

Interactive Second Language Learning from News Websites

Abstract

We propose a web browser extension that allows readers to learn a second language vocabulary while reading news online. Injected tooltips allow readers to look up selected vocabulary and give interactive tests to assess vocabulary mastery.

We discover that two key system components needed improvement, both which stem from the need to model context. These two issues are in practical word sense disambiguation (WSD) to aid translation quality and constructing the interactive tests. We start with Microsoft’s Bing translation API but employ additional dictionary based heuristics that significantly improve translation quality over a baseline in both coverage and accuracy. We also propose techniques for generating appropriate distractors for multiple-choice word mastery tests. Our preliminary user survey confirms the need and viability of such a language learning platform.

1 Introduction

Learning a new language from language learning websites is time consuming. Research shows that regular practice, guessing, memorization (Rubin, 1975) as well as immersion into real scenarios (Naiman, 1978) hastens language learning process. To make second language learning attractive and efficient, we seek to interleave language learning with a popular daily activity: online news reading.

Most existing language learning software are either instruction-driven or user-driven. Duolingo¹ is a popular instruction-driven sys-

tem that teaches through structured lessons. Instruction driven systems demand dedicated learner time on a daily basis and are limited by learning materials as lesson curation is often labor-intensive.

In contrast, many people informally use Google Translate² to learn vocabulary, making it a prominent example of a user-driven system. Translate, however, lacks the rigor of a learning platform as it lacks tests to allow learners to demonstrate mastery. In our work, we merge learning and assessment within the single activity of news reading. Our system also adapts to the learner’s skill during assessment.

We propose a system to enable online news readers to efficiently learn a new language. Our prototype targets Chinese language learning while reading English language news. Learners are provided translations of open-domain words for learning from an English news page. In the same environment – for words that the system deems mastered by the learner – learners are assessed by replacing the original English text in the article with their Chinese translations and asked to translate them back given a choice of possible translations. The system, deployed as a Chrome web browser extension, is triggered when readers visit a preconfigured list of news websites.

A key design property of our language learning extension is only active on certain news websites. This is important as news articles typically are classified with respect to a news category, such as *finance*, *world news*, and *sports*. If we know which category of news the learner is viewing, we can leverage this contextual knowledge to improve the learning experience.

In the development of the system, we discov-

¹<https://www.duolingo.com/>

²<https://translate.google.com/>

ered two key components that can be affected by this context modeling. We report on these developments here. In specific, we propose algorithms: (i) for translating English words to Chinese from news articles, (ii) for generating distractor translations for learner assessment.

2 The SystemA Chrome Extension

We give a running scenario to illustrate the use of our language learning platform, SystemA. When a learner browses to an English webpage on a news website, our extension selectively replaces certain original English words with their Chinese translation (Figure 1, middle). While the meaning of the Chinese word is often apparent in context, the learner can choose to learn more about the replaced word, by mousing over the translation to reveal a definition tooltip (Figure 1, left) to aid mastery of the Chinese word. Once the learner has encountered the replaced word a few times, SystemA will assess the learner’s mastery by generating a multiple choice translation test on the target word (Figure 1, right). Our learning platform thus can be viewed as have three logical components: *translating*, *learning* and *testing*.

Translating. We pass the main content of the webpage from the extension client to our server for candidate selection and translation. As certain words are polysemous, the server must select the most appropriate translation among all possible meanings. Our initial selection method replaces any instance of words stored in our dictionary. For translation, we check the word’s stored meanings against the machine translation of each sentence obtained from the Microsoft Bing Translation API (hereafter, “Bing”). Matches are deemed as correct translations and are pushed back to the Chrome client for rendering.

Learning. Hovering the mouse over the replacement Chinese word causes a tooltip to appear, which gives the translation, pronunciation, simplified and traditional written form, and a `More` link that loads additional contextual example sentences (that were previously translated by the backend) for the learner to study. The `More` link must be clicked for activation, as we find this two-click architecture helps to minimize latency and the loading of unnecessary data. The server

keeps record of the learning tooltip activations, logging the enclosing webpage URL, the target word and the user identity.

Testing. After the learner encounters the same word a pre-defined number $t = \text{BUG}$ times, SystemA generates a MCQ test to assess mastery. When the learner hovers over the replaced word, the test is shown for the learner to select the correct answer. When an option is clicked, the server logs the selection, and the correct answer is revealed by the client extension. Statistics on the user’s test history are also updated.

2.1 News Categories

As our learning platform is active only on certain news websites, we model the news category of both individual words and webpages. Of particular import to SystemA is the association of words to a news category, which is used downstream in both word sense disambiguation (Section 3) and the generation of distractors in the interactive tests (Section 4). Here, our goal is to automatically find highly relevant words to a particular news category – e.g., “what are typical *finance* words?”.

We first obtain a large sample of categorized English news webpages, by creating custom crawlers for specific news websites. We use a seed list of words that are matched against a target webpage’s URL. If any match, the webpage is deemed to be of that category. For example, a webpage that has the seed word “football” in its URL is deemed of category “Sports”. After a survey of a number of news websites, we decided on seven categories: namely, “World”, “Technology”, “Sports”, “Entertainment”, “Finance”, “Health” and “Travel”.

We tokenize and part-of-speech tag the main body text of the categorized articles, discarding punctuation and stopwords. The remaining words are classified to a news category based on document frequency. A word w is classified to a category c if it appears a tunable threshold $\delta = \text{BUG}$ more often than its average category document frequency. Note that a word can be categorized to multiple categories under this scheme.

3 Word Sense Disambiguation System

As we all know, one word often have multiple translations in another language, and our exten-

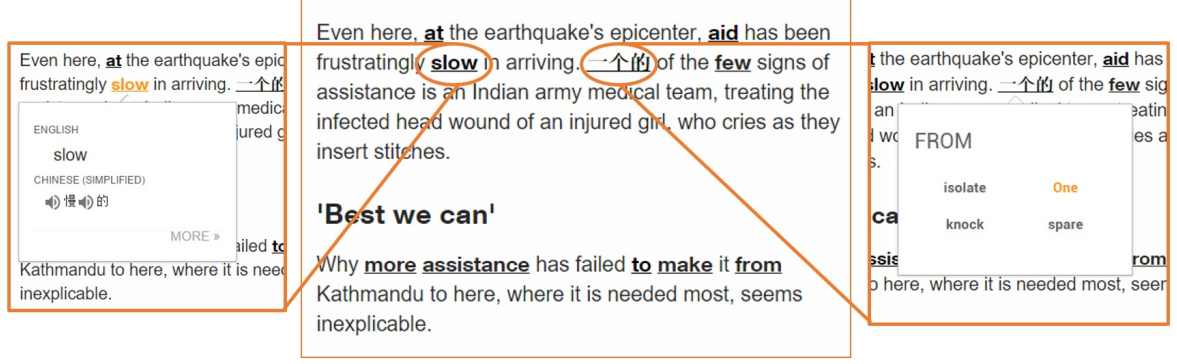


Figure 1: Merged screenshots of our Chrome extension on the CNN English article *Treacherous journey to epicenter of deadly Nepal earthquake*. Underlined components are clickable to yield tooltips of two different forms: (l) a definition for learning, (r) a multiple-choice interactive test.

Table 1: Experimental results.

Chinese Category	English Category	Example Words
Entertainment	Entertainment	饰演
Military	World	服役
Finance	Finance	存款
Sports	Sports	比分
International	World	—
Social	World	—
Technology	Tech	芯片
Lady	Fashion	瘦
Auto	Travel	奔驰
Game	—	—
Education	—	—
—	Health	—

sion is expected to show the most appropriate one based on the context. We call such translation selection as word sense disambiguation (WSD). WSD is an open task in natural language processing, aiming at identifying the proper sense (*i.e.*, meaning) of a word in a context, when the word has multiple meanings (Navigli, 2009). Traditionally, WSD system identifies the proper sense in the same language, while we show the proper sense in the form of another language.

In WSD, context information is the key to disambiguate word sense. We, therefore, make use of different granularity of context, *i.e.*, the category of the news, the word class, and the sentence, to select proper translations from our bilingual dictionary.

3.1 News Category

Topic information have been shown useful in WSD (Boyd-Graber et al., 2007). Take English word “interest” as an example. In finance related articles, “interest” is more likely to be “a share, right, or title in the ownership of property” (“利息” in Chinese), than ‘the feeling of a person whose attention, concern, or curiosity is particularly engaged by something’ (“兴趣”). Therefore, analysing the topic of the original article and selecting the translation with the same topic label might help disambiguate the word sense. We leverage the algorithm described in Section 2.1 to obtain the category for news and candidate Chinese translations.

3.2 Part-of-Speech Tagger

The word class, *i.e.*, the Part-of-Speech (POS) tag is believed to be beneficial for WSD (Wilks and Stevenson, 1998) and Machine Translation (Toutanova et al., 2002; Ueffing and Ney, 2003). For example, the English word “book” has two major classes, verb and noun, meaning “reserve” (“预定” in Chinese) and “printed work” (“书”), respectively. The two Chinese translations have the same POS tag as their corresponding English counterpart. Therefore, once knowing the POS tag for the English word in the context, we are able to pick up the Chinese translation with the same POS tag from the dictionary. In our system, we employ Stanford Log-linear Part-of-Speech tagger (Toutanova et al., 2003) to obtain the POS tag for English word, and POS tag for Chinese

Table 2: Example translations of WSD approaches. The target words are italicized and the proper translations are bolded. We omit the results for Category-based method, due to its poor performance.

English Sentence	Dictionary	Baseline	POST	Bing	Bing+	Bing++
... a very <i>close</i> friend of ...	verb: 关闭, 合, 关 ... adj: 密切, 紧密, 闭合 ... 亲密 ...	关闭	密切	亲密	亲密	亲密
... kids cant <i>stop</i> singing ...	verb: 停止, 站, 阻止, 停 ...	停止	阻止	停止	停止	停止
... it was about elsa being happy and <i>free</i> ...	adj: 免费, 自由, 游离, 畅, 空闲的...	免费	免费	自由	自由	自由
... why obama's <i>trip</i> to my homeland is meaningful ...	noun: 旅, 旅程 ... 旅游 ...	旅	旅	旅	旅行	旅行
... winning more points in the <i>match</i> ...	noun: 匹配, 比赛, 赛, 对手, 对手, 火柴 ...	匹配	匹配	比赛	比赛	比赛
... <i>state</i> department spokeswoman jen psaki said that the allies had a long history of cooperation ...	noun: 态, 国, 州, 状况 ... verb: 声明, 陈述, 述, 申明 ... 发言 ... adj: 国家的 ...	态	态	发言	发言人	国家

words are contained in our dictionary. In some cases, after applying this rule, there is still multiple candidate Chinese translations and we will choose the most frequently used one.

3.3 Machine Translation

A richer context can be exploited is the neighboring word. We regard the whole sentence as the context for the target word and send the sentence to Microsoft Bing Translator³, an online machine translation system with a limited free usage. The returned Chinese translation, however, does not have explicit word alignment to the original English sentence. Therefore, we need additional processing on the Chinese sentence, in order to find the Chinese word that is aligned to the English target.

Bing. As potential Chinese translations are available in our dictionary, the most intuitive processing is to perform a substring match, *i.e.*, check whether the candidate Chinese translation is a substring of the Bing translation. If more than one candidate is matched, we pick up the longest one as the final output. If none is matched, our system will not show translation for the target English

word.

Bing+. The previous method is limited by the coverage of our dictionary. As language is flexible, it is likely our dictionary does not capture all the possible Chinese translations. To alleviate this, we relax the substring restriction, allowing the Bing translation to be a super-string of a candidate translation in our dictionary. To this end, we first segment the Bing translation with Stanford Chinese Word Segmenter (Chang et al., 2008), and then use the matching rule to find the proper Chinese word.

Bing++. In the previous method, it is possible that one Chinese candidate translation in our dictionary matches multiple Chinese words in Bing translation. However, we do not know which Chinese word is corresponded to the target English word. This suggests word alignment information will be useful to resolve this issue. To obtain the alignment, We send the original English sentence and Chinese translation to Bing Word Alignment API⁴, and then apply the same matching rule as Bing+.

³<https://www.bing.com/translator/>

⁴<https://msdn.microsoft.com/en-us/library/dn198370.aspx>

3.4 Evaluation

To evaluate the effectiveness of our proposed methods, we randomly sampled 707 words and their sentences from recent CNN news articles, and manually annotated the ground truth translation for each target English word. We report both the **coverage** (*i.e.*, the chances that a system is able to return a translation) and **accuracy** (*i.e.*, the chances that a translation is appropriate). For comparison purpose, we also report the performance for the baseline method – always select the most frequently used Chinese translation.

Table 3: Experimental results.

	Coverage	Accuracy
Baseline	100%	57.3%
News Category	2.0%	7.1%
POSTagger	94.5%	55.2%
Bing	78.5%	79.8%
Bing+	75.7%	80.9%
Bing++	76.9%	97.4%

Table 3 shows the experimental results for the six methods. As expected, frequency-based baseline achieves 100% of coverage, but a low accuracy (57.3%). POS tagger method shows the same trend. News category based method is the worst among the methods, which suggests using category alone is not sufficient for WSD. On one hand, news category only provides a high-level context. On the other hand, not all of word senses have a strong topic tendency. The three Bing methods improve the accuracy iteratively and all have a reasonable coverage. Among all the methods, Bing++ is the best in terms of accuracy (97.4%), significantly better the others. This suggests the sentence-level context is the most beneficial for our WSD task.

4 Distractors Generation Algorithm

Vocabulary testing is a key functionality in our extension. In this section, we investigate a way to automatically generate suitable distractors (in English form) for a target word. We postulate "a set of suitable distractors" as: 1) being the same form as the target word, 2) fitting the reading context, and 3) having proper difficulty level according to user's knowledge skill.

By applying Part-of-Speech tagger, we obtain the POS tag for the target word, and then restrict the candidate distractors to be selected from the same word class. To make the distractors fitting the context, we identify news category (approach is detailed in Section 2.1), and select the distractors from the same category.

The difficulty of a distractor is measured by its **semantic distance** to the target word: a closer is, a more difficult distractor is. To quantify the semantic distance, we employ Lin Distance (Lin, 1998) to measure the distance between two words in WordNet (Miller, 1995) and define distractors to be difficult if the Lin Distance score is below some threshold. By observing the generated distractors, we empirically set 0.1 as the threshold.

4.1 User knowledge Aware Approach

As mentioned, our extension logs user's detailed learning history, and we categorize user's knowledge on a certain word into three levels, based on the number of times that he/she has encountered the word. Then we adopt different strategies to generate distractors for users in different knowledge level.

Knowledge Level 1 (K1): This indicates user is tested on this word for the first time. Considering this, our system prefers to generate simple distractors, and thus randomly select three words from the same news category.

Knowledge Level 2 (K2): This indicates that the user has known this word for three times. Therefore, the testing is expected to be harder. Our system first randomly selects two words from the same news category. For the third distractor, the system keeps randomly selecting distractor from the same category, computing its semantic distance to the target word, and stops until meeting a difficulty one.

Knowledge Level 3 (K3): This indicates that the user already has a good understanding of the word, *i.e.*, passing the test for six times. As such, we make the test even harder, and choose three difficulty distractors from the new category.

4.2 Evaluation

To compare with our proposed method, we further implemented an existing distractor generation method used in WordGap system (Knoop and

Wilske, 2013). WordGap still adopts knowledge-based approach, selecting the synonyms of synonyms (computed in WordNet) as the distractors. That is, we select the most frequently used word (referred as w_1), from the target word’s synonym set. Then we select the synonyms of w_1 and call this set as s_1 . Synset s_1 contain all the words that are synonyms of synonyms of the target word. Finally we select three most frequently used words from s_1 as distractors for the baseline approach.

For our proposed method, we adopt three different strategies to generate distractors, according to user’s knowledge level. In our evaluation, we focus on assessing distractors generated for two extreme cases, *i.e.*, knowledge level 1, and knowledge level 2. Therefore, we conduct a pairwise comparison – K1 vs. Baseline, and K3 vs. Baseline, using the same testing dataset.

4.2.1 User Study

To compare the two approaches in generating distractors, we turn to users, asking users to compare the plausibility of distractors. We randomly selected 50 sentences from recent news articles and then chose a noun or adjective from the sentence as the target word. In our survey, each question looks like a real MCQ quiz: we show the original sentence (leaving the target word as blank) as the context, and randomly display the six distractors and the target word as choices. Users are required to read the sentence and select the correct answer (that they think) as rating 1, and rank the other choices from 2 (most plausible) to 7 (least plausible) based on their plausibility. Figure 2 shows an example survey question.

We have two tests (K1 vs. Baseline, and K3 vs. Baseline) and each contains 50 questions. We further group 25 questions as one session, and give users the freedom to participate one or more sessions. Each question will be answered by at least five different users. Finally, we recruited 15 users from our university, and half of them are native English speakers.

4.2.2 Results

As each question is answered by five different users, we compute the average rating for each choice. A lower rating means a more plausible (harder) distractor. Unsurprisingly, the rating for all the target words is low (BUG in average), as

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 2: A sample survey question

they are the ground truth. This implies that users were answering the questions seriously, and the evaluation quality

Table 4: Results: Baseline vs. Knowledge Level 1 Algorithm

	Number of winning questions	Average score
Baseline	27	3.84
K1	23	4.10

Table 5: Results: Baseline vs. Knowledge Level 3 Algorithm

	Number of winning questions	Average score
Baseline	21	4.16
K3	29	3.49

Table 4 and Table 5 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 2 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 4 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.

4.2.3 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 Platform Viability and Usability Survey

We have thus far described and evaluated two critical components that can benefit from capturing the learner’s news article context. In the larger context, we also need to check the viability of second language learning intertwined with news reading. In a requirements survey prior to the prototype development, two-thirds of the respondents indicated that although they have used language learning software, they use it infrequently (less than once per week), giving us motivation for our development.

Post-prototype, we conducted a summative survey to assess whether our prototype product satisfied the target niche, in terms of interest, usability and possible interference with normal reading activities. We gathered 16 respondents, 15 of which were between the ages of 18–24. 11 (the majority) also claimed native Chinese language proficiency.

The respondents felt that the extension platform was a viable language learning platform (3.4 of 5; on a scale of 1 “disagreement” to 5 “agreement”) and that they would like to try it when available for their language pair (3 of 5).

In our original prototype, we replaced the original English word with the Chinese translation. While most felt that replacing the original English with the Chinese translation would not ham-

per their reading, they still felt a bit uncomfortable (3.7 of 5). This finding prompted us to review and change the default setting of the learning tooltip to simply add an underline to hint at the tooltip presence.

6 Conclusion

We have described *SystemA*, a software extension and server backend to transform the web browser into a second language learning platform. Leveraging web-based machine translation APIs and a static dictionary, it offers a viable user-driven language learning experience by pairing an improved, context-sensitive tooltip definition capability with the generation of context-sensitive multiple choice questions.

SystemA is potentially not confined to use in news websites; one respondent noted that they would like to use it on arbitrary websites, but currently we feel usable word sense disambiguation is difficult enough even in the restricted domain. We also note that respondents are more willing to use a mobile client for news reading, such that our future development work may be geared towards an independent mobile application, rather than a browser extension.

References

- Jordan L Boyd-Graber, David M Blei, and Xiaojin Zhu. 2007. A topic model for word sense disambiguation. In *EMNLP-CoNLL*, pages 1024–1033.
- Pi-Chuan Chang, Michel Galley, and Christopher D. Manning. 2008. Optimizing chinese word segmentation for machine translation performance. In *Proceedings of the Third Workshop on Statistical Machine Translation*, StatMT ’08, pages 224–232, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Susanne Knoop and Sabrina Wilske. 2013. Wordgap-automatic generation of gap-filling vocabulary exercises for mobile learning. In *Proceedings of Second Workshop NLP Computer-Assisted Language Learning at NODALIDA*, pages 39–47.
- Dekang Lin. 1998. An information-theoretic definition of similarity.
- George A. Miller. 1995. Wordnet: A lexical database for english. *Commun. ACM*, 38(11):39–41, November.

- Neil Naiman. 1978. *The good language learner*, volume 4. Multilingual Matters.
- Roberto Navigli. 2009. Word sense disambiguation: A survey. *ACM Computing Surveys*, 41(2):10:1–10:69, February.
- Joan Rubin. 1975. What the “good language learner” can teach us. *TESOL quarterly*, pages 41–51.
- Knoop; Sabrina Wilske Susanne. 2013. Wordgap - automatic generation of gap-filling vocabulary exercises for mobile learning.
- Kristina Toutanova, H Tolga Ilhan, and Christopher D Manning. 2002. Extensions to hmm-based statistical word alignment models. In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10*, pages 87–94. Association for Computational Linguistics.
- 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.
- Nicola Ueffing and Hermann Ney. 2003. Using pos information for statistical machine translation into morphologically rich languages. In *Proceedings of the tenth conference on European chapter of the Association for Computational Linguistics-Volume 1*, pages 347–354. Association for Computational Linguistics.
- Yorick Wilks and Mark Stevenson. 1998. The grammar of sense: Using part-of-speech tags as a first step in semantic disambiguation. *Natural Language Engineering*, 4(2):135–143, June.