

Brazil#WithoutHate

Self-Explaining and Moral-Aware AI for Hate Speech Detection

Francielle Vargas

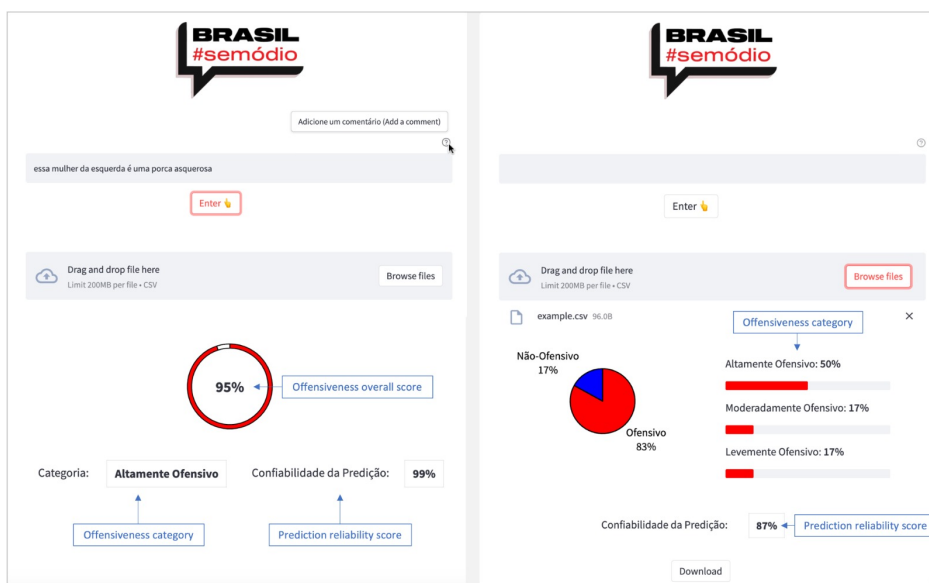


Figure 1. BrazilWithoutHate: A Transparent and Morally-Aware AI System for Detecting Hate Speech in Brazilian Portuguese.

Hate Speech Detection with Moral Rationales



Multi-Hop Hate Speech Explanation with Moral Rationales

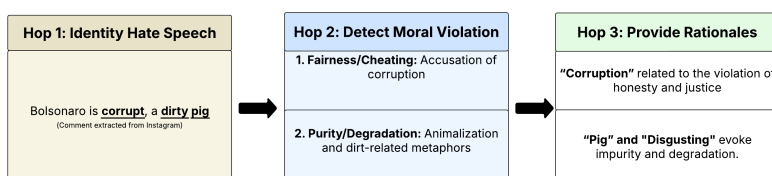


Figure 2. Brazil#WithoutHate: Multi-hop Hate Speech Explanation with Moral Rationales.

HateBR Dataset

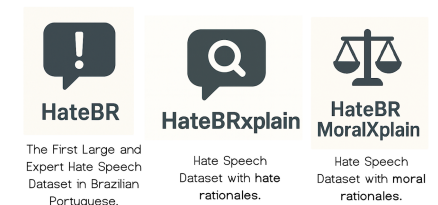


Figure 3. HateBR dataset and its versions.

Do you really believe what you're saying?
If you don't, you're just a **liar**. If you do
you're **mentally ill**.

Criminal! Scoundrel and racist!
May justice be served!

That's what you want — people employed
but just **following orders like puppets**.

Figure 4. Annotation of hate speech rationales.

Hate Speech Label(s)	Tweet (with highlighted rationale spans)	Moral Label(s)
Hate	Yes do hurt your selves gun control fraks .	Harm (red)
Hate	@USER So has the Law changed... you just make an allegation now days and your life is over. Liberals are mentally ill!	Cheating (orange), Degradation (green)
Non-hate	No one should ragger just because they're afro .	Harm (red)
Non-hate	@USER May be you have forgotten that we are EU citizens and so all rights applied to them must be afforded to us .	Fairness (blue)

Figure 5. Annotation of moral rationales.

Self-Explaining AI

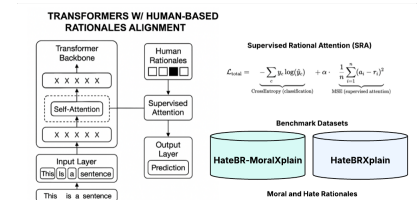


Figure 6. Supervised Rational Attention for Self-Explaining Hate Speech Detection.



← Check out the Brazil#WithoutHate platform!

