

Rethinking Offensive Text Detection as a Multi-Hop Reasoning Problem

Anonymous ACL submission

Abstract

We introduce the task of implicit offensive text detection in dialogues, where a statement may have either an offensive or non-offensive interpretation, depending on the listener and context. We argue that reasoning is crucial for understanding this broader class of offensive utterances, and create Mh-RIOT (Multi-hop Reasoning Implicitly Offensive Text Dataset), to support research on this task. Experiments using the dataset show that state-of-the-art methods of offense detection perform poorly when asked to detect implicitly offensive statements, achieving only ~ 0.11 accuracy.

In contrast to existing offensive text detection datasets, Mh-RIOT features human-annotated chains of reasoning which describe the mental process by which an offensive interpretation can be reached from each ambiguous statement. We explore the potential for a multi-hop reasoning approach by utilizing existing entailment models to score the transitions of these chains, and show that even naive reasoning models can result in improved performance in most situations. Analysis of the chains provides insight into the human interpretation process and emphasizes the importance of incorporating additional commonsense knowledge.

1 Introduction

With the development and popularity of online forums and social media platforms, the world is becoming an increasingly connected place to share information and opinions. However, the benefit these platforms provide to society is often marred by the creation of an unprecedented amount of bullying, hate, and other abusive speech¹. Such toxic speech has detrimental effects on online communities, and can cause great personal harm. Some efforts by the NLP community to address this

¹Disclaimer: due to the nature of this work, data and examples may contain content which is offensive to the reader.

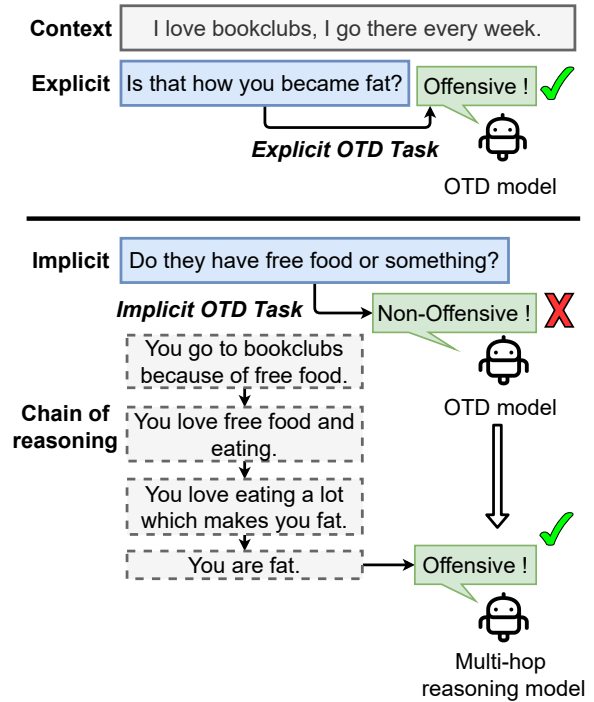


Figure 1: An instance illustrating Explicit OTD, Implicit OTD and our multi-hop reasoning approach.

problem have achieved high accuracies in classifying toxic speech in specific domains, such as sexist (Golbeck et al., 2017), racist (Waseem, 2016), or otherwise hateful text (Ross et al., 2016; Gao and Huang, 2017; Davidson et al., 2017).

While many instances of toxic speech are blatant and easily identified with sentence-level classifiers, not all offensive text contains obvious indicators. Waseem et al. (2017) argues for the classification of offensive text into two categories, (1) **explicit offensive text**², which is unambiguous in its potential to be offensive and often includes overtly offensive terms, such as slurs, and (2) **implicit offensive text**, which is more ambiguous, and may use sarcasm, innuendo, or other rhetorical

²Waseem et al.(2017) originally defined these terms as "explicit/implicit abusive text", but we adopt the phrase "offensive text" as used by the OTD community.

054 devices to hide the intended nature of the statement. 105
055 In this work we argue that there exists a direct 106
056 relationship between these tasks, and that each 107
057 implicitly offensive statement corresponds to an 108
058 explicitly offensive statement which is realized 109
059 through the interpretation process. This explicitly 110
060 offensive statement is closer to the sentiment the 111
061 listener feels when interpreting the statement as 112
062 offensive. Consider the example in Figure 1, a 113
063 dialogue between two speakers, S1 and S2: 114
064

065 S1: "I love bookclubs, I go every week"

066 S2: "Do they have free food or something?" 117
067

068 By itself, the statement by S2 is innocuous and 118
069 could be interpreted as a simple prompt for more 119
070 information about the bookclub. However, other 120
071 interpretations of this statement could lead S1 to 121
072 arrive at a number of explicitly offensive statements, 122
073 such as (1) "*You are poor*", (2) "*You are fat*", (3) 123
074 "*You are not smart/sophisticated*". Thus we con- 124
075 sider the chain of reasoning which constitutes the 125
076 interpretation to be a crucial part of recognizing 126
077 implicitly offensive statements.

078 The importance of more complex reasoning 127
079 when resolving such ambiguities in offensive con- 128
080 tent is not new. The Hateful Memes dataset (Kiel- 129
081 et al., 2021) pairs images with unrelated text cap- 130
082 tions. Both of these components are benign when 131
083 considered independently, but combining them can 132
084 occasionally create memes with offensive interpre- 133
085 tations. Consequently, approaches which jointly 134
086 reason over a combined representations of each 135
087 modality outperform those which treat each modal- 136
088 ity independently, hindering the system’s ability to 137
089 perform more complex reasoning.

090 To study this phenomenon purely in the text 138
091 domain, we use human annotators to construct a 139
092 dataset consisting of (1) an implicitly offensive 140
093 statement, (2) a corresponding explicitly offensive 141
094 statement, and (3) a chain of reasoning mapping 142
095 (1) to (2). When evaluated on the explicitly of- 143
096 fensive examples, state-of-the-art models perform 144
097 well, achieving > 90% accuracy. However, when 145
098 applied to the implicit OTD samples, the accu- 146
099 racy of the models drops to an average of about 147
100 < 11%. We then explore the use of a multi-hop 148
101 reasoning-based approach by utilizing a pre-trained 149
102 entailment model to score the transitions along each 150
103 "hop" of the reasoning chain. When incorporating 151
104 additional knowledge (from human annotations)

105 into the premises of each entailment, we achieve 106
107 higher accuracy than comparable methods which 108
109 do not utilize the reasoning chain. We present this 110
111 as evidence that a multi-hop reasoning-based ap- 112
113 proach is a promising solution to this problem, and 114
115 release our data to support further research into this 116
117 problem.

Our contributions in this work are threefold: 118

- We propose the task of implicit offensive 119
text detection (Implicit OTD), and construct a 120
dataset to research on this topic. The dataset 121
contains annotations of reasoning chains to 122
support study into multi-hop approaches. 123
- We conduct experiments using existing state- 124
of-the-art OTD models, and show they per- 125
form poorly on Implicit OTD task. 126
- We examine the use of entailment models as 127
part of a multi-hop reasoning approach for 128
Implicit OTD, showing improved accuracy in 129
most cases. We provide an analysis of which 130
types of reasoning are most challenging, and 131
which types of external knowledge is required. 132

2 Related Works 133

OTD in Text Classification Early approaches to 134
OTD relied primarily upon dictionaries like hate- 135
base³ to lookup offensive words and phrases. The 136
creation of OTD datasets enabled the development 137
of ML-based approaches utilizing simple features, 138
such as bag-of-word representations (Davidson 139
et al., 2017). With the advent of social media plat- 140
forms, many resources have been developed for 141
identifying toxic comments in web text (Waseem 142
and Hovy, 2016; Davidson et al., 2017), includ- 143
ing a number of deep learning-based methods (Pit- 144
silis et al., 2018; Zhang et al., 2018b; Casula et al., 145
2020; Yasaswini et al., 2021; Djandji et al., 2020). 146
Notably, all of these methods can be described as 147
building a contextual representation of a sentence 148
(whether trained end-to-end or on top of existing 149
pre-trained language models), and making a classi- 150
fication based on this representation.

OTD in Dialogue Systems As user-facing tech- 151
nologies, preventing dialogue systems from pro- 152
ducing offensive statements is crucial for their role 153
in society. As noted in Dinan et al. (2020), toxicity 154
in generated dialogue may begin with biases and 155

³www.hatebase.org

offensive content in the training data, and debiasing techniques focused on gender can reduce the amount of sexist comments generated by the resulting system. Similar outcomes can be obtained through adjustments to the model or training procedure, for instance, toxic words can be masked during training to reduce their role in model predictions (Dale et al., 2021). GeDi (Krause et al., 2020) proposed using class-conditional LMs as discriminators to reduce the toxicity produced by large pre-trained LMs (GPT-2). Additionally it may also be important to identify offensive statements made *to* a dialogue system, as it has been shown that dialogue systems can react with counter-aggression (Cercas Curry and Rieser, 2018), and systems which continuously learn during deployment may incorporate toxic user responses into future generations.

Subjectivity in OTD Previous work has hit upon the role that an individual’s own perspective may play when determining offensiveness. For instance, in the Offensive Language Identification Dataset (OLID), a widely used OTD dataset (Zampieri et al., 2019a,b, 2020), annotations exist on a hierarchy. Each level dictates the target of the offensive text, in terms of their identity as a group, individual, or entity. But to our knowledge, a person’s identity or attributes have not played a critical role in existing OTD research. OLID was also augmented with labels for capturing the degree of explicitness (Caselli et al., 2020), and may also support research into resolving implicitly offensive statements. However, implicitness in OLID is defined primarily as the lack of an overtly offensive word or slur, and the aforementioned personal attributes or subjectivity of interpretation are not considered. Our dataset differs in this respect, as we consider not just if a statement is offensive, but *how* it can be considered offensive, by defining the interpretation process as a chain of reasoning towards a subjective experience. In this sense, a more similar approach comes from normative reasoning in moral stories (Emelin et al., 2020), where a short chain of reasoning is used to assess morality of actions and consequences.

3 Data

We propose Mh-RIOT as a dataset for the study of Implicit OTD as a multi-hop reasoning problem, and for use as a diagnostic to test models’ ability to identify implicitly offensive statements.

Each example in the dataset consists of three parts:

1. A personal attribute of the reader/listener.
2. An implicitly offensive statement, its corresponding explicitly offensive statement, and a non-offensive statement.
3. A chain of reasoning, describing the iterative process of how the ambiguity of the implicitly offensive statement can be resolved into the corresponding explicitly offensive statement. Appendix A lists some sample chains in Mh-RIOT.

We collect annotations for Mh-RIOT using Amazon Mechanical Turk (AMT). Four pilot experiments were conducted to select qualified annotators for the final annotation. The instructions provided to the annotators can be found in Appendix C.

3.1 Annotation Scheme

Personal Attribute As we have defined in Section 1, we argue that the context in which a statement occurs is crucial to understanding its potential in creating an offensive interpretation, and therefore the context should play an important role in the annotation task. However, providing an overly specific context can increase the difficulty of providing a relevant implicitly offensive statement. To make the annotation task more feasible we reduce the context to a single feature: a personal attribute of the reader/listener.

The set of attributes is obtained from the personas in the PERSON-CHAT corpus (Zhang et al., 2018a), of the form “*I like sweets.*”, or “*I work as a stand up comedian.*” Attributes related to ethnicity, gender, and other protected classes are manually removed, leaving 5334 distinct attributes. We divide the attributes into several categories (detailed category information can be found in Appendix B) before randomly sampling a subset of 920 attributes, uniformly across categories, in order to increase the number of workers assigned to each attribute.

Implicit, Explicit and Non-offensive Text For each example, workers were provided 3 diverse attributes and asked to choose one as writing prompt. The workers are then instructed to provide annotation in the form of example sentences, including: *Implicitly offensive statement* Utterances that do not express an overt intention to cause offense and often require complicated reasoning or external

knowledge to be fully recognized as offensive contents.

Explicitly offensive statement Utterances which contain an obvious and direct intention to cause offense without external knowledge or reasoning processes.

Non-offensive statement Utterances that do not cause offense under the context initiated with the attribute.

Both explicit and implicit offensive statements should share the same meaning in terms of how they are offensive. Non-offensive statements are collected to construct a balanced dataset and to evaluate the accuracy of existing OTD models.

Chain of Reasoning A distinguishing characteristic of our work is the collection of chains of reasoning to explain the interpretation process for implicitly offensive text. We represent the chain of reasoning as a series of sentence-to-sentence rewrites, similar to natural logic (MacCartney and Manning, 2014). One practical advantage of using a sentence-based representation for reasoning steps (in comparison to a structured representation like predicate-argument tuples) is that it allows the use of powerful text-to-text (T5) (Raffel et al., 2019) and entailment models (Liu et al., 2019; He et al., 2021), which are trained on sentence-level input.

Formally each chain begins with an implicitly offensive statement (0-th step, denoted as s_0) and ends with an explicit offense (s_l), making the length of the chain the number of steps between s_0 and s_l , inclusive.

3.2 Post-processing

We were able to collect 2657 examples from the AMT and performed post-processing to ensure the quality of the data. We define three processes to edit the collected annotations in order to standardize the format of the reasoning steps, listed below. Examples with steps that can not be handled by any of the processes are removed from the dataset. To reduce biases in post-processing, we assign 3 workers to each task.

Attribute Insertion Rule (AIR) We insert the attribute statement into the first reasoning step (s_1) to make this information accessible to any model taking the sentence as input. For instance, for an example with the attribute, “*I am colorblind.*” and the implicit offensive statement, “*Oh, that would explain your wardrobe!*”, the reasoning step “*Oh,*

Knowledge

Only the best can win contests.
Classic things are usually old.
Grown-ups don’t play with dolls.
Parents want children to be independent.
Overworking makes people exhausted.

Table 1: Samples of the knowledge used to construct chains of reasoning.

your color blindness would explain your wardrobe! 298
 generated by the worker is tagged as AIR. 299

Knowledge Insertion Rule (KIR) Steps that are used to introduce external commonsense knowledge are tagged as KIR. For instance, to support the reasoning process from step “*You are a grown-up who can’t afford to rent a house.*” to “*You are poor.*”, the knowledge of “*Poor people can’t afford to rent a house.*” is introduced. The following step “*You are poor.*” is then tagged as KIR. To better understand the effectiveness of external knowledge, we also extract the commonsense knowledge during the post-processing (Table 1). 300-310

Rephrasing Rule (RR). Steps that have equivalent meaning to previous steps but can be simplified by rephrasing are tagged as RR. For instance, to express more explicit offensive meaning, an reasoning step written as a question “*Do you like meat too much, or just food in general?*” is rephrased as a declarative sentence step “*You must love food too much in general.*” and tagged as RR. 311-318

3.3 Post-processing Results 319

Of the initially collected 2657 examples, 1050 remained after the post-processing. The high task rejection rate (60.5%) also conveys the difficulty of this content generation task. In the dataset, the average length of a reasoning chain is 4.84 steps, with a minimum length of 3 (60 examples) and a maximum of 6 (39 examples). Among all three tags, RR is most frequently applied (59.6%), followed by KIR (21.5%) and AIR (18.9%). 320-328

4 Experiments 329

We evaluate the difficulty of the Implicit OTD task using existing state-of-the-art models, before exploring a multi-hop approach to Implicit OTD using existing entailment models to score transitions in the reasoning chains. 330-334

Models	Accuracy						
	Mh-RIOT				Twitter	OffensEval	Toxicity
	Implicit	Explicit	Non	All	All	All	All
RoBERTa-Twitter	1.7	79.0	99.7	59.5	85.9	85.8	89.1
BERT-OffensEval	15.9	93.2	99.2	62.8	82.2	82.4	84.2
ALBERT-OffensEval	9.7	88.6	94.5	65.2	82.4	82.7	85.2
BERT-Toxicity	14.8	96.6	98.5	61.9	81.2	81.9	83.6
ALBERT-Toxicity	11.4	91.5	94.9	62.8	79.4	80.3	82.6
Avg.	10.7	89.8	97.4	62.5	82.2	82.6	84.9

Table 2: Performance of SOTA OTD models on the classification task. *Non*: Non-offensive.

4.1 Sentence Classification

We begin by evaluating existing state-of-the-art OTD models on both the Implicit-OTD and Explicit-OTD task. These include BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and ALBERT (Lan et al., 2020), three pretrained large scale language models fine-tuned on existing OTD datasets, which produce the highest accuracy reported on the explicit OTD task.

These models are fine-tuned on three OTD datasets, including (1) the OLID/OffensEval2019 dataset (Zampieri et al., 2019a), discussed in Section 2, which contains 14,200 labeled tweets and includes implicit offensive statements, (2) the TWEETEVAL (Barbieri et al., 2020) multi-task offensive Twitter set for detecting irony, hate speech and offensive language, and (3) the Google Jigsaw Toxic Comments dataset⁴ which contains 159,571 samples in the training set. In the subsequent sections we refer to these datasets as OffensEval, Twitter, and Toxicity, respectively.

Table 2 shows the results of the baseline models on correctly classifying the implicitly and explicitly offensive text as offensive/non-offensive (systems are denoted as a hyphenated combination of pretrained model and dataset). In every situation, the performance on the implicit task is significantly lower. The overall trend is perhaps unsurprising, as implicit examples lack clear indicators of offensiveness, such as highly offensive words. However, the degree to which these models underperform in the Implicit-OTD task illustrates the extent to which these tasks differ, and highlights the risk of deploying such models to perform this task in real-world situations.

An underlying assumption of this work and the

⁴Google Jigsaw Toxic Comments

motivation for reasoning chains is the expectation that as the reasoning process is applied, the interpretation of the implicitly offensive utterance becomes increasingly (explicitly) offensive. We evaluate the extent to which this holds true in the dataset, using the baseline systems to predict the offensiveness of each rewrite across the reasoning chain. Appendix D shows that this is indeed the case, that moving down the reasoning chain correlates with higher accuracy, and implying that each step gradually reveals more of the offensive connotations in implicit offense. It also verifies that the collected/annotated chains have the property of being orderly.

4.2 Reasoning by Entailment

The results of Section 4.1 indicate two things: current OTD systems perform poorly on the implicit OTD task, and the difficulty of using existing models decreases as each successive step of the reasoning chain is applied. This insight hints at a potential approach to implicit OTD: apply a reasoning model to map initial statements to their simplest and most explicit corresponding offensive statement (and score the likelihood of it being entailed by the original statement), and then score the resulting statement with a dedicated OTD model. In essence, this decomposes a difficult inference into a series of smaller inferences which may be tackled with higher accuracy by current models. We explore the possibility using this approach with existing models, assuming the human-annotated chains as gold proof paths.

We treat the problem of scoring reasoning chains as a multi-hop textual entailment problem as in Figure 2. Using an existing state-of-the-art textual entailment model, we score the transition from each step s_i to the next, s_{i+1} . Such models take as input

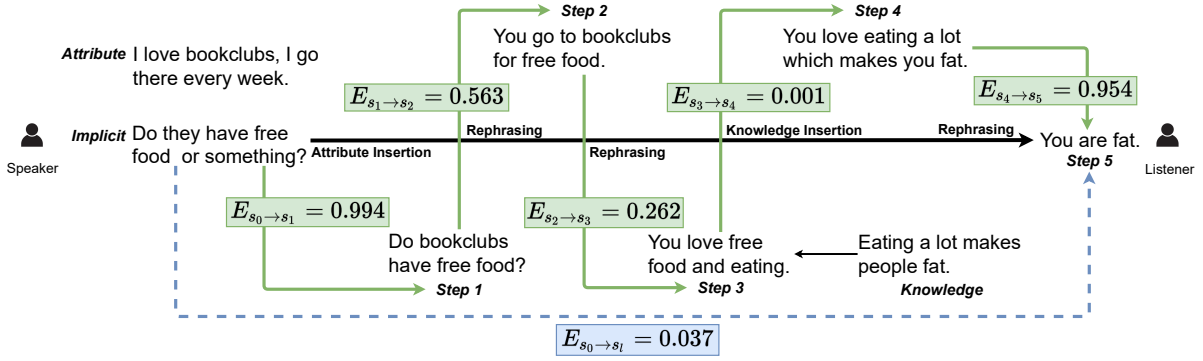


Figure 2: An example demonstrating the entailment experiment. Entailment scores between adjacent steps are given by the text entailment models. Arrows represent the entailment processes. $E_{s_i \rightarrow s_j}$ represents the entailment score from step i to step j , where s_0 represents the implicit offense and s_l represents the last step (step 4 in this example) of the chain.

a pair of texts, <premise, hypothesis>, and output scores for a set of labels indicating “entailment” ($E_{p \rightarrow h}$), “neutral” and “contradiction” ($C_{p \rightarrow h}$). An example reasoning step, the premise “*You look like someone who could use more exercise.*” entails the hypothesis “*You are fat.*”.

A naive approach to multi-hop reasoning is to treat each transition as an independent event, and model the probability of a reasoning chain as a product of transition scores. In the context of reasoning chains, we define the probability of a chain c as:

$$E(c) = \prod_{i=0}^{l-1} E_{s_i \rightarrow s_{i+1}} \quad (1)$$

We refer to this as *MUL*, the product model approach to multi-hop reasoning. For the entailment model scoring each transition in the chain, we consider two systems, one derived from **DeBERTa-base** (He et al., 2021) and one from **RoBERTa-large** (Liu et al., 2019). Both systems were finetuned on the MNLI corpus (Nangia et al., 2017), a standard corpus for textual entailment.

In our experiments we are most interested in comparing the scores of *MUL* to those of methods which ignore the reasoning chain, either by scoring the entailment of the explicitly offensive statement given the implicit one ($s_0 \rightarrow s_l$), or by using one of the current state-of-the-art approaches to classify the implicit statement directly (Table 2). While *MUL* is a naive model, any advantage of a model with such strong independence assumptions suggests areas where future multi-hop reasoning models could significantly improve over non-reasoning “single hop” counterparts.

The results of the multi-hop experiments are presented in Table 3. We observe that under most conditions, *MUL* outperforms $E_{s_0 \rightarrow s_L}$ by a modest margin. The performance of *MUL* does suffer on the longest reasoning chains as a result of an increasing number of < 1.0 multiplications (a consequence of the independence assumptions), negating the margins between the two systems. The detailed results can be found in Appendix G.

In terms of the types of reasoning which are most beneficial, we observe large changes in the transition scores before and after knowledge is integrated into the reasoning process, i.e., around KIR steps. We examine this behavior further, analyzing the performance of OTD models on predicting the final layer at points s_{k-1} and s_k , before and after knowledge integration (Table 5). We observe significant (2-3 fold) improvements when predicting after knowledge is integrated. Similar results can also be observed on textual inference models as shown in Appendix E.

To explore the effectiveness of the external knowledge, we utilize the extracted knowledge mentioned in Section 3.2 and perform an additional set of experiments (denoted k+) where the external knowledge acquired in data annotation is added to each statement as a conjunction, until after a KIR step occurs. For instance, if the knowledge in s_k is “*Eating too much can make people fat.*”, this knowledge will then be connected to all steps in $\{s_i | i = 0, 1, \dots, k - 1\}$ to form “ $\langle s_i \rangle$ and *eating too much can make people fat.*” As shown in Table 3, adding knowledge increases scores for both models, but notably resulting in a significant advantage to the RoBERTa product model, which

Entailment Scores										
Step	RoBERTa					DeBERTa				
	Chain Length					Chain Length				
	3	4	5	6	ALL	3	4	5	6	ALL
$s_0 \rightarrow s_1$	64.7	84.4	89.9	90.0	-	68.4	78.2	86.5	90.7	-
$s_1 \rightarrow s_2$	37.1	58.0	46.9	57.4	-	29.7	46.1	41.2	45.0	-
$s_2 \rightarrow s_3$	73.6	55.1	42.5	50.2	-	64.4	50.5	35.5	44.3	-
$s_3 \rightarrow s_4$		58.2	61.6	40.6	-		51.0	55.6	37.5	-
$s_4 \rightarrow s_5$			60.9	65.9	-			50.0	63.3	-
$s_5 \rightarrow s_6$				67.5	-				57.8	-
MUL_{s_0, \dots, s_l}	14.3	13.1	4.6	5.4	11.5	12.1	7.7	1.8	3.3	6.8
$E_{s_0 \rightarrow s_l}$	17.2	9.1	4.4	5.6	7.6	8.3	5.9	2.4	3.6	4.5
$MUL_{s_0, \dots, s_l}(k+)$	38.1	32.0	17.9	16.5	23.5	30.2	20.3	7.6	4.0	14.1
$E_{s_0 \rightarrow s_l}(k+)$	35.9	15.9	10.8	8.6	15.0	25.3	11.9	7.5	6.6	10.9

Table 3: Entailment scores between various steps of the reasoning chain, and the scores of a product model processing each step sequentially (MUL). Column headers indicate subsets of the data, where all chains are of 3, 4, 5, or 6 steps respectively. $k+$: scores indicate those where external knowledge is concatenated to all statements prior to a KIR step.

now outperforms direct prediction, and all previous baseline models, in all scenarios. The resulting system is also more robust to long reasoning chains. We even observe that the performance margins over direct prediction in the 6-step chains exceeds that of 3-step setting.

5 Discussion

We introduced this work based on a hypothesis of multi-hop approach as having a conceptual advantage over existing approaches to offensive text detection, in that humans must each be performing some reasoning process in order to find statements either offensive or unoffensive in different situations. We then showed that this conceptual advantage could translate to an empirical one, and showed performance gains over current approaches. However, we do so under strong assumptions and with access to additional information. How realistic is our experimental setup?

5.1 A Perfect Reasoning Model?

A concern in our initial entailment experiment is the naive product reasoning model. As mentioned in Section 4.2, an ideal reasoning model for this multi-hop approach should be a generative model that outputs explicitly offensive statements directly from implicitly offensive statements with reasoning process handled internally. In this sense, existing

contextual paraphrase generation models (Kazemnejad et al., 2020; Niu et al., 2021) can be promising candidates for a generative reasoning model. Such models aim at paraphrasing sentences while incorporating knowledge and external reasoning process and thus possess the potential to handle the reasoning process underlying implicit offensive statements after trained on large amount of data. But to what extent can a perfect reasoning model benefit the performance?

To answer that, we can assume a perfect paraphrasing model and the task reduces to whether we can predict the first transition from the implicit statement to the next step in the chain. This is akin to moving from an observed statement to a hypothetical knowledge base, upon which reasoning can occur to produce the explicitly offensive analog, which can be classified with high accuracy. As shown in Table 4a, the initial transition, $E_{s_0 \rightarrow s_1}$, can be predicted with much higher score than the direct prediction, $E_{s_1 \rightarrow s_l}$. On one hand, This result shows that even if the model is aware of the corresponding explicitly offensive rewrite, it has difficulty directly understanding the inference relationship between them. Lower $C_{implicit \rightarrow non}$ also shows the difficulty distinguishing implicit and non-offensive statements as shown in Table 4b. On the other hand, these results show the possibility of getting better inference by grounding the implicit

Entailment Scores		
Steps	RoBERTa	DeBERTa
$s_0 \rightarrow s_1$	86.1	83.1
$s_0 \rightarrow s_l$	6.7	3.9

(a)

Contradiction Scores		
Steps	RoBERTa	DeBERTa
<i>implicit</i> \rightarrow <i>non</i>	13.7	17.9
<i>explicit</i> \rightarrow <i>non</i>	94.6	97.0

(b)

Table 4: The entailment scores (a) and contradiction scores (b) from implicit statements to non-offensive statements versus explicit statements to non-offensive statements.

statement in a knowledgebase that follows the general structure of the reasoning chains, which can finally result in improvements in the overall classification accuracy. In other words, with such a paraphrasing model (rather than our naive product model), we should be able to improve the accuracy ultimately close to $E_{s_0 \rightarrow s_1}$.

5.2 What Knowledge is Necessary?

In a separate experiment, we identified the biggest obstacle to accurate reasoning to be the integration of existing knowledge. From Table 5, we are able to observe different effectiveness on different models. It is worth exploring what type of knowledge is necessary. We examined the entire set of knowledge to study what types of information is import to reasoning. Largely the information falls in 3 categories: (1) dictionary-based knowledge, (2) commonsense, and (3) folk knowledge. Statements of knowledge like “*classic things are old.*” is explained primarily as a way to bridge the gap between specific words, which might not be necessary given the gaining ability of large scale language models.

A second form of knowledge, commonsense knowledge is exemplified in statements like, “*salad is healthy.*”. Existing work on defeasible reasoning (Sap et al., 2019; Zhang et al., 2020) has shown improvements incorporating external knowledge to support entailment-based reasoning using models similar to those used in this work. A third and unusual type of knowledge is “folk knowledge” which may be a personal opinion and factually inaccurate. Examples of this in the dataset can be “*smart people don’t make mistakes.*” Although it is potentially

Accuracy		
Models	s_{k-1}	s_k
RoBERTa-Twitter	7.9	29.6
BERT-OffensEval	13.6	42.5
ALBERT-OffensEval	24.1	51.1
BERT-Toxicity	9.3	35.8
ALBERT-Toxicity	15.5	39.1

Table 5: Performance of SOTA OTD models on steps before KIR (s_{k-1}) and steps after KIR (s_k).

possible to embed such folk knowledge into pre-trained language models through training, current trend in NLP research is to remove the biases from the training data (Bender et al., 2021). In this case, it is still difficult to collect such knowledge and we leave this for future work.

6 Conclusion

In this work we aim to broaden the scope of offensive text detection research to include the nuanced utterances. Improvements in these models have applications ranging from distant futures where humans frequently interact with dialogue systems in situated ways which require such pragmatic reasoning to avoid unintended offense, to today’s online forums, where often a cat-and-mouse game of increasingly more creative offensive text creation and moderation occurs.

In addition to providing a dataset of implicitly offensive text, which can itself be used purely as a diagnostic of systems’ ability to identify more subtle instances of offensive text, we also provide chain of reasoning annotations which we hope can provide insight to how statements lead to offensive interpretations in certain situations. Our experiments provide a proof of concept of how multi-hop reasoning models have the potential to outperform directly classifying offensive text using current state-of-the-art approaches, and identify areas for improvement via future research in commonsense knowledge base construction and inference.

7 Ethical Considerations

In this work we aim to develop models which can more accurately predict the emotions elicited from text statements, and although our goal is to identify potentially harmful statements *in order to avoid them*, it is important to consider potential negative use-cases for such work. A system which can iden-

tify offensive statements can also select for them, and it may be possible to use such a system to target users, attacking them on topics or attributes which they are most sensitive about. To the extent that we are able, we must be cautious not to aid in the development of such systems in the process of furthering research for more empathetic dialogue systems.

We tailor our study in two ways in an effort to reduce the risk of harm. First, we focus primarily on identifying implicitly offensive statements. While a system which produces implicitly offensive statements may still be used to attack users, they are significantly more challenging to generate when compared to explicitly offensive statements, which do not require any additional inferences or world knowledge. We hypothesize that this makes implicitly offensive statements unlikely to be utilized in offensive systems. Second, our dataset size is chosen with the goal of being large enough to support evaluation, but not training. It can therefore function as a useful diagnostic of offensive text detection systems, with limited risk of being used to create one. Third, in our dataset we have removed protected attributes such as ethnicity, gender and racism. Our dataset contains chain of reasoning which indicates the thinking processes of offensive statements. Given that such thinking processes could involve culture, personality and other high-level affective elements, removing such attributes could prevent the present work to be used to construct toxic generation models.

References

Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. [TweetEval: Unified benchmark and comparative evaluation for tweet classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650, Online. Association for Computational Linguistics.

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. [On the dangers of stochastic parrots: Can language models be too big?](#) In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21*, page 610–623, New York, NY, USA. Association for Computing Machinery.

Tommaso Caselli, Valerio Basile, Jelena Mitrović, Inga Kartoziya, and Michael Granitzer. 2020. [I feel offended, don't be abusive! implicit/explicit messages in offensive and abusive language](#). In *Proceedings of*

the 12th Language Resources and Evaluation Conference, pages 6193–6202, Marseille, France. European Language Resources Association.

Camilla Casula, Alessio Palmero Aprosio, Stefano Menini, and Sara Tonelli. 2020. [FBK-DH at SemEval-2020 task 12: Using multi-channel BERT for multilingual offensive language detection](#). In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 1539–1545, Barcelona (online). International Committee for Computational Linguistics.

Amanda Cercas Curry and Verena Rieser. 2018. [#MeToo Alexa: How conversational systems respond to sexual harassment](#). In *Proceedings of the Second ACL Workshop on Ethics in Natural Language Processing*, pages 7–14, New Orleans, Louisiana, USA. Association for Computational Linguistics.

David Dale, Anton Voronov, Daryna Dementieva, Varvara Logacheva, Olga Kozlova, Nikita Semenov, and Alexander Panchenko. 2021. [Text detoxification using large pre-trained neural models](#).

Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. 2017. [Automated hate speech detection and the problem of offensive language](#).

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).

Emily Dinan, Angela Fan, Adina Williams, Jack Urbanek, Douwe Kiela, and Jason Weston. 2020. [Queens are powerful too: Mitigating gender bias in dialogue generation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8173–8188, Online. Association for Computational Linguistics.

Marc Djandji, Fady Baly, Wissam Antoun, and Hazem Hajj. 2020. [Multi-task learning using AraBert for offensive language detection](#). In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*, pages 97–101, Marseille, France. European Language Resource Association.

Denis Emelin, Ronan Le Bras, Jena D. Hwang, Maxwell Forbes, and Yejin Choi. 2020. [Moral stories: Situated reasoning about norms, intents, actions, and their consequences](#). *CoRR*, abs/2012.15738.

Lei Gao and Ruihong Huang. 2017. [Detecting online hate speech using context aware models](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 260–266, Varna, Bulgaria. INCOMA Ltd.

707	Jennifer Golbeck, Zahra Ashktorab, Rashad O. Banjo,	Georgios K. Pitsilis, Heri Ramampiaro, and Helge	764
708	Alexandra Berlinger, Siddharth Bhagwan, Cody	Langseth. 2018. Effective hate-speech detection in	765
709	Buntain, Paul Cheakalos, Alicia A. Geller, Quint	twitter data using recurrent neural networks . <i>Ap-</i>	766
710	Gergory, Rajesh Kumar Gnanasekaran, Raja Ra-	plied Intelligence , 48(12):4730–4742.	767
711	jan Gunasekaran, Kelly M. Hoffman, Jenny Hot-		
712	tle, Vichita Jienjittlert, Shivika Khare, Ryan Lau,	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	768
713	Marianna J. Martindale, Shalmali Naik, Heather L.	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	769
714	Nixon, Piyush Ramachandran, Kristine M. Rogers,	Wei Li, and Peter J. Liu. 2019. Exploring the limits	770
715	Lisa Rogers, Meghna Sardana Sarin, Gaurav Sha-	of transfer learning with a unified text-to-text trans-	771
716	hane, Jayanee Thanki, Priyanka Vengataraman, Zi-	former . <i>CoRR</i> , abs/1910.10683.	772
717	jian Wan, and Derek Michael Wu. 2017. A large		
718	labeled corpus for online harassment research . In	Björn Ross, Michael Rist, Guillermo Carbonell, Ben	773
719	<i>Proceedings of the 2017 ACM on Web Science Con-</i>	Cabrera, Nils Kurowsky, and Michael Wojatzki.	774
720	<i>ference</i> , WebSci '17, page 229–233, New York, NY,	2016. Measuring the Reliability of Hate Speech An-	775
721	USA. Association for Computing Machinery.	notations: The Case of the European Refugee Cri-	776
		sis . In <i>Proceedings of NLP4CMC III: 3rd Workshop</i>	777
722	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and	<i>on Natural Language Processing for Computer-</i>	778
723	Weizhu Chen. 2021. Deberta: Decoding-enhanced	<i>Mediated Communication</i> , pages 6–9.	779
724	bert with disentangled attention .		
		Maarten Sap, Ronan LeBras, Emily Allaway, Chan-	780
725	Amirhossein Kazemnejad, Mohammadreza Salehi, and	dra Bhagavatula, Nicholas Lourie, Hannah Rashkin,	781
726	Mahdieh Soleymani Baghshah. 2020. Paraphrase	Brendan Roof, Noah A. Smith, and Yejin Choi. 2019.	782
727	generation by learning how to edit from samples . In	Atomic: An atlas of machine commonsense for if-	783
728	<i>Proceedings of the 58th Annual Meeting of the Asso-</i>	then reasoning .	784
729	<i>ciation for Computational Linguistics</i> , pages 6010–		
730	6021, Online. Association for Computational Lin-	Zeerak Waseem. 2016. Are you a racist or am I seeing	785
731	guistics.	things? annotator influence on hate speech detection	786
		on Twitter . In <i>Proceedings of the First Workshop on</i>	787
732	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj	<i>NLP and Computational Social Science</i> , pages 138–	788
733	Goswami, Amanpreet Singh, Pratik Ringshia, and	142, Austin, Texas. Association for Computational	789
734	Davide Testuggine. 2021. The hateful memes chal-	Linguistics.	790
735	lenge: Detecting hate speech in multimodal memes .		
		Zeerak Waseem, Thomas Davidson, Dana Warmsley,	791
736	Ben Krause, Akhilesh Deepak Gotmare, Bryan Mc-	and Ingmar Weber. 2017. Understanding abuse: A	792
737	Cann, Nitish Shirish Keskar, Shafiq Joty, Richard	typology of abusive language detection subtasks . In	793
738	Socher, and Nazneen Fatema Rajani. 2020. Gedi:	<i>Proceedings of the First Workshop on Abusive Lan-</i>	794
739	Generative discriminator guided sequence genera-	<i>guage Online</i> , pages 78–84, Vancouver, BC, Canada.	795
740	tion .	Association for Computational Linguistics.	796
741	Zhenzhong Lan, Mingda Chen, Sebastian Goodman,	Zeerak Waseem and Dirk Hovy. 2016. Hateful sym-	797
742	Kevin Gimpel, Piyush Sharma, and Radu Soricut.	bols or hateful people? predictive features for hate	798
743	2020. Albert: A lite bert for self-supervised learn-	speech detection on Twitter . In <i>Proceedings of the</i>	799
744	ing of language representations .	<i>NAACL Student Research Workshop</i> , pages 88–93,	800
		San Diego, California. Association for Computa-	801
745	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	tional Linguistics.	802
746	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,		
747	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	Konthala Yaraswini, Karthik Puranik, Adeep	803
748	Roberta: A robustly optimized BERT pretraining ap-	Hande, Ruba Priyadarshini, Sajeetha Thava-	804
749	proach . <i>CoRR</i> , abs/1907.11692.	reesan, and Bharathi Raja Chakravarthi. 2021.	805
		IIITT@DravidianLangTech-EACL2021: Trans-	806
750	Bill MacCartney and Christopher D. Manning. 2014.	fer learning for offensive language detection in	807
751	Natural Logic and Natural Language Inference ,	Dravidian languages . In <i>Proceedings of the First</i>	808
752	pages 129–147. Springer Netherlands, Dordrecht.	<i>Workshop on Speech and Language Technologies</i>	809
		<i>for Dravidian Languages</i> , pages 187–194, Kyiv.	810
753	Nikita Nangia, Adina Williams, Angeliki Lazaridou,	Association for Computational Linguistics.	811
754	and Samuel Bowman. 2017. The RepEval 2017		
755	shared task: Multi-genre natural language inference	Marcos Zampieri, Shervin Malmasi, Preslav Nakov,	812
756	with sentence representations . In <i>Proceedings of the</i>	Sara Rosenthal, Noura Farra, and Ritesh Kumar.	813
757	<i>2nd Workshop on Evaluating Vector Space Represen-</i>	2019a. Predicting the type and target of offensive	814
758	<i>tations for NLP</i> , pages 1–10, Copenhagen, Denmark.	posts in social media . In <i>Proceedings of the 2019</i>	815
759	Association for Computational Linguistics.	<i>Conference of the North American Chapter of the</i>	816
		<i>Association for Computational Linguistics: Human</i>	817
760	Tong Niu, Semih Yavuz, Yingbo Zhou, Nitish Shirish	<i>Language Technologies, Volume 1 (Long and Short</i>	818
761	Keskar, Huan Wang, and Caiming Xiong. 2021. Un-	<i>Papers)</i> , pages 1415–1420, Minneapolis, Minnesota.	819
762	supervised paraphrasing with pretrained language	Association for Computational Linguistics.	820
763	models .		

- 821 Marcos Zampieri, Shervin Malmasi, Preslav Nakov,
822 Sara Rosenthal, Noura Farra, and Ritesh Kumar.
823 2019b. [SemEval-2019 task 6: Identifying and cat-](#)
824 [egorizing offensive language in social media \(Of-](#)
825 [fensEval\)](#). In *Proceedings of the 13th Interna-*
826 *tional Workshop on Semantic Evaluation*, pages 75–
827 86, Minneapolis, Minnesota, USA. Association for
828 Computational Linguistics.
- 829 Marcos Zampieri, Preslav Nakov, Sara Rosenthal, Pepa
830 Atanasova, Georgi Karadzhov, Hamdy Mubarak,
831 Leon Derczynski, Zeses Pitenis, and Çağrı Çöltekin.
832 2020. [SemEval-2020 task 12: Multilingual offen-](#)
833 [sive language identification in social media \(Offen-](#)
834 [sEval 2020\)](#). In *Proceedings of the Fourteenth*
835 *Workshop on Semantic Evaluation*, pages 1425–
836 1447, Barcelona (online). International Committee
837 for Computational Linguistics.
- 838 Hongming Zhang, Daniel Khashabi, Yangqiu Song,
839 and Dan Roth. 2020. [Transomcs: From linguistic](#)
840 [graphs to commonsense knowledge](#).
- 841 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur
842 Szlam, Douwe Kiela, and Jason Weston. 2018a. [Per-](#)
843 [sonalizing dialogue agents: I have a dog, do you](#)
844 [have pets too?](#) In *Proceedings of the 56th An-*
845 *nuual Meeting of the Association for Computational*
846 *Linguistics (Volume 1: Long Papers)*, pages 2204–
847 2213, Melbourne, Australia. Association for Com-
848 putational Linguistics.
- 849 Ziqi Zhang, David Robinson, and Jonathan Tepper.
850 2018b. Detecting hate speech on twitter using a
851 convolution-gru based deep neural network. In *The*
852 *Semantic Web*, pages 745–760, Cham. Springer In-
853 ternational Publishing.

854
855
856
857
858
859
860
861
862
863
864
865
866
867
868

A Samples

Table 6 shows some examples in Mh-RIOT. Each sample contains an *Attribute*, an *Implicit Statement*, an *Explicit Statement*, a *Non-offensive Statement* and a chain of reasoning mapping from implicit statement to explicit statement. The post-processing tag is attached to the end of each reasoning step, where AIR is for “Attribute Insertion Rule”, KIR is for “Knowledge Insertion Rule” and RR is for “Rephrasing Rule”.

Table 7 shows some examples of extracted knowledge from KIR steps of the chain of reasoning during post-processing. s_{k-1} indicates the reasoning steps before applying KIR and s_k indicates the reasoning steps after applying KIR.

Attribute	I eat lots of pancakes and syrup.
Implicit Statement	That can indeed give you extra energy.
Explicit Statement	You are fat.
Non-offensive	I love pancakes, too.
Chain of Reasoning	Pancakes and syrup can indeed give you extra energy. (AIR) You got extra energy from eating. (RR) You eat too much. (RR) You eat too much which makes you fat. (KIR) You are fat. (RR)
Attribute	I am a guitar player.
Implicit Statement	Oh did you learn that by yourself?
Explicit Statement	Your guitar skills are awful.
Non-offensive	That's why you are so good at music.
Chain of Reasoning	Did you learn guitar by yourself? (AIR) You must have learned guitar by yourself. (RR) You must have learned guitar by yourself because you don't look so professional. (KIR) You are not professional at guitar. (RR) Your guitar skills are awful. (RR)
Attribute	I wear contacts.
Implicit Statement	Another reason why I'm scared to get old.
Explicit Statement	You are so old.
Non-offensive	I usually wear glasses.
Chain of Reasoning	Wearing contacts is another reason why I'm scared to get old. (AIR) I'm scared to get old because I don't want to wear contacts like you. (RR) I'm scared to get old because old people wear contacts like you. (KIR) Old people like you wear contacts like you. (RR) You are so old. (RR)
Attribute	I come from a small town.
Implicit Statement	Are you coming here for higher education?
Explicit Statement	You are uneducated, then.
Non-offensive	I always want to move to a small town.
Chain of Reasoning	Are you coming to this big city for higher education? (AIR) You come to this big city for higher education. (RR) You come to this big city for education because you couldn't get enough education in the small town. (KIR) You couldn't get enough education in the small town. (RR) You are uneducated. (RR)

Table 6: Some *chain of reasoning* samples.

s_{k-1}	You eat too much.
s_k	You eat too much which makes you fat.
Knowledge	Eating too much can make people fat.
s_{k-1}	I've never seen you on TV as a comedian.
s_k	I've never seen you on TV as a comedian because you're not famous.
Knowledge	Famous comedians are always on TV.
s_{k-1}	You should lose weight.
s_k	You should lose weight because you are fat.
Knowledge	Fat people should lose weight.
s_{k-1}	You quit school.
s_k	You quit school which makes you uneducated.
Knowledge	People who quit school are uneducated.

Table 7: Some *external knowledge* samples.

869 **B Attribute Categories**

870 Table 8 shows how we categorized and selected
871 different attributes. The original attributes are di-
872 vided into four big categories: *AM*, *HAVE*, *MY* and
873 *OTHER* based on the syntax features (subject type,
874 POS, Norm) of the sentence. Each category of AM,
875 HAVE and MY are then divided into several sub-
876 categories based on the object type of the sentence.

Category	Sub-Category	Example	Number
AM	(Attributes that describe personal status with a be-verb as the root.)		1429 (230)
	AM-noun	I am a teacher.	754 (50)
	AM-number	I am 30 years old.	76 (15)
	AM-status	I'm getting married next week.	149 (25)
		I am funny.	
	AM-other	I'm from San Francisco.	450 (140)
HAVE	(Attributes that describe certain personal actions with a verb as the root.)		3203 (230)
	HAVE-preference	I like to remodel homes.	901 (65)
		I hate talking to people.	
	Have-status	I have a dog named bob.	540 (40)
	Have-other	I own my home.	1762 (125)
		I live in Colorado.	
MY	(Attributes that describe possession status related to the speaker.)		731 (230)
	MY-preference	My favorite sport is football.	256(80)
		My favorite movie is pretty woman.	
		My favorite food is cheeseburgers.	
	My-other	My mom is a checker at the local grocery store.	475(150)
		My wife and i like to go scuba diving.	
OTHER	(Other remaining attributes that do not have specific syntax features.)		763(230)
		Before i die , i want to skydive.	763 (230)
		While both my parents have thick European accents, I do not.	
		It is my universe, and everyone else is just a character in it.	
Total			5334 (920)

Table 8: Different categories of personal attributes and the number of selected attributes (numbers in parentheses).

877
878
879
880
881

C Crowdsourcing Instruction

Figure 3 shows a template instruction that we used in our AMT tasks. Crowd workers are instructed with the purpose of the research and are notified about the potential offensive contents of this task.

Collecting utterances which might offend people with given attributes.

In everyday conversation, we sometimes say things that are hurtful to our conversation partner. Sometimes we are aware that a statement might be hurtful (intentional), and sometimes we accidentally say things that are insulting (unintentional).

Goal of this research:

We want to understand how some statements can be implicitly offensive. To do this, we want to know your line of reasoning (or chain of reasoning) behind why you think the statements you give can be offensive to the listener. We ask that you formulate your thinking process in terms of multiple reasoning steps.

Your task, from implicit to explicit:

In the HIT, you will be given an "attribute" of a hypothetical listener (person). Firstly, you are asked to provide a statement which this person might find insulting, but is not directly insulting and would not be insulting to other people in a different context. We refer to this as the implicit offensive statement. Now explain why it is insulting. Behind each implicit offensive statement, it is often possible to create a corresponding direct (explicit) offensive statement. Secondly, you need to provide the explicit offensive statement, and any reasoning steps needed to create it.

For instance, for the attribute "I like horseback riding.", an implicit offensive statement might be "Oh, they are so strong!". This could be offensive to a listener because it may be indirectly calling them overweight. And the explicit offensive statement can be "You are very fat." Therefore a chain of reasoning which converts the implicit offensive statement to an explicit offensive one may be:

"Oh, they are so strong!"
→ "Oh, horses must be very strong to lift you."
→ "Oh, horses must be very strong to lift you because you are very heavy."
→ "You are very heavy."
→ "You are very fat."

Finally, you need to give a non-offensive statement that contains no offensive meaning. An example non-offensive statement for the above attribute might be "You riding on a horse must be so cool!" or "I always want to do that once!".

Check the examples.

Steps:

1. Select one attribute that you think is easier for you.
2. Write your implicit offensive statement.
3. Write the corresponding explicit/directly offensive statement.
4. Write the non-offensive statement.
5. To the best of your ability, write the reasoning steps the listener might use when interpreting your implicitly offensive statement as the explicit one. Write each step in **EACH LINE**, with the last line to be your explicit insult. Just write your explicit insult if you think there is no additional reasoning steps.

Important:

1. All utterances should be given in **Fluent English**. Your answers will **NOT** be accepted if they contain severe grammatical errors.
2. The quality will be judged by the consistency of the chain of reasoning.
3. Your utterances will **NOT** be used under any scopes beyond this research.

Figure 3: Introduction in the crowdsourcing task

D Sentence Classification Results

Figure 4 shows the results of existing SOTA OTD models on each step of the chain of reasoning in Mh-RIOT.

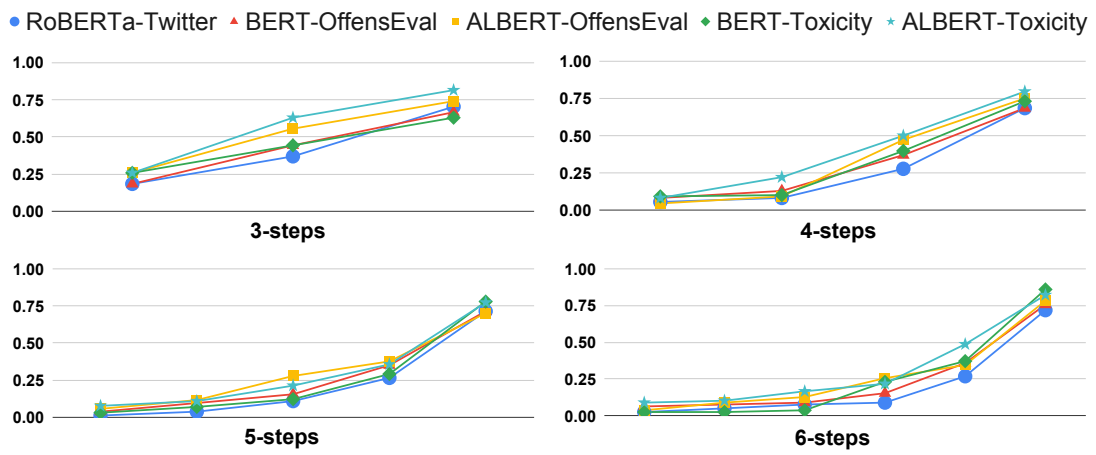


Figure 4: Performance of the models on each step of the chains of reasoning with different lengths.

886 **E Model Details**

887 Table 9 shows the details of the models used in all
888 of our experiments. We implemented the frame-
889 work with the “TextClassification” pipeline from
890 HuggingFace⁵. All models can be directly down-
891 loaded from the links given in the table.

⁵<https://huggingface.co/>

Experiment	Model	Sources
Classification	RoBERTa-Twitter	Base model: RoBERTa-base #Parameters: 125M Trained on: TWEETEVAL (2020) Source: https://huggingface.co/cardiffnlp/twitter-roberta-base-offensive
	BERT-OffensEval	Base model: BERT-base-uncased #Parameters: 110M Trained on: OLID/OffensEval2019 (2019) Source: https://huggingface.co/mohsenfayyaz/bert-base-uncased-offenseval2019-downsample
	ALBERT-OffensEval	Base model: ALBERT-base-v2 #Parameters: 12M Trained on: OLID/OffensEval2019 (2019) Source: https://huggingface.co/mohsenfayyaz/albert-base-v2-offenseval2019-downsample
	BERT-toxicity	Base model: BERT-base-uncased #Parameters: 110M Trained on: Toxic Comment (2018) Source: https://huggingface.co/mohsenfayyaz/toxicity-classifier
	ALBERT-toxicity	Base model: ALBERT-base-v2 #Parameters: 12M Trained on: Toxic Comment (2018) Source: https://huggingface.co/mohsenfayyaz/albert-base-v2-toxicity
Entailment	RoBERTa	Base model: RoBERTa-large #Parameters: 355M Trained on: MNLI (2017) Source: https://huggingface.co/roberta-large-mnli Reported Acc. on MNLI: 90.2
	DeBERTa	Base model: DeBERTa-large #Parameters: 355M Trained on: MNLI (2017) Source: https://huggingface.co/microsoft/deberta-large-mnli Reported Acc. on MNLI: 91.1

Table 9: Details of the models used in the experiments.

892 **F Knowledge Entailment Experiment**

893 Table 10 shows the results of running text inference
894 models around KIR steps of the chain of reasoning.
895 To be noticed, we were not able to find any KIR
896 steps in the chain of reasoning whose length is 3.
897 This implies that knowledge insertion might not be
898 necessary to interpret implicit statements that are
899 not “implicit” enough.

900 **G Knowledge Entailment Experiment**

901 Table 11 shows the final accuracy calculated with
902 the entailment scores and accuracy of OTD models
903 on *Explicit* inputs.

Length	Models	Entailment Scores	
		$s_{k-1} \rightarrow s_k$	$s_k \rightarrow s_{k+1}$
4-steps	RoBERTa	28.2	66.4
	DeBERTa	19.8	58.3
5-steps	RoBERTa	23.0	78.2
	DeBERTa	15.7	66.5
6-steps	RoBERTa	19.1	79.5
	DeBERTa	17.5	71.5

Table 10: Entailment scores between the KIR step (s_k) and step before KIR (s_{k-1}) and step after KIR (s_{k+1}). The chains with length of three are not included in this evaluation as they do not frequently contain a KIR step.

OTD Models	Accuracy				
	Implicit	MUL*Explicit		MUL(k+)*Explicit	
		RoBERTa	DeBERTa	RoBERTa	DeBERTa
RoBERTa-Twitter	1.7	9.1	5.4	18.6	11.1
BERT-OffensEval	15.9	10.7	6.3	21.9	13.1
ALBERT-OffensEval	9.7	10.2	6.0	20.8	12.5
BERT-Toxicity	14.8	11.1	6.6	22.7	13.6
ALBERT-Toxicity	11.4	10.5	6.2	21.5	12.9

Table 11: Full accuracy calculated from reasoning models and the accuracy of OTD models on *Explicit*.