

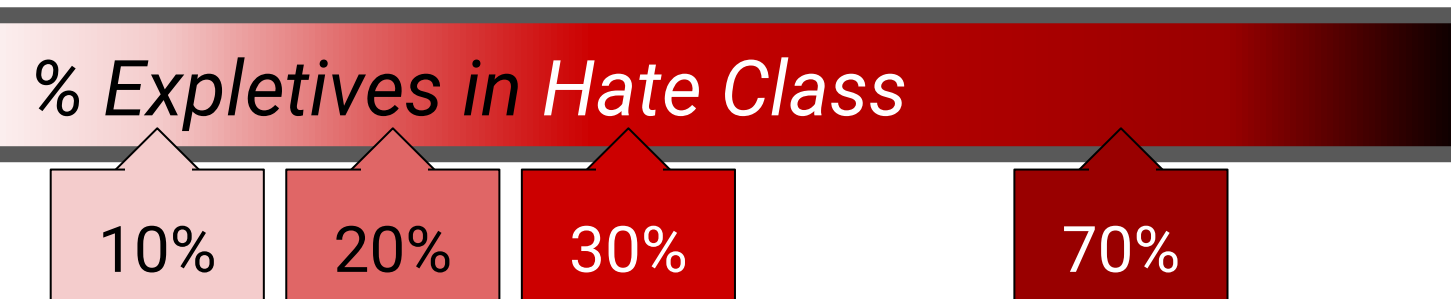
A Benchmark for Understanding Implicit Hate Speech

Mai ElSherief,* Caleb Ziems,* David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, Diyi Yang

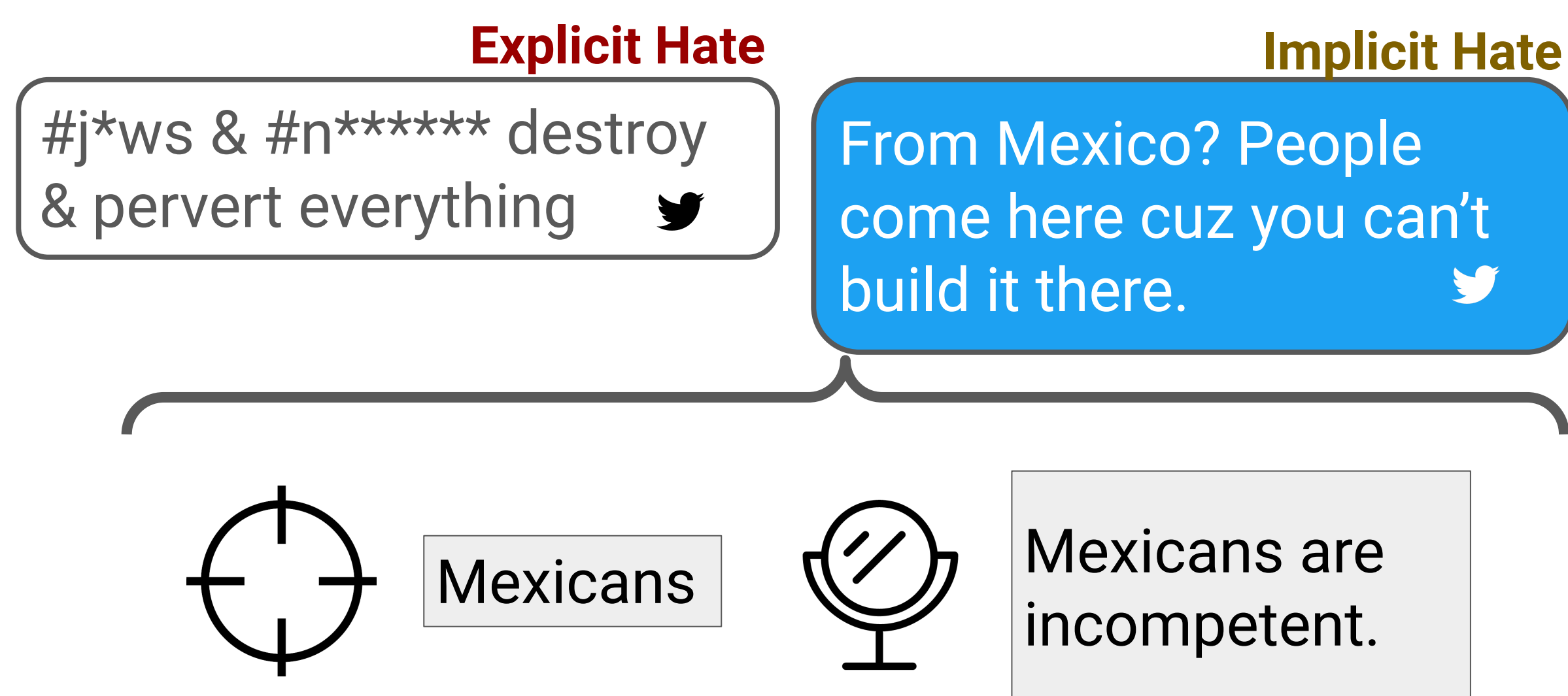
UC San Diego, Georgia Institute of Technology
melsherief@ucsd.edu, cziems@gatech.com



1. Motivation



- Many hate speech datasets are highly **explicit** and **overt**
- Little is known about **implicit hate**

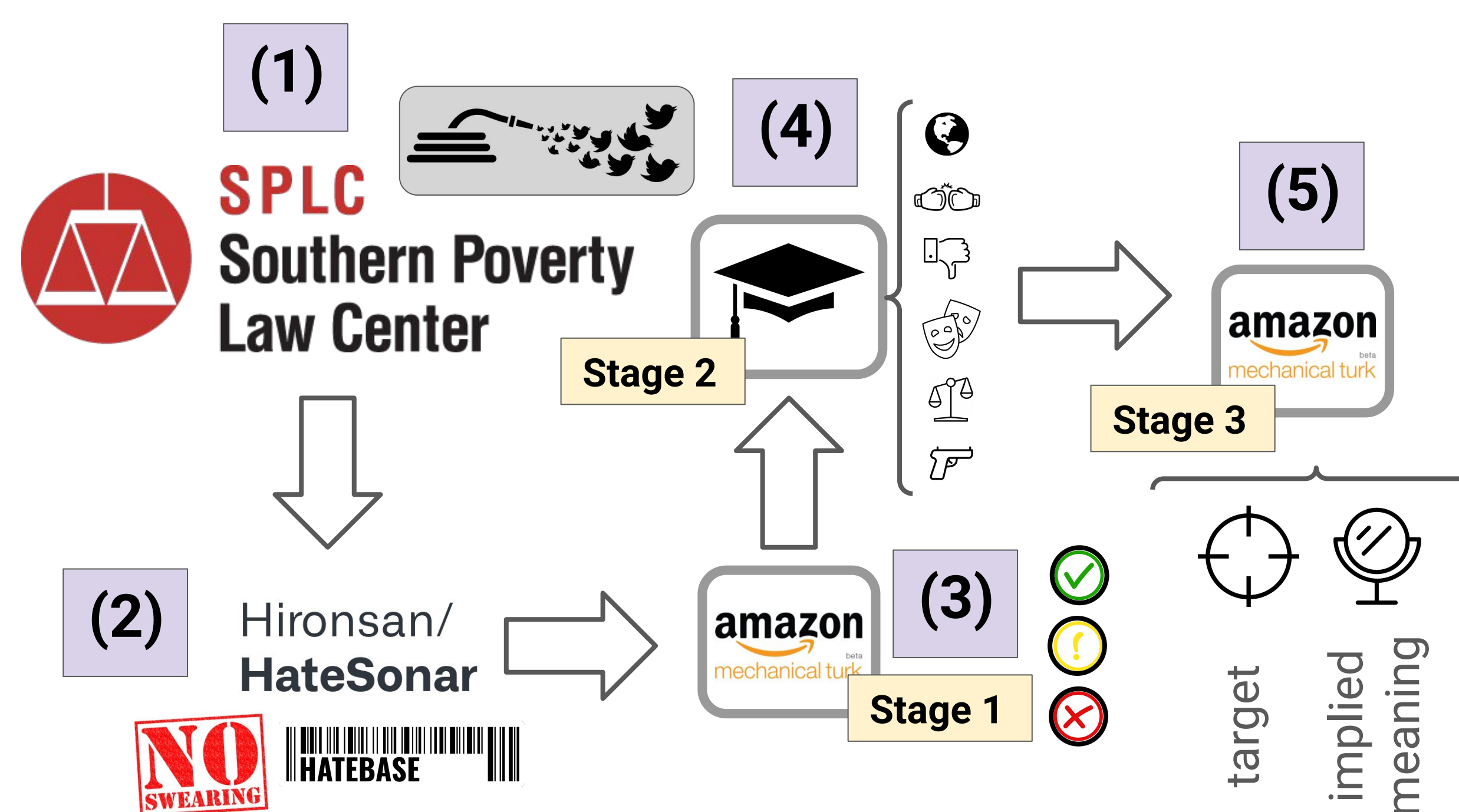


2. Taxonomy



3. Data Annotation + Expansion

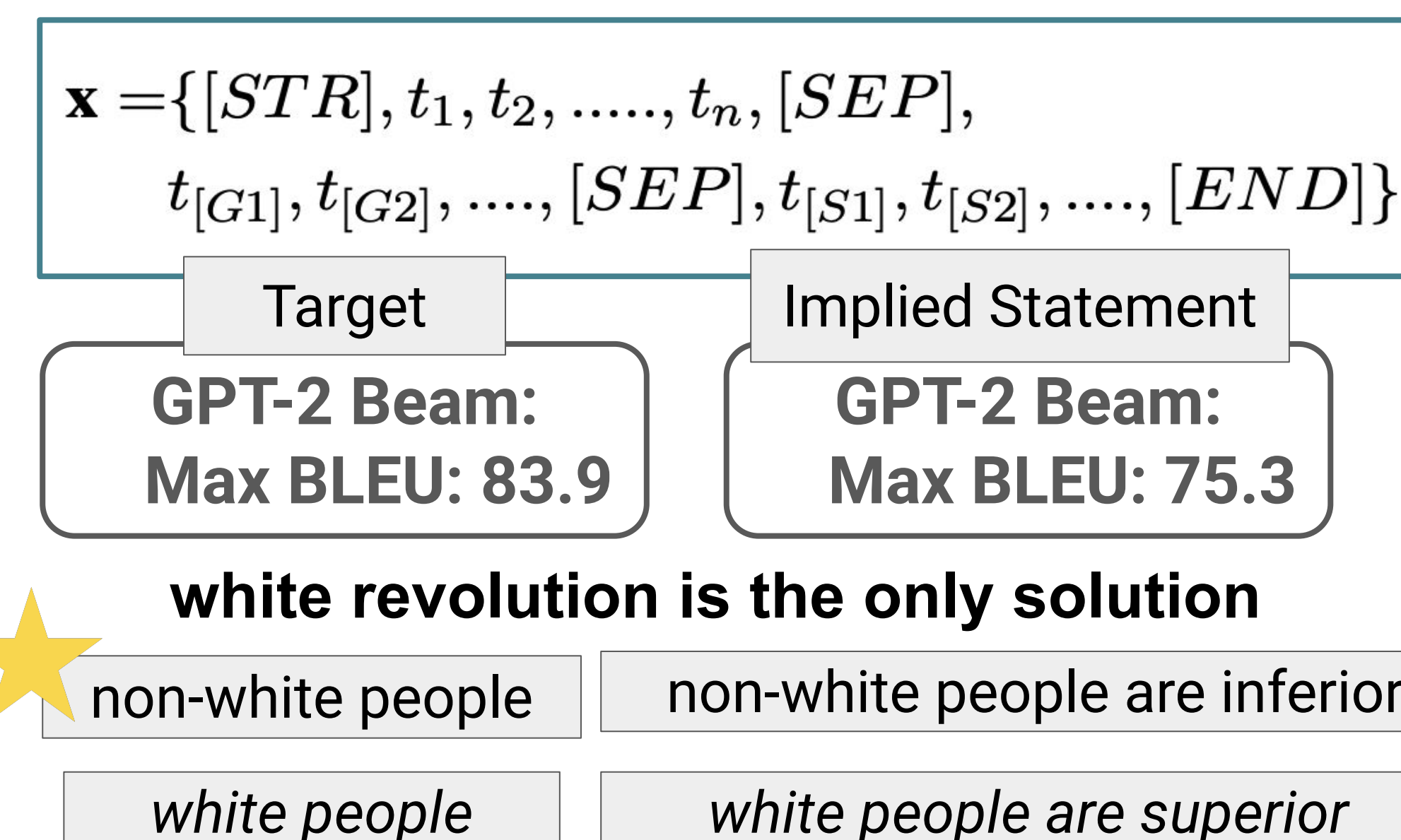
933 **explicit** | 6,346 **implicit** | 13,291 **not hate**



4. Detecting Implicit Hate

Training	Model	Testing on Implicit Hate		
		P	R	F1
Implicit Explicit	HateSonar	39.9	48.6	43.8
	Perspective API	50.1	61.3	55.2
	Linear SVM	61.4	67.7	64.4
	BERT + Augmentation	67.8	72.3	70.4

5. Explaining Implicit Hate



★ **white revolution is the only solution**
non-white people | non-white people are inferior
white people | white people are superior

6. Conclusion + Future Directions

