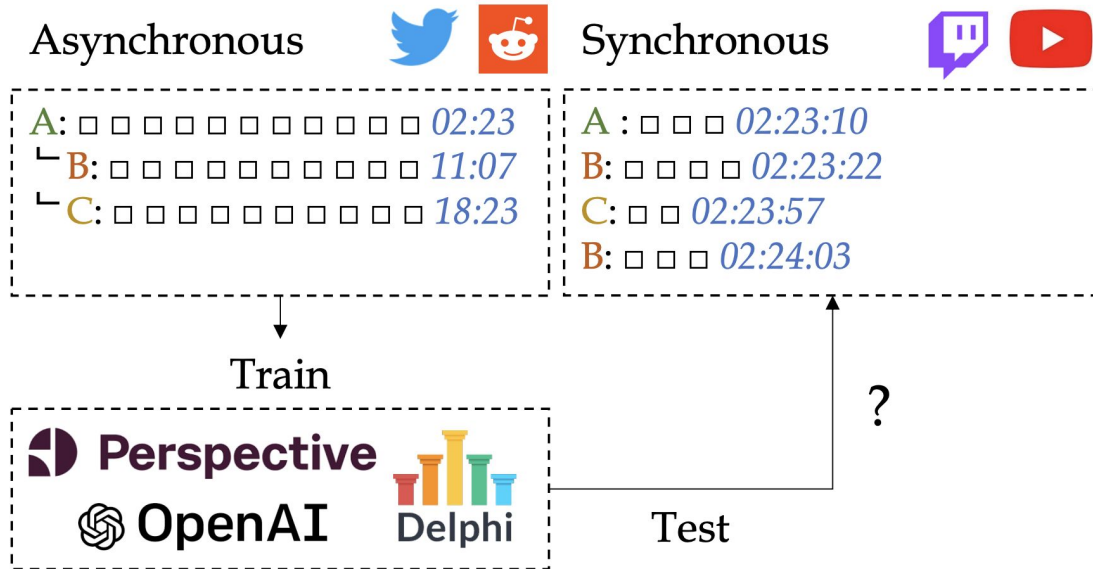- Incivility 46.8%
- Spam 11.0%
- Meta-Rules 2.8%

**Similarity**
(1) **Harassment** (HIB) and **Incivility** take up a large portion.
(2) **Off-topic, Inappropriate contents, Privacy** exist but less enforced.

**Difference**
(1) **Spam**, **Meta-Rules** cover significantly higher portion in Twitch.
(2) Twitch streamers are less concerned about **contents**.

17

# Performance of existing model endpoints.

Recap



Can these models perform well in
detecting violations also in the **synchronous** domain?

# Performance of existing model endpoints.

*Binary classification on (Discrimination & HIB)*

| Model | Precision | Recall | F1 |
|---|---|---|---|
| ToxiGen | 0.31 | 0.91 | 0.46 |
| Perspective API | 0.39 | 0.95 | 0.56 |
| OpenAI moderation | 0.11 | 0.94 | 0.20 |
| OpenAI content filter | 0.55 | 0.86 | 0.67 |

Existing models do not frequently produce false positives (**high recall**).
However, they perform poorly in detecting toxic messages (**low precision**).

# Norm Classification

Binary classification for each category

**+** *Chat history*

*You are XXX*

Violating **'discrimination'** or not?

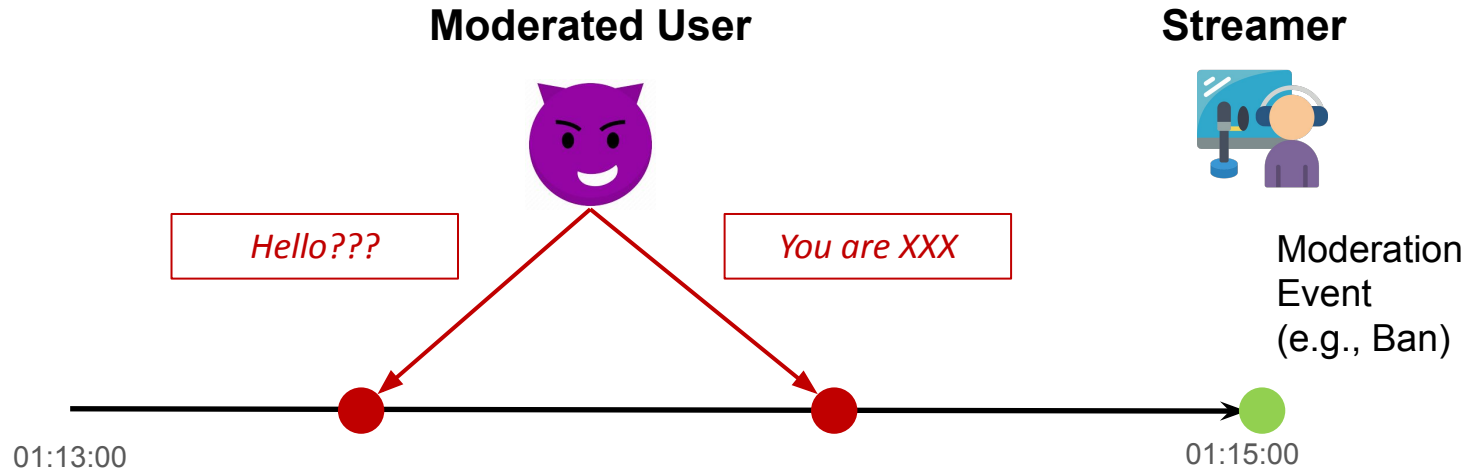Violating **'harassment'** or not?

⋮

Violating **'spam'** or not?

# Norm Classification
Context

# Norm Classification
Single-user Context



**Moderated User**

**Streamer**

**Single-user context** → *Hello???* → *You are XXX*

Moderation Event (e.g., Ban)

01:13:00                                                              01:15:00

# Norm Classification
Multi-user Context

# Norm Classification

Performance

**All**: Any norm violated? or not?

| Context | All | HIB | Off Topic | Spam | Meta-Rules | Incivility |
|---|---|---|---|---|---|---|
| - | 0.70 | 0.52 | 0.07 | 0.63 | **0.65** | 0.28 |
| Single-user context | 0.78 | 0.50 | 0.05 | 0.67 | 0.58 | 0.28 |
| Multi-user context (event) | 0.75 | 0.44 | 0.05 | 0.66 | 0.60 | 0.17 |
| Multi-user context (utterance) | 0.91 | **0.61** | 0.10 | 0.66 | **0.65** | 0.24 |
| Multi-user context (first) | **0.95** | **0.61** | 0.08 | **0.70** | 0.62 | **0.45** |

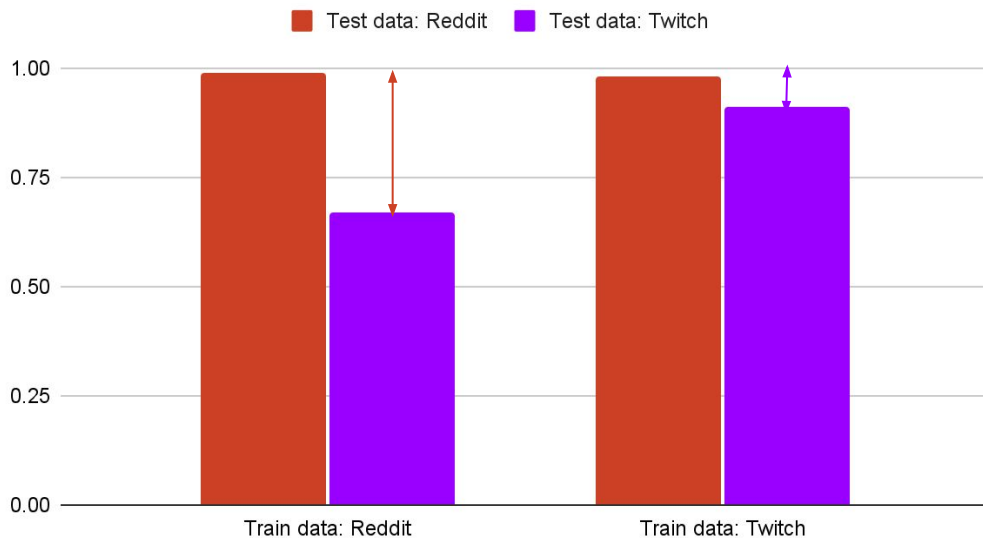**Context matters a lot for training better norm classification model!**

- Chats from other users (multi-user context) help determine the toxicity
- However, temporal gap between event and the actual offending chat may be substantial

# Norm Classification

Train on Reddit → Test on Twitch / Train on Twitch → Test on Reddit



Binary classification on "ALL" category

Model trained on Twitch **shows better generalizability** than model trained on Reddit

# Conclusion

1.  Norm violations vary between synchronous platforms like Twitch and asynchronous ones like Reddit

2.  **Context is crucial** for improving models that detect chat norm violations.

3.  A model trained on Twitch data may be more effective in identifying real-world norm violations.
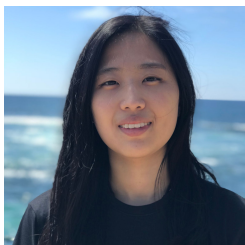
# Q&A


Dong-Ho Lee


Jihyung Moon


Sungjoon Park


Hyundong Cho


Woojeong Jin


Chanyoung Park


Jonathan May


Jay Pujara

**Data labeling sponsor:** SELECTSTAR