# Conversation—AI

Lucy Vasserman

# Contents

We're a unit within Alphabet that builds technology to make the world safer. Our team tackles a range of global security issues including defending against digital attacks, mitigating the rise of online hate and harassment, countering online extremism, and fighting censorship.
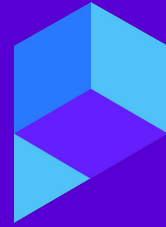
# Conversation AI Effort

# Mission
*Protect voices in conversation*

# Our work
API, tools, and research

Perspective

# Abuse and toxicity have led people to give up on conversations.
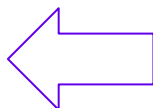
# Voices
# are silenced

# People
# are siloed
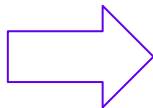
People stop expressing themselves and the
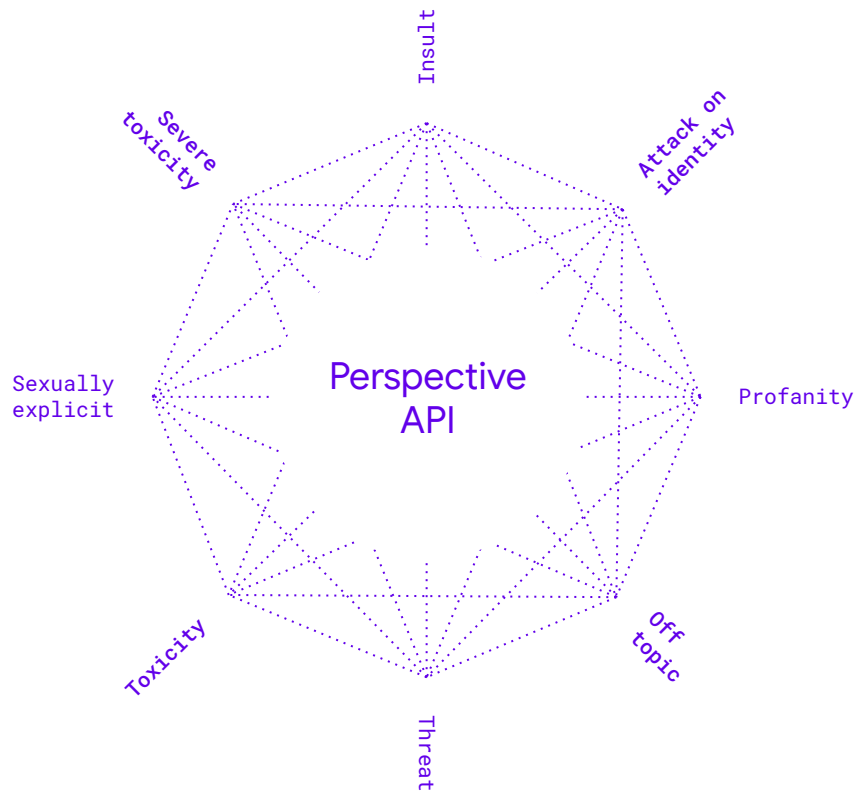loudest voices shout other out of the room .

By optimizing for likes/shares platforms create filter
bubbles so that people who disagree don't interact, or
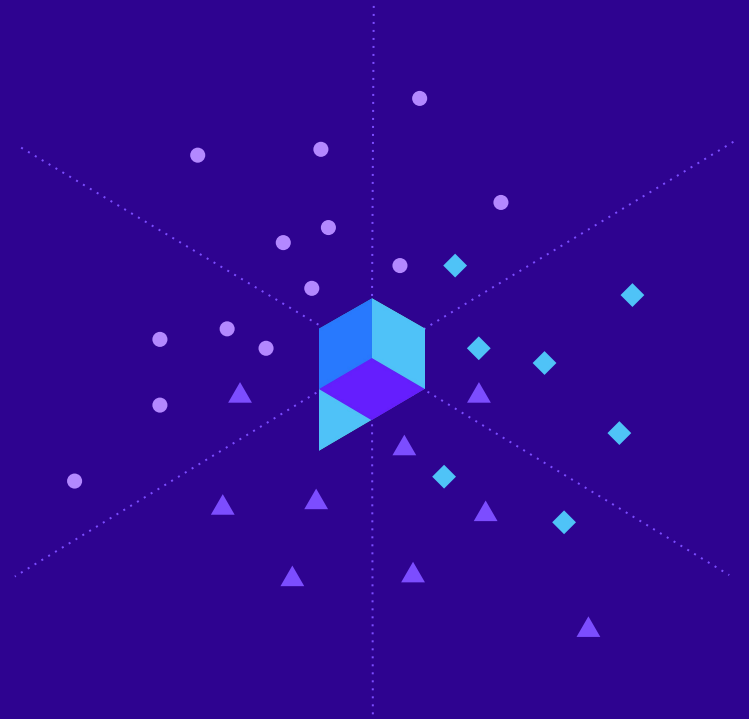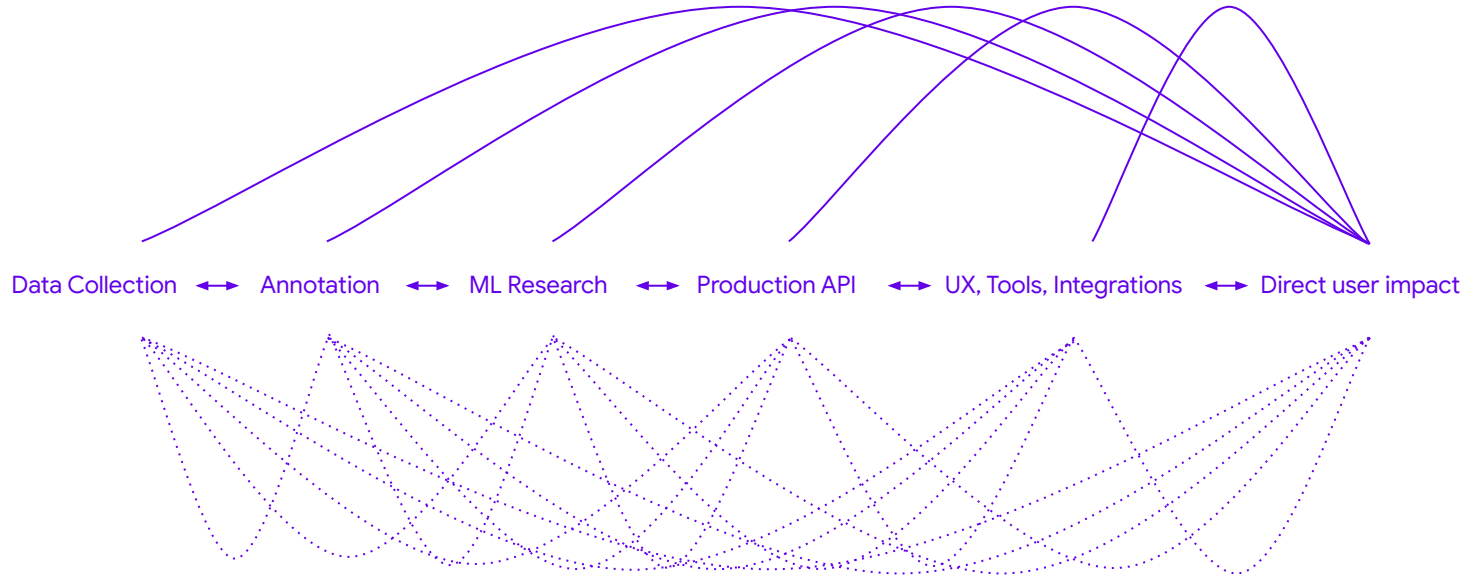they shut down comments and discussion all together.

# Perspective aims to classify the emotional impact of language.

*Is this a rude, disrespectful, or unreasonable comment that is likely to make you leave a discussion?*

# Outputs



Data Collection ⟷ Annotation ⟷ ML Research ⟷ Production API ⟷ UX, Tools, Integrations ⟷ Direct user impact

# Success Metrics

↗

**Participation**

Increase in the number of voices in a discussion

↗

**Diversity**

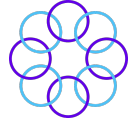Increase diversity of voices in a discussion

↘

**Toxicity**
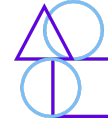
Reduce the prevalence of toxic comments online

↗

**Action**

Increase in action against toxicity across ecosystem

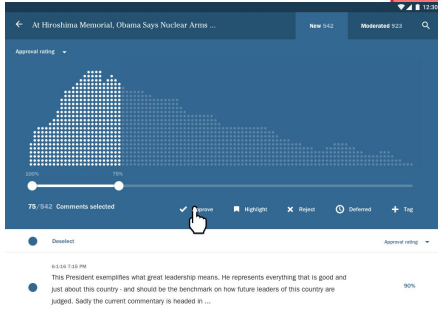How we work
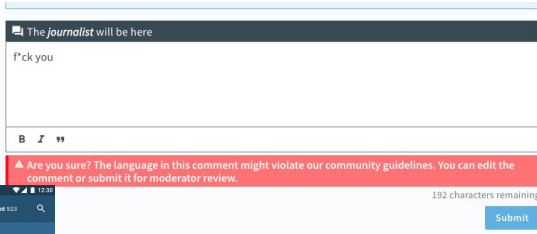
# Values

Community

Topic Neutrality

Transparency

Privacy

Inclusivity

# Experiences



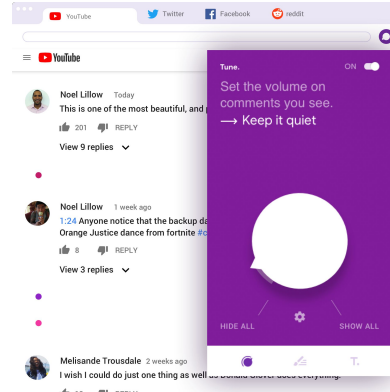## Moderation

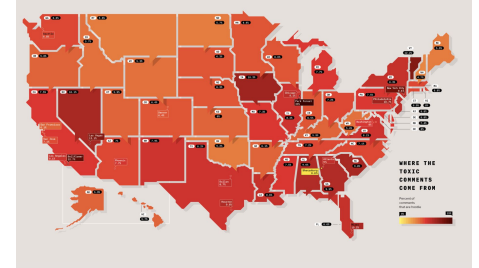Help community managers set rules and review comments faster.

## Authorship

Help people understand the impact of what they are writing.

## Readership

Help people discover the conversations that interest them.

## Visual trends

Help creators build data visualizations to better understand conversations at scale.

# Transparency

# Public Demo

Having an easy to use public demo has enabled us to find and fix problems

# Model Cards

A **Model Card** is a documentation framework that outlines:

- Evaluation results
- Intended usage
- Insight into training processes

The False Positive - Medium Blog
Conversation AI - Jigsaw

# Unintended Bias

# False "toxic" positives

A naively trained model will have some strong unintended biases illustrated by these false-positive examples...

| Comment | Toxicity score |
| --- | --- |
| The Gay and Lesbian Film Festival starts today. | 0.82 |
| Being transgender is independent of sexual orientation. | 0.52 |
| A Muslim is someone who follows or practices Islam. | 0.46 |

# Bias Mitigation

*Bias caused by dataset imbalance*
- Frequently attacked identities are overrepresented in toxic comments
- Length matters

Add *assumed non-toxic data* from Wikipedia articles to fix the imbalance.
- Original dataset had 127,820 examples
- 4,620 non-toxic examples added

| Term | Comment Length | | | | |
|---|---|---|---|---|---|
| | 20-59 | 60-179 | 180-539 | 540-1619 | 1620-4859 |
| ALL | 17% | 12% | 7% | 5% | 5% |
| gay | 88% | 77% | 51% | 30% | 19% |
| queer | 75% | 83% | 45% | 56% | 0% |
| homosexual | 78% | 72% | 43% | 16% | 15% |
| black | 50% | 30% | 12% | 8% | 4% |
| white | 20% | 24% | 16% | 12% | 2% |
| wikipedia | 39% | 20% | 14% | 11% | 7% |
| atheist | 0% | 20% | 9% | 6% | 0% |
| lesbian | 33% | 50% | 42% | 21% | 0% |
| feminist | 0% | 20% | 25% | 0% | 0% |
| islam | 50% | 43% | 12% | 12% | 0% |
| muslim | 0% | 25% | 21% | 12% | 17% |
| race | 20% | 25% | 12% | 10% | 6% |
| news | 0% | 1% | 4% | 3% | 3% |
| daughter | 0% | 7% | 0% | 7% | 0% |

# How can we measure unintended bias?

Definitions

- **Unintended bias** exists if the model performance varies across different subgroups

- **Subgroups** are the identities mentioned in the text (not the identities of the author or recipient)

Metrics

- Metrics should be threshold independent

# Measuring Overall Model Performance - AUC

*How good is the model at distinguishing toxic from non-toxic examples? (ROC-AUC)*

AUC (for a given test set) = Given two randomly chosen examples, one in-class (e.g. one is toxic and the other is not), AUC is the probability that the model will give the in-class example the higher score.

# Measuring Overall Model Performance - AUC

*How good is the model at distinguishing toxic from non-toxic examples? (ROC-AUC)*

AUC (for a given test set) = Given two randomly chosen examples, one in-class (e.g. one is toxic and the other is not), AUC is the probability that the model will give the in-class example the higher score.



Toxic Comments

Non-toxic Comments

# Subgroup AUC

Measures low subgroup performance.

Detects if the model performs worse on subgroup comments than it does on comments overall.



background

subgroup

model_score

Toxic Comments

Non-toxic Comments

# Subgroup AUC

Measures low subgroup performance.

Detects if the model performs worse on subgroup comments than it does on comments overall.

# Background Positive Subgroup Negative (BPSN) AUC

Measures subgroup shifts to the right

Detects if the model systematically scores comments from the subgroup higher.

# Background Positive Subgroup Negative (BPSN) AUC

Measures subgroup shifts to the right

Detects if the model systematically scores comments from the subgroup higher.

# Background Positive Subgroup Negative (BPSN) AUC

Measures subgroup shifts to the left.

Detects if the model systematically scores comments from the subgroup lower.



background

subgroup

model_score

Toxic Comments

Non-toxic Comments

# Background Positive Subgroup Negative (BPSN) AUC

Measures subgroup shifts to the left.

Detects if the model systematically scores comments from the subgroup lower.

# Evaluation on synthetic data

*Synthetic data shows real improvement!*

Comments are generated using simple templates

text: *"I am <identity>"*
label: *non-toxic*

text: *"I hate <identity>"*
label: *toxic*

# Public dataset for bias research

**~2 million comments** released by Civil Comments platform

Annotated for toxicity (all)
>    *Is this a rude, disrespectful, or unreasonable comment that is likely to make you leave a discussion?*

Annotated for identity content (~360k)
>    *What genders are mentioned in this comment?*
>    *What races or ethnicities are mentioned in this comment?*
>    *etc...*

# Evaluation on real data

*Real data shows mixed results*

TOXICITY@1

| | subgroup_auc | bpsn_auc | bnsp_auc |
|---|---|---|---|
| male | 0.88 | 0.9 | 0.92 |
| female | 0.89 | 0.91 | 0.92 |
| transgender | 0.82 | 0.87 | 0.9 |
| homosexual_gay_or_lesbian | 0.8 | 0.79 | 0.94 |
| christian | 0.9 | 0.94 | 0.9 |
| jewish | 0.84 | 0.87 | 0.92 |
| muslim | 0.82 | 0.86 | 0.92 |
| atheist | 0.88 | 0.91 | 0.91 |
| black | 0.81 | 0.85 | 0.92 |
| white | 0.81 | 0.84 | 0.92 |
| asian | 0.9 | 0.92 | 0.92 |
| latino | 0.85 | 0.89 | 0.9 |
| psychiatric_or_mental_illness | 0.89 | 0.89 | 0.93 |

TOXICITY@6

| | subgroup_auc | bpsn_auc | bnsp_auc |
|---|---|---|---|
| male | 0.89 | 0.91 | 0.93 |
| female | 0.89 | 0.91 | 0.93 |
| transgender | 0.8 | 0.91 | 0.87 |
| homosexual_gay_or_lesbian | 0.77 | 0.88 | 0.86 |
| christian | 0.89 | 0.94 | 0.9 |
| jewish | 0.84 | 0.86 | 0.93 |
| muslim | 0.81 | 0.88 | 0.91 |
| atheist | 0.89 | 0.91 | 0.92 |
| black | 0.82 | 0.79 | 0.96 |
| white | 0.82 | 0.82 | 0.95 |
| asian | 0.9 | 0.9 | 0.94 |
| latino | 0.86 | 0.87 | 0.94 |
| psychiatric_or_mental_illness | 0.88 | 0.89 | 0.94 |

# Evaluation on real data - short comments only

*The unintended bias was worse for short comments.*

*Bias mitigation brought performance on short comments closer to overall performance, but bias still exists.*
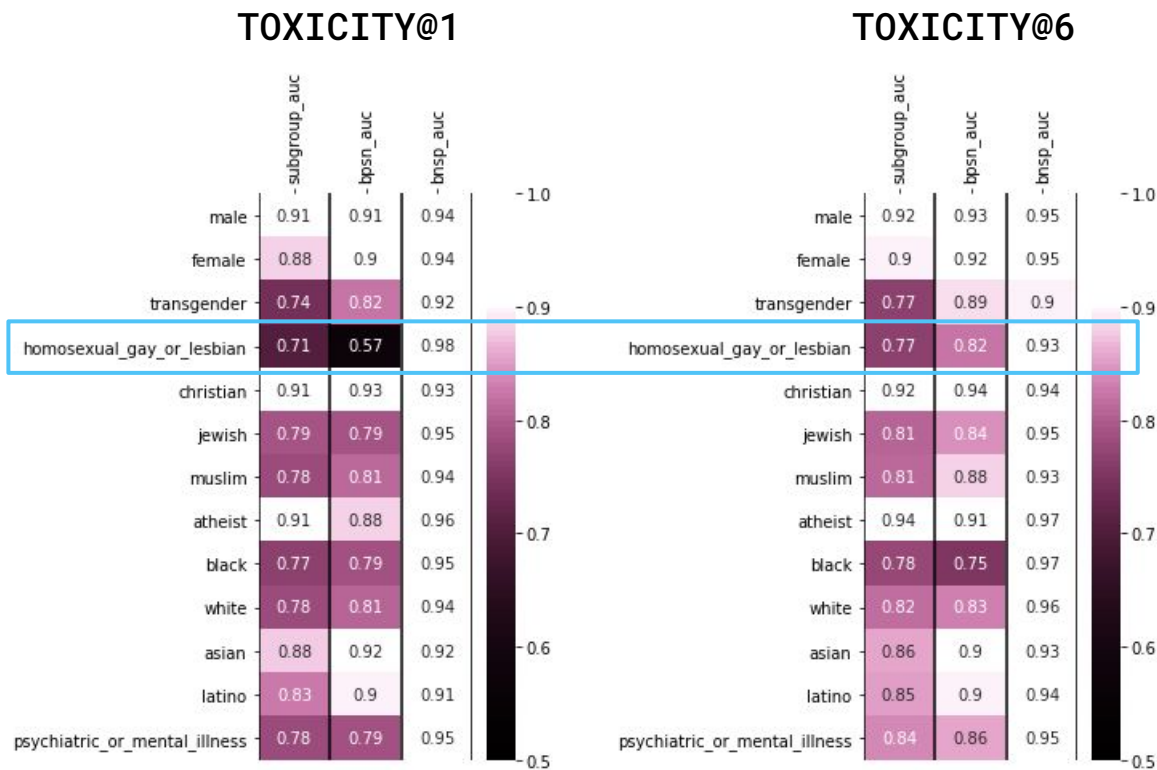
## TOXICITY@1

| | subgroup_auc | bpsn_auc | bnsp_auc |
|---|---|---|---|
| male | 0.91 | 0.91 | 0.94 |
| female | 0.88 | 0.9 | 0.94 |
| transgender | 0.74 | 0.82 | 0.92 |
| homosexual_gay_or_lesbian | 0.71 | 0.57 | 0.98 |
| christian | 0.91 | 0.93 | 0.93 |
| jewish | 0.79 | 0.79 | 0.95 |
| muslim | 0.78 | 0.81 | 0.94 |
| atheist | 0.91 | 0.88 | 0.96 |
| black | 0.77 | 0.79 | 0.95 |
| white | 0.78 | 0.81 | 0.94 |
| asian | 0.88 | 0.92 | 0.92 |
| latino | 0.83 | 0.9 | 0.91 |
| psychiatric_or_mental_illness | 0.78 | 0.79 | 0.95 |

## TOXICITY@6

| | subgroup_auc | bpsn_auc | bnsp_auc |
|---|---|---|---|
| male | 0.92 | 0.93 | 0.95 |
| female | 0.9 | 0.92 | 0.95 |
| transgender | 0.77 | 0.89 | 0.9 |
| homosexual_gay_or_lesbian | 0.77 | 0.82 | 0.93 |
| christian | 0.92 | 0.94 | 0.94 |
| jewish | 0.81 | 0.84 | 0.95 |
| muslim | 0.81 | 0.88 | 0.93 |
| atheist | 0.94 | 0.91 | 0.97 |
| black | 0.78 | 0.75 | 0.97 |
| white | 0.82 | 0.83 | 0.96 |
| asian | 0.86 | 0.9 | 0.93 |
| latino | 0.85 | 0.9 | 0.94 |
| psychiatric_or_mental_illness | 0.84 | 0.86 | 0.95 |

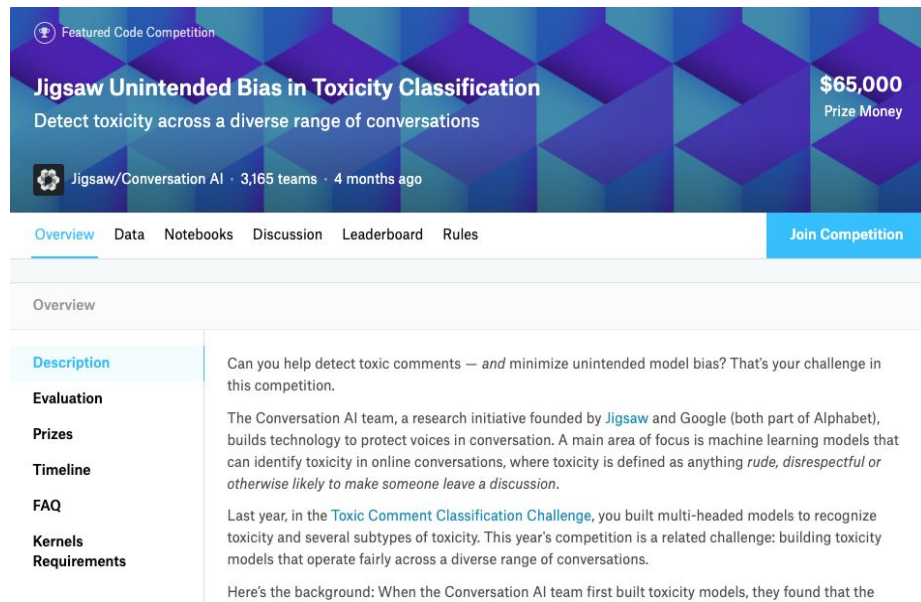Kaggle Competition

# Kaggle Competition

**Data**
2 million comments set from Civil Comments

**Evaluation**
Generalized mean of three bias AUCs for all identities and overall AUC

**Results**
3k+ teams researching bias mitigation techniques
Winners used BERT models and identity-aware data weighting

Questions?