

# Bert-Based Promotional Words Detection Classifier Development

Zikun Lin (supervised by: John R. Kender)

*Columbia University, New York, NY, United States*

---

## Abstract

From the analysis last semester, we can see the necessity for us to find a way to get rid of advertising and promotion words from our corpora. After manually picking out promotion sentences, we get a list of "blackwords" from video descriptions corpora, which can be used as a dataset in machine learning. In this report, I use the traditional model Word2Vec and the most advanced NLP model "BERT" combining with other machine learning methods to help us picking out the promotion sentences from video descriptions.

---

## 1. Introduction

As discussed in the previous report, there are a lot of promotion words and sentences appearing in the video descriptions. Typical categories of these words contains:

- (a) Media promotion and subscription request. Such as: news, videos, subscribe, visit.
- (b) URL links and social network account. Such as: twitter, facebook.
- (c) Media name. Such as: vice, 'arirang. In fact, the most "vice"s appear in the descriptions are not acting as "vice president", but "VICE News" instead.

Before applying NLP and ML methods to video descriptions and transcripts dataset, we hope to get rid of these words as many as possible. However, applying manual work to all of these corpora is neither efficient nor doable. Therefore, we hope to develop a "blackwords-detection classifier" to assist us complete this task.

Several traditional NLP methods can take essential roles in our research. Especially in recent years, with the help of deep learning, models like LSTM, Word2Vec, BERT become popular in the NLP field. It enables us to convert words into vectors, making measurements and operations like clustering and classification possible. We can use these well-developed models and methods in our research to make the analysis more quickly and efficiently.

17 The main contributions of this research can be summarized as follows:

- 18 • I utilized the most recent context-based natural language processing model “BERT”  
19 to get embedding vectors for each word instance in our corpora.
- 20 • I utilized the “blackwords” dataset I collected last semester to train a “blackwords-  
21 detection classifier” based on SVM classifier. After grid search and fine-tuning, I  
22 achieved the accuracy of more than 60%. This classifier is able to be used in Stanley’s  
23 deep cross-cultural system.
- 24 • I used video descriptions in AlphaGo and Florida shooting events as training set to  
25 train a final classifier. This classifier obtained a nearly 90% accuracy in the testing  
26 set, i.e. video descriptions in lunar rover and Thailand cave rescue events.

27 To help readers better understand the models used in experiments, I describe the datasets  
28 and some related work in Section 2. Then I will focus on my experiments in Section 4. And  
29 in Section 5, I will conclude my research and plan for the future.

## 30 2. Related Work

### 31 2.1. News Descriptions

32 As described in my previous report, in order to carry out the NLP part of our experiment,  
33 having clean data based on our task is necessary. In my experiments, I mainly use the news  
34 description data collected by Andy beforehand. These data contain paragraphs from the  
35 video description on several topics (AlphaGo, Florida shooting, lunar rover, Thailand cave  
36 rescue) in both YouTube and CGTN. For promotion sentences, I use the dataset manually  
37 picked by myself last semester, which contains two topics (AlphaGo and Florida shooting).

### 38 2.2. Word2Vec

39 Representing words in a vector space is an efficient way to group similar words and analyze  
40 the distribution of a set of words. [Rumelhart et al. \(1988\)](#), [Mikolov et al. \(2013a\)](#) and  
41 [Mikolov et al. \(2013b\)](#)'s papers described methods and improvements to represent word and  
42 phrases and their compositionality on a vector space. Particularly, [Mikolov et al. \(2013a\)](#)  
43 introduced the Skip-gram model, which is an efficient method for learning high-quality vector  
44 representations of words from large amounts of unstructured text data, and it is one of the  
45 most popular ways to train word vectors.

46 In order to carry out our experiments quickly, I use Google's pre-trained word and phrase  
47 vectors<sup>1</sup>, so that we do not need to take much time training from massive datasets. Instead,  
48 with the help of [Řehůřek and Sojka \(2010\)](#)'s Gensim library, we only need to call

```
49 model = gensim.models.KeyedVectors.load_word2vec_format()
```

50 function to load the model and get the vector representation we need.

### 51 2.3. BERT

52 In Nov.2018, [Devlin et al. \(2018\)](#) from Google AI Google AI open sourced a new technique  
53 for NLP pre-training called **B**idirectional **E**ncoder **R**epresentations from **T**ransformers, or  
54 BERT.

55 BERT is the first deeply bidirectional, unsupervised language representation, pre-trained  
56 using only a plain text corpus. Especially, with the "bidirectional" strategy, BERT is able  
57 to embed words according to its context. In order to utilize this strategy, they use the

---

<sup>1</sup><https://code.google.com/archive/p/word2vec/>

58 straightforward technique of masking out some of the words in the input and then condition  
59 each word bidirectionally to predict the masked words.

60 With the help of their pretrained model, we are able to train models on our corpora in  
61 several hours just on our own machines. In my experiments, I used their most recently  
62 released “BERT-Large, Uncased (Whole Word Masking)” model<sup>2</sup>, with 24 layers and 1024  
63 hidden dimensions on their Github repo<sup>3</sup>. After minor modifications to their code (shown  
64 in Appendix 1), we can easily get a 1024-dimension embedding for every word occurrence in  
65 our training corpus.

## 66 2.4. SVM

67 Support Vector Machine (SVM) is a supervised learning model. The original SVM algorithm  
68 was invented by Vladimir N. Vapnik and Alexey Ya. Chervonenkis in 1963. After about 30  
69 years, [Boser et al. \(1992\)](#) introduced kernel SVM to create nonlinear margins. The current  
70 standard incarnation (soft margin) was proposed by [Cortes and Vapnik \(1995\)](#).

71 SVM can greatly help our task by providing a lightweight but effective classifier to classify  
72 two different types of words (“black” and “white”). In my experiments, I mainly use the  
73 SVM model implemented in sklearn Python package. In this model, we are able to adjust  
74 the kernel coefficient for radial-based-function (gamma) and penalty parameter of the error  
75 term (C).

---

<sup>2</sup>[https://storage.googleapis.com/bert\\_models/2019\\_05\\_30/wwm\\_uncased\\_L-24\\_H-1024\\_A-16.zip](https://storage.googleapis.com/bert_models/2019_05_30/wwm_uncased_L-24_H-1024_A-16.zip)

<sup>3</sup><https://github.com/google-research/bert>

## 76 3. Hand-labeled blackwords

### 77 3.1. Video Descriptions on AlphaGo

#### 78 3.1.1. Blacklist Words for Video Descriptions in English

79 Sentences containing advertisement and promotion are treated as “blacklist sentences”. The  
80 typical “blacklist sentence” is like this one<sup>4</sup>:

81 *Subscribe to VICE News here: <http://bit.ly/Subscribe-to-VICE-News>*  
*Check out VICE News for more: <http://vicenews.com>*  
*Follow VICE News here:*  
*Facebook: <https://www.facebook.com/vicenews>*  
*Twitter: <https://twitter.com/vicenews>*  
*Tumblr: <http://vicenews.tumblr.com/>*  
*Instagram: <http://instagram.com/vicenews>*  
*More videos from the VICE network: <https://www.fb.com/vicevideo>*

82 The top 60 words in the blacklist after eliminating punctuation are shown in Table 1.

Word	Frequency	Word	Frequency
news	110	with	22
to	98	is	22
the	90	out	22
on	88	in	21
	77	a	21
subscribe	70	channel	20
vice	60	your	19
us	53	this	17
official	53	cnbc	17
for	52	website	16
and	50	arirang	16
here	46	by	16
<a href="http://www.facebook.com/arirangtvtwitter">httpwwwfacebookcomarirangtvtwitter</a>	39	my	16
<a href="http://twitter.com/arirangworldinstagram">httptwittercomarirangworldinstagram</a>	39	<a href="http://bit.ly/subscribe/vice/news/check">httpbitlysubscribetovicenewscheck</a>	15
visit	36	<a href="http://vicenews.com/follow">httpvicenewscomfollow</a>	15
of	35	herefacebook	15
‘arirang	34	<a href="https://www.facebook.com/vicenewstwitter">httpswwwfacebookcomvicenewstwitter</a>	15
news’	34	<a href="http://stwitter.com/vicenewstumblr">httpstwittercomvicenewstumblr</a>	15
pagesfacebooknews	34	<a href="http://vicenewstumblr.com/instagram">httpvicenewstumblrcominstagram</a>	15
<a href="http://www.facebook.com/news/arirang/homepage">httpwwwfacebookcomnewsariranghomepage</a>	34	<a href="http://instagram.com/vicenewsmore">httpinstagramcomvicenewsmore</a>	15
<a href="http://www.arirang.com/facebook">httpwwwarirangcomfacebook</a>	34	network	15
<a href="http://instagram.com/arirangworld">httpinstagramcomarirangworld</a>	34	<a href="https://www.fb.com/vicevideo">httpswwwfbcomvicevideo</a>	15
more	31	like	14
facebook	29	tv	14
from	27	intel	14
youtube	26	please	13
you	26	at	13
videos	25	cbs	13
our	25	app	12
twitter	24	software	12

Table 1: Top 60 Words in English Blacklist on AlphaGo

83 From this blacklist, we can see several typical categories of words:

<sup>4</sup><https://www.youtube.com/watch?v=8dMFJpEgNLQ>

84 (a) Media promotion and subscription request. Such as: news (110), videos (25), subscribe  
85 (70), visit (36), please (13).

86 (b) URL links and social network account. Such as: <http://www.facebook.com/arirangtv> twitter  
87 (39), facebook (29).

88 (c) Media name. Such as: vice (60), ‘arirang (34). In fact, the most “vice”s appear in the  
89 descriptions are not acting as “vice president”, but “VICE News” instead.

90 This list will help us a lot in future work dealing with other video descriptions and transcripts  
91 since there are many similar patterns in videos from other topics.

### 92 3.1.2. Blacklist and Whitelist Words for Video Descriptions in Chinese

93 Sentences containing advertisement and promotion are treated as “blacklist sentences”. The  
94 typical “blacklist sentence” in Chinese is like this one<sup>5</sup>:

螃蟹科技微信公众号：螃蟹科技 (*pangxiekeji*) 螃蟹科技 QQ 群 419859745

如果对我们的栏目有什么建议或者对智能数码有什么需要了解的，在公众号中回复你想了解的，我们来帮你解答。

95 **Translation:** “Crab Technologies” Wechat Official Account: Crab Technologies (*pangxiekeji*) “Crab Technologies”  
QQ Group 419859745. If you have any suggestions for our column or need to know about smart digital, reply  
96 what you want to know in the Wechat official account, we will answer.

97 And the typical “whitelist sentence” in Chinese is like this one<sup>6</sup>:

柯洁将在下月迎战谷歌旗下的著名人工智能围棋软件 *AlphaGo*。

98 **Translation:** *Ke Jie will be battling with Google’s famous AI Go software AlphaGo next month.*

99 As described in Section ??, the cultural differences on descriptions make it more difficult  
100 for analysis in Chinese. In this particular topic on AlphaGo, as we can expect, besides  
101 the sentences most relevant to AlphaGo and Ke Jie (blacklist sentences) and the promotion  
102 and advertising sentences (whitelist sentences), there are also other sentences that don’t  
103 belong to any of these two lists, which are referred as “irrelevant sentences”. Although time-  
104 consuming, it is quite interesting to read through all of these Chinese descriptions. Some  
105 typical irrelevant sentences are shown below.

106 The following one<sup>7</sup> is collected from a “technology news weekly digest” video. AlphaGo  
107 only serves as a small part in this video. So most contents in this video are irrelevant with  
108 AlphaGo and Ke Jie.

---

<sup>5</sup><https://www.bilibili.com/video/av4116312>

<sup>6</sup>[https://www.iqiyi.com/v\\_19rrbttzwb.html](https://www.iqiyi.com/v_19rrbttzwb.html)

<sup>7</sup><http://v.qq.com/page/u/m/2/u0305zm41m2.html>

三星 Note6/7 工程图曝光联想 Moto Z 真机图泄露柯洁 AlphaGo 即将开战  
10 万块军工级手机发布全球首款带夜视仪的手机发布

**Translation:** *Exposure of Samsung Note6/7 Engineering Drawings, Lenovo Moto Z Real Machine Map Leakage, Ke Jie and AlphaGo are about to battle, 100,000 military grade mobile phones released, first mobile phone with night vision in the world*

The following one<sup>8</sup> is collected from a funny video imagining AlphaGo playing League of Legends game. These kinds of videos are not from the news, but there are several such kinds of videos on these Chinese websites.

153. 如果 AlphaGo 来玩英雄联盟

**Translation:** *If AlphaGo plays LOL*

The following one<sup>9</sup> is collected from an industry introduction sentence. It used “AlphaGo” to express that they are using the modern techniques and they are among the first tier.

英飞凌德累斯顿智能工厂，工业 4.0 的 “AlphaGo”

**Translation:** *Infineon Dresden Intelligent Factory, the “AlphaGo” of Industrial 4.0*

The following one<sup>10</sup> is a bit special. This is a self-edited video with no informative content, and there are several similar videos like this on Chinese video websites. The uploaders of these videos want to express their fondness for somebody or something, so they made these videos using the existing video footage. In this video, the content is mainly collected and edited from news video clips, so most scenes are relevant to the AlphaGo topic. Also, there are many keywords on this topic in the description. Therefore, it will be easily recognized as “related news” if using blacklists and whitelists only.

这个视频的构思想了一年多（是的没写错）从去年小李人机的时候开始想，直到今年才在小十一的古力……  
啊不是，鼓励之下开始动手

一个 AI 爱上了人类，最终他们在一起了故事 \# 严肃

第一次做剧情向，剧情比较凌乱，希望能看懂

送给小十一！希望喜欢！！注 1: AlphaGo 来自于 *Ex Machina-Domhnall Gleeson*

注 2: 主 CP 为 AlphaGo/柯洁，副 CP 为木谷实/吴清源，古力/李世石

注 3: 2017 年 6 月 2 日更新微调版本。具体剧情见回复

**Translation:** *This video has been conceived for more than a year (yes, correctly written) since Lee Sedol's battle last year, and it was not until this year that Gu Li was in eleventh ranking.*

*This is a story. An AI falls in love with a human being and eventually they get together \# seriously*

*This is the first time that I make a story video, and the plot is messy. I hope you can understand it.*

*It's a present for the Eleventh! Hope you like it! Note 1: AlphaGo comes from Ex Machina-Domhnall Gleeson*

*Note 2: The main couple is AlphaGo / Ke Jie, secondary couples are Minoru Kitani / Wu Qingyuan, Gu Li / Lee Sedol.*

*Note 3: Updated fine-tuned version on June 2, 2017. See the reply for the specific plot.*

<sup>8</sup>[https://www.iqiyi.com/w\\_19rub12smp.html](https://www.iqiyi.com/w_19rub12smp.html)

<sup>9</sup><https://v.qq.com/x/page/i0188drze8u.html>

<sup>10</sup><https://www.bilibili.com/video/av10975529/>

126 In conclusion, Chinese descriptions are much more complicated, so it is challenging to carry  
 127 out a two-class classification for Chinese descriptions.  
 128 After eliminating these irrelevant sentences, the top 60 words in the blacklist are shown in  
 129 Table 2, the top 60 words in the whitelist are shown in Table 3.

Word	Translation	Frequency	Word	Translation	Frequency
，		87	科技	science and technology	10
的	's	64	更多	more	9
：		61	！		9
。		25	我	I	9
碧蓝	Azur	23	玩家	player	9
在	at	21	加入	enter	9
航线	Lane	20	了	have done ...	9
中途岛	Midway	16	qq		8
群	group	15	《		8
游戏	game	14	》		8
如果	if	14	com		8
微信	Wechat	14	...		8
集	episode	13	您	you	8
主	main	13	可以	can	8
alphago		13	喂	Hello	8
公众	public	12	id		8
号	account	12	服务器	server	8
交流	communicate	12	服	server	8
欢迎	welcome	12	644132397		8
都	all	11	请	please	7
allen		11	粉丝	fans	7
关注	follow	11	也	also	7
更	more	11	详细	detailed	7
up		11	攻略	strategy	7
有	have	10	尽	use all	7
围棋	Go	10	wiki		7
是	is	10	你	you	6
视频	video	10	对	to	6
、		10	和	and	6
大家	everyone	10	【		6

Table 2: Top 60 Words in Chinese Blacklist on AlphaGo

130 From the lists shown, we can see that the whitelist for Chinese is much more reliable than  
 131 blacklist: For whitelist, there are about 10 words that appear more than 100 times, most  
 132 of which are highly relevant to the topic. However, the top blacklist that has real meaning  
 133 is “碧蓝 (Azur)”, which only has a frequency of 23. This phenomenon has shown that the  
 134 blacklist in Chinese is much messy than whitelist, thus much less reliable.

<sup>11</sup>The word “dog” has the same pronunciation as the word “Go” in Chinese, so “Alpha Go” will sometimes be referred as “Alpha Dog” in Chinese news.



Word	Translation	Frequency	Word	Translation	Frequency
alphago		418	战胜	defeat	29
,		366	将	will do	29
的	's	302	对弈	play chess with	29
柯洁	Ke Jie	177	狗	dog <sup>11</sup>	28
。		170	0		25
围棋	Go	131	棋手	chess player	25
大战	battle	117	deepmind		24
人机	human and computer	113	上	up	23
了	have done ...	106	4		23
“		89	比赛	game, competition	23
李世石	Lee Sedol	85	中	middle	23
”		84	谷歌	Google	22
人类	human beings	83	vs		22
在	at	82	第	-th	22
是	be	70	master		21
:		63	阿尔法	Alpha	21
人工智能	artificial intelligence	59	《		21
月	month	57	》		21
5		51	1		21
日	date	51	被	be done	20
3		44	用	use	20
与	and	43	不	no	20
战	battle	42	人	human	20
和	and	38	手	hand	19
中国	China	37	乌镇	Wuzhen (a place in China)	19
?		36	进行	be in progress	19
对	to	35	团队	team	19
,		31	我们	we	19
ai		30	你	you	18
!		30	马云	Jack Ma	18

Table 3: Top 60 Words in Chinese Whitelist on AlphaGo

## 135 3.2. Video Descriptions on Florida Shooting

136 After a discussion, we found that the AlphaGo event is a little general, which means several  
137 irrelevant events appeared under the search result of the “AlphaGo” keyword. In order to  
138 sort out better “blacklist” and “whitelist”, we turn our eyes to another event, the Florida  
139 Shooting tragedy.

### 140 3.2.1. Blacklist Words for Video Descriptions in English

141 The top 60 words in the blacklist after eliminating punctuation are shown in Table 4.  
142 From the original data, we can observe that most of the videos come from different sources  
143 comparing to AlphaGo videos. However, they share many common blacklist words. For  
144 example, the English Blacklist on Florida Shooting shown in Table 4 and the English Blacklist  
145 on AlphaGo shown in Table 1 share 5 words in top 10, 8 words in top 20, 22 words in top  
146 40. This means different media also use similar words in advertising and promotion. This  
147 observation makes our future work much easier since we can use the blacklists above to  
148 discover most of the targets in new topics.

Word	Frequency	Word	Frequency
news	1074	full	83
on	657	is	83
the	619	fox	78
and	511	local	77
cbs	501	episodes	74
to	400	google	74
here	362	all	74
of	324	cbc	72
nbc	234	our	69
subscribe	196	it	68
evening	182	as	67
you	164	broadcast	66
with	156	access	63
a	132	devices	61
morning	126	day	59
watch	121	stories	59
twitter	118	business	59
your	113	original	58
today	110	coverage	56
facebook	109	guardian	54
instagram	108	entertainment	53
	106	new	52
channel	105	video	52
for	105	digital	52
in	105	source	50
this	104	mobile	48
more	98	shows	47
live	97	breaking	46
latest	95	apps	45
from	91	across	45

Table 4: Top 60 Words in English Blacklist on Florida Shooting

## 149 4. Experiments

150 In this section, I will focus on my experiments leading to my final classifier and the results  
151 of these experiments.

### 152 4.1. Data Preprocessing

153 My previous report and Section 2.1 described several properties of promotional sentences in  
154 our corpora. Before using BERT model to get word embedding vectors, Several preprocessing  
155 steps are necessary. With the help of BERT vocabulary, I modified and added some of the  
156 steps from last semester’s approach, which improves the quality of our corpora.

157 (1) Cleanup: For every word occurs in sentence, I used regular expression matching to  
158 separate words by spaces and several particular punctuations (!?,.:’”());). In order to get  
159 rid of the disturbances of multiple dots in urls, I also added another rule to detect and  
160 cleanup these urls.

161 (2) Tokenize: Not every word in the corpora is in the dictionary. Therefore, in order to  
162 carry out the word embedding successfully, we need to tokenize these words. There are  
163 several cases that we need to consider about:

- 164 • The word is in the dictionary: This is the simplest case. We can directly use the  
165 original word as the token.
- 166 • The word is a combination of two other words in the dictionary: First several  
167 models of BERT carefully considered this case. They used a greedy algorithm to  
168 split up these kind of words. This will make our model much more complicated  
169 without performance improvement. Therefore, I did not apply this strategy in my  
170 experiments.
- 171 • The word is not in the dictionary: we use “[UNK]” token for every word not in the  
172 dictionary.

173 (3) (Optional) Remove stopwords and re-tokenize: Stopwords are those words that can be  
174 ignored in search engine. These words are necessary components in sentences, but they  
175 don’t contribute to the overall topics or styles. We can remove or replace these words  
176 with a special token to ignore the effect of these words. In my following experiments, I  
177 replace these stopwords with a special token ‘#’.

178 As shown in Table 5, after preprocessing, we are able to get the number of words and unique  
179 words for each topic. For AlphaGo and Florida shooting topics, we are also able to get the  
180 number of promotional words and unique promotional words.

Topics	# Word Instances	#Promotional Word Instances
AlphaGo	23934	5534
Florida shooting	19826	6486
Lunar Rover	10280	?
Thailand Cave Rescue	21406	?

Table 5: Word Count

## 4.2. Word Embedding

In the previous report, I used word2vec model to embed every word to a 300-dimension vector space. Also, in order to plot the results on a two-dimensional figure, I also use the t-SNE method described in [Maaten and Hinton \(2008\)](#) to visualizes high-dimensional data by giving each data point a location in a two-dimensional map.

As described in 2.2 and 2.3, we can take advantage of two different word embedding methods developed by Google. Here are several main differences between them:

- (1) Word2vec is a simpler model, and we can directly download the “word to vector” dictionary provided by Google and have a quick reference to get the vector. BERT is a much more complicated model, and we need to calculate the vectors again for every word in new corpora.
- (2) Word2vec is simply based on word. Regardless of its context, every instance of the same word share exactly the same word vector. BERT is based on word and its context. The deep neural network in its model will calculate the word vector based on a vocabulary list as well as the “masked” context. This is beneficial for words with more than one meaning (e.g. bank).

Besides these reasons, BERT is more suitable for this projects in the following ways:

- (1) Since our corpora only contain 200 video descriptions for each topic, it’s obvious that we have very limited training examples. If using word2vec, the size of the training set will be limited to around 2,000, since the same word is treated as only one single sample. However, if using BERT, each word instance is a sample, which enlarges the size of the training set to more than 10,000. This can avoid over-fitting problem if we apply machine learning methods to it.
- (2) We are training a binary classifier for each word, and there are many words that appear in both “blackword” list and “whiteword” list if we take every different word as a single sample. However, if we use BERT, every word instance is a single sample,

In recent years, BERT has been proved to be one of the best model in NLP field. It has

208 already been applied to many NLP tasks like reading comprehension, cloze, etc. Although  
 209 fine-tuning the model requires huge amount of computing resources and time, the original  
 210 model itself is sufficient for our project. In the following sections, I will be using their  
 211 latest model “BERT-Large, Uncased (Whole Word Masking)” published in May, 2019, which  
 212 produces an 1024-dimension vector for each word instance.

213 Here is the comparison of word embeddings using word2vec (shown in Figure 1) and BERT  
 214 (shown in Figure 2). Red dots represent “blackwords” and blue dots represent “non-blackwords”.

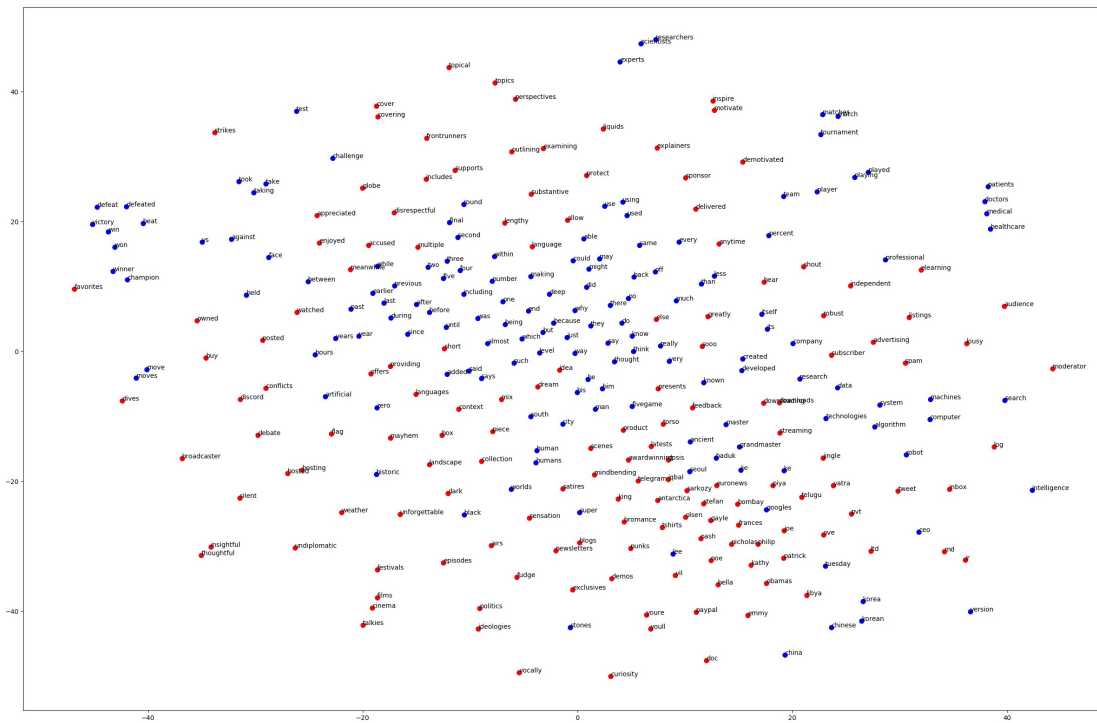


Figure 1: Word2vec embedding

215

216 From the figures shown, we can make several observations:

- 217 (1) Using word embeddings generated by word2vec, it’s difficult to observe a clear boundary  
 218 based on the t-SNE figure. Since the overall boundary of Figure 1 is not so clear, I do  
 219 these steps to help generate the “boundaries” shown in the Figure 3: I first use t-SNE  
 220 as before to plot the words on a 2-d vector space. After that, I use SVM (support vector  
 221 machine) to do the “classification” step. This step does not mean to “train a classifier”  
 222 for future words. It just serves as a method to find the boundary between the two classes.
- 223 (2) Using word embeddings generated by BERT, the boundary becomes much clearer. We  
 224 can also observe many clusters. In fact, every cluster usually represent the same word

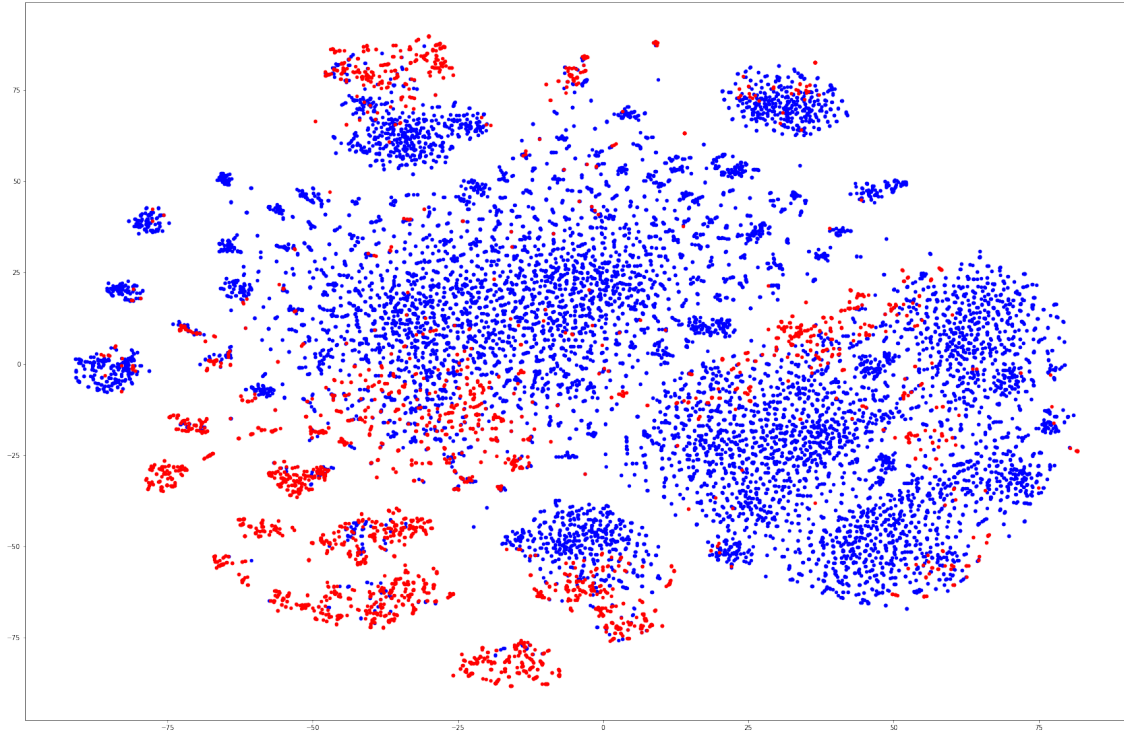


Figure 2: BERT embedding

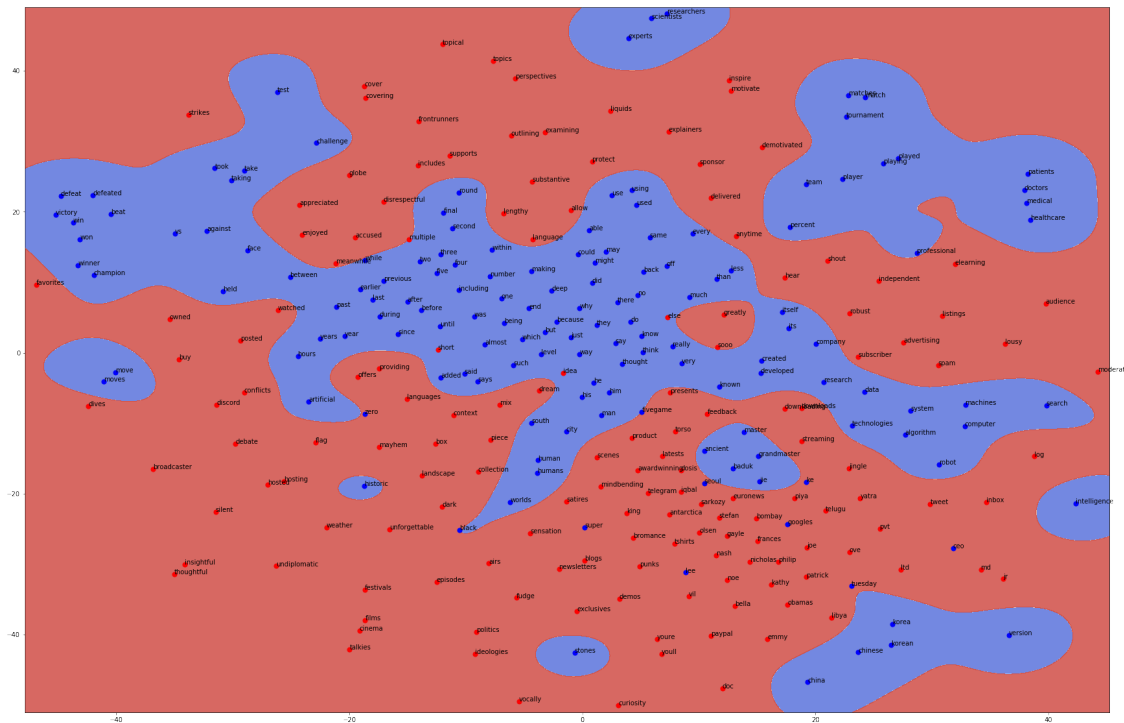


Figure 3: Word2vec embedding with boundaries generated by SVM

225 appearing in different places. If the two instances of the same word have similar usages  
226 or context, they tend to have similar word embedding.

### 227 4.3. Classifier Training

228 The BERT embedding shown in Figure 2 suggests that it's possible to train a classifier based  
229 on the word vectors. In this section, I will introduce the procedures of fine-tuning the model  
230 and the results I get from it. Since the word2vec embeddings are not able to provide us  
231 with good results, I will mainly focus on BERT embeddings in the following experiments  
232 and results.

#### 233 4.3.1. Feasibility

234 In order to quantify the results and conclusions shown above, I used SVM with default C  
235 value and  $\gamma=0.001$  to train the very first model based on 70% of AlphaGo descriptions.  
236 The training error rate is 2.87% and 8.20%, which are both really good for this task. This  
237 result suggests that it's feasible to apply SVM model to train and fine-tune our classifier.

#### 238 4.3.2. Grid Search and Fine-Tuning

239 In order to get the best performance, I used grid search to fine-tune the gamma and C  
240 values. In order to carry out the grid search, I trained 100 different models based on 10  
241 gamma values and 10 C values, and then plot the results in a 3-dimensional figure.

242 The true positive rates based on different gamma-C pairs are shown in Figures 4 and 5.

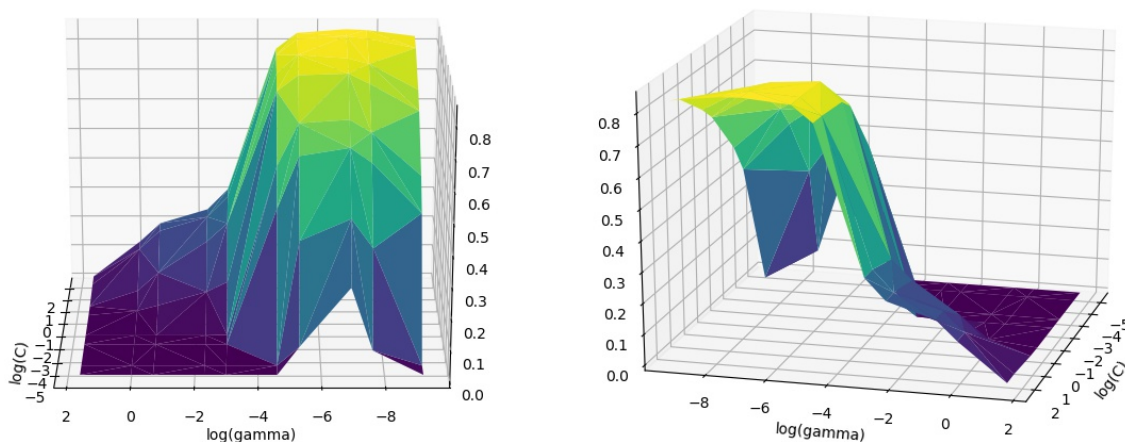


Figure 4: The true positive rates based on different gamma-C pairs on validation set (30% AlphaGo data)

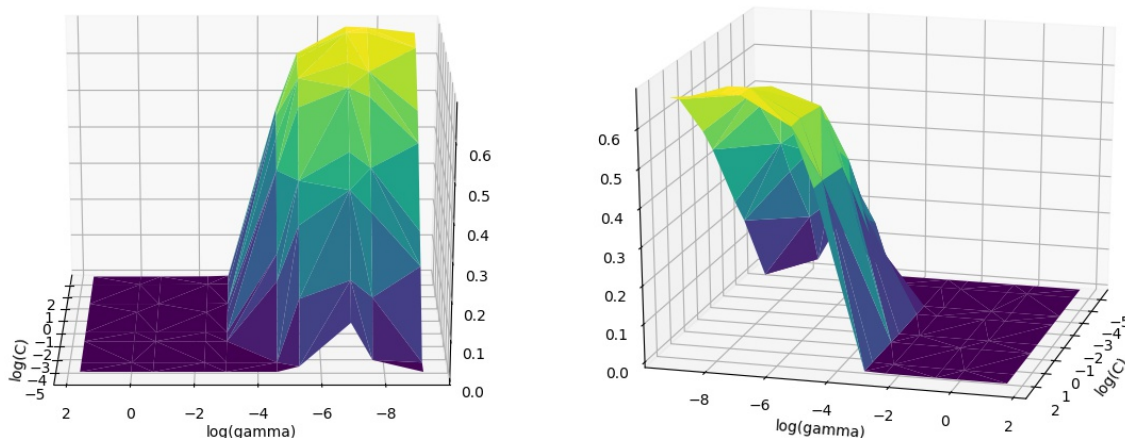


Figure 5: The true positive rates based on different gamma-C pairs on test set (Florida data)

243 From the two graphs, we can observe that with a lower gamma and higher C, we tend to get  
 244 better performance. The basic reasons are<sup>12</sup>:

245 • Intuitively, the gamma parameter defines how far the influence of a single training  
 246 example reaches, with low values meaning “far” and high values meaning “close”. The  
 247 gamma parameters can be seen as the inverse of the radius of influence of samples  
 248 selected by the model as support vectors.

249 Therefore, in our model, if using a lower gamma value, we are having higher  $\sigma$  in the  
 250 original SVM model, leading to more support vectors. This is beneficial to our model  
 251 with more than 15,000 training examples in only 1024 dimensions.

252 • The C parameter trades off correct classification of training examples against maxi-  
 253 mization of the decision function’ s margin. For larger values of C, a smaller margin  
 254 will be accepted if the decision function is better at classifying all training points cor-  
 255 rectly. A lower C will encourage a larger margin, therefore a simpler decision function,  
 256 at the cost of training accuracy. In other words “C” behaves as a regularization pa-  
 257 rameter in the SVM.

258 Therefore, in our model, since we have much more training examples than dimensions,  
 259 we’d like to use a higher C to avoid “underfitting” problem.

260 Finally, I decided to use C=10 and gamma=0.0005 in my final model. This is based on  
 261 the results above as well as reasonable analysis. In the following sections, I will use these  
 262 parameters to train the final models.

<sup>12</sup>[https://scikit-learn.org/stable/auto\\_examples/svm/plot\\_rbf\\_parameters.html](https://scikit-learn.org/stable/auto_examples/svm/plot_rbf_parameters.html)



263 **4.3.3. Results on AlphaGo**

264 For the AlphaGo dataset, I used 70% of sentences as training set, 30% of sentences as  
265 validation set, all Florida sentences as testing set. The results are shown in Table 6.

	Train	Validation	Test
Error Rate	1.10%	7.23%	25.10%
True Positive Rate	96.91%	84.61%	68.31%
True Positive Count	3542	1578	13544
False Positive Count	72	229	323
False Negative Count	113	287	6282
True Negative Count	13059	5040	6163

Table 6: Results based on 70% of AlphaGo description sentences as training set

266 From this table, we can also observe that the performance on testing set is much worse than  
267 training set and validation set. This is because we are training on AlphaGo dataset but  
268 testing on Florida dataset. In order to generalize our model to different topics, I trained on  
269 both AlphaGo and Florida dataset in the next section.

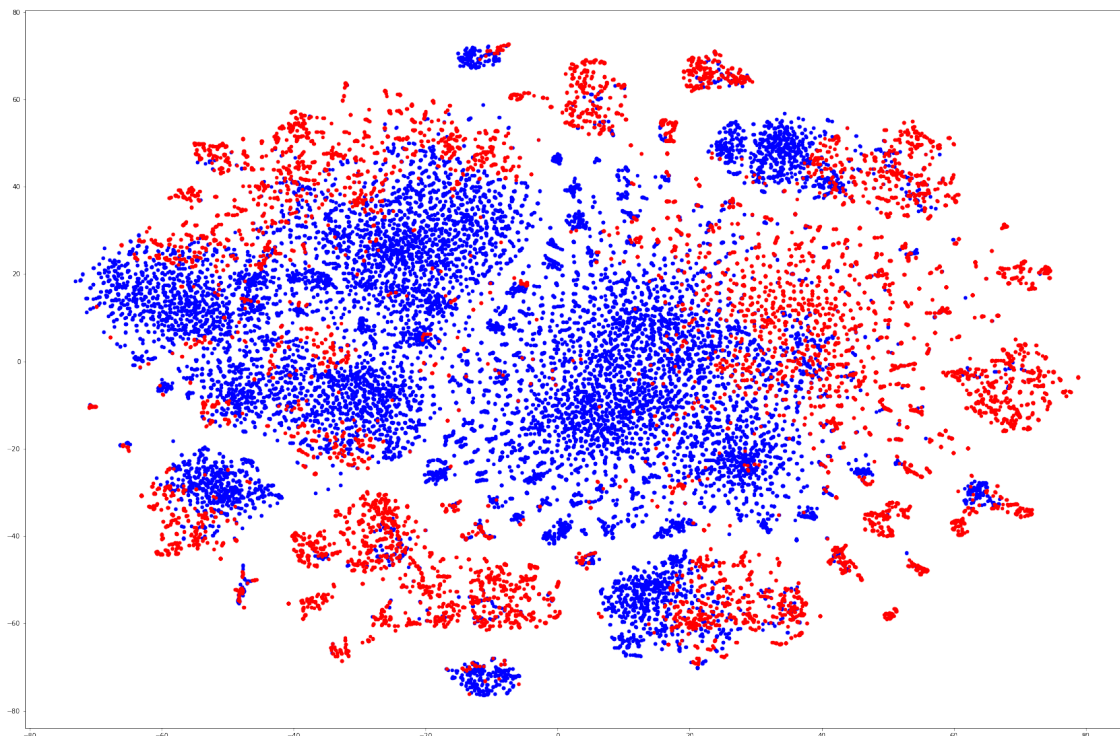


Figure 6: BERT embedding

#### 270 4.3.4. Results on AlphaGo and Florida Shooting

271 Since we are mixing all AlphaGo and Florida description sentences together, we need to  
272 make sure that it’s still possible to train an SVM model based on the dataset. Figure 6  
273 shows the new BERT embedding after t-SNE operation.

274 Similar as previous experiments, I used 70% of sentences as training set, 30% of sentences  
275 as validation set. The prediction results are shown in Table 7.

	Train	Validation
Error Rate	1.96%	3.71%
True Positive Rate	97.96%	97.38%
True Positive Count	14037	10727
False Positive Count	405	252
False Negative Count	293	289
True Negative Count	20901	3328

Table 7: Results based on 70% of (AlphaGo+Florida) description sentences as training set

276 Comparing to the previous model, it’s obvious that both accuracy and true positive rate  
277 improved a lot. And it shows that it’s possible to use this model on other topics.

#### 278 4.4. Final Results

279 After showing the possibility to train a model with high accuracy, I started use the model  
280 trained in the previous section to predict the other two topics, Thailand cave rescue and  
281 lunar rover.

282 After some observation, I found that “Thailand” topic contains a much higher ratio of  
283 blackwords than “lunar rover” topic. Therefore, I put the first 200 lines of the prediction in  
284 Appendix 2. We can observe that most of “blackwords” are predicted out, which means our  
285 model have a reasonable accuracy and sensitivity for promotional words detection.

## 5. Conclusion and Future Work

From the work described above, we can make several conclusions.

- (a) We can manually collect “blacklist” and “whitelist” sentences and words from video descriptions. Due to the cultural differences, people treat descriptions differently on English and Chinese platforms. Also, video descriptions collected from Chinese video websites has more variety.
- (b) As the most recent word embedding model, BERT has a much higher performance comparing to the traditional word2vec embedding. This allows us to train a model with more training samples based on the meaning in the contents.
- (c) Using t-SNE, we are able to visualize the word vectors and observe the feasibility for training a classifier. Using fine-tuned SVM model, I finally trained a classifier with relatively high accuracy and true positive rate. This means we managed to use the hand-labeled data we have and BERT vectors to train a classifier for promotional words detection.

These are possible future steps for this research.

- (a) Use BERT to train a new model for video descriptions in Chinese based on hand-labeled Chinese video descriptions on AlphaGo topic.
- (b) Consider how to generalize this model to video transcript, which may have a different features and distributions comparing to video descriptions.

## References

- 305
- 306 D. E. Rumelhart, G. E. Hinton, R. J. Williams, et al., Learning representations by back-  
307 propagating errors, *Cognitive modeling* 5 (1988) 1.
- 308 T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in  
309 vector space, arXiv preprint arXiv:1301.3781 (2013a).
- 310 T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations  
311 of words and phrases and their compositionality, in: *Advances in neural information*  
312 *processing systems*, pp. 3111–3119.
- 313 R. Řehůřek, P. Sojka, Software Framework for Topic Modelling with Large Corpora, in:  
314 *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, ELRA,  
315 Valletta, Malta, 2010, pp. 45–50. <http://is.muni.cz/publication/884893/en>.
- 316 J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional  
317 transformers for language understanding, arXiv preprint arXiv:1810.04805 (2018).
- 318 B. E. Boser, I. M. Guyon, V. N. Vapnik, A training algorithm for optimal margin classifiers,  
319 in: *Proceedings of the fifth annual workshop on Computational learning theory*, ACM, pp.  
320 144–152.
- 321 C. Cortes, V. Vapnik, Support-vector networks, *Machine learning* 20 (1995) 273–297.
- 322 L. v. d. Maaten, G. Hinton, Visualizing data using t-sne, *Journal of machine learning*  
323 *research* 9 (2008) 2579–2605.

## 324 Appendix 1: Modification to the original BERT code

### 325 5.1. extract\_features.py

```
326
327 @@ -31,6 +31,8 @@ flags = tf.flags
328
329 FLAGS = flags.FLAGS
330
331 +flags.DEFINE_string("token_dir", None, "")
332 +
333 flags.DEFINE_string("input_file", None, "")
334
335 flags.DEFINE_string("output_file", None, "")
336 @@ -60,6 +62,10 @@ flags.DEFINE_bool(
337     "Whether to lower case the input text. Should be True for uncased "
338     "models and False for cased models.")
339
340 +flags.DEFINE_bool(
341 + "do_whole_word_mask", False,
342 + "Whether to use whole word masking rather than per-WordPiece masking.")
343 +
344 flags.DEFINE_integer("batch_size", 32, "Batch size for predictions.")
345
346 flags.DEFINE_bool("use_tpu", False, "Whether to use TPU or GPU/CPU.")
347 @@ -211,6 +217,8 @@ def convert_examples_to_features(examples, seq_length, tokenizer):
348     """Loads a data file into a list of `InputBatch`s."""
349
350     features = []
351 +
352 + f = open(FLAGS.token_dir, 'w')
353     for (ex_index, example) in enumerate(examples):
354         tokens_a = tokenizer.tokenize(example.text_a)
355
356 @@ -279,15 +287,11 @@ def convert_examples_to_features(examples, seq_length, tokenizer):
357     assert len(input_mask) == seq_length
358     assert len(input_type_ids) == seq_length
359
360 - if ex_index < 5:
361 - tf.logging.info("*** Example ***")
362 + if ex_index < 1000 and ex_index % 2 == 0:
363     tf.logging.info("unique_id: %s" % (example.unique_id))
364 - tf.logging.info("tokens: %s" % " ".join(
365 - [tokenization.printable_text(x) for x in tokens]))
366 - tf.logging.info("input_ids: %s" % " ".join([str(x) for x in input_ids]))
367 - tf.logging.info("input_mask: %s" % " ".join([str(x) for x in input_mask]))
```

```

368 - tf.logging.info(
369 - "input_type_ids: %s" % " ".join([str(x) for x in input_type_ids]))
370 + f.write(" ".join([tokenization.printable_text(x) for x in tokens]) + "\n")
371 + f.write("\n")
372 + tf.logging.info("*****")
373
374     features.append(
375         InputFeatures(
376

```

## 377 5.2. tokenization.py

```

378
379 @@ -172,7 +172,7 @@ class FullTokenizer(object):
380     for token in self.basic_tokenizer.tokenize(text):
381         for sub_token in self.wordpiece_tokenizer.tokenize(token):
382             split_tokens.append(sub_token)
383 -
384 + print("\t split tokens: " + ', '.join(split_tokens))
385     return split_tokens
386
387     def convert_tokens_to_ids(self, tokens):
388 @@ -327,35 +327,14 @@ class WordpieceTokenizer(object):
389
390     output_tokens = []
391     for token in whitespace_tokenize(text):
392 - chars = list(token)
393 - if len(chars) > self.max_input_chars_per_word:
394 - output_tokens.append(self.unk_token)
395 - continue
396
397 - is_bad = False
398 - start = 0
399 - sub_tokens = []
400 - while start < len(chars):
401 - end = len(chars)
402 - cur_substr = None
403 - while start < end:
404 - substr = "".join(chars[start:end])
405 - if start > 0:
406 - substr = "##" + substr
407 - if substr in self.vocab:
408 - cur_substr = substr
409 - break
410 - end -= 1
411 - if cur_substr is None:

```

```

412 - is_bad = True
413 - break
414 - sub_tokens.append(cur_substr)
415 - start = end
416 -
417 - if is_bad:
418 - output_tokens.append(self.unk_token)
419 + #
420 + # Zikun's new code
421 + #
422 + if token in self.vocab:
423 + output_tokens.extend([token])
424     else:
425 - output_tokens.extend(sub_tokens)
426 + output_tokens.append(self.unk_token)
427     return output_tokens
428
429
430 @@ -378,7 +357,7 @@ def _is_control(char):
431     if char == "\t" or char == "\n" or char == "\r":
432         return False
433     cat = unicodedata.category(char)
434 - if cat.startswith("C"):
435 + if cat in ("Cc", "Cf"):
436         return True
437     return False
438
439 @@ -390,8 +369,10 @@ def _is_punctuation(char):
440     # Characters such as "^", "$", and "`" are not in the Unicode
441     # Punctuation class but we treat them as punctuation anyways, for
442     # consistency.
443 - if ((cp >= 33 and cp <= 47) or (cp >= 58 and cp <= 64) or
444 - (cp >= 91 and cp <= 96) or (cp >= 123 and cp <= 126)):
445 + if cp == 91 or cp == 93:
446 + return False
447 + if (cp >= 33 and cp <= 47) or (cp >= 58 and cp <= 64) or (cp == 92) or \
448 + (cp >= 94 and cp <= 96) or (cp >= 123 and cp <= 126):
449     return True
450     cat = unicodedata.category(char)
451     if cat.startswith("P"):
452

```

453 **Appendix 2: First 100 Result Segments on Thailand Rescue Topic**  
454 **(Bold words are predicted “blackwords”)**

455 hi friends , today would like share rescue operation plan cave incident ( thailand ) .  
456 idea slightly different tesla ceo ' idea . think much possible rescue 12 boys coach .  
457 idea insert **pipe** system cave rescue children help rescue pods ( like pods used chile mine rescue operation )  
458 .  
459 please watch video . also **welcome new** ideas suggestions helping take part rescue operation .  
460 **thank !**  
461 **please share . http : google . http : youtube .**  
462 leader thailand ' rescue mission save 12 boys soccer coach flooded cave says , ” limited amount time . ”  
463 spoke former thai navy seal , , died inside cave complex lack oxygen .  
464 ben tracy reports chiang rai .  
465 ” **cbs morning** ” channel : **http : ” cbs morning ” : http**  
466 **: latest installment** ” note self , ” ” **cbs morning** , ” : **http : ” cbs morning ” : http**  
467 **: ” cbs morning ” facebook : http : . ” cbs morning ” twitter : http : ” cbs morning ” :**  
468 **http : latest news best original reporting cbs news delivered . : http : news go ! download cbs**  
469 **news mobile apps : http : new episodes shows**  
470 **love across devices next day , stream local news live , watch full seasons cbs fan favorites anytime**  
471 **, anywhere cbs access . try free ! http : king , ”**  
472 **cbs morning** ” offers thoughtful , substantive source news information daily audience 3 million  
473 viewers . emmy broadcast presents mix daily news , coverage developing stories national global  
474 significance , interviews  
475 **leading figures politics , business entertainment . check local listings ” cbs morning ” broadcast**  
476 **times .**  
477 one boys rescued caves thailand said ” shocked ” found .  
478 12 boys football coach making first public appearance following ordeal caves .  
479 please **http :**  
480 thailand erupted celebration 12 youth football players coach trapped flooded cave northern chiang rai  
481 province two weeks rescued , following astonishing mission world .  
482 final four school boys coach , trapped darkness cave complex 18 days , today carried operation came end .  
483 residents chiang rai , city closest caves , took streets celebrate ,  
484 drivers car horns pedestrians dancing outside hospital wild boar fc players recovering .  
485 original article : **http : . co . original video : http : . co . daily**  
486 mail **facebook : http : mail : http : mail snap : https : . daily mail twitter : http : mail : http :**  
487 **co . mail :**  
488 **https : google . free daily mail mobile app : http : co .**  
489 last member rescue team leave thai cave , australian doctor richard harris , lost father .  
490 **time http : get closer world entertainment celebrity news time gives access insight people**  
491 **make watch , read share . https : youtube . ? list . money helps learn spend invest money .**  
492 **find advice guidance count negotiate ,**  
493 **save everything . https : youtube . ? list . find latest developments science technology access**  
494 **brings ideas people changing world . https : youtube .**  
495 **? list . let time show everything need know drones , autonomous cars , smart devices latest**  
496 **inventions shaping industries way : youtube . ? list . stay**



497 **date breaking news around world trusted reporting , insight : youtube . ? list . connect : http**  
 498 **: : https : : https : facebook .**  
 499 **: https : google . : https : . ? : http : : time . brings insight , access authority news . news**  
 500 **publication nearly century experience , coverage shapes understand world . daily news , inter-**  
 501 **views , science , technology , politics , health , entertainment , business updates ,**  
 502 **well exclusive videos person year , time 100 created acclaimed writers , producers editors .**  
 503 **father** australian doctor helped rescue trapped thai soccer team dies : youtube .  
 504 channel ' teresa tang looks challenges divers face rescuing remaining 9 people trapped cave complex . meet  
 505 wild heroes helped bring safety : **https : us : https :**  
 506 **. com https : facebook . https : . https : twitter . https :**  
 507 team divers resumed daring rescue mission free group boys flooded cave chiang rai , thailand july 9 . timeline  
 508 happened far .  
 509 12 thai boys coach found alive inside cave complex chiang rai province , nine days went missing . ' bird '  
 510 view cave system . read full coverage  
 511 search rescue **operation : https : . follow us : https : . com https : facebook . https : . https :**  
 512 **twitter . https :**  
 513 spokesperson us army team helping thailand rescue 12 boys assistant football coach flooded cave talk chal-  
 514 lenges met operation .  
 515 new details emerge 12 boys soccer coach survived two weeks flooded cave thailand . cbs **2** ' anna werner  
 516 reports .  
 517 thai official heading cave rescue says next phase operation start hours . joy benedict reports .  
 518 former elite diver navy died bringing oxygen flooded cave network youth football team coach still trapped .  
 519 **: http : . top stories today ? click watch :**  
 520 **https : youtube . ? list . : watched news channel ! http : youtube . available 13 languages :**  
 521 **https : youtube . english :**  
 522 **website : http : . : https : facebook . : http : : http : : http :**  
 523 cave rescue operation trapped thai soccer team hit wall appears easy solution getting 12 boys coach safely .  
 524 **welcome national , flagship nightly newscast cbc national watch** videos  
 525 **: https : youtube . voice opinion connect us online : national updates facebook : https :**  
 526 **facebook . national updates twitter : https : national cbc television ' flagship**  
 527 **news program . airing six days week , show delivers news , feature documentaries analysis**  
 528 **canada ' leading journalists .**  
 529 indian played key role saving children trapped thailand ' cave . prasad , designing engineer brothers limited  
 530 company district , played key role saving lives 12 children coach  
 531 cave thailand . 12 . **hindi channel latest updates movies related videos . tube : https : youtube .**  
 532 **follow us twitter : https : us facebook : https :**  
 533 **facebook . circle google plus : https : google . download app : https : google .**  
 534 patrick decker says ' " still long way go " ongoing effort rescue members youth soccer team coach cave  
 535 thailand .  
 536 members wild boar soccer team described moment found , **news** conference chiang rai , thailand . **nbc**  
 537 **news : http : watch nbc video : http : news leading**  
 538 **source global news information . find clips nbc nightly news , meet press , original digital videos**  
 539 **. channel news stories , technology , politics , health , entertainment**  
 540 **, science , business , exclusive nbc investigations . connect nbc news online ! visit . com : http**  
 541 **: nbc news facebook : http : nbc news twitter : http**

542 : **nbc news : http : nbc news : http : nbc news : http** : soccer boys speak dramatic **rescue flooded**  
543 **cave nbc news**  
544 cave system takes experienced diver five hours boys back full strength suffering exhaustion starvation found  
545 . original article : http : . **co** . original **video** : http  
546 : . **co** . **daily mail facebook : http : mail : http : mail snap : https : . daily mail twitter : http :**  
547 **mail : http**  
548 : **co** . mail : **https : google . free daily mail mobile app : http : co** .  
549 thailand ' art bridge chiang rai , created giant painting commemorate rescue operation 12 boys coach stuck  
550 cave . report nadia .  
551 australian federal police divers , dr harris , whole australian team thank . term hero gets used lot . examples  
552 personnel worked cave rescue thailand .  
553 pleasure morning . **australia** looks forward welcoming home safely **later week** .  
554 listen tale race time rescue 12 boys soccer coach trapped cave two weeks . smarter . faster . colorful get  
555 **story http** : . even  
556 ? ! usa **youtube channel : http : usa today facebook : https : facebook . usa today twitter :**  
557 **https :**  
558 thai navigate difficult terrain underground find alternative ways extract 12 boys football coach trapped cave  
559 complex since jun 23 . latest updates rescue efforts : **https : . follow us**  
560 : **https : . com https : facebook . https : . https : twitter . https :**  
561 captain dan brown discusses thailand cave rescue  
562 divers working free 12 boys coach trapped cave northern thailand must navigate dark , flooded tunnels six  
563 hours reach . takes another five hours return . details extraordinary operation underway non  
564 emerged thursday , pushed ahead multiple plans free boys trapped underground almost two weeks . 11 .  
565 **hindi channel latest updates movies related videos . tube : https : youtube . follow**  
566 **us twitter : https : us facebook : https : facebook . circle google plus : https : google . download**  
567 **app : https : google .**  
568 authorities thailand say twelve boys coach trapped inside cave ready dangerous dive flooded narrow passage  
569 . ' ve trapped almost two weeks . comments follow death former thai navy seal part  
570 rescue team . died lack oxygen . john joe regan reports . : **http : : http : : http : : http : : http**  
571 : **website : http : world**  
572 two british divers , john jason honoured part rescuing 12 boys football coach cave complex thailand ' chiang  
573 rai province . meet wild heroes helped bring safety :  
574 **https : full coverage rescue : https : us : https : . com https : facebook . https : . https : twitter**  
575 **. https :**  
576 pressure mounting thai authorities bring forward rescue plan 12 boys coach trapped deep inside flooded cave  
577 northern thailand , death former navy diver drop oxygen levels underground . read : https :  
578 **rt . http : live http : rt ! http : youtube . like us facebook http : facebook . us telegram https**  
579 **: us https : us twitter http**  
580 : **us http : us http : google . ( russia today ) global news network broadcasting moscow washington**  
581 **studios . rt first news channel break 1 billion youtube views .**  
582 john , bristol , one number foreign expert divers drafted rescue 12 boys football coach trapped thai cave nine  
583 days . tells bbc points west moment fellow  
584 divers first discovered children alive . mr , member south mid wales cave rescue team , said knew found due  
585 smell cave . please **http** :  
586 rescue ! ' wireless equipment used rescue 12 boys trapped caves thailand . **find : https : . follow us :**

587 **facebook : https :**  
588 **facebook . : https : : https : .**  
589 dive teams thailand rescued four boys flooded jungle cave monday . **watch** full episode ' **world news**  
590 **tonight ' : https : full episodes world news tonight : http : go .**  
591 northern thailand raced time ominous monsoon season youth soccer team trapped partially flooded cave  
592 advance heavy rains forecast later week .  
593 two weeks trapped cave , 12 members wild football team coach rescued . risky operation led thai navy seals  
594 , international team managed get boys complicated often narrow  
595 exit route . , coordinating operation , said ' thailand ' mission impossible ' **guardian news youtube http**  
596 **: cave rescue : boys get ? https**  
597 **: . support guardian https : guardian https : . guardian youtube network : guardian www .**  
598 **youtube . jones talks http : football http : sport http : culture http**  
599 **: science tech http :**  
600 elephant calf tumbled well eastern thailand reunited **mother** rescued group villagers . **us youtube : https**  
601 **: app apple store ( ios ) : https**  
602 **: apple . download app google play ( android ) : https : google . follow us : facebook : https :**  
603 **facebook . : https :**  
604 **. ? : https : : https : . : http : . : http :**  
605 thirteen members thai football team stranded cave rescued ending operation save . june 23rd , boys coach  
606 went explore cave complex , heavy rain flooded tunnels , leaving trapped  
607 . john joe regan reports international rescue mission . : **http : : http : : http : : http : : http : website**  
608 **: http :**  
609 world  
610 may secret passage cave thailand youth soccer team ( bottom right cave ) trapped 11 days , emerged today  
611 . boys aged 11 16 told ( left top  
612 right ) heard dogs barking , children playing despite 800 metres underground . led officials think may another  
613 way ' chimney hole ' surface . Chiang Rai provincial governor  
614 , overseeing rescue , said 30 teams searching . believes must one boys able breathe long . left : trapped  
615 coach , 25