

*Latent Hatred:*  
A Benchmark for Understanding  
Implicit Hate Speech

**Mai ElSherief,\* Caleb Ziems,\***

David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt,  
Munmun De Choudhury, Diyi Yang

**UC San Diego**  
JACOBS SCHOOL OF ENGINEERING

**Georgia Tech**  **College of Computing**

# Introduction: A Benchmark for Implicit Hate Speech

 **Content Warning:** may contain upsetting examples.

# Hate Speech Datasets

Basile et al. (2019)

Davidson et al. (2017)

Djuric et al. (2015)

Founta et al. (2018)

Burnap and Williams  
(2014)

Gao and Huang  
(2017)

Warner and  
Hirschberg (2012)

Waseem and Hovy  
(2016)

de Gibert et al. (2018)

Kennedy et al. (2018)

Sap et al. (2020)

Zampieri et al. (2019)

# Hate Speech Datasets

● Basile et al. (2019)

● Davidson et al. (2017)

Djuric et al. (2015)

● Founta et al. (2018)

Burnap and Williams  
(2014)

Gao and Huang  
(2017)

Warner and  
Hirschberg (2012)

● Waseem and Hovy  
(2016)

de Gibert et al. (2018)

Kennedy et al. (2018)

Sap et al. (2020)

Zampieri et al. (2019)

● Seeded w/ Hate Lexicon

# Hate Speech Datasets Are:



Explicit

● Basile et al. (2019)

● Davidson et al. (2017)

Djuric et al. (2015)

● Founta et al. (2018)

Burnap and Williams  
(2014)

Gao and Huang  
(2017)

Warner and  
Hirschberg (2012)

● Waseem and Hovy  
(2016)

de Gibert et al. (2018)

Kennedy et al. (2018)

Sap et al. (2020)

Zampieri et al. (2019)

● Seeded w/ Hate Lexicon

# Hate Speech Datasets Are:

 Explicit



 Seeded w/ Hate Lexicon



# Hate Speech Datasets Are:



Explicit

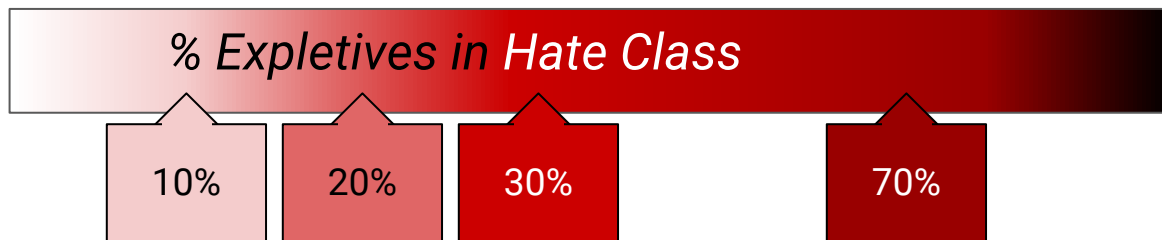


Targeted



● Seeded w/ Hate Lexicon

● Seeded w/ Racial Identifiers

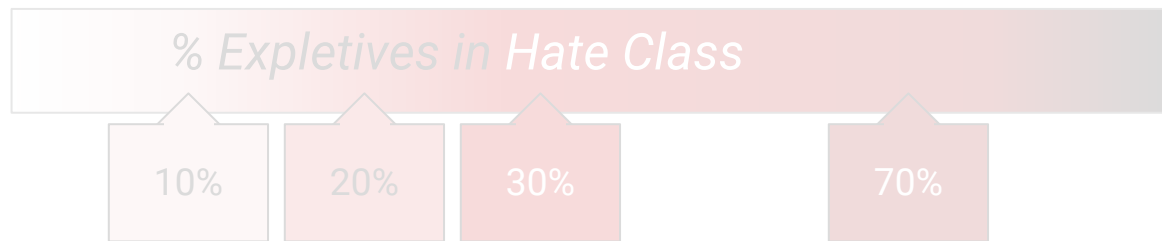


# Hate Speech Datasets Are:

Explicit  Targeted



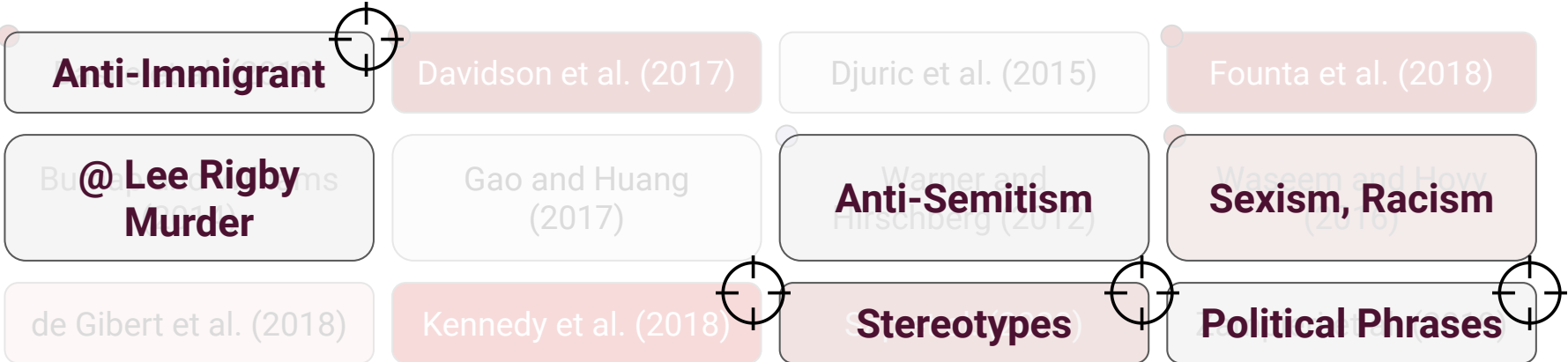
● Seeded w/ Hate Lexicon  
● Seeded w/ Racial Identifiers





# Hate Speech Datasets Are:

Explicit     Targeted



● Seeded w/ Hate Lexicon  
● Seeded w/ Racial Identifiers



# Hate Speech Datasets

## Hate Speech

### Explicit Hate Speech

Basile et al. (2019)

Davidson et al. (2017)

Djuric et al. (2015)

Founta et al. (2018)

Burnap and Williams  
(2014)

Gao and Huang (2017)

Warner and Hirschberg  
(2012)

Waseem and Hovy  
(2016)

de Gibert et al. (2018)

Kennedy et al. (2018)

Sap et al. (2020)

Zampieri et al. (2019)

# Hate Speech Datasets

## Hate Speech

### Explicit Hate Speech

Basile et al. (2019)

Davidson et al. (2017)

Djuric et al. (2015)

Founta et al. (2018)

Burnap and Williams  
(2014)

Gao and Huang (2017)

Warner and Hirschberg  
(2012)

Waseem and Hovy  
(2016)

de Gibert et al. (2018)

Kennedy et al. (2018)

Sap et al. (2020)

Zampieri et al. (2019)

### Implicit Hate Speech

???

# Hate Speech Datasets

## Hate Speech

### Explicit Hate Speech

Basile et al. (2019)

Davidson et al. (2017)

Djuric et al. (2015)

Founta et al. (2018)

Burnap and Williams  
(2014)

Gao and Huang (2017)

Warner and Hirschberg  
(2012)

Waseem and Hovy  
(2016)

de Gibert et al. (2018)

Kennedy et al. (2018)

Sap et al. (2020)

Zampieri et al. (2019)

### Implicit Hate Speech

Kennedy et al. (2018)

# Hate Speech Datasets

## Hate Speech

### **Explicit** Hate Speech

*(intentional)*  
*stereotyping,*  
*racism, sexism,*  
*targeted threats*

### **Implicit** Hate Speech

# Hate Speech Datasets

## Hate Speech

### **Explicit** Hate Speech

*(intentional)*

*stereotyping,  
racism, sexism,  
targeted threats*

### **Implicit** Hate Speech

*(intentionally harmful)*

*circumlocution, coded language, colloquialisms, connotations, dog  
whistles, entity framing, euphemisms, hidden threats, idioms, inferiority  
assumptions, irony, metaphors, presuppositions, symbolic language, ...*

# Hate Speech Datasets

## Hate Speech

### **Explicit** Hate Speech

*(intentional)*  
stereotyping,  
racism, sexism,  
targeted threats

### **Implicit** Hate Speech

*(intentionally harmful)*  
circumlocution, coded language, colloquialisms, connotations, dog  
whistles, entity framing, euphemisms, hidden threats, idioms, inferiority  
assumptions, irony, metaphors, presuppositions, symbolic language, ...

Bias (e.g. [Social Bias Frames](#), [PowerTransformer](#))

Microaggressions ([Breitfeller et al. 2019](#))

# Hate Speech Datasets

## Hate Speech

### Explicit Hate Speech

*(intentional)*  
stereotyping,  
racism, sexism,  
targeted threats

### Implicit Hate Speech

*(intentionally harmful)*  
circumlocution, coded language, inequalities, connotations, dog  
whistles, entity frames, animosities, hidden threats, idioms, inferiority  
assumptions, metaphors, presuppositions, symbolic language, ...

**[ MISSING FROM THE LITERATURE ]**

Bias (e.g. Social Bias Frames, PowerTransformer)

Microaggressions (Breitfeller et al. 2019)



# Hate Speech Datasets

## Hate Speech

### Explicit Hate Speech

*(intentional)*  
stereotyping,  
racism, sexism,  
targeted threats

### Implicit

*(intentionally not)*  
circumlocution, coded language, metaphors, presuppositions, dog  
whistles, entity frames, animosities, inferiority  
assumptions, metaphors, presuppositions, page, ...

Can we {understand, detect, automatically explain}  
this complex subset of hate speech?

[LITERATURE]







[MISSING]

Bias (e.g. Social Bias Frames, PowerTransformer)

Microaggressions (Breitfeller et al. 2019)

# Our Contributions


## 1. Implicit Hate Taxonomy

- a.  Incitement of Violence
- b.  Inferiority Language
- c.  Irony
- d.  Stereotypes and Misinformation
- e.  Threats and Intimidation
- f.  White Grievance

# Our Contributions

1. Implicit Hate Taxonomy
2. Implicit Hate Benchmark Dataset


# Our Contributions

1. Implicit Hate Taxonomy
2. Implicit Hate Benchmark Dataset
  - a.  hate target\*




# Our Contributions

1. Implicit Hate Taxonomy
2. Implicit Hate Benchmark Dataset

a.  hate target\*

b.  implied meaning\*

# Our Contributions


1. Implicit Hate Taxonomy
2. Implicit Hate Benchmark Dataset
  - a.  hate target\*
  - b.  implied meaning\*
  - c.  fine-grained implicit hate category


# Our Contributions

1. Implicit Hate Taxonomy

2. Implicit Hate Benchmark Dataset

a.  hate target\*

b.  implied meaning\*




c.  fine-grained implicit hate category

3. Baseline Classifiers for Detecting Implicit Hate



transformers

# Our Contributions

1. Implicit Hate Taxonomy
2. Implicit Hate Benchmark Dataset
  - a.  hate target\*
  - b.  implied meaning\*
  - c.  fine-grained implicit hate category
3. Baseline Classifiers for Detecting Implicit Hate
4. Generative LMs for Explaining Implicit Hate



GPT / GPT-2



# 1. Implicit Hate Taxonomy

# 1. Implicit Hate Taxonomy

Incitement of Violence

Inferiority Language

Irony

Stereotypes and Misinformation

Threats and Intimidation

White Grievance

# 1. Implicit Hate Taxonomy

## > Incitement of Violence

*"Hitler was Germany –Germans shall rise again!"*

Inferiority Language

Irony

Stereotypes and Misinformation

Threats and Intimidation

White Grievance

# 1. Implicit Hate Taxonomy

Incitement of Violence 

> Inferiority Language 

*"It's not a coincidence the best places to live are majority white."*

Irony

Stereotypes and Misinformation

Threats and Intimidation

White Grievance

# 1. Implicit Hate Taxonomy

Incitement of Violence 

Inferiority Language 

> Irony 

*“Horrors... Disney will be forced into hiring Americans”*

Stereotypes and Misinformation

Threats and Intimidation

White Grievance

# 1. Implicit Hate Taxonomy

Incitement of Violence 

Inferiority Language 

Irony 

> Stereotypes and Misinformation 

*“Can someone tell the black people in Chicago to stop killing one another before it becomes Detroit?”*

Threats and Intimidation

White Grievance

# 1. Implicit Hate Taxonomy

Incitement of Violence 

Inferiority Language 

Irony 

Stereotypes and Misinformation 

> Threats and Intimidation 

*"It won't be long before whitey is walking thru the slums, hoping and preying for a white victory."*

White Grievance

# 1. Implicit Hate Taxonomy

Incitement of Violence 

Inferiority Language 

Irony 

Stereotypes and Misinformation 

Threats and Intimidation 

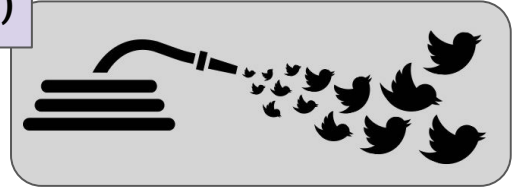
> White Grievance 

*"Black lives matter and white lives don't? Sounds racist."*

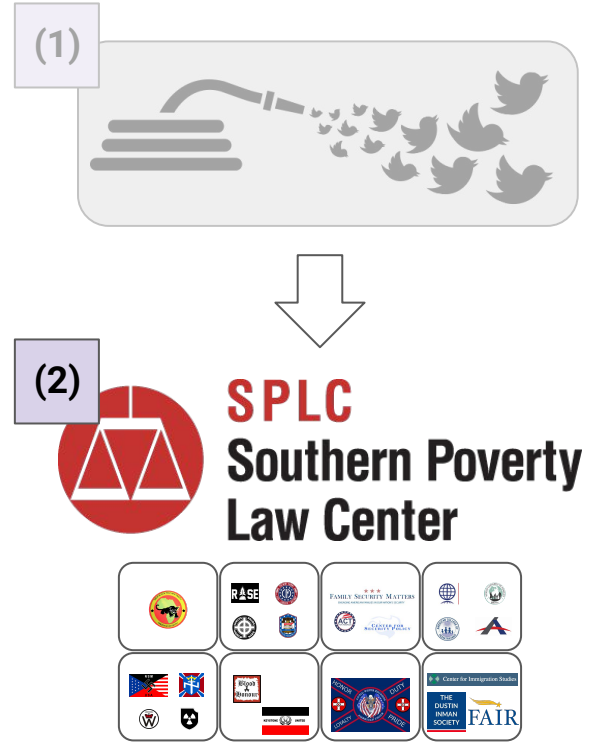


## 2. Implicit Hate Benchmark Dataset

(1)



## 2. Implicit Hate Benchmark Dataset



# 2. Implicit Hate Benchmark Dataset

**Black Separatist (27.1%)**



**White Nationalist (16.4%)**



**Anti-Muslim (8.9%)**



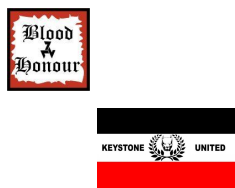
**Anti-LGBT (7.4%)**



**Neo-Nazi (6.2%)**



**Racist Skinhead (5.1%)**



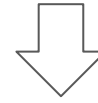
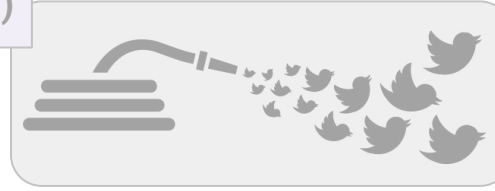
**KKK (5.0%)**



**Anti-Immigrant (2.1%)**



(1)



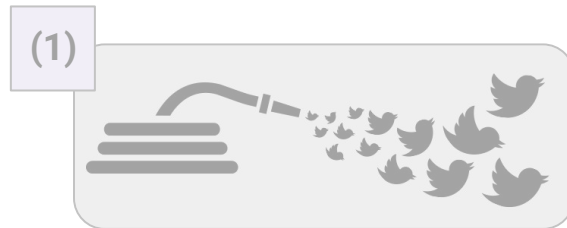
(2)



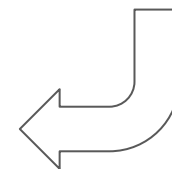
**SPLC**  
**Southern Poverty**  
**Law Center**



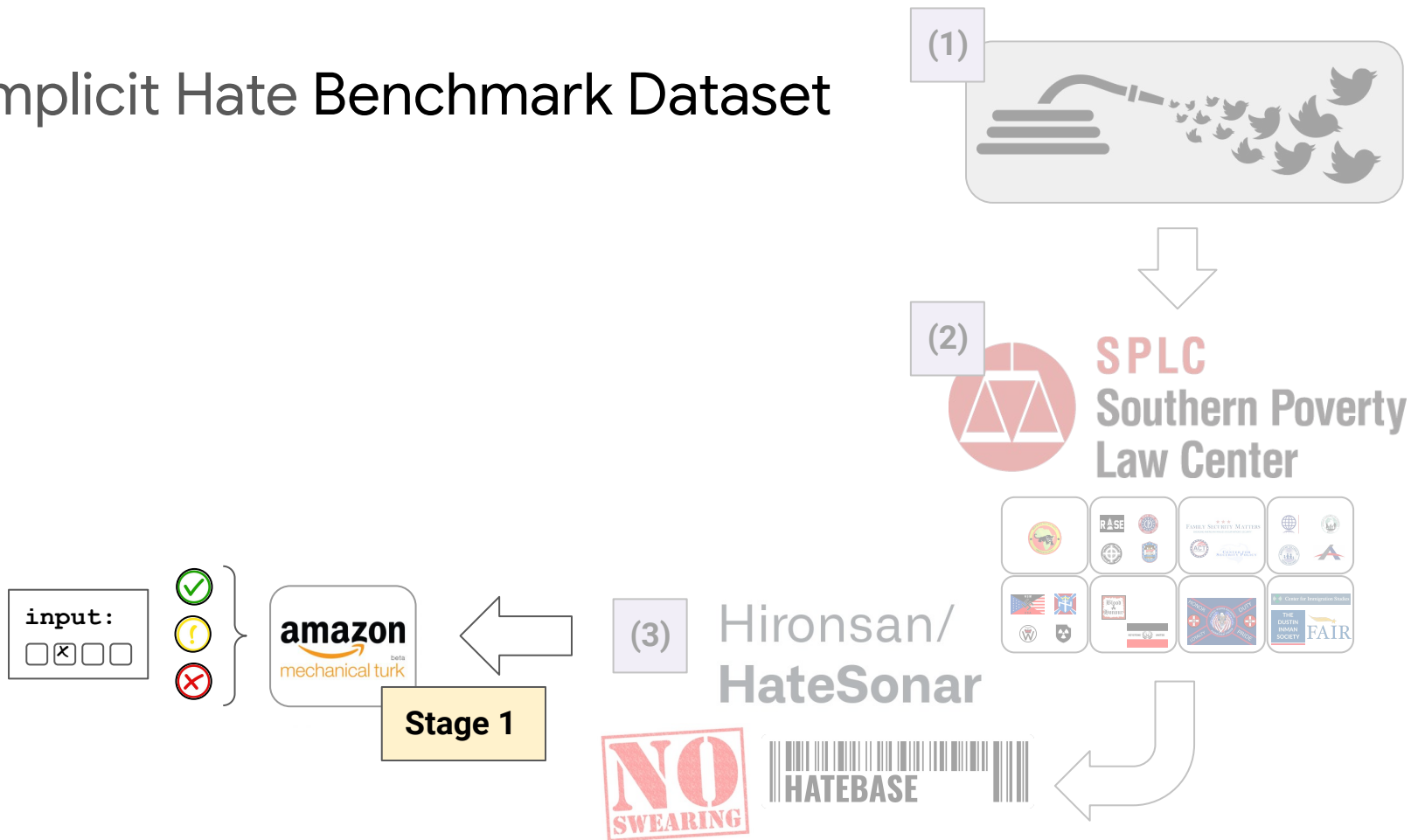
## 2. Implicit Hate Benchmark Dataset



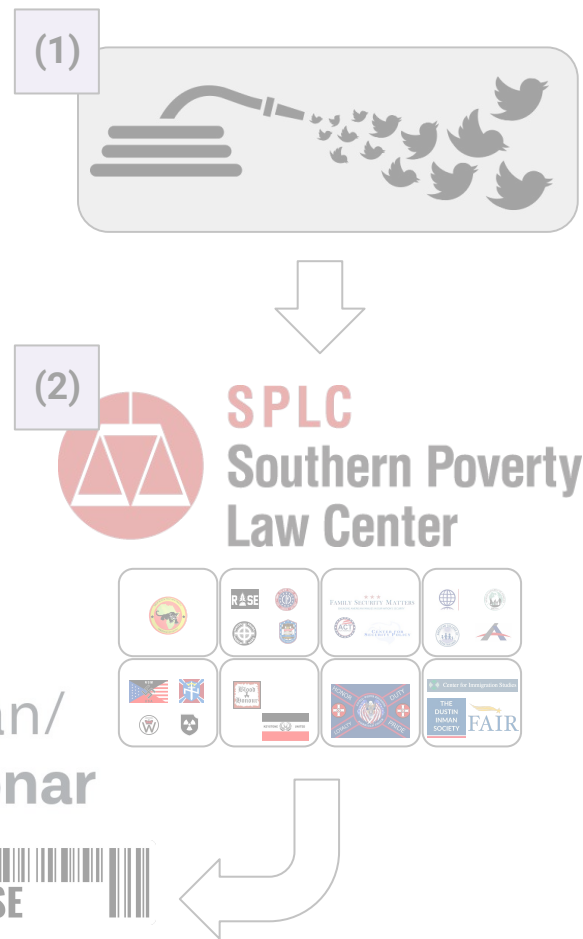
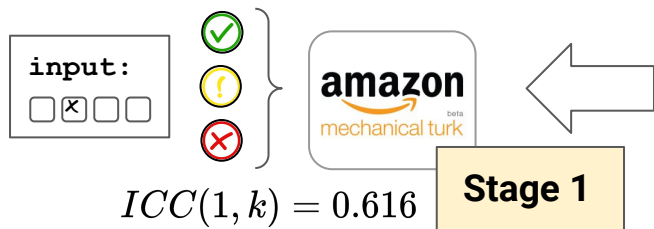
(3) Hironson/  
**HateSonar**



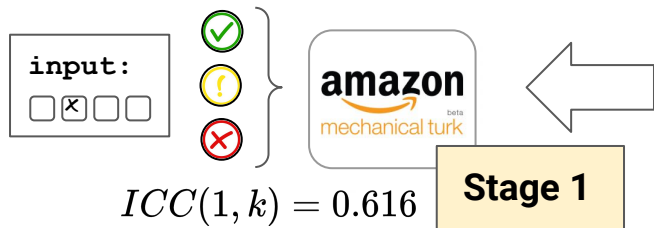
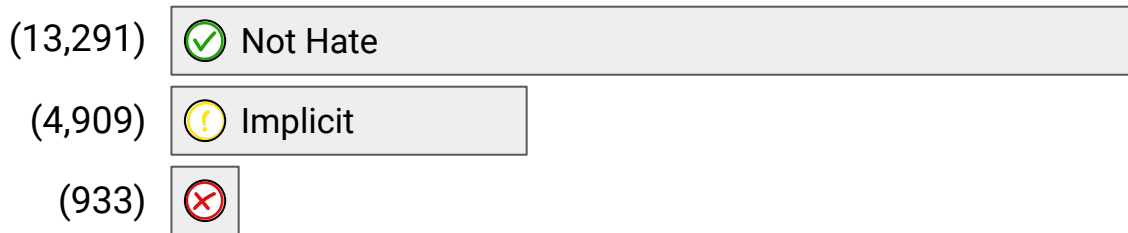
## 2. Implicit Hate Benchmark Dataset



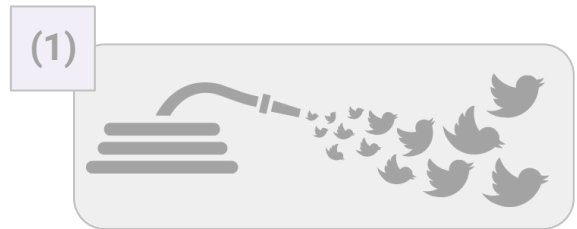
## 2. Implicit Hate Benchmark Dataset



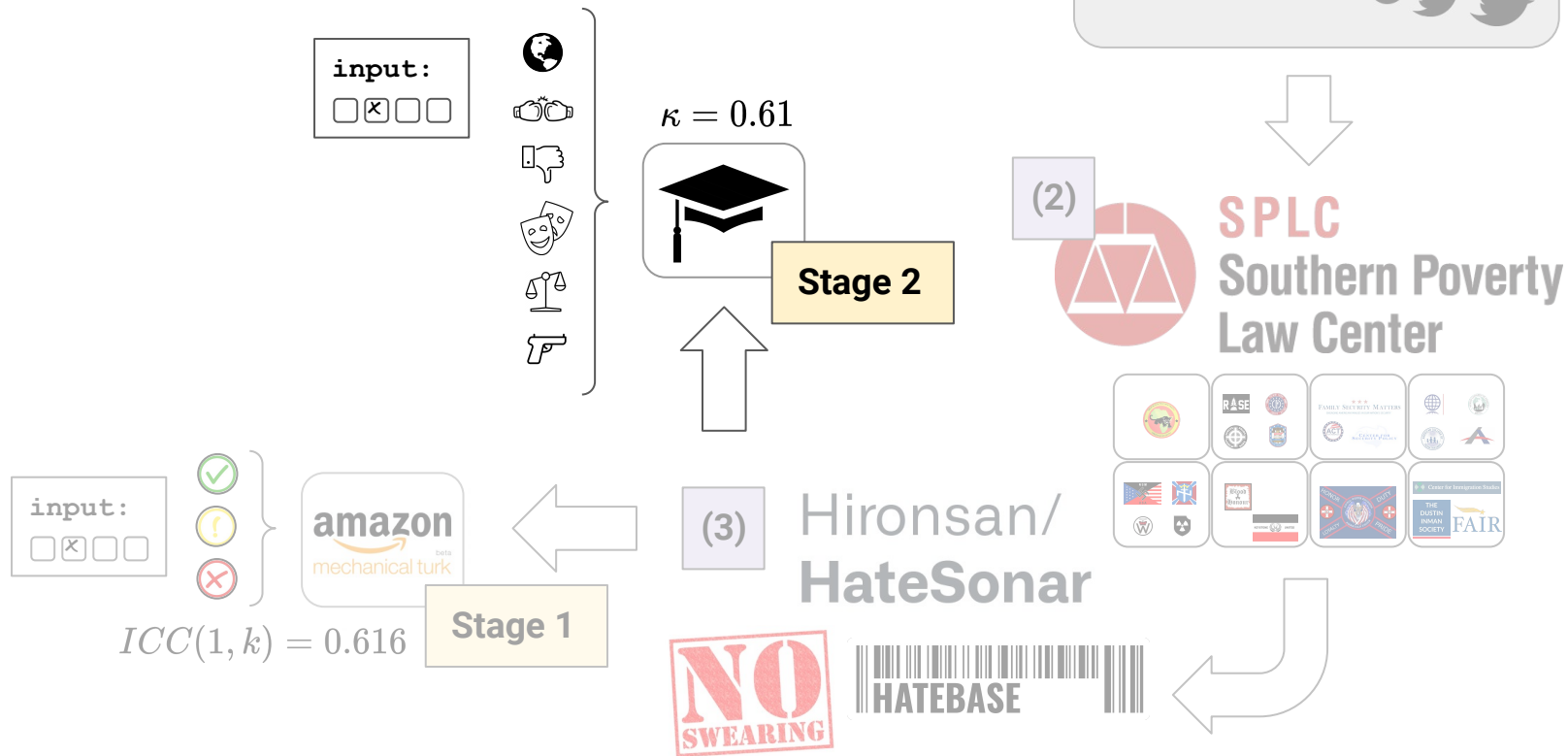
## 2. Implicit Hate Benchmark Dataset



(3) Hironson/  
HateSonar

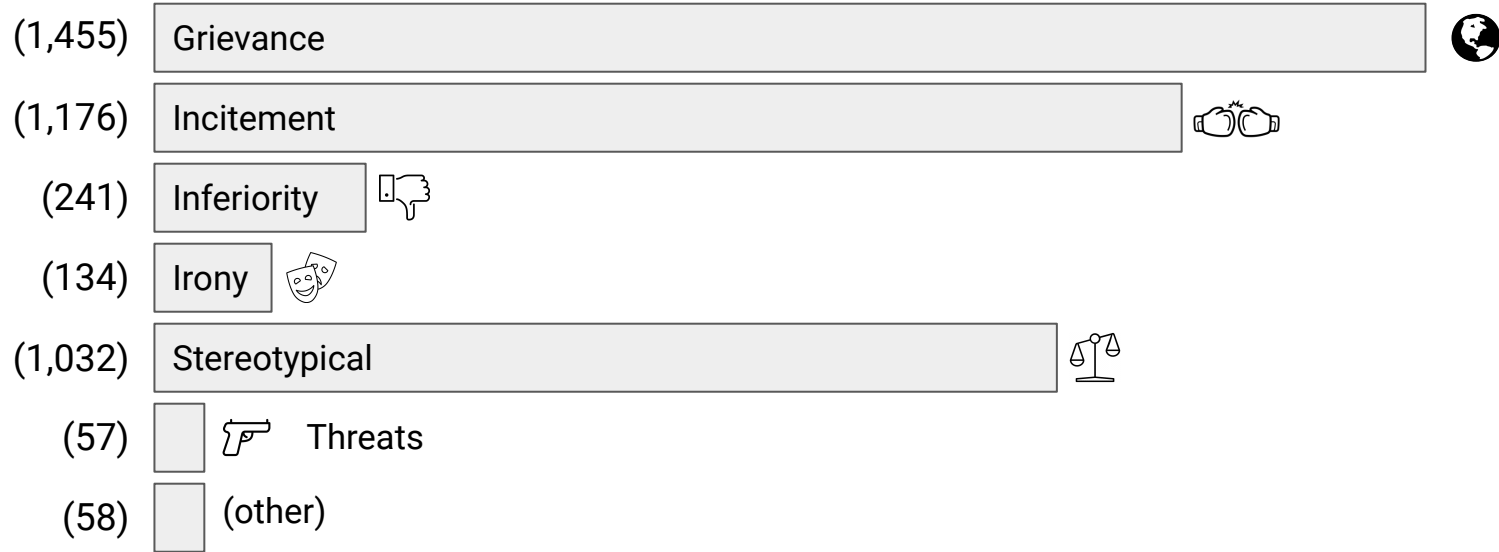


## 2. Implicit Hate Benchmark Dataset

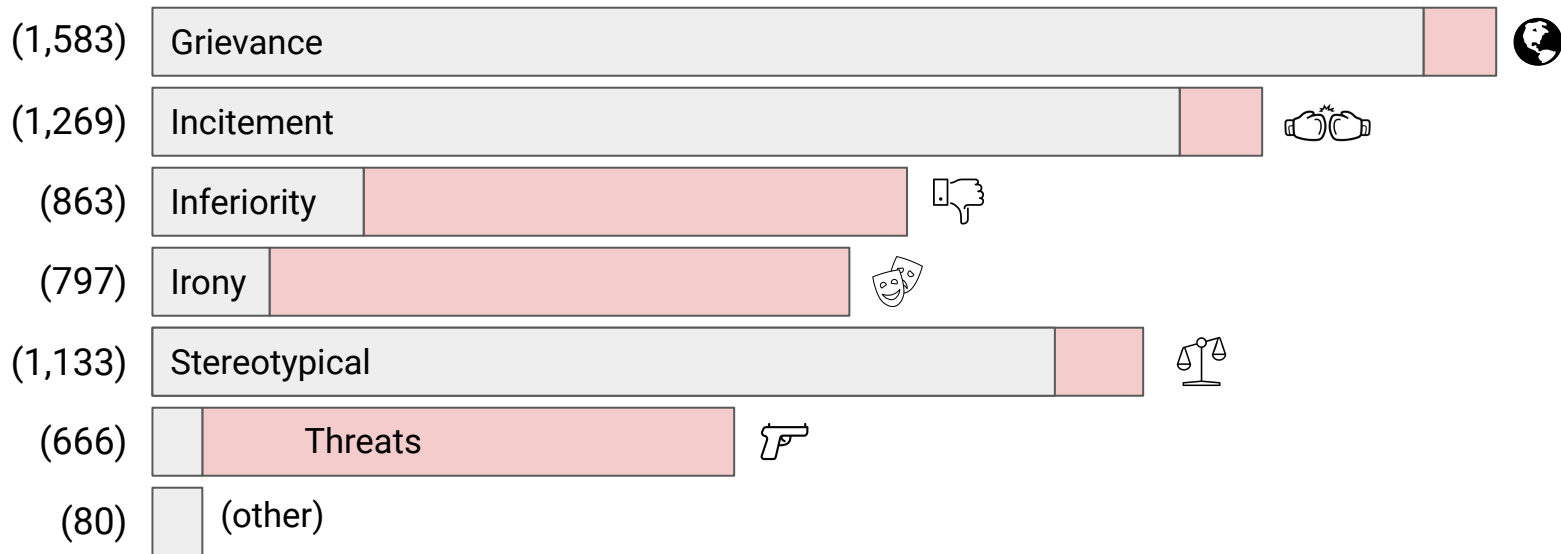




## 2. Implicit Hate Benchmark Dataset

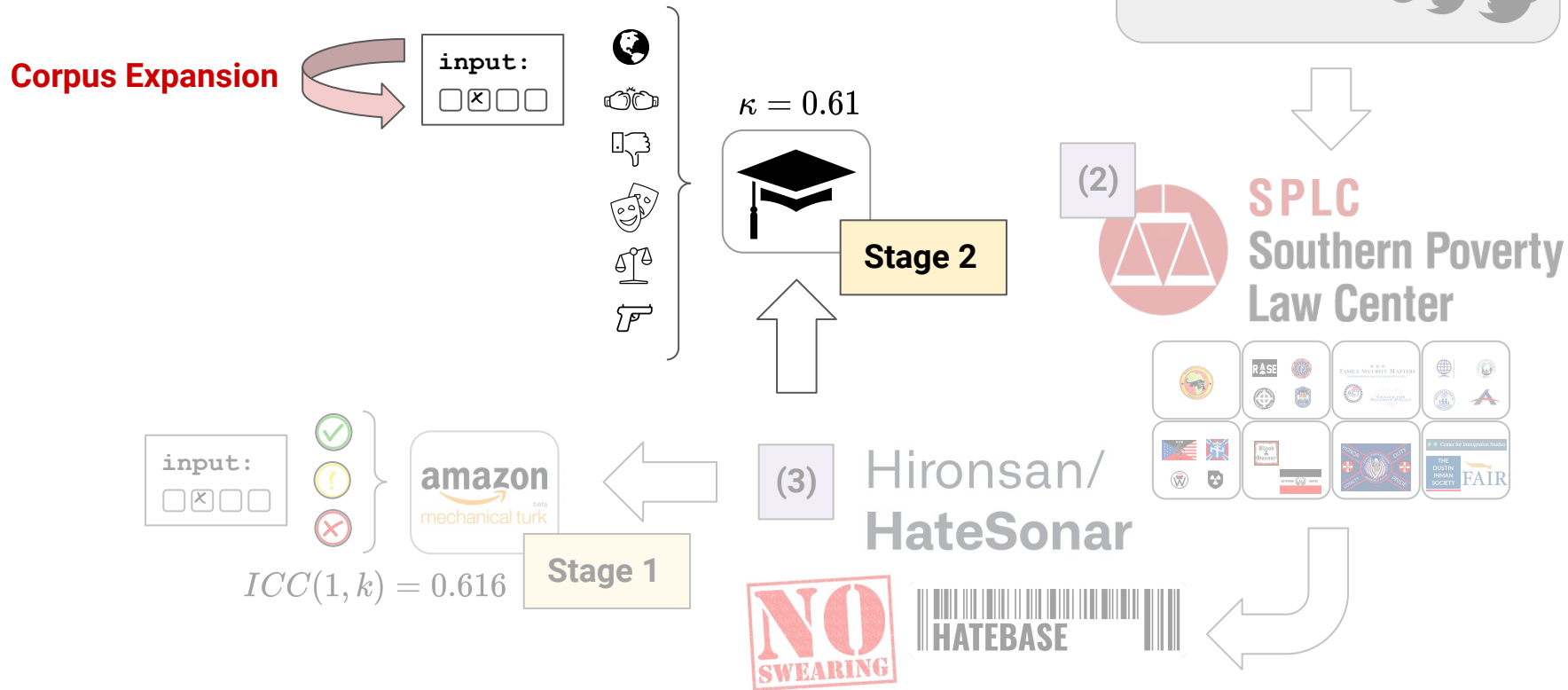


## 2. Implicit Hate Benchmark Dataset

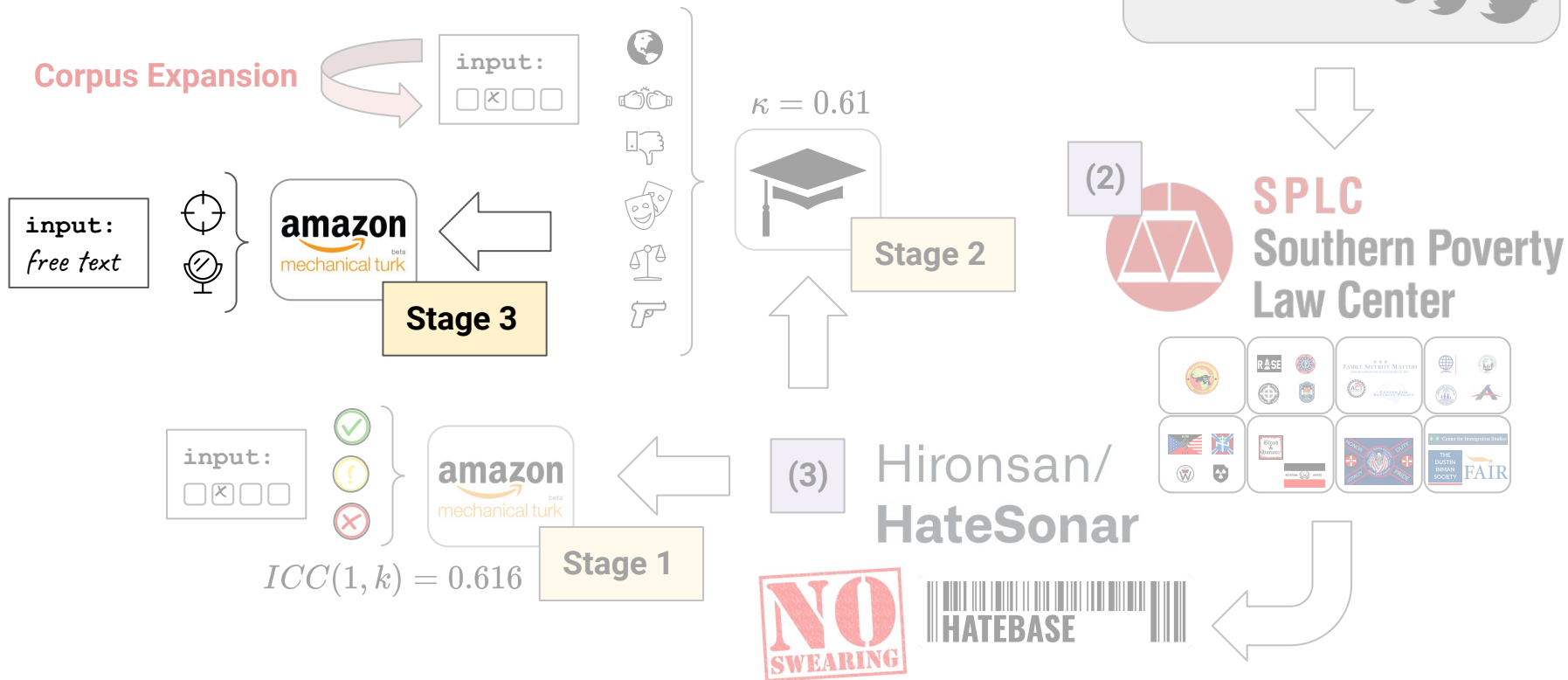


**6,346** Total After **Corpus Expansion**

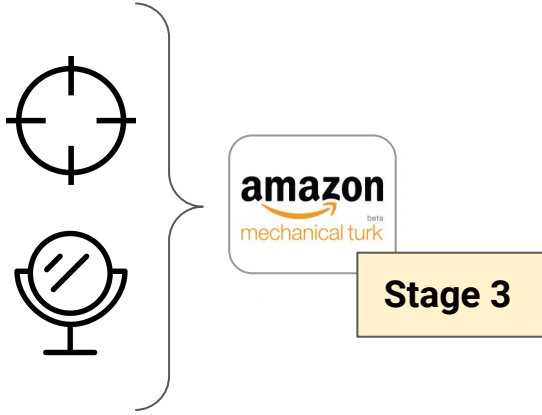
## 2. Implicit Hate Benchmark Dataset



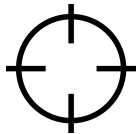
## 2. Implicit Hate Benchmark Dataset



## 2. Implicit Hate Benchmark Dataset



## 2. Implicit Hate Benchmark Dataset



**Target:** {Black folks, Muslim folks, Non-whites}



**Implied Statement:** <targets> {do, are, commit} <predicate>

e.g. *"Mexicans are incompetent"*

### 3. Baseline Classifiers for Detecting Implicit Hate

---

SVM (n-grams)

SVM (TF-IDF)

SVM (GloVe)

BERT

BERT + Aug

BERT + Aug + Wikidata

BERT + Aug + ConceptNet

---

### 3. Baseline Classifiers for Detecting Implicit Hate

#### Augmentation

SVM (n-grams)

SVM (TF-IDF)

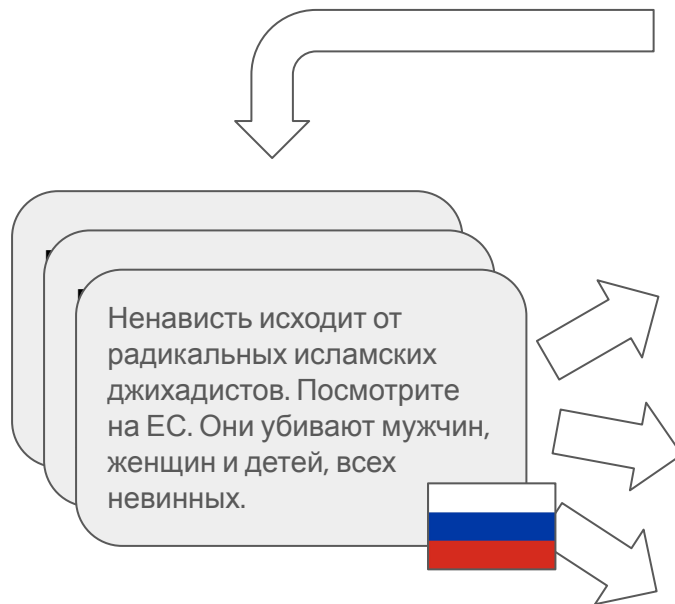
SVM (GloVe)

BERT

BERT + Aug

BERT + Aug + Wikidata

BERT + Aug + ConceptNet



*Hate is from radical Islamic jihadists. Look at the EU. They kill men women and children all innocents.*

The hatred of radical Islamic jihadists is directed at men, women, and children, all of whom are innocent.

Hatred comes from radical jihadists who kill innocent men, women, and children.

The hatred of radical Islamic jihadists goes so far that they kill men, women and children, all innocents.



### 3. Baseline Classifiers for Detecting Implicit Hate

#### WikiData

---

SVM (n-grams)

SVM (TF-IDF)

SVM (GloVe)

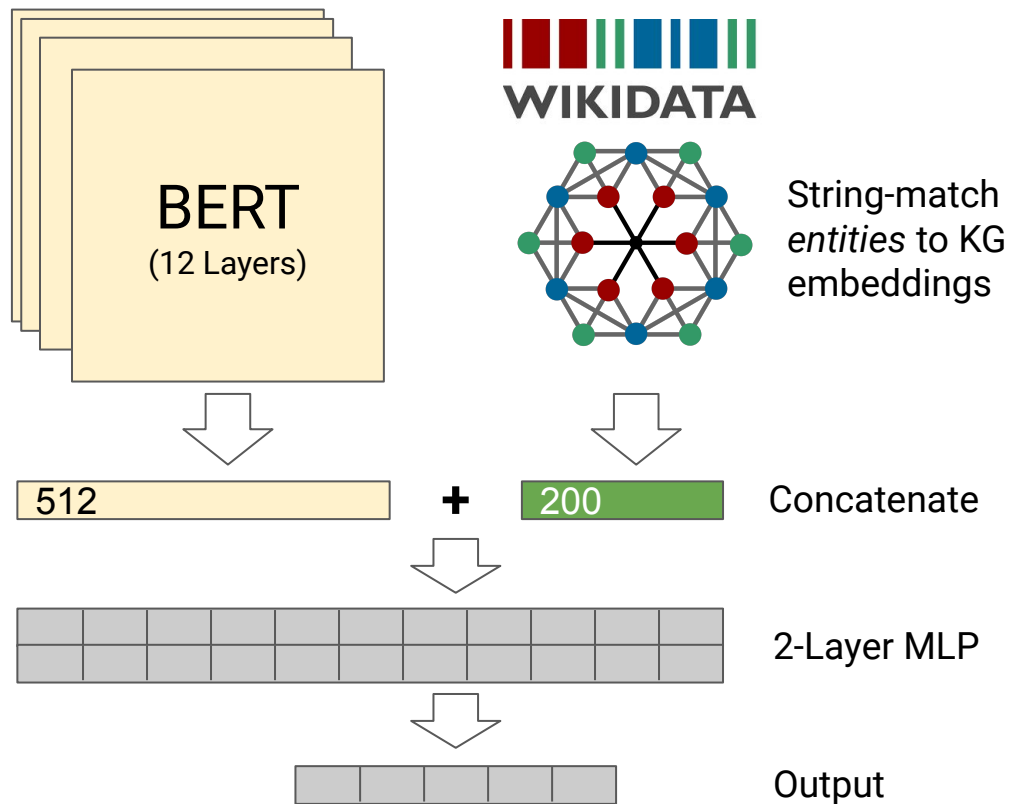
BERT

BERT + Aug

BERT + Aug + Wikidata

BERT + Aug + ConceptNet

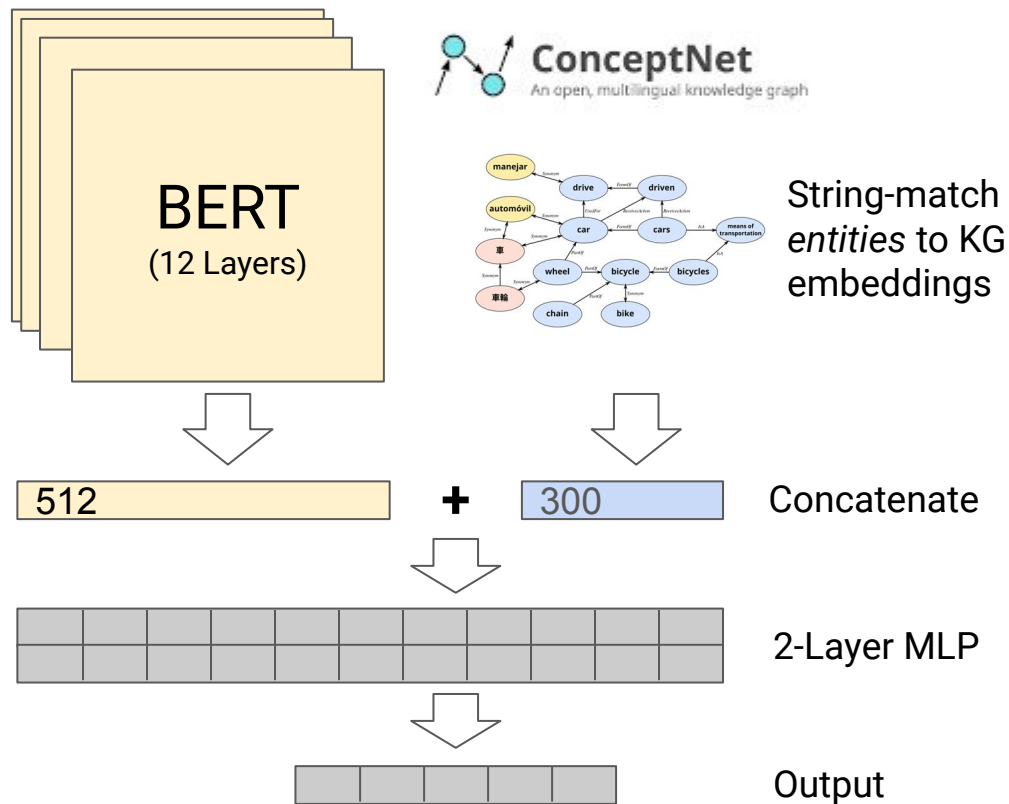
---



### 3. Baseline Classifiers for Detecting Implicit Hate

#### ConceptNet

- SVM (n-grams)
- SVM (TF-IDF)
- SVM (GloVe)
- BERT
- BERT + Aug
- BERT + Aug + Wikidata
- BERT + Aug + ConceptNet**



### 3. Baseline Classifiers for Detecting Implicit Hate

---

#### **Models**

SVM (n-grams)

SVM (TF-IDF)

SVM (GloVe)

BERT

BERT + Aug

BERT + Aug + Wikidata

BERT + Aug + ConceptNet

---

### 3. Baseline Classifiers for Detecting Implicit Hate

Models	Binary Classification				Implicit Hate Categories			
	P	R	F	Acc	P	R	F	Acc
Hate Sonar	39.9	48.6	43.8	54.6	-	-	-	-
Perspective API	50.1	61.3	55.2	63.7	-	-	-	-
SVM (n-grams)	61.4	67.7	64.4	72.7	48.8	49.2	48.4	54.2
SVM (TF-IDF)	59.5	68.8	63.9	71.6	53.0	51.7	51.5	56.5
SVM (GloVe)	56.5	65.3	60.6	69.0	46.8	48.9	46.3	51.3
BERT	<b>72.1</b>	66.0	68.9	<b>78.3</b>	<b>59.1</b>	57.9	58.0	62.9
BERT + Aug	67.8	<b>73.2</b>	<b>70.4</b>	77.5	58.6	<b>59.1</b>	<b>58.6</b>	<b>63.8</b>
BERT + Aug + Wikidata	67.6	72.3	69.9	77.3	53.9	55.3	54.4	62.8
BERT + Aug + ConceptNet	68.6	70.0	69.3	77.4	54.0	55.4	54.3	62.5

### 3. Baseline Classifiers for Detecting Implicit Hate

#### Error Analysis



## 4. Generative LMs for Explaining Implicit Hate

## 4. Generative LMs for Explaining Implicit Hate



Message

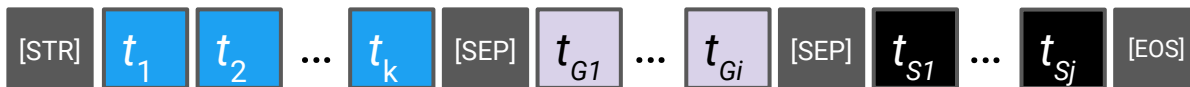


Implied Statement

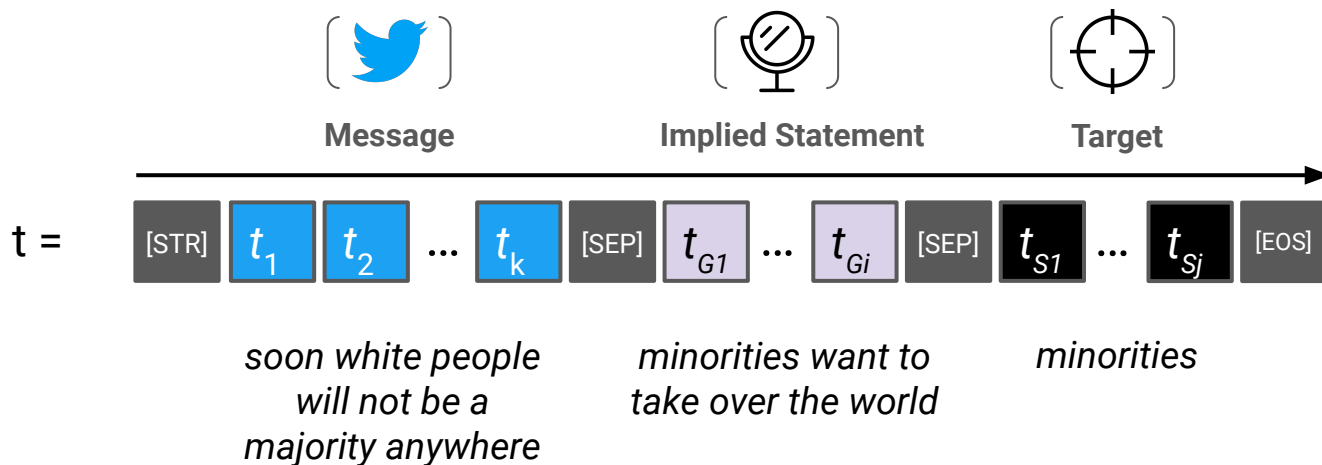


Target

t =



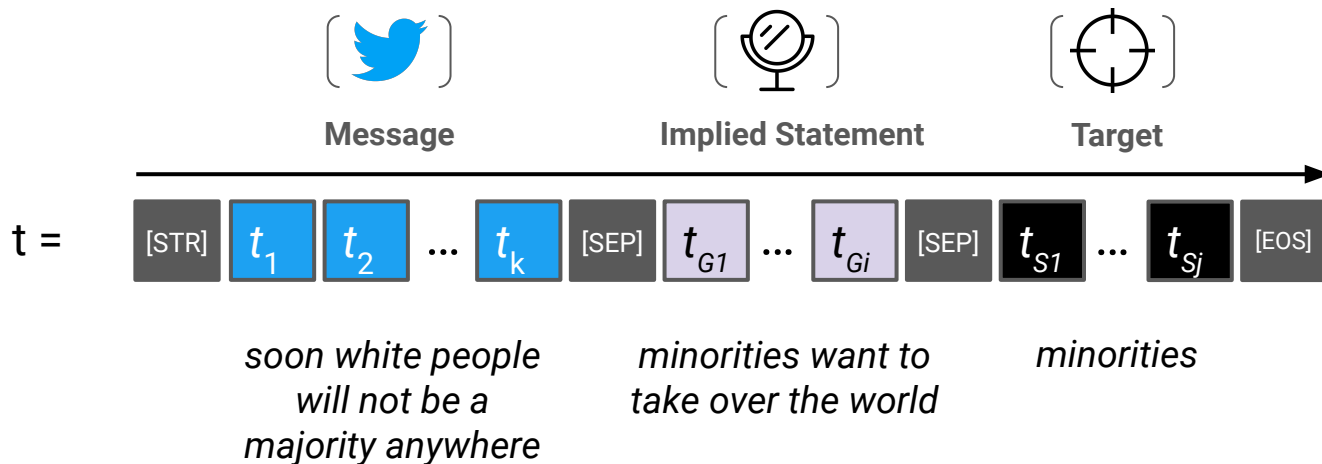
## 4. Generative LMs for Explaining Implicit Hate



Cross-Entropy Loss: 
$$-\sum_l \log P(\tilde{t}_l | t_{<l})$$





## 4. Generative LMs for Explaining Implicit Hate



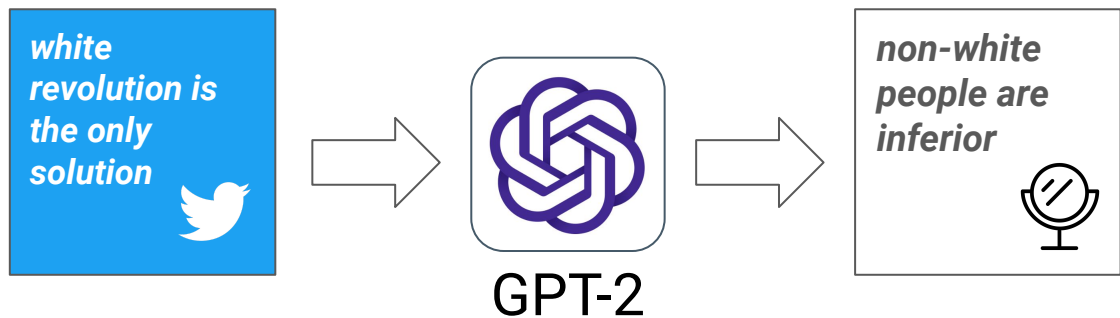
Cross-Entropy Loss: 
$$-\sum_l \log P(\tilde{t}_l | t_{<l})$$

Models:  GPT  GPT-2

## 4. Generative LMs for Explaining Implicit Hate

Models	Target Group 				Implied Statement 			
	BLEU	BLEU*	Rouge-L	Rouge-L*	BLEU	BLEU*	Rouge-L	Rouge-L*
GPT-gdy	43.7	65.2	42.9	63.3	41.1	58.2	31	45.3
GPT-top-p	57.7	76.8	55.8	74.6	55.2	69.4	40	53.9
GPT-beam	59.3	81	57.3	78.6	57.8	73.8	46.5	63.4
GPT-2-gdy	45.3	67.6	44.6	66	42.3	59.3	32.7	47.4
GPT-2-top-p	58.0	76.9	56.2	74.8	55.1	69.3	39.6	53.1
GPT-2-beam	<b>61.3</b>	<b>83.9</b>	<b>59.6</b>	<b>81.8</b>	<b>58.9</b>	<b>75.3</b>	<b>48.3</b>	<b>65.9</b>

## 4. Generative LMs for Explaining Implicit Hate



---

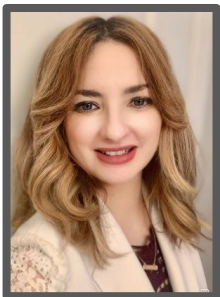
**ground truth:**  
*white people are superior to other races*

# Conclusion

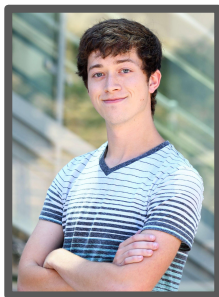
1. Implicit Hate Taxonomy
2. Implicit Hate Benchmark Dataset
3. Baseline Classifiers for Detecting Implicit Hate
4. Generative LMs for Explaining Implicit Hate



GT-SALT/  
[implicit-hate](https://github.com/GT-SALT/implicit-hate)



 @mai\_elsherief



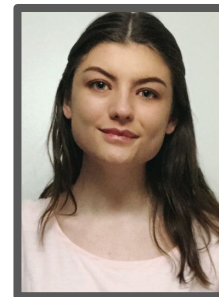
 @cziems



 @DMuchlinski



Vaishnavi Anupindi



Jordyn Seybolt

# Thanks!

Please reach out



GT-SALT/  
[implicit-hate](#)



 @munmun10



 @Diyi\_Yang

Post	Type	Target Group	Implied Statement
(a) white revolution is the only solution	incitement	non-white people <i>white people</i>	non-white people are inferior <i>white people are superior to other races</i>
(b) yes you are fine in a white majority country. how is mexico doing these days? people come here because you couldn't build it.	inferiority	mexican people <i>mexicans</i>	mexican people do not build things <i>mexicans are incompetent</i>
(c) black people only seem to thrive in white countries. how strange	irony	black people <i>black people</i>	black people are inferior <i>black people can't survive without white people</i>
(d) soon white people will not be a majority anywhere. you believe	grievance	whites <i>minorities</i>	whites have been driven out of the country <i>minorities want to take over the world</i>
(e) muslims are always wanting to kill someone! its their duty to kill the infidel.	stereotypical	muslims <i>muslims</i>	muslims are killers <i>muslims are violent</i>
(f) oh how i wish a man with a gun was there! wake up white people! they hate us	threatening	white people <i>non-whites</i>	white people are hated <i>non-whites hate whites</i>

Table 5: Example posts from our dataset along with their implicit category labels, the GPT-2 generated target and implied statements (first row of each block), and the ground truth target and implied statements (final row of each block, in *italics*). Generated implied statements are semantically similar to the ground truth statements.

[View instructions](#)

\$(tweet\_text)

Does this text attack a person or a group of people, explicitly or implicitly, based on their protected characteristics?

- Yes, explicit hate speech
- Yes, implicit hate speech
- Not hate speech

Tip: Explicit hate speech uses explicit hatred expressions or offensive keywords such as n\*gger, c\*nt, etc. to attack a person or a group of people based on their protected characteristics. Implicit hate speech is a more complex attack on a person or a group of people without relying on explicit keywords. Beware that implicit discriminatory speech may be subtle. Protected characteristics include ethnicity, race, national origin, religion, sex, gender, and sexual orientation. Note that pointing out racism should not be considered as hate speech.

---

**Submit**

Figure 2: Amazon Mechanical Turk interface used to collect ternary annotations (explicit hate, implicit hate, and not hate) for our first stage.

\$(tweet\_text)

**The following tweet has been categorized as "implicit hate speech" in a prior labeling stage; a more complex and subtle attack on a person or a group of people based on their protected characteristics without relying on explicit keywords.**

The goal of the task is to infer both the targeted group (GROUP) and what the post is actually implying about that group.

Step 1: The targeted group might be ethnicity, religion, class, or sexually oriented-related among other characteristics such as immigration.

Step 2: The second step in this task would be to determine what is really implied by the post. For this section, we ask you to write structured language, using the group identified in the prior step, such as (GROUP do/does \_\_\_\_, GROUP are \_\_\_\_, GROUP kill \_\_\_\_, GROUP have \_\_\_\_, GROUP commit \_\_\_\_)

Q1) Which group of people does this post refer to? (GROUP)

Example of answers are: black folks, asian folks, muslims, jews, latino/latina folks, immigrants, etc.

---

Q2) What aspect/stereotype/characteristic of this group is referenced or implied by the post? -- Use simple phrases and do not copy paste from the post.

Use the GROUP identified in the previous question to form a simple phrase and DO NOT COPY PASTE from the post. Examples of simple phrases include but are not limited to: GROUP do/does \_\_\_\_, GROUP are \_\_\_\_, GROUP kill \_\_\_\_, GROUP have \_\_\_\_, GROUP commit \_\_\_\_

Examples of common stereotypes include: Women are \*\*\*, Immigrants take \*\*\*, Muslims kill \*\*\*, Liberals are \*\*\*

---

Figure 3: Amazon Mechanical Turk interface used to collect the hate target and the implied statement per implicit hate speech post.



	Macro				Grievance			Incitement			Inferiority		
	P	R	F	Acc	P	R	F	P	R	F	P	R	F
SVM (n-grams)	48.8	49.2	48.4	54.2	65.6	53.6	59.0	53.7	55.8	54.7	49.7	46.4	48.0
SVM (TF-IDF)	53.0	51.7	51.5	56.5	66.9	56.7	61.4	60.4	56.2	58.2	46.0	45.3	45.6
SVM (GloVe)	46.8	48.9	46.3	51.3	63.7	48.6	55.1	55.2	46.7	50.6	45.8	39.7	42.5
BERT	<b>59.1</b>	57.9	58.0	62.9	65.4	63.9	64.6	62.4	<b>56.6</b>	59.4	<b>65.4</b>	57.9	61.4
BERT + Aug	58.6	<b>59.1</b>	<b>58.6</b>	<b>63.8</b>	67.6	<b>65.7</b>	<b>66.6</b>	<b>66.8</b>	56.5	<b>61.2</b>	61.0	59.0	59.9
BERT + Aug + Wikidata	53.9	55.3	54.4	62.8	<b>68.8</b>	63.0	65.8	62.7	55.9	59.1	60.3	<b>60.8</b>	<b>60.4</b>
BERT + Aug + ConceptNet	54.0	55.4	54.3	62.5	67.6	64.9	66.2	63.8	52.7	57.7	62.1	57.7	59.7

	Irony			Stereotypical			Threatening		
	P	R	F	P	R	F	P	R	F
SVM (n-grams)	41.4	51.8	46.0	60.7	52.7	56.4	52.0	72.2	60.5
SVM (TF-IDF)	43.9	55.4	48.9	60.9	58.8	59.8	55.3	72.2	62.7
SVM (GloVe)	48.7	55.4	51.8	59.3	53.9	56.5	50.2	74.3	59.9
BERT	<b>62.3</b>	<b>63.8</b>	<b>63.0</b>	58.5	69.3	63.4	<b>67.2</b>	71.5	69.3
BERT + Aug	62.0	62.3	62.1	<b>62.0</b>	<b>70.1</b>	<b>65.8</b>	65.0	<b>75.6</b>	<b>69.8</b>
BERT + Aug + Wikidata	60.0	63.1	61.4	60.7	69.3	64.7	64.2	73.8	68.6
BERT + Aug + Conceptnet	61.5	63.3	62.3	59.1	70.0	64.0	62.4	74.7	67.9

Table 6: Fine-grained implicit hate classification performance, averaged across five random seeds. Macro scores are further broken down into category-level scores for each of the six main implicit categories, and we omit scores for *other*. Again, the BERT-based models beat the linear SVMs on  $F_1$  performance across all categories. Generally, augmentation improves recall, especially for two of the minority classes, *inferiority* and *threatening*, as expected. Knowledge graph integration (Wikidata, Conceptnet) does not appear to improve the performance.

	<b>White Nationalist</b>	<b>Neo-Nazi</b>	<b>A-Immgr</b>	<b>A-MUS</b>	<b>A-LGBTQ</b>	<b>KKK</b>
Nouns (N)	identity evropa activists alt-right whites	adolf bjp india modi invaders	immigration sanctuary aliens border cities	islam jihad islamic muslim(s) sharia	potus democrats trump abortion dumbocrats	ku klux hood niggas brother
Adjectives (A)	white hispanic anti-white third racial	more non-white german national-socialist white	illegal immigrant dangerous ice criminal	muslim political islamic migrant moderate	black crooked confederate fake racist	alive edgy white outed anonymous
Hashtags (#)	#projectsiege #antifa #berkrally #altright #endimmigration	#swrm #workingclass #hitler #freedom #wpww	#noamnesty #immigration #afire #fairblog #stopsanctuarycities	#billwarnerphd #stopislam #makedclisten #bansharia #cspi	#defundpp #pjnet #unbornlivesmatter #religiousfreedom #prolife	#opkkk #hoodsoff #mantears #kkk #anonymous

Table 7: Top five salient nouns, adjectives, and hashtags identified by measuring the log odds ratio informative Dirichlet prior (Monroe et al., 2008) for the following ideologies: White Nationalist, Neo-Nazi, Anti-Immigrant (A-Immgr), Anti-Muslim (A-MUS), Anti-LGBTQ (A-LGBTQ), and Ku Klux Klan (KKK).

## 2. Implicit Hate Benchmark Dataset

