

# Cross-lingual Transfer Learning for Irony Detection

Chen Hu

University of Minnesota  
huxxx853@umn.edu

Jiaqi Liu

University of Minnesota  
liu00687@umn.edu

Keyang Xuan

University of Minnesota  
xuan0008@umn.edu

## Abstract

Irony is a literary technique that is widely used across languages. A text's ironic intent is defined by its context incongruity. Accurate irony detection is crucial to effective sentiment analysis as well as harassment and hate speech detection in social media. However, detecting ironic statements is tough for a machine as irony enables one speaker or writer to conceal their true intention of negativity under the guise of overt positive representation. In this project, we aim to study this common feature of context incongruity in ironic sentences among different languages and formulate a universal multilingual model which is of paramount importance to increase the overall performance of irony detection. The preliminary result showed that irony detection benefited from mixed language datasets and multilingual models. In order to enhance the model's recognition of context incongruity, we proposed to use prompt tuning as our major technique. By inserting a learnable soft prompt at the beginning of each sentence, the fine-tuning is considered to be more directed to its downstream task. However, our prompt tuning results did not improve the performance significantly. The number of tokens appended and mask construction might have a big impact on our results. Future works should focus on the mask construction so that the soft prompt tuning could function as a hint for model to train.

## 1 Introduction

The development of the social website has been a rich source of non-literal language such as irony and sarcasm. As a result, research in automatic irony detection has thrived in recent years, for the purpose of better understanding and producing human language. The irony is defined in various ways, but one common agreement on irony identifies it as a figurative language whose actual meaning is

different from its superficial meaning (Kong and Qiu, 2011). Accurate irony detection could have a border impact on different research aspects. For example, in order to detect irony accurately, advanced text-mining techniques need to be applied. Besides, failure in irony detection would influence accurate sentimental analysis, and might further affect online harassment detection which has realistic usage in social media.

One of the challenges in irony detection comes from the inconsistency between contextual meaning and literal meaning. Sentiment polarity contrast is a common feature used to determine ironic language (Joshi et al., 2015). For example, in the sentence "I love being ignored", incongruity comes from the contradiction between the positive polarity word "love" and the negative polarity word "ignored". Other features including lexical factors (Kreuz and Caucci, 2007), punctuation marks, and syntactic patterns (Davidov et al., 2010) were investigated in English-based irony detection.

While studies based on English have provided a relatively comprehensive understanding of irony detection, there is still a lack of a systematic map about how irony detection could be applied to non-English language. Even though irony is a common linguistic phenomenon that appears in almost every language, some of its features still vary in different cultures and in structural properties of a specific language (Xing and Xu, 2015; Calvo et al., 2020; Cignarella et al., 2018). For example, words with all capitalized characters typically would be recognized to have non-literal meaning in English (Karoui et al., 2019) while other languages such as Chinese or Japanese do not have such capitalized characters. Chinese irony detection is more challenging because it is either composed of short statements in social media (Li et al., 2019), or it contains widely-used emoji that have unique meanings in the Chinese linguistic environment.

Although there are some papers investigating the

039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079

models and features specific to Chinese ironic detection, our interest mainly focuses on investigating the pattern similarity between irony sentences of different languages. We proposed a novel application of the multilingual model to learn English and Chinese irony. The rest of the paper is organized as follows: we discussed some related works in Section 2. In Section 3, we describe our overall objective, data resources, and proposed methods in detail. In Section 4, we showed our results in which the performance of monolingual and multilingual models was compared. We also showed our soft prompt tuning results. In Section 5, we concluded our work and discussed some potential failure reasons as well as future directions.

## 2 Related Work

### 2.1 Tokenization

One challenge in our project is that the tokenization methods for Chinese and English are different. Unlike English which naturally has space as a sign for tokenization, Chinese sentence as a collection of characters is more ambiguous to segment. The most popular type of tokenization adopted by pretrained-language models (PLM) including Chinese is sub-word tokenization, such as byte pair encoding (BPE) (Sennrich et al., 2016), Word-Piece (Schuster and Nakajima, 2012) and unigram language model segmentation (Kudo, 2018). Apart from sub-word tokenization, a simple character tokenizer is also common to use (Sun et al., 2019). Another difference is that Chinese words do not need to do lemmatization and stemming. That is to say, Chinese characters do not need to deal with tense or plural.

### 2.2 Irony Detection

The task of irony detection is to classify a piece of text as ironic or non-ironic. Current approaches to irony detection can be classified into three classes, namely rule-based approaches, classical feature-based machine learning methods, and deep neural network models (Zhang et al., 2019). Rule-based approaches generally rely on linguistic features such as sentiment lexicon or hashtags to detect irony (Sulis et al., 2016) while classical feature-based machine learning approaches use hand-crafted features for irony detection, such as sentiment lexicon, subjectivity lexicon, emotional category features, emotional dimension features or structural features (Farias and Rosso, 2017). In

this project, we more focus on deep learning-based approaches where (deep) features are automatically derived from texts using neural network models (Zhang et al., 2019).

Within deep learning-based models, some researchers use a pre-trained convolutional neural network for extracting sentiment, emotion and personality features for irony detection (Poria et al., 2016) while other researchers use CNN-LSTM structures for irony detection (Ghosh and Veale, 2017). However, context-based models utilize both content and contextual information required for irony detection, which leads to a growing interest in using neural language models for pre-training for various tasks in natural language processing. Given the highlighted importance of context to capture figurative language phenomena and the difficulties of data annotation, transfer learning approaches such as transformers are gaining attention in various domain adaptation problems. People (Potamias et al., 2020) propose Recurrent CNN Roberta (RCNN-RoBERTa), a hybrid neural architecture building on RoBERTa architecture, which is further enhanced with the employment and devise of a recurrent convolutional neural network. They report performance with an accuracy of 79% on the SARC dataset (Khodak et al., 2017). Similarly, an ensemble of Roberta and Albert on *GetitOffMyChest* dataset (Jaidka et al., 2020) achieve a performance of 85% accuracy with an F1 score of 0.55 (Dadu and Pant, 2020). BERT is used along with aspect-based sentiment analysis to extract the relation between context dialogue sequence and response. They obtain an F1 score of 0.73 on the Twitter dataset and 0.73 on the Reddit dataset (Javdan et al., 2020). However, among all sentimental polarity tests, irony classification shows its hardness, as the overall performance is over 10% lower (Zhang et al., 2019).

Few researchers perform the task of cross-lingual irony detection. It is shown that monolingual models trained separately in different languages using multilingual word representation can open the door to irony detection in different languages (Ghanem et al., 2020). The effectiveness of dependency-based syntactic features is also found in irony detection in a multilingual perspective (Cignarella et al., 2020). However, among all these approaches, modern cross-lingual transformer-based models have seldom been applied and overall performance barely reaches 70%.

180 It is our objective to study the pattern similarities  
181 between irony sentences of different languages and  
182 formulate a novel model that increases the accuracy  
183 of cross-lingual irony detection.

## 184 2.3 Prompt Tuning

185 Due to the rich knowledge obtained by Pretrained  
186 Language Models (PLMs), prompt tuning was pro-  
187 posed by a series of studies to bridge the gap  
188 between pre-training objectives and down-stream  
189 tasks (Hu et al., 2021; Schick and Schütze, 2020;  
190 Liu et al., 2021). Prompt turning shows excel-  
191 lent performance in few-shot learning and zero-  
192 shot learning. Among all prompt tuning mecha-  
193 nisms, "p-tuning" or "prompt tuning" learns "soft  
194 prompts" to condition frozen language models to  
195 perform specific downstream tasks (Lester et al.,  
196 2021). These mechanics improved model gener-  
197 alization and avoid over-fitting to a specific task  
198 domain. In our cases, we will use prompt tuning  
199 as it may explicitly reveal the context incongruity  
200 in an ironic sentence and thus improve the model  
201 performance. In this case, the multilingual irony  
202 patterns can be emphasized and learned with a lim-  
203 ited amount of data inputs. The detail is explained  
204 in the following section.

## 205 3 Problem Formulation

206 In this section, we present our training objective,  
207 general pipeline, and the data resources we use. For  
208 the whole project, we follow our pipeline tightly  
209 to explore the transfer-ability of irony property be-  
210 tween English and Chinese and we aim to improve  
211 the performance on the task compared to the base-  
212 line model by enhancement techniques such as  
213 prompt tuning.

### 214 3.1 Training Objective

215 Even though there is no previous work explor-  
216 ing the transferability of irony detection between  
217 two languages, many XLM approaches can be  
218 used for this task. Among those, we use XLM-  
219 RoBERTa (*XLM-R*) as the base model since it was  
220 pretrained on CommonCrawl data which contained  
221 more amount of data compared to the Wiki-100  
222 (Wenzek et al., 2020) and *XLM-R* has been proved  
223 to have better performance on many cross-lingual  
224 tasks comparing to other XLM approaches (Con-  
225 neau et al., 2020). Our goal is to create a bet-  
226 ter model architecture, which can beat the perfor-  
227 mance of *XLM-R* in our cross-lingual irony detec-

tion task.

Prompt tuning is recognized as an effective tool  
for few-shot learning and zero-shot learning tasks.  
It adapts the downstream tasks by inserting text  
pieces i.e. template, to the input and transforms  
a classification problem into a masked language  
modeling problem. In our case, since the training  
datasets only have a few thousand examples and  
we generally think irony pattern such as context  
incongruity has been learned by large multilingual  
models, we will utilize prompt tuning to explicitly  
direct the model to our specific task domain.

### 240 3.2 Data Resources

241 In this case, we pick English and Chinese as the  
242 linguistic basis for irony resources. For the En-  
243 glish irony resource, we combine the Reddit ironic  
244 corpus which is composed of about 2000 Reddit  
245 comments (Wallace et al., 2014) and Twitter ironic  
246 dataset from SemEval-2018 Task 3 (Van Hee et al.,  
247 2018) which contained about 4000 tweets. For  
248 the Chinese irony resource, we make a combined  
249 irony dataset by both Ciron which is collected from  
250 Weibo for irony annotation in simplified-style (Xi-  
251 ang et al., 2020) and NTU irony corpus which con-  
252 sists of messages in traditional Chinese version  
253 from the microblogging platform based on emoti-  
254 cons (Tang and Chen, 2014). Also, to keep the  
255 same format, we first convert all posts in the NTU  
256 irony corpus from the traditional version to the sim-  
257 plified Chinese version. Then, since the original  
258 scale in Ciron is from 1 to 5 which corresponds to  
259 not irony to strongest irony, we manually re-scale  
260 the label from which we convert the original 1,2 la-  
261 bels to -1 in the new dataset representing not irony  
262 and we convert the original 4,5 labels to 1 represent  
263 irony. For the rest sentences which are labeled as  
264 3, we manually delete those due to their neutral  
265 meaning.

### 266 3.3 General Pipeline

267 Fig 1. showed the general pipeline for our train-  
268 ing process. For our task, the general pipeline is  
269 composed of these steps: Firstly, we train a multi-  
270 lingual model based on either the Chinese training  
271 set, English training set, or Mix training set and at  
272 the same time train a monolingual model by the  
273 same training data. Then the fine-tuning model of  
274 the multilingual model will be separately tested on  
275 the Chinese testing set and English testing set and  
276 its performance will compare to the result from  
277 the monolingual fine-tuning model. The perfor-

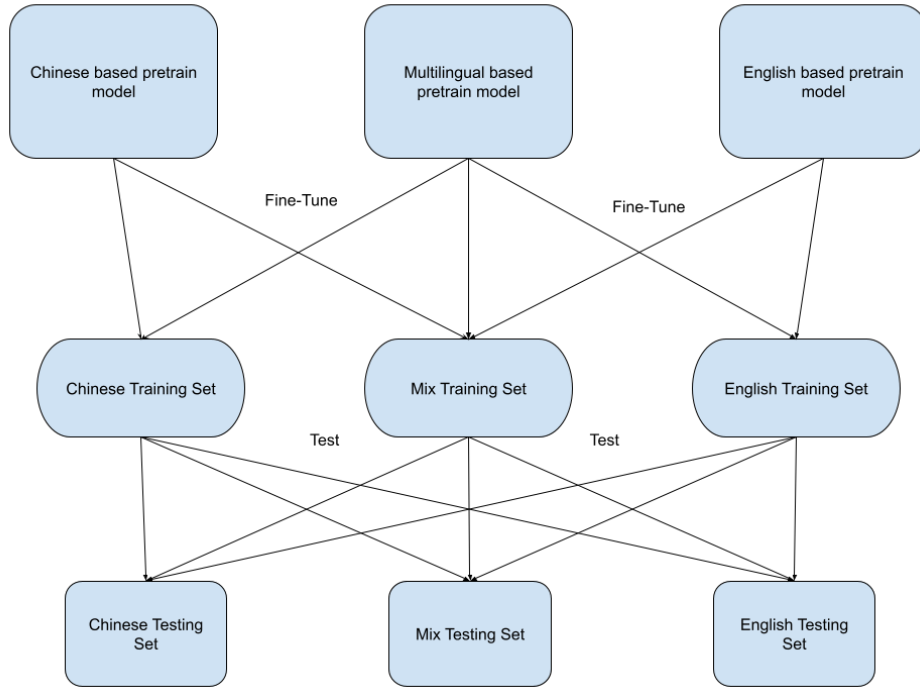


Figure 1: General Pipeline

278 mance of the monolingual model and *XLM-R* are  
 279 the benchmark for us to evaluate our self-defined  
 280 approach’s performance.

281 Fig 2. showed our general pipeline for prompt  
 282 tuning. We generally followed the algorithm stated  
 283 in the paper (Lester et al., 2021) and placed a 20-  
 284 token-lengths prompt in front of each text embed-  
 285 ding. The prompt weight was randomly initialized,  
 286 tuned during training steps, and universally the  
 287 same for all input sentences.

## 288 4 Results

### 289 4.1 Preliminary results - Comparison between 290 Monolingual and Multilingual Models

291 We first started with the experiment via two mono-  
 292 lingual models and one multilingual model. The  
 293 idea is to compare the performance of the mono-  
 294 lingual and multilingual models, and whether the  
 295 models mentioned above could benefit from a mix-  
 296 ture use of Chinese and English compared with  
 297 only using the single dataset. Table 1. summarized  
 298 the datasets that we used in this experiment.

299 GTP-2 (Radford et al., 2019) was used to test  
 300 how a model pretrained by English could detect  
 301 Chinese irony, English irony, and Mix dataset. We  
 302 trained GTP-2 on either the English training set or  
 303 the Mixture training set since it is not reasonable to  
 304 fine-tune it on the Chinese training set. Similarly,

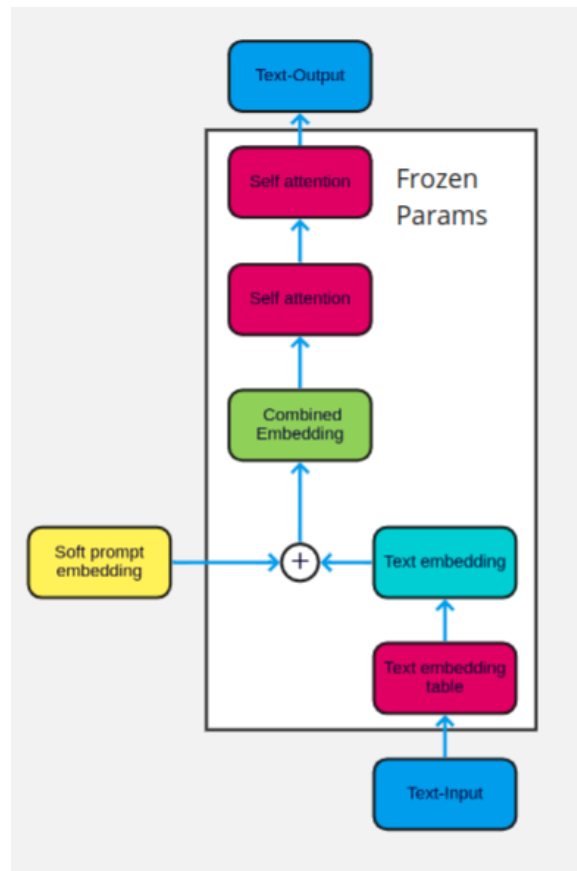


Figure 2: Soft Tuning Flowchart (Lester et al., 2021)

Training set			
	Irony	Non-irony	Total
English	2308	3068	5376
Chinese	1487	5582	7069
Mix	3795	8650	12445
Testing set			
	Irony	Non-irony	Total
English	441	733	1174
Chinese	389	1379	1768
Mix	830	1112	2942

Table 1: Training and testing dataset size

we used Bert Chinese Based model (CPT) (Shao et al., 2021) to test how a model pretrained by Chinese could detect irony in all training sets. We only fine-tuned the CPT model on the Chinese and Mixture datasets. XLM-RoBERTa (Conneau et al., 2019) was the multilingual model we selected for irony detection. It was trained on three datasets and tested on three datasets respectively as well. We used accuracy and F1 score as the evaluation metric. It was notable that due to the application of irony detection, missing an irony case was more harmful than misclassifying a non-irony case. That is to say, the recall was more important in such cases. Thus, we would value the F1 score heavily instead of accuracy.

#### 4.1.1 Discussion

Fig 3. showed the accuracy and F1 score for Chinese irony detection via GPT-2, CPT, and XRoberta models trained on either the English dataset (orange line) or the Mixture dataset (blue line). As mentioned above, there was no result for CPT trained in English. Based on the plot, one thing to notice was that models trained on a mixture dataset outperformed the models trained on the English dataset. The accuracy of Chinese irony detection by XRoberta trained in English slightly decreased compared to GTP-2 models while the F1 score increased on the other hand from 0.12 to 0.50. Since accuracy is the sum of correctly identifying true positive cases and true negative cases, one possible reason could be due to the inability of XRoberta to classify non-irony cases while the recall and precision remained relatively high. Besides, accuracy and F1 score did not change dramatically across different models when the models were trained on mixture data. One possible reason might be due to the information contained in the mixture dataset is sufficient for models to make predictions on Chi-

GPT-2			
	English	Chinese	Mix
English	0.66/0.34	0.78/0.13	0.78/0.13
Mix	0.65/0.28	0.89/0.70	0.79/0.50
Bert-Chinese			
	English	Chinese	Mix
Chinese	0.58/0.13	0.74/0.12	0.70/0.41
Mix	0.64/0.54	0.90/0.77	0.81/0.65
XLM-Roberta			
	English	Chinese	Mix
English	0.70/0.61	0.66/0.41	0.67/0.48
Mix	0.71/0.65	0.89/0.74	0.82/0.68

Table 2: Monolingual and multilingual model results (column - training sets, row - testing sets)

nese irony detection.

Fig 4. showed the accuracy and F1 score for English irony detection via GPT-2, CPT, and XRoberta models trained on either the Chinese dataset (orange line) or the mixture dataset (blue line). On the contrary with results from models tested on Chinese irony, there is a significant increase in accuracy and F1 scores from monolingual models to multilingual models no matter what training set was used. Similarly, models trained on the mixture dataset had better accuracy and F1 scores compared to the ones trained on the single dataset.

We summarized all our monolingual and multilingual model results in Table 2. In conclusion, we found that the multilingual model would have higher performance than the monolingual model (in the condition of the same training and testing dataset). Besides, even though Chinese and English are two languages with quite different language patterns, when we incorporated more information in the training (either by using the multilingual model or training with mix dataset), the performance was always better. This can indicate that there is a common pattern between irony in those two quite different languages, such as contextual incongruity.

#### 4.1.2 False Cases and Error Analysis

From the testing results, we generated some mistakenly classified examples to see whether there was any potential pattern that could make the model confuse. Table 3. showed a group of falsely classified sentences with their ground label and predicted results. Based on these samples, several patterns are worth mentioning: Firstly, in the sentence "I just love being ignored [smile]", the emoji mark contains a large portion of the ironic sentiment of the whole sentence. However, such kind of expres-

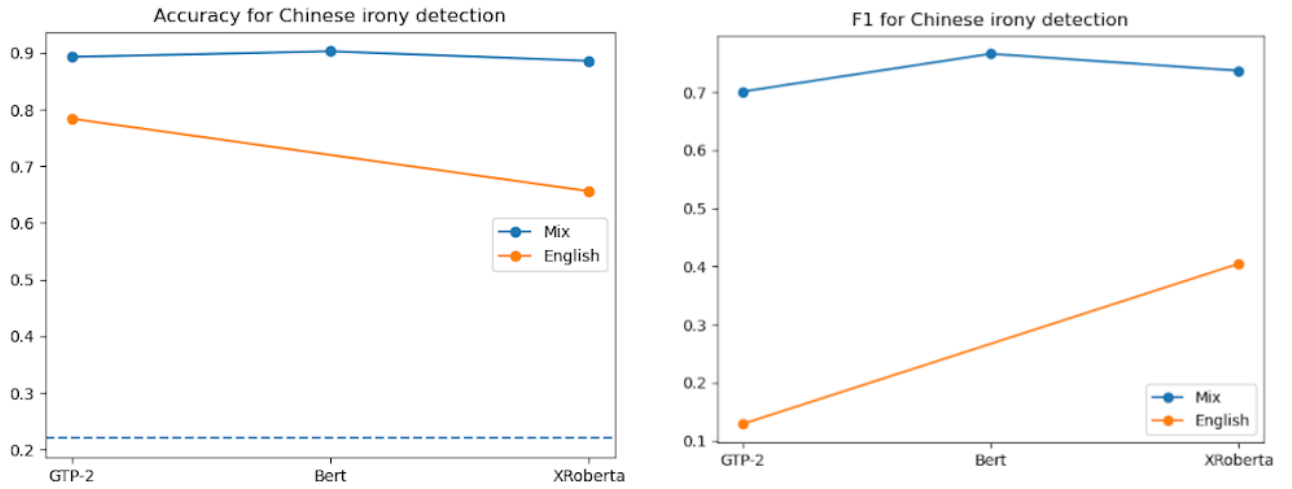


Figure 3: Accuracy and F1 for chinese irony detection

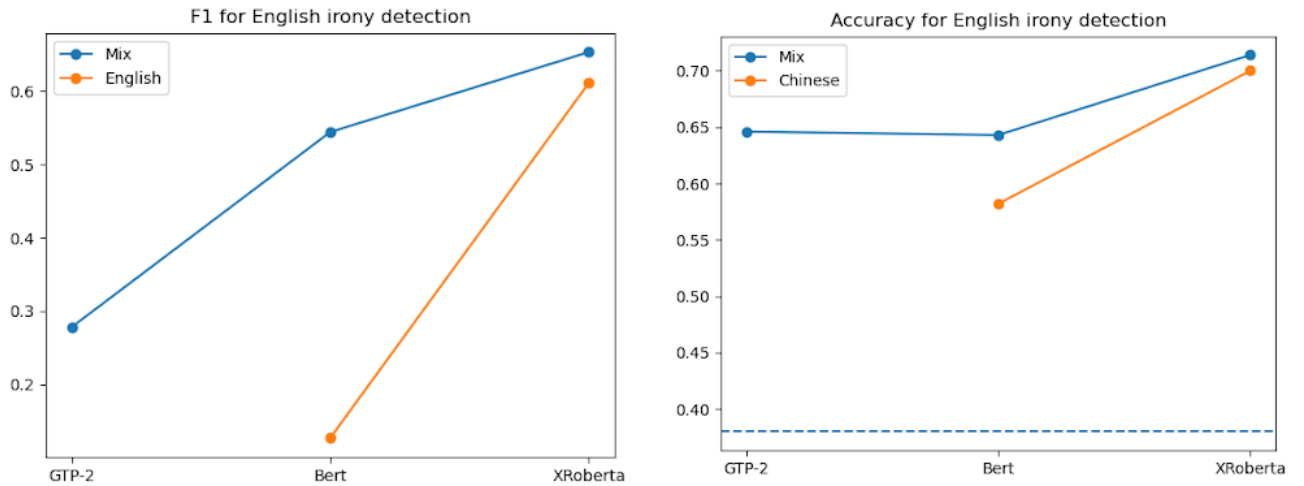


Figure 4: Accuracy and F1 for English irony detection

### Error Analysis

Content	True Label	Predicted Label
I just love being ignored [smile]	1	0
I just drank a healthy, homemade, all-fruit smoothie..in a @Budweiser beer glass#irony	1	0
I am so ready for Monday. #sarcasm	1	0
@GalloSays this game is pathetic. How are they losing this game?	0	1
@robinhosking where did THAT come from?!	0	1
天气可以再热一点没关系 :-&	1	0
这个酒柜的设计如何? 真的很棒是吧?	0	1
晚上的电视节目 可以再无聊一点点! :-&	1	0

Table 3: Failure cases

sion is difficult for the model to learn and capture the meaning which probably is the reason for false prediction. Secondly, post such as "I just drank a healthy, homemade, all-fruit smoothie..in a @Budweiser beer glass" really needs more supplementary evidence or details to help make a correct classification. In detail, the original context of such kind of post delivers trivial patterns related to irony semantics, and most time these posts are affiliated with other elements including pictures or material so we can only detect its meaning correctly only by analyzing these elements together. Thirdly, since the rhetorical question is one of the commonest formats of irony sentiment, our model might be overly sensitive to the question mark. As a result, some general questions without any irony semantic may also be classified into irony.

**4.1.3 Limitations**

Although the multilingual model trained on a mixture dataset had relatively high accuracy and F1 score, the performance was not what we expected. One limitation is that in the English testing dataset, some hashtags including irony were not deleted. Based on the failure cases shown, it seemed that having irony as a hashtag did not improve the performance. We would still try to avoid such situations in any datasets. Another limitation is that the models still relied heavily on superficial meaning instead of contextual meaning. An example could be the usage of emojis in the testing data mentioned above. One way to improve the understanding of contextual meaning could be grouping together some combinations of words. This could be either done by tokenization or by adding a soft prompt so that the model had the freedom to learn patterns with a small subset of training samples.

**4.2 Enhancement by Prompt Tuning**

**4.2.1 Evaluation in GPT-3 Interface**

Due to the limitation and the attempt to use soft tuning, We first explored several prompt tuning techniques including the aforementioned one in GPT-3 interface to verify whether prompt tuning is appropriate for our specific task. To be notifiable, we decided to use failure cases in the preliminary result and checked whether a proportion of them could be recognized by the model when a prompt is added.

According to Fig.5 and Fig.6, we could find that when the input was a single ironic sentence, GTP-3

Check whether the sentence is ironic or not  
 Input: A member of PETA wears leather shoes.

The sentence is not ironic.

Figure 5: Example of irony detection via single prompt

Check whether the sentence is ironic or not  
 Example: A marriage counselor files for divorce.  
 Output: Yes  
 Example: The police station gets robbed.  
 Output: Yes

Example: A fire station burns down.  
 Output: Yes

Input: A member of PETA wears leather shoes.  
 Yes

Figure 6: Example of irony detection via several prompts

was unable to correctly identify the property. However, after giving several examples of ironic sentences, there was a certain chance for GTP-3 to identify whether the input sentences are ironic or not. Thus, prompt tuning might be appropriate for our ironic detection tasks.

**4.2.2 Soft Prompt Tuning**

In this step, We utilized the "soft prompt tuning" technique, in which a trainable length of embedding will be added to the front of each input sentence embedding, as depicted in Fig. 2. The prompt weight was randomly initialized, tuned during training steps, and universally the same for all input sentences. We compared the results with a pre-trained fine-tuning model to analyze the efficiency of soft prompt tuning techniques.

In table 4, we showed results on soft prompt tuning with GPT-2 and CPT models. Due to the time limitation and the incompatibility of the model, we did not accomplish testing prompt tuning on the multilingual XLM-Roberta model. Overall, compared to the preliminary results shown in Table 2,

428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449

GPT-2			
	English	Chinese	Mix
English	0.66/0.04	0.78/0.21	0.75/0.14
Mix	0.45/0.31	0.69/0.14	0.67/0.11
Bert-Chinese			
Chinese	0.56/0.14	0.79/0.09	0.69/0.08
Mix	0.47/0.2	0.7/0.11	0.64/0.12
XLM-Roberta			
English	/	/	/
Mix	/	/	/

Table 4: Prompt tuning results (column - training sets, row - testing sets)

we found that the general performance with prompt tuning was not better than without it. Among the 12 trials, only CPT model trained on the Chinese set and tested on the Chinese set had a 0.05 improvement in accuracy. We attributed this to insufficient prompt length tuning and inadequate prompt initialization and the detailed analysis is explained in the next section.

## 5 Future Directions and Conclusions

Unexpectedly, our soft prompting tuning results did not have a significant improvement compared to the preliminary results. One possible explanation could be the number of tokens appended to the original prompts. Based on the previous paper, the number of tokens was a variable that need to be tuned and had a huge impact on model performances. Unfortunately, we spent lots of time on model tuning and training and only had limited time to tune the model hyperparameters. Besides, the way that model initialized the appended tokens would also affect the model performance. The current model randomly assigned token values to the prompt which did not introduce additional useful information to the models in the first place. From Fig.6, we could find that the example prompts need to have the same language structure and patterns as the input prompt in order to have GTP-3 work.

Thus, one future direction that needs to pay more attention to is the initial prompt construction. For example, one prompt tuning technique called composition converts the task into several sub-tasks. For the irony classification task with this technique, we can aggregate the sentence with the prompt such as "This sentence is [Mask] from context. This sentence is [Mask] from meaning. So this sentence is [Mask]". The first and second masks can choose

from "positive" and "negative" while the third mask can choose from "irony" and "non-irony". Even though this initialized prompt will change during the training process, based on research in (Lester et al., 2021), the final prompt would still have a high chance to be localized around the initial prompt. In this way, we explicitly formulate the mechanism of detecting irony and expose that to the model. We expect this type of prompt tuning will help generalize the model and improve performance as the soft prompting functions as a hint for our model to train in the process.

To conclude, we found that the multilingual model would generally have higher performance than the monolingual model. Even though Chinese and English are two languages with quite different language patterns, the better performance with mixed datasets indicates that there is a common pattern between irony and different languages, such as contextual incongruity. Despite our prompt tuning techniques did not work as expected, we still think it is a strong model enhancement technique and we want to investigate more in the future.

## References

- Hiram Calvo, Omar J Gambino, and Consuelo Varinia García Mendoza. 2020. Irony detection using emotion cues. *Computación y Sistemas*, 24(3):1281–1287.
- Alessandra Teresa Cignarella, Valerio Basile, Manuela Sanguinetti, Cristina Bosco, Paolo Rosso, and Farah Benamara. 2020. Multilingual irony detection with dependency syntax and neural models. *arXiv preprint arXiv:2011.05706*.
- Alessandra Teresa Cignarella, Simona Frenda, Valerio Basile, Cristina Bosco, Viviana Patti, Paolo Rosso, et al. 2018. Overview of the evalita 2018 task on irony detection in italian tweets (ironita). In *Sixth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian (EVALITA 2018)*, volume 2263, pages 1–6. CEUR-WS.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Asso-*



538			
539		<i>ciation for Computational Linguistics</i> , pages 8440–	
540		8451, Online. Association for Computational Lin-	
		guistics.	
541	Tanvi Dadu and Kartikey Pant. 2020. Sarcasm detec-		
542		tion using context separators in online discourse. In	
543		<i>Proceedings of the Second Workshop on Figurative</i>	
544		<i>Language Processing</i> , pages 51–55.	
545	Dmitry Davidov, Oren Tsur, and Ari Rappoport. 2010.		
546		Semi-supervised recognition of sarcasm in twitter	
547		and amazon. In <i>Proceedings of the fourteenth con-</i>	
548		<i>ference on computational natural language learning</i> ,	
549		pages 107–116.	
550	DI Hernández Farias and Paolo Rosso. 2017. Irony, sar-		
551		casm, and sentiment analysis. In <i>Sentiment Analysis</i>	
552		<i>in Social Networks</i> , pages 113–128. Elsevier.	
553	Bilal Ghanem, Jihen Karoui, Farah Benamara, Paolo		
554		Rosso, and Véronique Moriceau. 2020. Irony de-	
555		tection in a multilingual context. In <i>European Con-</i>	
556		<i>ference on Information Retrieval</i> , pages 141–149.	
557		Springer.	
558	Aniruddha Ghosh and Tony Veale. 2017. Magnets for		
559		sarcasm: Making sarcasm detection timely, contex-	
560		tual and very personal. In <i>Proceedings of the 2017</i>	
561		<i>Conference on Empirical Methods in Natural Lan-</i>	
562		<i>guage Processing</i> , pages 482–491.	
563	Shengding Hu, Ning Ding, Huadong Wang, Zhiyuan		
564		Liu, Juanzi Li, and Maosong Sun. 2021. Knowl-	
565		edgeable prompt-tuning: Incorporating knowledge	
566		into prompt verbalizer for text classification. <i>arXiv</i>	
567		<i>preprint arXiv:2108.02035</i> .	
568	Kokil Jaidka, Iknoor Singh, Ni Chhaya, Lyle Ungar, and		
569		Jiahui Lu. 2020. A report of the cl-aff offmychest	
570		shared task: Modeling supportiveness and disclosure.	
571	Soroush Javdan, Behrouz Minaei-Bidgoli, et al. 2020.		
572		Applying transformers and aspect-based sentiment	
573		analysis approaches on sarcasm detection. In <i>Pro-</i>	
574		<i>ceedings of the second workshop on figurative lan-</i>	
575		<i>guage processing</i> , pages 67–71.	
576	Aditya Joshi, Vinita Sharma, and Pushpak Bhat-		
577		tacharyya. 2015. Harnessing context incongruity for	
578		sarcasm detection. In <i>Proceedings of the 53rd An-</i>	
579		<i>nuual Meeting of the Association for Computational</i>	
580		<i>Linguistics and the 7th International Joint Confer-</i>	
581		<i>ence on Natural Language Processing (Volume 2:</i>	
582		<i>Short Papers)</i> , pages 757–762.	
583	Jihen Karoui, Farah Benamara, and Veronique Moriceau.		
584		2019. <i>Automatic Detection of Irony: Opinion Mining</i>	
585		<i>in Microblogs and Social Media</i> . John Wiley & Sons.	
586	Mikhail Khodak, Nikunj Saunshi, and Kiran Vodrahalli.		
587		2017. A large self-annotated corpus for sarcasm.	
588		<i>arXiv preprint arXiv:1704.05579</i> .	
589	Lingpeng Kong and Likun Qiu. 2011. Formalization		
590		and rules for recognition of satirical irony. In <i>2011</i>	
591		<i>International Conference on Asian Language Pro-</i>	
592		<i>cessing</i> , pages 135–138. IEEE.	
	Roger Kreuz and Gina Caucci. 2007. Lexical influences		593
	on the perception of sarcasm. In <i>Proceedings of the</i>		594
	<i>Workshop on computational approaches to Figurative</i>		595
	<i>Language</i> , pages 1–4.		596
	Taku Kudo. 2018. <a href="#">Subword regularization: Improv-</a>		597
	<a href="#">ing neural network translation models with multiple</a>		598
	<a href="#">subword candidates</a> . In <i>Proceedings of the 56th An-</i>		599
	<i>nuual Meeting of the Association for Computational</i>		600
	<i>Linguistics (Volume 1: Long Papers)</i> , pages 66–75,		601
	Melbourne, Australia. Association for Computational		602
	Linguistics.		603
	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.		604
	The power of scale for parameter-efficient prompt		605
	tuning. <i>arXiv preprint arXiv:2104.08691</i> .		606
	An-Ran Li, Emmanuele Chersoni, Rong Xiang, Chu-		607
	Ren Huang, Qin Lu, et al. 2019. On the “easy” task		608
	of evaluating chinese irony detection.		609
	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding,		610
	Yujie Qian, Zhilin Yang, and Jie Tang. 2021. Gpt		611
	understands, too. <i>arXiv preprint arXiv:2103.10385</i> .		612
	Soujanya Poria, Erik Cambria, Devamanyu Hazarika,		613
	and Prateek Vij. 2016. A deeper look into sarcastic		614
	tweets using deep convolutional neural networks.		615
	<i>arXiv preprint arXiv:1610.08815</i> .		616
	Rolandos Alexandros Potamias, Georgios Siolas, and		617
	Andreas-Georgios Stafylopatis. 2020. A transformer-		618
	based approach to irony and sarcasm detection.		619
	<i>Neural Computing and Applications</i> , 32(23):17309–		620
	17320.		621
	Alec Radford, Jeff Wu, Rewon Child, David Luan,		622
	Dario Amodei, and Ilya Sutskever. 2019. Language		623
	models are unsupervised multitask learners.		624
	Timo Schick and Hinrich Schütze. 2020. Exploit-		625
	ing cloze questions for few shot text classification		626
	and natural language inference. <i>arXiv preprint</i>		627
	<i>arXiv:2001.07676</i> .		628
	Mike Schuster and Kaisuke Nakajima. 2012. <a href="#">Japanese</a>		629
	<a href="#">and korean voice search</a> . In <i>2012 IEEE International</i>		630
	<i>Conference on Acoustics, Speech and Signal Process-</i>		631
	<i>ing (ICASSP)</i> , pages 5149–5152.		632
	Rico Sennrich, Barry Haddow, and Alexandra Birch.		633
	2016. <a href="#">Neural machine translation of rare words with</a>		634
	<a href="#">subword units</a> . In <i>Proceedings of the 54th Annual</i>		635
	<i>Meeting of the Association for Computational Lin-</i>		636
	<i>guistics (Volume 1: Long Papers)</i> , pages 1715–1725,		637
	Berlin, Germany. Association for Computational Lin-		638
	guistics.		639
	Yunfan Shao, Zhichao Geng, Yitao Liu, Junqi Dai,		640
	Fei Yang, Li Zhe, Hujun Bao, and Xipeng Qiu.		641
	2021. Cpt: A pre-trained unbalanced transformer		642
	for both chinese language understanding and genera-		643
	tion. <i>arXiv preprint arXiv:2109.05729</i> .		644

- 645 Emilio Sulis, Delia Irazú Hernández Farías, Paolo  
646 Rosso, Viviana Patti, and Giancarlo Ruffo. 2016. Fig-  
647 urative messages and affect in twitter: Differences be-  
648 tween# irony,# sarcasm and# not. *Knowledge-Based*  
649 *Systems*, 108:132–143.
- 650 Yu Sun, Shuohuan Wang, Yu-Kun Li, Shikun Feng,  
651 Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu,  
652 Hao Tian, and Hua Wu. 2019. [ERNIE: enhanced](#)  
653 [representation through knowledge integration](#). *CoRR*,  
654 abs/1904.09223.
- 655 Yi-jie Tang and Hsin-Hsi Chen. 2014. [Chinese irony](#)  
656 [corpus construction and ironic structure analysis](#). In  
657 *Proceedings of COLING 2014, the 25th International*  
658 *Conference on Computational Linguistics: Technical*  
659 *Papers*, pages 1269–1278, Dublin, Ireland. Dublin  
660 City University and Association for Computational  
661 Linguistics.
- 662 Cynthia Van Hee, Els Lefever, and Véronique Hoste.  
663 2018. Semeval-2018 task 3: Irony detection in en-  
664 glish tweets. In *Proceedings of The 12th Interna-*  
665 *tional Workshop on Semantic Evaluation*, pages 39–  
666 50.
- 667 Byron C. Wallace, Do Kook Choe, Laura Kertz, and  
668 Eugene Charniak. 2014. [Humans require context to](#)  
669 [infer ironic intent \(so computers probably do, too\)](#).  
670 In *Proceedings of the 52nd Annual Meeting of the*  
671 *Association for Computational Linguistics (Volume 2:*  
672 *Short Papers)*, pages 512–516, Baltimore, Maryland.  
673 Association for Computational Linguistics.
- 674 Guillaume Wenzek, Marie-Anne Lachaux, Alexis Con-  
675 neau, Vishrav Chaudhary, Francisco Guzmán, Ar-  
676 mand Joulin, and Edouard Grave. 2020. [CCNet:](#)  
677 [Extracting high quality monolingual datasets from](#)  
678 [web crawl data](#). In *Proceedings of the 12th Lan-*  
679 *guage Resources and Evaluation Conference*, pages  
680 4003–4012, Marseille, France. European Language  
681 Resources Association.
- 682 Rong Xiang, Xuefeng Gao, Yunfei Long, Anran Li,  
683 Emmanuele Chersoni, Qin Lu, and Chu-Ren Huang.  
684 2020. [Ciron: a new benchmark dataset for Chinese](#)  
685 [irony detection](#). In *Proceedings of the 12th Lan-*  
686 *guage Resources and Evaluation Conference*, pages  
687 5714–5720, Marseille, France. European Language  
688 Resources Association.
- 689 Frank Z Xing and Yang Xu. 2015. A logistic regression  
690 model of irony detection in chinese internet texts.  
691 *Res. Comput. Sci.*, 90:239–249.
- 692 Shiwei Zhang, Xiuzhen Zhang, Jeffrey Chan, and Paolo  
693 Rosso. 2019. Irony detection via sentiment-based  
694 transfer learning. *Information Processing & Manage-*  
695 *ment*, 56(5):1633–1644.