

SECOND LANGUAGE LEARNING FROM NEWS WEBSITES

First Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Second Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Abstract

Learning a second language is difficult and requires constant revision and immersion. Fortunately, many of us take the time to update ourselves through reading news on a daily basis. In this paper, we merge both of these goals into a Web browser extension that allows a reader to learn and master vocabulary items. We conducted a user survey to evaluate our system against user requirements collected through an earlier survey. Since we find a word's context to be useful in learning a vocabulary, we further adopt word sense disambiguation (WSD) technique to show the best translation for each word in the context. Our proposed WSD method, leveraging the extension of standard machine translation system, significantly better baseline methods in both coverage and accuracy. We also elaborate on the issues of determining appropriate distractors for multiple-choice word mastery quizzes within the scope of the project.

1 Introduction

Learning a new language from language learning websites is time consuming. It is, therefore, necessary to make second language learning attractive and efficient. Further, since habitual learning is effective, we seek to interleave language learning with a popular daily activity. Reading news online is once such activity. Further recent increase in the popularity of portable devices has made online news reading popular than ever before (Kathryn Zickuhr and Brenner, 2012). We leverage on this culture to provide users of news

websites with an opportunity to learn a second language.

We propose a system to enable online news readers to efficiently learn a new language while they are reading news on news websites. We propose a Chrome extension which would run on the client (Chrome web browser) when readers visit news websites on a preconfigured list.

Learning a new vocabulary is the most time consuming and boring part of language learning¹. Perhaps, this justifies the poor adoption of current second language learning systems. We, therefore, focus on enabling language learners build their vocabulary efficiently while providing them with an enjoyable user experience.

2 Methods

There are many existing language learning software, which, fall into two categories, learning by lessons and learning vocabularies. In the first category, lessons are purposefully designed to help users easily learn a foreign language. Duolingo² is a popular websites in this category. For the second category users are guided to recite lists of words, or provided with a translation for their input word in the foreign language. Google Translate³ stands out in this category. The service is available as desktop / mobile / web software including a chrome extension. We mainly compare our system with the aforementioned two softwares. Table 1 summarises important differences between our system and all these existing tools. Each difference serves as a motivation for developing our extension.

“Duolingo is a free language-learning and crowd-

¹<https://neltachoutari.wordpress.com/tag/vocabulary/>

²<https://www.duolingo.com/>

³<https://translate.google.com/>

sourced text translation platform”⁴. Most people start to use Duolingo when they know a little or nothing about the new language. They starting from some basic lessons and improve step by step. However, our target audience is a mix novice and intermediate level learners of the foreign language. We can not only help beginners learn a new language but also help them continue their learning by allowing them to practice their foreign language. There are also a lot articles with their translations in Duolingo, but all the articles and their translations are manually added by Duolingo or users from Duolingo. Therefore, parallel articles in Duolingo are old and limited. However, our chrome extension is always working even for those up to the minute news and our user can just practice their foreign language in their daily readings.

Google Translate: “Highlight or right-click on a section of text and click on Translate icon next to it to translate it to your language”⁵. Google Translate is a chrome extension that displays only the translation when user select a section, which can be a word, a phrase, a sentence or even a whole page. Our chrome extension will translate a single word only, and display the translation, following with the pronunciations and example sentences to help user understand and remember this word. Compared with our extension, Google Translate is more like an extension to help user understand the content of the page. Furthermore, our extension will display the most appropriate translation as it will refer to the context of the word.

2.1 Software Design

Based on the user study, we divided the extension into three components: translating, learning and testing. After user opened a news website, some words in the main content will be replaced by their translation from user’s preferred foreign language, and this is what our translating component is doing. If user want to know more about the replaced word, he can simply move his mouse over the translation and a window will pop over to help user learn this word, and this is learning component. If user have encountered some word for a

⁴<http://en.wikipedia.org/wiki/Duolingo>

⁵http://en.wikipedia.org/wiki/Google_Chrome_Extensions

Table 1: Summary of the differences

	Duolingo	Google Translate	Chrome Extension
Lessons	Yes	No	No
User’s foreign language level	Low	Low-High	Low-High
Time consuming	Yes	No	No
Resource	Limited	Infinite	Infinite
Customizable	Yes	No	Yes
Link to External Dictionary	No	No	Yes

few times, we will generate some quiz for him and this is testing component.

2.1.1 Translating

Yahoo Sports’ Charles Robinson highlighted **two** of the **more controversial** calls, **or** non-calls, **that** went **against** the Badgers. The **first** was Justise Winslow possibly stepping out of bounds before dishing the ball to Jahlil Okafor for a Duke bucket. The other was a close out-of-bounds decision in which Winslow looked to have touched the ball last:

Figure 1: Screen shot of Translating Component

After the original web page, our chrome extension will fetch the content of the news and pass them to the server paragraph by paragraph. After receiving the content, server will compare every word in the paragraph with the words in our vocabulary. If there are some matches, which simply means there are some words that need to be replaced. As every English word might have a few Chinese meanings, our server must select the most appropriate translation among all the meanings. The way that we are trying to solve this problem so far is to compare all the Chinese meanings with the translation of the whole sentence from Bing Translate. If any of the Chinese meanings is the substring of the translation of the sentence, our server will choose that meaning (This is not a proper and accurate way to solve this problem, but it is much better than randomly choose one Chinese meanings. Also, this would be my main research problem that I need to solve next semester).

Then, server will pass a JSON string that contains all the words that need to be replaced, their Chinese meanings as well as their pronunciations back to front end. Then, front end will replace the content of the news paragraph by paragraph, in which some words have been replaced. Figure 1 is the screen shot of this component.

2.1.2 Learning



Figure 2: Screen shot of popover with highlighted English word



Figure 3: Screen shot of popover with highlighted Chinese word

By moving mouse on the Chinese word for one second, a window with its English meaning and pronunciation will pop over. Figure 2 is the screen shot of the pop over without its example sentence. If user want to know how to use this word, he can just click the button next to the pronunciation to get the example sentences of this word. After user click the button to get example sentence, our extension will send a request to server and wait for server's response. There is another way of doing this, which is simply get the example sentences together with the words in the Translating component. However, the example sentences contains much more characters comparing with the pronunciation. We want to maximize the loading

speed and minimize the data transferred between front end and server, so we decided to split the pop over content into two request. Figure 3 is the screen shot of the pop over with its example sentences.

2.1.3 Testing

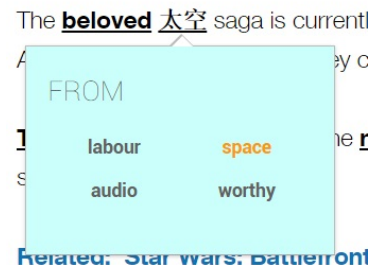


Figure 4: Screenshot of English test popover



Figure 5: Screen shot of Chinese test popover

When user has encountered the same word for a few times, our system will generate a quiz about this word for him. If user move mouse over this replaced word, a window with a quiz will pop over. After user select one option, this window will tell user the correct answer and sent whether the answer is correct to server. Figure 4 is the screen shot of our testing popover in English. Figure 5 is the screen shot of our testing popover in Chinese.

3 Word Sense Disambiguation System

As we all know, one word may have multiple translations in another language, and our extension is expected to select the most appropriate one based on the context. We call such translation selection as cross-lingual word sense disambiguation (WSD).

WSD is an open problem in natural language processing and ontology, aiming at identifying the

proper sense of a word (i.e. meaning) in a context, when the word has multiple meanings. However, the system that I have implemented is different from the traditional WSD. Normally, WSD is to identify its sense in its original language, but my system is to identify its sense in another language.

It has been proved by a lot of studies that context is the key to learning a new language. This further suggests that a good word sense disambiguation system is necessary for our Chrome Extension. It can help user remember vocabulary efficiently and is a key feature that differentiate our software from other language learning tools.

In this following, I describe four approaches that I have tried to accomplish WSD system, which is also my main progress in the second semester. The four approaches are:

- Frequency based: always selecting the most frequent translation (the baseline),
- Part-of-Speech Tag based: selecting the translation based on the Part-of-Speech Tag of the English word
- Translation based: Selecting the translation based on the result from existing Machine Translation systems
- Category based: Selecting the translation based on the category of the news article

3.1 Baseline

The simplest way to select a translation from the candidates is by random. However, the correctness of this method is very low, probably less than 20%, and is not a good baseline for other methods to compete with. Another simple idea is to always select the most commonly used translation. Luckily, when I crawled the dictionary, Google Translate does provide usage frequency of each Chinese Translation. This turns out to be a much better result, and thus serves as a fair baseline method.

3.2 Part-of-Speech Tagger

As we all know, many English words have more than one Part-of-Speech (POS) tags and their Chinese translations in different POS may differ a lot.

For example, the word “book” has two POS tags, noun and verb. If it is used as a noun, mostly it means a handwritten or printed work of fiction or nonfiction, which should be translated as “书”, and mostly means to reserve if used as a verb, which should be translated as “预定”. Therefore, if I can get the POS tag of the English word, it might help me identify its sense or the Chinese translation.

Table 2: Sample input/output of POSTagger

Input	Output	Tag Description
They are treating me like family now	They—PRP are—VBP treating—VBG me—PRP like—IN family—NN now—RB	IN — preposition or conjunction

Among all possible on-line sources, Stanford Log-linear Part-of-Speech Tagger (Toutanova et al., 2003) is the most stable and well performed Part-of-Speech Tagger, which is developed by The Stanford Natural Language Processing Group. Column one and column two of Table 2 is one sample input and output. However, the tag from Part-of-Speech Tagger is not the one we normally used. I need one more step to match the tag to its Part-of-Speech. Therefore, I followed the Part-of-Speech guidelines from the University of Pennsylvania (Penn) Treebank Tag-set. The last column of Table 2 is one example of the matching.

Algorithm 1 is the approach of using POSTagger. From Algorithm 1, firstly, if the word “like” need to be translated, the algorithm will fetch all the Chinese translations as well as their Part-of-Speech tag from our dictionary. Secondly, the algorithm will send the original English sentence to Part-of-Speech Tagger, which is a Java package and has been wrapped into a server. After the client has got the output from the server, it will fetch the corresponding tag and match it to Part-of-Speech tag based on the guidelines mentioned above. Lastly, it will select the translations based on the POS.

In most cases, this algorithm is only a filter. There

Algorithm 1 Part-of-Speech Tagger

Require: POS Tagger, Dictionary $\langle \text{English, Chinese, POS} \rangle$, Input English
Pairs $\leftarrow \langle \text{word, POS} \rangle \leftarrow \text{POSTagger} \leftarrow \text{English}$
if *EnglishWord* \in *Dictionary.English* **then**
 TranslationList \leftarrow
 Dictionary.Chinese + *Dictionary.POS*
 for *word* \in *Pairs.word* **do**
 if *word* = *EnglishWord* **then**
 POSResult \leftarrow *Pairs.POS*
 Break
 end if
 end for
 for *POS* \in *TranslationList.POS* **do**
 if *POS* = *POSResult* **then**
 FinalResult \leftarrow
 TranslationList.Chinese
 Break
 end if
 end for
end if
return *FinalResult*

might be a few Part-of-Speech tags matched from the output and one POS tag might have a few corresponding Chinese translations. In this case, I will choose the translation with the highest frequency of use, which is actually a combination of POSTagger and baseline, so that this system is testable independently.

3.3 Machine Translation

Since our target is to select the most appropriate translation based on the context, using existing Machine Translation (MT) systems is also a good approach, as all of them will certainly translate words based on the context.

There are mainly two kinds of MT systems. One is off-line Machine Translation systems, which are mostly not available as mostly they are build for internal usage. Luckily, NUS NLP group built one MT system before, and it has been wrapped into a server, so that I can use it as an on-line service. The other one is on-line MT system, which are wrapped as a server and open to public, such as Bing Translator or Google Translate. As

Table 3: Example input/output of POSTagger

Input English	they are treating me like family now
English Word	like
Translation List	verb : 喜欢, 爱, 爱好, 待见, 好, 看上, 喜, 喜爱, 喜好 adjective : 一样, 似, 同, 相似 conjunction : 如同 noun : 类 preposition : 好像, 好比, 好以 adverb : 不啻, 若
Pairs $\langle \text{word, POS} \rangle$	they—PRP are—VBP treating—VBG me—PRP like—IN family—NN now—RB
Final Result	好像

only Bing Translator is free, I decide to try Bing Translator as well.

The first priority of choosing a MT system is its translation quality, if it can give me a result that nearly as good as a result from human translation, then the Chinese word that I generated from the MT system will have a high chance to be correct as well. However, after I tried both MT systems, the performance of the one from NLP group is worse than Bing Translator and the server is very unstable, I decide to use Bing Translator as my Machine Translation system.

3.3.1 Bing

Table 4: Example input/output of Bing Translator

Input	Output
including a 45-caliber pistol a pump shotgun and an ar-15 rifle	包括 45 口径手枪唧筒式猎枪和 ar-15 步枪
they are asking for privacy at this time	在这个时候他们正在寻求隐私
possessing cartridges used exclusively by the military and carrying a firearm without a license	拥有只供军方使用的墨盒和携带火器的许可证

Bing Translator, also called Microsoft Translator, is a on-line Machine Translation system that developed by Microsoft team with a cloud-based API that is conveniently integrated into multiple products, tools, and solutions. Table 4 is the sample input and output of Bing Translator.

Algorithm 2 Bing Translator

Require: Bing Translator, Dictionary $\langle \text{English}, \text{Chinese} \rangle$, Input English
ChineseTranslation \leftarrow BingTranslator \leftarrow English
if *EnglishWord* \subset *Dictionary.English* **then**
 TranslationList \leftarrow *Dictionary.Chinese*
 MaxLength = 0
 for *ChineseWord* \subset *TranslationList.Chinese* **do**
 if (*ChineseWord* \subset *ChineseTranslation*) \cap
 (*MaxLength* $<$ *ChineseWord.Length*) **then**
 FinalResult \leftarrow *ChineseWord*
 end if
 end for
end if
return *FinalResult*

Table 5: Example input/output of Bing Translator

Input English	including a 45-caliber pistol a pump shotgun and an ar-15 rifle
English Word	pump
Translation List	verb:抽, 抽水, 打气, 唧, 唧筒, 套 noun:抽水机, 唧筒
Chinese Translation	包括 45 口径手枪唧筒式猎枪和 ar-15 步枪
Final Result	唧筒

Algorithm 2 is the steps of using Bing Translator. The original English sentence is “including a 45-caliber pistol a pump shotgun and an ar-15 rifle” and “pump” is the word that we want to translate. Firstly, this algorithm will fetch all the Chinese translations from the database. Next, it will

send the original English sentence to Bing Translator using the API provided by Microsoft and get the result that returned from Bing Translator. After that, for each Chinese translation, I will check whether this translation is a substring of the Bing Translator result. If there are a few translations that can match with the Bing Translator result, I will select the longest translation. If there are a few translations with the same length and all of them can match with the Bing Translator result, I will select the translation with the highest frequency of use. In this example, both “唧” and “唧筒” are the substrings of Bing Translator result. As “唧筒” have two characters and “唧” only have one character, this algorithm will take “唧筒” as the final result.

3.3.2 Bing+

Table 6: Another example of Bing Translator

Input English	as a result of the latest premier league broadcast rights deal all of its teams have made it into the worlds top 40 clubs
English Word	top
Translation List	顶部, 顶端, 顶, 颠, 盖, 极, 尖, 尖峰, 面, 上身, 头, 上面的, 最大的, 最高的, 盖
Chinese Translation	由于最新英超联赛转播的权交易所有其团队已经进入世界顶级40名俱乐部
Final Result	顶

Table 6 is another example of Bing Translator. In this example, “top” is the word that need to be translated. If we manually translated the word “top” based on the Bing Translator, I would say the correct translation should be “顶级”. However, the word “顶级” is not covered by our dictionary, so the final result generated by Bing algorithm is “顶” as this is the only word that matches Bing Translator. Although the meaning of the word “顶” is almost the same as that of the word “顶级” and I, personally, would prefer to make “顶” as one of the correct result when I tried to evaluate this example, the word “顶级” is more

accurate in this context. Therefore, is there any way to solve this problem and make the result more accurate? Yes, solving this problem requires our system to generate Chinese words that are not covered by our dictionary and the only way is to use Word Segmenter.

Table 7: Example of Stanford Word Segmenter

Input	Output
问题是你将要运行直赫顿到他说的第一修正案	问题, 是, 你, 将要, 运行, 直赫顿, 到, 他, 说的, 第一, 修正案
博士苏斯博士知道这将是他的最后一本书	博士, 苏, 斯, 博士, 知道, 这, 将, 是, 他, 出版, 的, 最后, 一, 本, 书
进入世界顶级40名俱乐部	进入, 世界, 顶级, 40, 名, 俱乐部

Text segmentation is the process of dividing written text into meaningful units, such as words, sentences, or topics. The term applies both to mental processes used by humans when reading text, and to artificial processes implemented in computers, which are the subject of natural language processing. I use Stanford Word Segmenter to do Chinese text segmentation. Stanford Word Segmenter is an open source Java package that developed by The Stanford Natural Language Processing Group. I wrapped it into a local server and Table 7 contains some example of the input and output of Stanford Word Segmenter.

Algorithm 3 is the approach of using Bing Translator together with Stanford Word Segmenter, and I would like to use Bing+ to represent this algorithm. Step one, two and three has been described in the previous section as it is exactly the same as Bing approach. From Bing approach, this algorithm will generate “顶” as the result. After that, Bing+ approach will send the Chinese sentence returned from Bing Translator to Stanford Word Segmenter. Then, this algorithm will use the segmented word that contains the Bing result as a substring or equals to the Bing result as the final result. In this example, the final result of Bing+ is “顶级” which is the best result that

Algorithm 3 Bing+

Require: Bing Translator, Word Segmenter, Dictionary $\langle \text{English}, \text{Chinese} \rangle$, Input English

```

ChineseTranslation  $\leftarrow$ 
BingTranslator  $\leftarrow$  English
SegmentedChineseTranslation  $\leftarrow$ 
WordSegmenter  $\leftarrow$  ChineseTranslation
if EnglishWord  $\subset$  Dictionary.English
then
    TranslationList  $\leftarrow$  Dictionary.Chinese
    MaxLength = 0
    for ChineseWord  $\subset$ 
    TranslationList.Chinese do
        if (ChineseWord  $\subset$ 
        ChineseTranslation)  $\cap$ 
        (MaxLength  $<$ 
        ChineseWord.Length) then
            FinalResult  $\leftarrow$  ChineseWord
        end if
    end for
    for SegmentedWord  $\subset$ 
    SegmentedChineseTranslation do
        if FinalResult  $\subset$  SegmentedWord
        then
            FinalResult  $\leftarrow$  SegmentedWord
        end if
    end for
end if
return FinalResult

```

can be generated from the result of Bing Translator and also a result that does not covered by our dictionary.

3.3.3 Bing++

Table 9 is another example of Bing+ Translator. In this example, “line” is the word that need to be translated and the correct translation based on Bing Translator should be “线上”. The Bing approach will get “线” as the result. After the Bing approach got the result, the next step should be finding the segmented word that contains “线” as a substring. However, both word “线上” and word “视线” contains word “线” as a substring and, obviously, only the word “线上” is the correct result and “视线” is actually the translation for English word “sight”. As both results generated by Bing+ approach are not covered by our dictionary, the algorithm cannot get the frequency of use informa-

Table 8: Example input/output of Bing+

Input English	as a result of the latest premier league ... have made it into the worlds top 40 clubs
English Word	top
Translation List	顶部, 顶端, 顶, 颠, 盖, 极, 尖, 尖峰, 面, 上身, 头, 上面的, 最大的, 最高的, 盖
Chinese Translation	由于, 最新, 英, 超联赛, 转播, 的, 权 ... 进入, 世界, 顶级, 40, 名, 俱乐部
Final Result	顶级

Table 9: Another example of Bing+ Translator

Input English	but the world no 46 played one of his best matches crucially not crumbling when the finish line was in sight
English Word	line
Translation List	线, 线路, 路线, 系, 行列, 划线于, 衲, 排, 诗句, 纹, 线条, 衬, 衲
Chinese Translation	但, 世界, 没有, 46, 起, 最, 重要, 的, 是, 不, 崩溃, 在, 终点, 线上, 视线, 的, 时候, 他, 最, 好, 的, 比赛, 之一
Final Result	线上 or 视线

tion. As a result, Bing+ approach will randomly choose one word, so there is only 50% to get the correct result. Is there any way to solve this problem and always choose the correct result in this case? Yes, as long as the algorithm could get the word alignment information, it will know exactly how to match the Chinese words with those English words.

Bitext word alignment or simply word alignment is the natural language processing task of identifying translation relationships among the words (or more rarely multiword units) in a bitext, resulting in a bipartite graph between the two sides of the bitext, with an arc between two words

Table 10: Example input/output of Word Alignment

Input 1	dr seuss knew it would be the last book he published
Output 1	博士苏斯博士知道这将是他的最后一本书
Output 1	0:1—0:1 3:7—2:5 9:12—6:7 14:15—8:8 17:21—9:9 23:24—10:10 26:28—14:14 30:33—15:17 35:38—18:19 40:41—11:11 43:51—12:13
Input 2	the problem is youre going to run straight headon into the first amendment he said
Output 2	问题是你将要运行直赫顿到他说的第一修正案
Output 2	4:10—0:1 12:13—2:2 15:19—3:3 21:28—4:5 30:32—6:7 34:41—8:8 43:48—9:10 50:53—11:11 59:63—15:16 65:73—17:19 75:76—12:12 78:81—13:13
Input 3	this is something preschoolers deal with all the time
Output 3	这是学龄前儿童处理所有的时间
Output 3	0:3—0:0 5:16—1:1 18:29—2:6 31:39—7:8 41:43—9:10 45:47—11:11 49:52—12:13

if and only if they are translations of one another. I use Bing Word Alignment API⁶ developed by Microsoft team to get the word alignment from English to Chinese Simplified. Luckily, although this API only support very few sets of language pairs, English to Chinese Simplified is one of the few supported sets. Table 10 has some examples of input and output from Bing Word Alignment. The left column is the original English sentence, the column in the middle is the translated Chinese sentence and the right column is the word alignment information. For word alignment information, the colon separates start and end index, the dash separates the languages, and space separates the words. For example, in the second

⁶<https://msdn.microsoft.com/en-us/library/dn198370.aspx>

column, “0:1—0:1” means the word “dr” should match with word “博士” and “9:12—6:7” means the word “knew” should match with word “知道”.

Algorithm 4 Bing++

Require: Bing Translator, Word Segmenter, Word Alignment, Dictionary <English, Chinese>, Input English

ChineseTranslation \leftarrow BingTranslator \leftarrow English

SegmentedChineseTranslation \leftarrow WordSegmenter \leftarrow *ChineseTranslation*

Pair < *EnglishWord*, *ChineseWord* > \leftarrow WordAlignment \leftarrow English

if *EnglishWord* \subset Dictionary.*English* **then**

TranslationList \leftarrow Dictionary.*Chinese*

MaxLength = 0

for *ChineseWord* \subset *TranslationList.Chinese* **do**

if (*ChineseWord* \subset *ChineseTranslation*) \cap (*MaxLength* < *ChineseWord.Length*) **then**

*FinalResult*₁ \leftarrow *ChineseWord*

end if

end for

for *SegmentedWord* \subset *SegmentedChineseTranslation* **do**

if *FinalResult*₁ \subset *SegmentedWord* **then**

*FinalResult*₁ \leftarrow *SegmentedWord*

end if

end for

for *Word* \subset *Pairs.EnglishWord* **do**

if *EnglishWord* = *Word* **then**

*FinalResult*₂ \leftarrow *Pairs.ChineseWord*

end if

end for

end if

return *FinalResult*₁ \cup *FinalResult*₂

Algorithm 4 is the approach of using Bing+ approach together with the Microsoft Bing Word Alignment. First few steps are exactly the same as Bing+ approach. In this example, “state” is the word that need to be translated. The result from Bing+ approach is “发言人”, which is the

Table 11: Example input/output of Bing++

Input English	state department spokeswoman jen psaki said that the al- lies had a long history of cooperation
English Word	state
Translation List	...陈, 陈说, 称, 称述, 发表, 发言...
Chinese Translation	国家, 部门, 的, 女, 发言人, jenpsaki, 说, 盟国, 有, 很, 长, 的, 合作, 历史
Bing+ Result	发言人
Word Alignment Result	国家

translation of “spokeswoman”, because the Chinese translation “发言” can be translated from both “state” and “spokeswoman”. Then step five will send the original English sentence to Bing Word Alignment. Now, there will be two final results, one from Bing+ approach and the other one from Bing Word Alignment and the algorithm will choose the correct one from these two results. In this example, “state” will match with “国家” and the algorithm will choose “国家” as the final result as well.

One general question about this approach is that, since I can get the official word alignment from Bing Word Alignment approach, is Bing+ approach still useful? Table 12 contains some examples of Bing+ approach and Bing++ approach. From left to right, the four columns are original English sentence, Chinese translation, result from Bing+ approach and result from Bing++ approach. It is very obvious that the result from Bing+ approach is the substring of the result from Bing++ approach, but which one is better? As the purpose of our Word Sense Disambiguation system is to select the most appropriate translation based on the context, but Bing Translator is a bit too smart comparing with our purpose. Bing Translator will generate some Chinese words that cannot be translated from any of the English word but can make this sentence clear and smooth. In this case, our system will choose the short answer instead of the long answer. That’s why in the

Table 12: Some examples of Bing+ and Bing++

Input English	oh the places youll go! rose to the bestseller list shortly after it was re- leased in 1990 and con- tinues to pop up there most every spring as high school and college grads transition to a new phase of life
English Word	spring
Bing+ Result	春天
Word Align- ment Result	年春天
Input English	the darkness of the book is what makes the optimism credible nel said
English Word	book
Bing+ Result	书
Word Align- ment Result	本书

Bing++ approach, I will keep the result both from Bing+ approach and Bing Word Alignment and choose the better one.

3.4 News Category

The word “interest” have two very different translations when it is used as a noun. One translation is “the feeling of a person whose attention, concern, or curiosity is particularly engaged by something”, which should be translated as “兴趣”. The other translation is “a share, right, or title in the ownership of property, in a commercial or financial undertaking, or the like”, which should be translated as “利益”. It is quite obvious that the second sense is mostly used in financial related topics. Therefore, if we can analyze the category of the original article and select the translation with the same category label, it might help disambiguate the word meaning.

Getting the category of the original news article is very simple. Most news websites have a manually assigned a category for each news article and in most cases, the category label is part of the URL. However, assigning a category for Chinese word

is not simple. As we are dealing with news, it is good to obtain such information from Chinese news domain. I crawled 100 Chinese news articles in each category from Baidu News⁷, making around 1000 news articles in total. After I got all the news articles, I send all the news articles to the Stanford Chinese Word Segmenter, and further calculate word document frequency under each category. For example, if word “interest” is found five times in article A and three times in article B, both article A and B are under “finance” category, then I will add two for category “finance” of word “interest” as it will be counted only once even it can be found multi times in one article. I will use “weight” to represent this value and “averageweight” is the average weight of all categories of one word. After that, I will normalize the weight and use Equation 1 and Equation 2 to assign categories for those Chinese words. Basically, the two equations means that, if this word can be found in at least ten different news articles and more than 80% of the articles are under the same category, then I will use this category for this word.

$$averageweight > 1 \quad (1)$$

$$threshold > 8 * averageweight \quad (2)$$

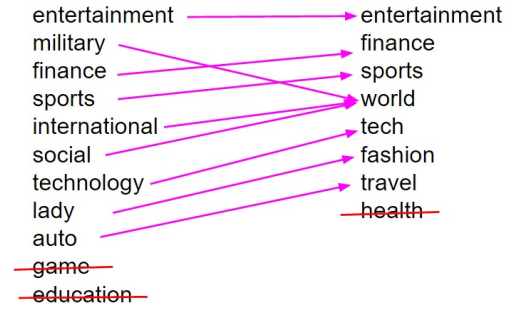


Figure 6: Alignment between categories

However, the categories in Chinese and English news are not the same. I manually aligned the categories and delete some categories when necessary. As shown in Figure 6, on the left side, the 11 categories are from Chinese news website and

⁷<http://news.baidu.com/>

on the right side, the eight categories are from English news website. During the alignment, I not only take the name of category into account, but also consider the semantics of the category.

Algorithm 5 News Category

Require: Dictionary $\langle \text{English, Chinese, category} \rangle$, Input English, News URL
if $\text{EnglishWord} \subset \text{Dictionary.English}$
then
 $\text{TranslationList} \leftarrow \text{Dictionary.Chinese}$
 $\text{Dictionary.category}$
 $\text{EnglishCategory} \leftarrow \text{URL.category}$
 for category \subset
 $\text{TranslationList.category}$ **do**
 if $\text{category} = \text{EnglishCategory}$ **then**
 $\text{FinalResult} \leftarrow \text{TranslationList.Chinese}$
 Break
 end if
 end for
end if
return FinalResult

Table 13: Example input/output of News Category

Input English	duncan told cnns don lemon hes just painting a picture of urban street life with his lyrics
English Word	picture
Translation List	...陈, 陈说, 称, 称述, 发表, 发言...
Category	entertainment
Final Result	影片

Algorithm 5 is the steps of using news category. The English sentence is the original sentence and “picture” is the word that need to be translated. Firstly, the algorithm will fetch all the Chinese translations for word “picture” and split them with comma. In this example, only the word “影片” has a category “entertainment”. Next, the algorithm will fetch the category of the English news article from the URL, which is also “entertainment”. In this case, the algorithm will use “影片” as the translation for word “picture”. If a few

words shares the same category, the algorithm will choose the translation with the highest frequency of use.

4 Distractors Generation Algorithm

The key research topic here is to investigate a way to automatically generate suitable distractors for a certain vocabulary test. The distractors are generated in English form.

4.1 Collecting category-related words

To generate good category-related distractors, it is essential to gather enough words that are more related in a certain category to serve as distractors candidates.

To find good “category-related” words, it is essential to get the words from those already classified news articles. The process involves 3 steps, crawling news contents from popular news website, preprocessing, and classifying word category.

4.1.1 Crawling news content

Several web crawlers are designed to get news content from different popular news websites. The crawler will detect URLs from each news website’s main page as and its sub-category pages. For example, there are sub-categories like “football”, “basketball” under main category “Sports”, and the crawler is able to crawl URLs from “football” page and “basketball” page as well.

After detailed comparison of most news websites, I divided news articles into seven categories, namely “World”, “Technology”, “Sports”, “Entertainment”, “Finance”, “Health” and “Travel”. Most news articles can be classified into one of the seven categories. The web crawler will store all paragraph tags from each websites and store them as one file under one category.

In this experiment, around 1400 news articles, i.e. 200 news articles from each category are crawled and stored locally. Post natural language processing is then applied to this category corpus.

4.1.2 Preprocessing

After storing all the news articles document into each category, the server uses Natural Language

Tool Kit (Edward, 2009) for word tokenizing and POS Tagging. The system will store the POS tag of each word. After elimination of all non-English words and those words that contain special symbols, like “O’Real”, “S\$40”, all words that contains only alphabetic letters are conserved. All stop words are also eliminated as well. They are stored as lower case for the ease of future process.

4.1.3 Classification

In statistical analysis step, the server counts the document frequency of each word in all those stored news articles, i.e. if word “scored” appeared 4 times in one article, it will only be counted as once. By following this approach we can successfully reduce the bias of some words only appear a lot of times in one article while don’t appear often in other article. As we are storing similar number of articles in each category, this approach will provide a fair comparison of each word’s popularity among different categories. After this step we will know the document frequency count of each word in different category. Assume C is the list of category names, and $f(w, C(i))=m$ means word w appeared in category $C(i)$ for m times, then the sum weight of word w as $sw(w)$ is calculated in Equation 3:

$$sw(w) = \sum_{i=1}^n f(w, C(i)) \quad (3)$$

The average weight of word w as $aw(w)$ is calculated in Equation 4::

$$aw(w) = sw(w)/n \quad (4)$$

A word w is classified into category $C(i)$ if it satisfies Equation 5::

$$f(w, C(i)) - aw(w) \geq \delta \quad (5)$$

The confidence factor δ can be a positive integer between 0 and the average number of articles in each category. It means on average, the word w must appear in a specific category $C(i)$ δ times more than it appear in other category before it can be classified into category $C(i)$.

In the example below in Figure 18, frequency counts for word “investment” in each category

are displayed. It is obvious that $sw(\text{“investment”}) = 2 + 1 + 2 + 10 + 3 + 2 + 1 = 21$, thus $aw(\text{“investment”}) = sw(\text{“investment”})/7 = 3$, in this case if we choose a confidence factor $\delta = 3$, word “investment” will be classified into category “Finance”, as $10-3 \geq \delta = 3$. However, if we choose a very big δ , for example $\delta = 8$, then word “investment” will not be classified into any category.

Table 14: Example of classification word into category

Category	Investment
Technology	2
World	1
Sports	2
Finance	10
Entertainment	3
Health	2
Travel	1

It is obvious that a higher confidence factor value will result in less number words get classified, but it will result in getting words that are more accurate. A lower confidence factor value will result in more number of words get classified, but less accurate in each category. In this experiment, after several round of tests and analysis, we chose a confidence factor value of 10, which is capable of producing enough number of classified words while maintaining the accuracy.

4.2 Generating distractors

My selection strategy in choosing distractors takes following parameters:

- News website URL
- News sentence
- Word to test
- User’s knowledge level of the word

4.2.1 Detect news category

After getting the news URL, our system needs to determine the category of the news. Based on the analysis from most popular news URLs, there is a

set of common identifiers that can identify the category of the news article. For example, technology news URL often contains “/tech”, “/science”, and if we find these strings in news URL, we will classify this news URL into “Technology” category. The algorithm will go through all category identifier in the list, and will return the category name the moment it finds a match. The current list of category provides reasonable accuracy for the purpose of detecting news category.

4.2.2 Detect Part-Of-Speech Tag

Given the target word and the target sentence, it is easy to run the NLTK POS tagger to get the correct POS tag of this word. This step is essential to help select distractors with similar forms, i.e. if the target word is adjective, it will be appropriate to choose three other adjectives, not verbs, as distractors.

4.2.3 Semantic Distance

Before we go to explain the next step, it is essential to introduce the semantic distance calculator we used in the server implementation.

The perspective of semantic relatedness or its inverse, semantic distance, is a concept that indicates the likeness of two words. It is more general than the concept of similarity as stated in WordNet’s synset relation. Similar entities in WordNet are classified into same synset based on their similarity. However, dissimilar entries may also have a close semantic connection by lexical relationships such as meronymy (car-wheel) and antonymy (hot-cold), or just by any kind of functional relationship or frequent association (pencil-paper, penguin-Antarctica) (Alexander, 2001). Semantic distance calculator aims to calculate the semantic relatedness score between two words.

There are many approaches to calculate semantic relatedness score. In this application, we are using Lin Distance (Lin, 1998) to calculate the semantic distance between two concepts. The detail of Lin Distance methodology is explained as follows.

Lin attempted to define a measure of semantic similarity that would be both universal and theoretically justified. There are three intuitions that he used as a basis:

- The similarity between arbitrary objects A

and B is related to their commonality; the more commonality they share, the more similar they are;

- The similarity between A and B is related to the differences between them; the more differences they have, the less similar they are.
- The maximum similarity between A and B is reached when A and B are identical, no matter how much commonality they share.

Based on the intuition above, Lin proposed his approach in measuring similarity between two concepts $c1$, $c2$ in Equation 6:

$$sim(c1, c2) = \frac{2 * \log_p(lso(c1, c2))}{\log_p(c1) + \log_p(c2)} \quad (6)$$

where $p(c)$ denotes the probability of encountering concept c , and $lso(c1, c2)$ denotes the lowest common subsumer, which is the lowest node in WordNet hierarchy that is a hypernym of $c1$ and $c2$.

The distance calculator will return a score from 0 to 1, as can be easily seen from the formula above. If the score is closer to 1, it means the two words are closer in semantic sense. This distance calculator will play an important role in the following algorithm.

4.2.4 Distractors Selection Algorithm

Based on the input parameters, at this stage the server has already got the current category of the news article and the correct POS tag of the target word to test. The server is going to generate distractors based on user’s knowledge level of the target word to test.

Knowledge level is 1: This indicates that the user has just learnt this word. The algorithm will randomly select three words from current category’s word list. The reason for using randomization is to avoid the situation that similar distractors are generated every time.

Knowledge level is 2: This indicates that the user has known this word for some times. The algorithm will randomly select two words from the current category’s word list as two distractors. Then the algorithm will randomly select word from the current category’s word list and calculated the semantic distance between the selected

word and the target word, once the score is above certain threshold, the selected word will be chose as the third distractor. The selection of threshold value will have a direct effect on the speed of distractors generation process. As a very high threshold value will result in more rounds of calculation in semantic distance calculator, and it will take a long time before the distractors are returned to the front end. After several rounds of analysis of each category's words and the results returned from semantic distance calculator, the threshold value of 0.1 is selected.

Knowledge level is 3: This indicates that the user has a good understanding of the word already; the algorithm will choose distractors solely based on results returned from semantic distance calculator. Similar to the approach when knowledge level is 2, the algorithm will randomly select word from current category's word list and calculate the semantic distance between the selected word and the target word. If the score is above certain threshold, the selected word is chosen as one of the distractors. The process is continued until the server can find three distractors.

5 Evaluation

This project has two main research parts, Distractors Generation algorithm and WSD system. To evaluate the distractors selection strategy, the best way to evaluate is to listen to users' voice. The WSD system is a standard research problem and can be evaluated with ground-truth, reporting its performance by coverage and accuracy.

5.1 Distractors Generation Algorithm

To evaluate the distractors selection strategy as described in this report, we chose the knowledge-based approach used by many other language learning systems, which is to utilize the WordNet data and selection distractors based on synonyms of synonyms. WordGap system uses this approach to generate vocabulary test for its android application.

In our implementation of the baseline algorithm, we will choose the most frequent used word w1 from the target word's synonym set, and select the most frequent used word w2 from word w1's synonym set. The selection process is continued until we can find 3 distractors to form a vocabu-

lary test. However, if the number of valid result we can get is less than 3, we will choose the word that shares the same antonym with the target word.

5.1.1 Designing Survey

To compare the two approaches in generating distractors, we designed several survey sets to ask users to compare the plausibility of distractors. We randomly selected 50 sentences from recent news articles and choose one noun or adjective inside the sentence as the target word to test. In the survey, participants are required to answer each question and rank the plausibility of all distractors from 1 to 7. The correct answer will be ranked as 1, and the least plausible distractor will be ranked as 7. A screenshot of one sample question is shown in Figure 7.

There are two evaluations to be done as follows:

41. The ranks of the opposition, civil society and labor movement have been decimated in the last 50 years through imprisonment without trial and _____ prosecution, and nearly every newspaper, TV channel and radio station is owned and run by the state *

	1	2	3	4	5	6	7
criminal	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
turn	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
outlaw	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
bend	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
terrorist	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
arrestment	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
young	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 7: A sample survey question

1. Compare Baseline with Knowledge Level 1 Algorithm

. Compare Baseline with Knowledge Level 3 Algorithm For each comparison, three distractors are generated from the baseline algorithm; three distractors are generated from the stated algorithm in this report. With the first comparison we will be able to see if the category information will help in selecting more suitable distractors. By comparing the results from the both evaluation, we will be able to see if semantic distance and category information will help improve the suitability of distractors.

5.1.2 Results

The evaluation contains 100 questions and is separated into 4 surveys, with each survey containing 25 questions. Each participant is free to choose one or more than one surveys. The purpose is to

reduce the workload in each survey to get better responses. The surveys are sent to Year 1 students from School of Computing, National University of Singapore. There are 15 valid responses with each participant ranking each distractor with a different weight from 1 to 7. Half of the participants are native English speakers.

Each participant’s rank will be the weight of the particular distractor in that question, i.e. if the user rank one distractor as rank “5”, the weight of this distractor in this user’s response will be 5. For each distractor of each question, the ranks of all users’ responses are summed. As the more plausible the distractor is, the higher rank it will have, thus if the sum is higher, the approach is not as plausible as the other from user’s point of view.

Table 15: Comparison 1 Baseline vs. Knowledge level 1 Algorithm

	Number of winning questions	Average score
Baseline	27	3.84
Level 1 Algorithm	23	4.10

Table 16: Comparison 2 Baseline vs. Knowledge level 3 Algorithm

	Number of winning questions	Average score
Baseline	21	4.16
Level 3 Algorithm	29	3.49

Table 15 and Table 16 showed the detailed result of each comparison. If for any question, the sum of weight from all participants for one approach is bigger than the other, then this approach is considered to have won this question. The “average score” is the average sum of weight from each approach for all questions. The lower the average score is, the better performance this approach has gained.

From Figure 7 we can see that in the first comparison, the baseline algorithm actually outscored

the knowledge level 1 generation algorithm by 4 questions, with a sum of weight lower than 0.26. From Table 15 we can see that in the second comparison, the knowledge level 3 generation algorithm surpassed the baseline algorithm by 8 questions, with the average weight of 3.49 vs 4.16.

5.1.3 Analysis

In knowledge level 1 generation algorithm, there is no semantic distance calculation involved. If the target word to test has no strong category indication, for example, words like “venue”, “week”, it is possible that the knowledge level 1 algorithm will select some distractors that are not as plausible as those coming from the target word’s synonym of synonym.

However, this problem is solved with the help of semantic distance calculator. In the knowledge level 3 generation algorithm, the distractors chosen are both semantic close and also category-related, which produced a relatively better experiment result.

Also in the baseline algorithm, it is possible that it will select words that are very rare in real life (Susanne, 2013), which may also have influence in the result.

5.2 WSD System

Our Word Sense Disambiguate System can be evaluated from two important aspects: coverage (i.e., is able to return a translation) and accuracy (i.e., the translation is proper). To this end, I manually annotate the ground truth. Each approach was evaluated right after I had implemented it, therefore, they was tested against a random but different set of recent news articles from CNN. Though the evaluation datasets are different, it is still fair to compare their results, as the size of all dataset is sufficiently large.

Firstly, we want our algorithm to return at least one result instead of blank. For POSTagger approach, if our dictionary do not cover the Part-of-Speech generated from Stanford POSTagger, the algorithm will return nothing. For News Category approach, as the algorithm will only assign categories for some of the Chinese translations and not all Chinese news categories can match with a English news category, so the algorithm sometimes will return nothing as well. For Bing+ and Bing++

approach, if none of the Chinese translations is the substring of the Bing result, the algorithm will return nothing. For Bing++ approach, if the word alignment information is phrase to phrase matching, for example, it may give a matching between “in order to” and its Chinese translation, the algorithm will return nothing. Alternatively, for all the listed algorithm listed above, they can always return the translation with the highest frequency of use, but in this case, we cannot know whether the result is generated from the algorithm itself or just the baseline. That’s why I choose to return a blank instead of the translation with the highest frequency of use.

Table 17: Coverage for different approaches

	Cover	Coverage
Baseline	707/707	100%
POSTagger	668/707	94.5%
News Category	14/707	2.0%
Bing	555/707	78.5%
Bing+	535/707	75.7%
Bing++	544/707	76.9%

Table 17 contains the coverage for different approaches. As the algorithm will try to translate some word only if it is covered by our dictionary, the coverage for Baseline is always 100%. The coverage for Bing, Bing+, Bing++ and POSTagger are roughly the same and all of them are acceptable. However, the coverage for News Category approach is only 1.9%. One reason is that when I set the threshold for assigning categories for Chinese word, I purposely make it very high to maximize the accuracy. If the accuracy is quite high, which means this approach is quite useful, then I will lower the threshold and find the balance point.

Secondly, we want our algorithm to be as accurate as possible, and the most ideal situation is that all the translation returned from the algorithm is the correct or the most appropriate translation in that context. When I evaluate the accuracy of these few approaches, I use a few news articles from CNN as the input data and manually select the most appropriate translation for all the output data. After that, I will compare the result from the algorithm and the result that I manually generated

and get the accuracy.

Table 18: Accuracy for different approaches

	Correct	Accuracy
Baseline	405/707	57.3%
POSTagger	369/668	55.2%
News Category	1/14	7.1%
Bing	443/555	79.8%
Bing+	433/535	80.9%
Bing++	530/544	97.4%

Figure 18 contains the accuracy of all the approaches. The last column is the accuracy for News Category approach and it is only 30%. As mentioned in above Chapter, since the accuracy is very low, there is no need to lower the threshold and try to allocate more categories for Chinese words. The accuracy for Baseline is 69%, which is already a fairly high accuracy. The accuracy for Bing and POSTagger is around 69% also, which is a bit lower than our expectation. The accuracy for Bing++ is 97% which I think is a very good result and it is already very hard to improve. Therefore, based on my test results, Bing++ is the best approach among these five approaches.

References

- Hirst Alexander, Budanitsky; Graeme. 2001. Semantic distance in wordnet: An experimental, application-oriented evaluation of five measures.
- Klein Edward, Loper; Ewan. 2009. Natural language tool kit.
- Kristen Purcell Mary Madden Kathryn Zickuhr, Lee Rainie and Joanna Brenner. 2012. Younger americans’ reading and library habits. *Pew Internet*.
- Dekang Lin. 1998. An information-theoretic definition of similarity.
- Knoop; Sabrina Wilske Susanne. 2013. Wordgap - automatic generation of gap-filling vocabulary exercises for mobile learning.
- Kristina Toutanova, Dan Klein, Christopher D. Manning, and Yoram Singer. 2003. Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology*

- *Volume 1*, NAACL '03, pages 173–180, Stroudsburg, PA, USA. Association for Computational Linguistics.