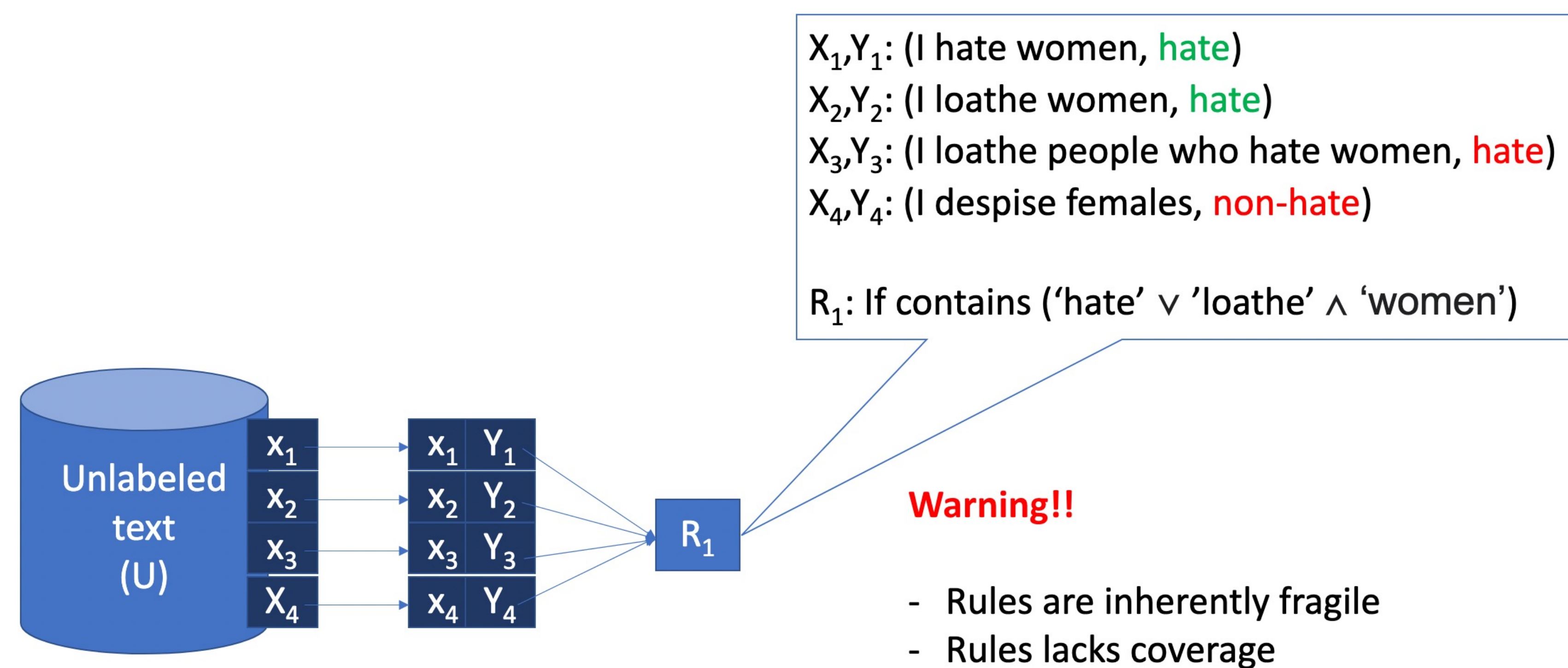
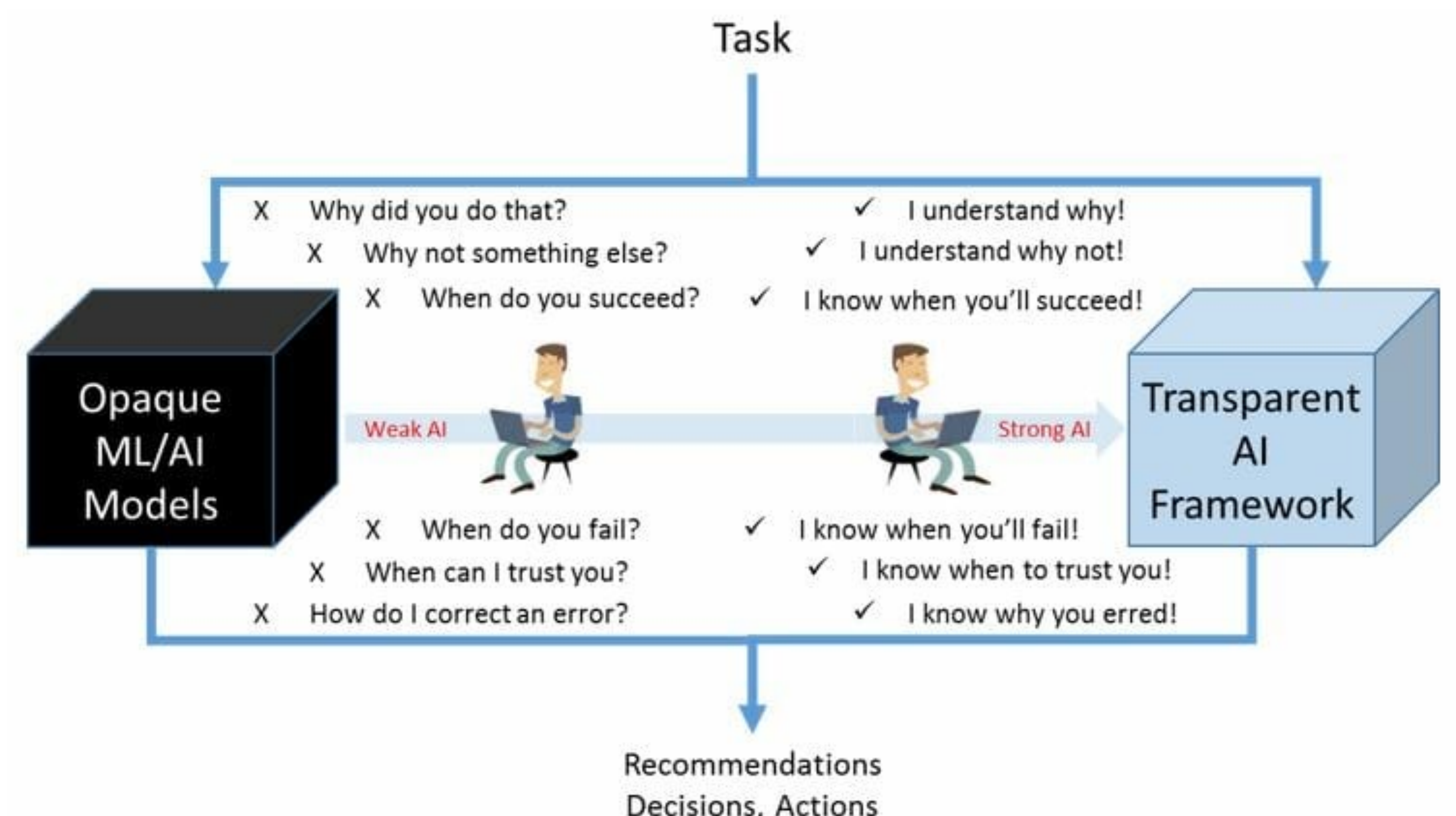


## Problem Statement

Classic approaches to content moderation typically apply a rule-based heuristic approach to flag content. Recent advances in deep learning have demonstrated the promise of using highly effective deep neural models to overcome these challenges. However, despite the improved performance, these data-driven models lack transparency and explainability, often leading to mistrust from everyday users and a lack of adoption by many platforms. In this paper, we present **Rule By Example (RBE)**: a novel exemplar-based contrastive learning approach for learning from logical rules for the task of textual content moderation. RBE is capable of providing rule-grounded predictions, allowing for more explainable and customizable predictions compared to typical deep learning-based approaches.

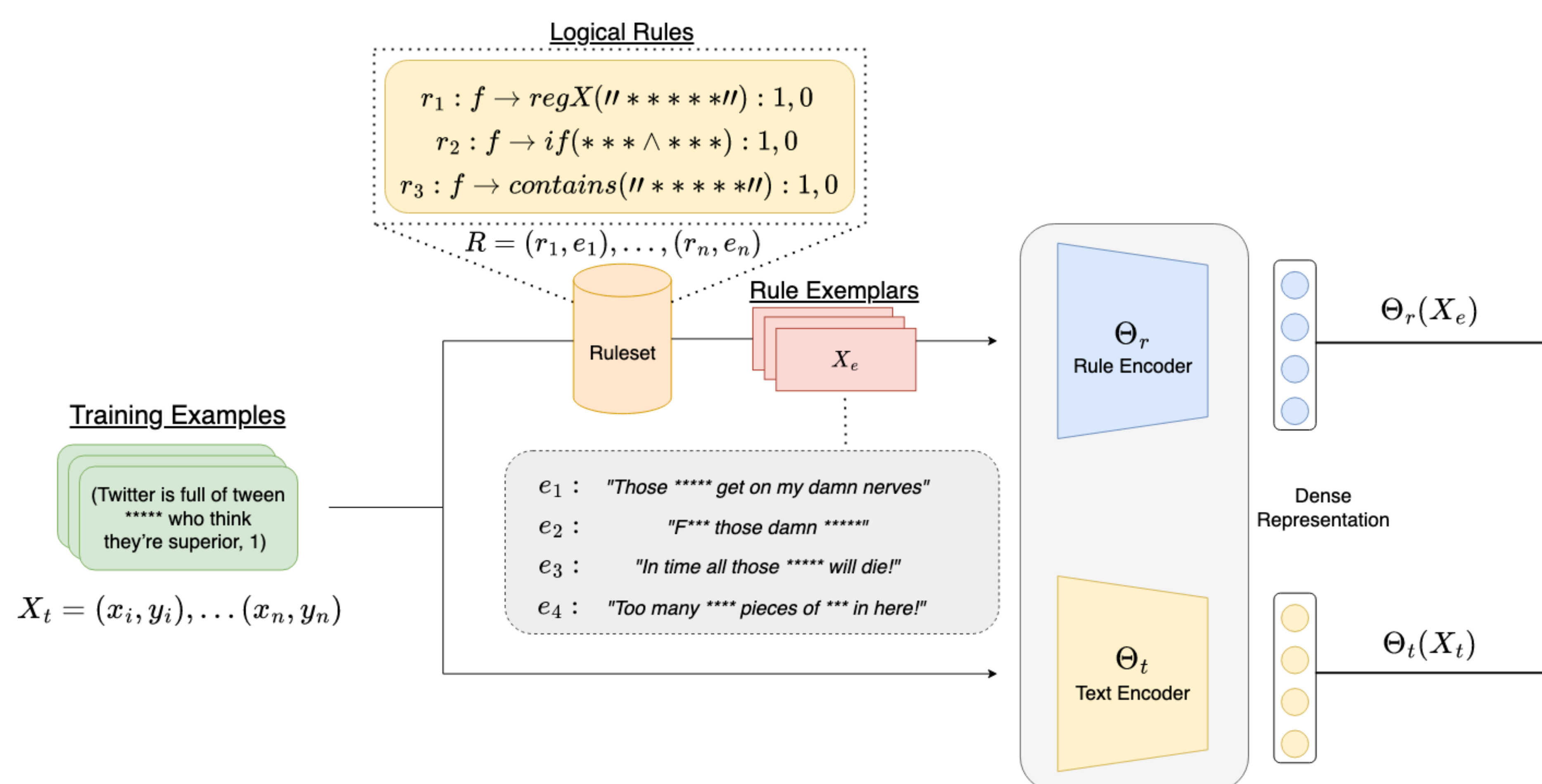


**Figure 1:** Generalization problem of rules. Logical rules, while easy to explain, are inherently fragile to the nuances of natural language.



**Figure 2:** Neural models while performant lack, transparency, customizability, predictability & explainability.

## Methods and Approaches



**Figure 3:** Rule By Example Framework (RBE) is comprised of two neural networks, a rule encoder and a text encoder, which jointly learn rich embedding representations for hateful content and the logical rules that govern them. Through Contrastive learning, RBE utilizes a semantic similarity objective that pairs hateful examples with clusters of rule exemplars that govern it.

## Results

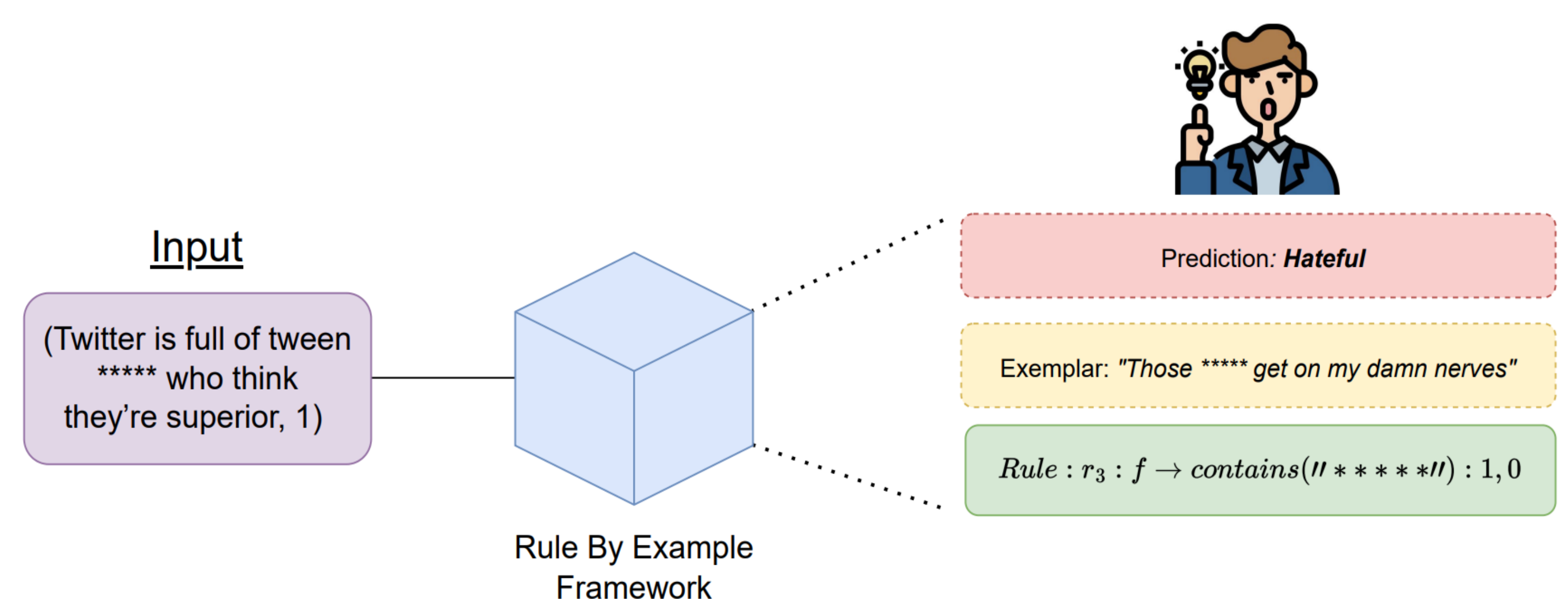
Content Moderation Using Rules (Fully Supervised)			
Model	HateXplain F1	Jigsaw F1	CAD F1
HateXplain Rationale Rules	0.752	-	-
Hate+Abuse Rules	0.719	0.226	0.290
CAD Rules	-	-	0.194
MPNet SeqCLS	0.823	0.581	0.463
BERT SeqCLS	0.824	0.563	0.433
Rule By Example (BERT   CAD Rules)	-	-	0.435
Rule By Example (MPNet   CAD Rules)	-	-	<b>0.476</b>
Rule By Example (BERT   Hate+Abuse Rules)	0.824	0.602	0.445
Rule By Example (MPNet   Hate+Abuse Rules)	<b>0.837</b>	<b>0.604</b>	0.476
Rule By Example (BERT   HateXplain Rationale)	0.816	-	-
Rule By Example (MPNet   HateXplain Rationale)	0.832	-	-

**Table 1:** Experiment Results in Fully Supervised Setting on hate speech classification datasets.

## Acknowledgements

We thank our anonymous reviewers for their feedback and suggestions. This work was conducted by the ROAR (Responsible & Open AI Research) team at Microsoft Cloud & AI.

## Rule Grounding



**Figure 4:** By displaying the rules and exemplars responsible, rule authors and users are better able to understand model predictions and can automatically adjust their ruleset to further improve model performance.

## Repo & Dataset

We release with RBE each of our derived rulesets for the HateXplain & Contextual Abuse datasets as well as a small subset of our internal Hate+Abuse list for public use! 🤗

