

FLOW: A First-Language-Oriented Writing Assistant System

Mei-Hua Chen*, Shih-Ting Huang+, Hung-Ting Hsieh*, Ting-Hui Kao+, Jason S. Chang+

*Institute of Information Systems and Applications

+Department of Computer Science

National Tsing Hua University

HsinChu, Taiwan, R.O.C. 30013

{chen.meihua, koromiko1104, vincent732, maxis1718, jason.jschang}@gmail.com

Abstract

Writing in English might be one of the most difficult tasks for EFL (English as a Foreign Language) learners. This paper presents FLOW, a writing assistance system. It is built based on first-language-oriented input function and context sensitive approach, aiming at providing immediate and appropriate suggestions including translations, paraphrases, and n-grams during composing and revising processes. FLOW is expected to help EFL writers achieve their writing flow without being interrupted by their insufficient lexical knowledge.

1. Introduction

Writing in a second language (L2) is a challenging and complex process for foreign language learners. Insufficient lexical knowledge and limited exposure to English might interrupt their writing flow (Silva, 1993). Numerous writing instructions have been proposed (Kroll, 1990) as well as writing handbooks have been available for learners. Studies have revealed that during the writing process, EFL learners show the inclination to rely on their native languages (Wolfersberger, 2003) to prevent a breakdown in the writing process (Arndt, 1987; Cumming, 1989). However, existing writing courses and instruction materials, almost second-language-oriented, seem unable to directly assist EFL writers while writing.

This paper presents FLOW¹ (Figure 1), an interactive system for assisting EFL writers in

composing and revising writing. Different from existing tools, its context-sensitive and first-language-oriented features enable EFL writers to concentrate on their ideas and thoughts without being hampered by the limited lexical resources. Based on the studies that first language use can positively affect second language composing, FLOW attempts to meet such needs. Given any L1 input, FLOW displays appropriate suggestions including translation, paraphrases, and n-grams during composing and revising processes. We use the following example sentences to illustrate these two functionalities.

Consider the sentence “*We propose a method to*”. During the composing stage, suppose a writer is unsure of the phrase “*solve the problem*”, he could write “*解決問題*”, a corresponding word in his native language, like “*We propose a method to 解決問題*”. The writer’s input in the writing area of FLOW actively triggers a set of translation suggestions such as “*solve the problem*” and “*tackle the problem*” for him/her to complete the sentence.

In the revising stage, the writer intends to improve or correct the content. He/She is likely to change the sentence illustrated above into “*We try all means to solve the problem.*” He would select the phrase “*propose a method*” in the original sentence and input a L1 phrase “*盡力*”, which specifies the meaning he prefers. The L1 input triggers a set of context-aware suggestions corresponding to the translations such as “*try our best*” and “*do our best*” rather than “*try your best*” and “*do your best*”. The system is able to do that mainly by taking a context-sensitive approach. FLOW then inserts the phrase the writer selects into the sentence.

¹ FLOW: <http://flowa1demo.appspot.com>

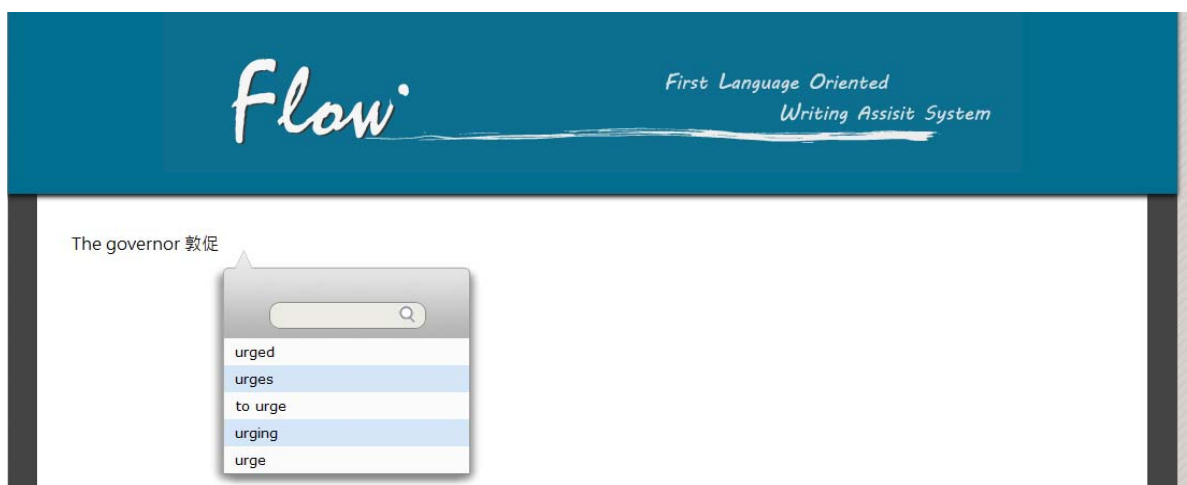


Figure 1. Screenshot of FLOW

In this paper, we propose a context-sensitive disambiguation model which aims to automatically choose the appropriate phrases in different contexts when performing n-gram prediction, paraphrase suggestion and translation tasks. As described in (Carpuat and Wu, 2007), the disambiguation model plays an important role in the machine translation task. Similar to their work, we further integrate the multi-word phrasal lexical disambiguation model to the n-gram prediction model, paraphrase model and translation model of our system. With the phrasal disambiguation model, the output of the system is sensitive to the context the writer is working on. The context-sensitive feature helps writers find the appropriate phrase while composing and revising.

This paper is organized as follows. We review the related work in the next section. In Section 3, we brief our system and method. Section 4 reports the evaluation results. We conclude this paper and point out future directions to research in Section 5.

2. Related Work

2.1 Sub-sentential paraphrases

A variety of data-driven paraphrase extraction techniques have been proposed in the literature. One of the most popular methods leveraging bilingual parallel corpora is proposed by Bannard and Callison-Burch (2005). They identify paraphrases using a phrase in another language as a pivot. Using bilingual parallel corpora for

paraphrasing demonstrates the strength of semantic equivalence. Another line of research further considers context information to improve the performance. Instead of addressing the issue of local paraphrase acquisition, Max (2009) utilizes the source and target contexts to extract sub-sentential paraphrases by using pivot SMT systems.

2.2 N-gram suggestions

After a survey of several existing writing tools, we focus on reviewing two systems closely related to our study.

PENS (Liu et al, 2000), a machine-aided English writing system, provides translations of the corresponding English words or phrases for writers' reference. Different from PENS, FLOW further suggests paraphrases to help writers revise their writing tasks. While revising, writers would alter the use of language to express their thoughts. The suggestions of paraphrases could meet their need, and they can reproduce their thoughts more fluently.

Another tool, TransType (Foster, 2002), a text editor, provides translators with appropriate translation suggestions utilizing trigram language model. The differences between our system and TransType lie in the purpose and the input. FLOW aims to assist EFL writers whereas TransType is a tool for skilled translators. On the other hand, in TransType, the human translator types translation of a given source text, whereas in FLOW the input,

either a word or a phrase, could be source or target languages.

2.3 Multi-word phrasal lexical disambiguation

In the study more closely related to our work, Carpuat and Wu (2007) propose a novel method to train a phrasal lexical disambiguation model to benefit translation candidates selection in machine translation. They find a way to integrate the state-of-the-art Word Sense Disambiguation (WSD) model into phrase-based statistical machine translation. Instead of using predefined senses drawn from manually constructed sense inventories, their model directly disambiguates between all phrasal translation candidates seen during SMT training. In this paper, we also use the phrasal lexical disambiguation model; however, apart from using disambiguation model to help machine translation, we extend the disambiguation model. With the help of the phrasal lexical disambiguation model, we build three models: a context-sensitive n-gram prediction model, a paraphrase suggestion model, and a translation model which are introduced in the following sections.

3. Overview of FLOW

The FLOW system helps language learners in two ways: predicting n-grams in the composing stage and suggesting paraphrases in the revising stage (Figure 2).

3.1 System architecture

Composing Stage

During the composing process, a user inputs S . FLOW first determines if the last few words of S is a L1 input. If not, FLOW takes the last k words to predict the best matching following n-grams. Otherwise, the system uses the last k words as the query to predict the corresponding n-gram translation. With a set of prediction (either translations or n-grams), the user could choose an appropriate suggestion to complete the sentence in the writing area.

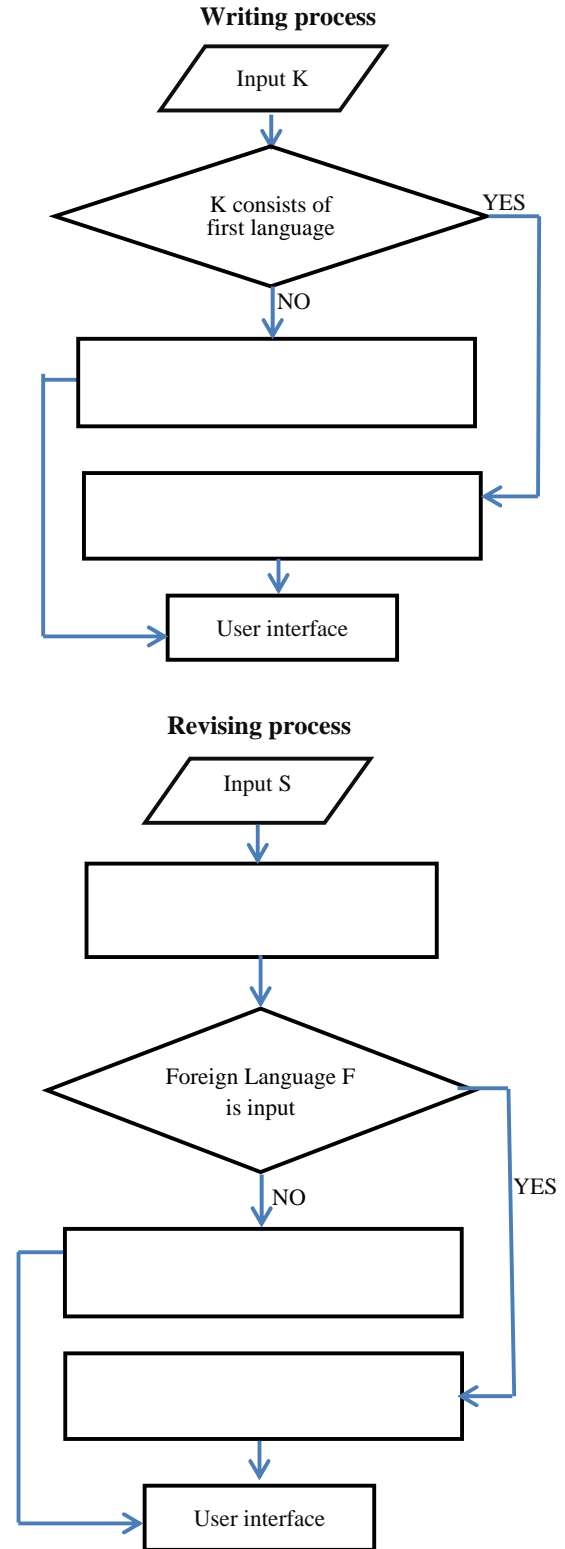


Figure 2. Overall Architecture of *FLOW* in writing and revising processes

Revising Stage

In the revising stage, given an input I and the user selected words K , *FLOW* obtains the word sequences L and R surrounding K as reference for prediction. Next, the system suggests sub-sentential paraphrases for K based on the information of L and R . The system then searches and ranks the translations.

3.2 N-gram prediction

In the n-gram prediction task, our model takes the last k words with m^2 English words and n foreign language words, $\{e_1, e_2, \dots, e_m, f_1, f_2 \dots, f_n\}$, of the source sentences S as the input. The output would be a set of n-gram predictions. These n-grams can be concatenated to the end of the user-composed sentence fluently.

Context-Sensitive N-gram Prediction (CS-NP)

The CS-NP model is triggered to predict a following n-gram when a user composes sentences consisted of only English words with no foreign language words, namely, n is equal to 0. The goal of the CS-NP model is to find the English phrase e that maximizes the language model probability of the word sequence, $\{e_1, e_2, \dots, e_m, e\}$:

$$e = \operatorname{argmax}_{e, \forall m \leq k} P(e|e_1, e_2, \dots, e_m)$$

$$P(e|e_1, e_2, \dots, e_m) = \frac{P(e_1, e_2, \dots, e_m, e)}{P(e_1, e_2, \dots, e_m)}$$

Translation-based N-gram Prediction (TB-NP)

When a user types a set of L1 expression $f = \{f_1, f_2 \dots, f_n\}$, following the English sentences S , the *FLOW* system will predict the possible translations of f . A simple way to predict the translations is to find the bilingual phrase alignments $T(f)$ using the method proposed by (Och and Ney, 2003). However, the $T(f)$ is ambiguous in different contexts. Thus, we use the context $\{e_1, e_2, \dots, e_m\}$ preceding f to fix the prediction of the translation. Predicting the translation e can be treated as a sub-sentential translation task:

$$e = \operatorname{argmax}_{e \in T(f)} P(e|e_1, e_2, \dots, e_m),$$

where we use the user-composed context $\{e_1, e_2, \dots, e_m\}$ to disambiguate the translation of f . Although there exist more sophisticated models which could make a better prediction, a simple naïve-Bayes model is shown to be accurate and efficient in the lexical disambiguation task according to (Yarowsky and Florian, 2002). Therefore, in this paper, a naïve-Bayes model is used to disambiguate the translation of f . In addition to the context-word feature, we also use the context-syntax feature, namely surrounding POS tag Pos , to constrain the syntactic structure of the prediction. The TB-NP model could be represented in the following equation:

$$e^* = \operatorname{argmax}_e P(e|e_1, e_2, \dots, e_m, p_1, p_2, \dots, p_m),$$

$$Pos = \{p_1, p_2, \dots, p_m\}$$

According to the Bayes theorem,

$$P(e|e_1, e_2, \dots, e_m, p_1, p_2, \dots, p_m)$$

$$= \prod_{e_i \in E} P(e_i|e) * \prod_{p_j \in P} P(p_j|e)$$

The probabilities can be estimated using a parallel corpus, which is also used to obtain bilingual phrase alignment.

3.3 Paraphrase Suggestion

Unlike the N-gram prediction, in the paraphrase suggestion task, the user selects k words, $\{e_1, e_2, \dots, e_k\}$, which he/she wants to paraphrase. The model takes the m words $\{r_1, r_2, \dots, r_m\}$ and n words $\{l_1, l_2, \dots, l_n\}$ in the right and left side of the user-selected k words respectively. The system also accepts an additional foreign language input, $\{f_1, f_2, \dots, f_l\}$, which helps limit the meaning of suggested paraphrases to what the user really wants. The output would be a set of paraphrase suggestions that the user-selected phrases can be replaced by those paraphrases precisely.

Context-Sensitive Paraphrase Suggestion (CS-PS)

The CS-PS model first finds a set of local paraphrases P of the input phrase K using the

² In this paper, $m = 5$.

pivot-based method proposed by Bannard and Callison-Burch (2005). Although the pivot-based method has been proved efficient and effective in finding local paraphrases, the local paraphrase suggestions may not fit different contexts. Similar to the previous n-gram prediction task, we use the naïve-Bayes approach to disambiguate these local paraphrases. The task is to find the best e such that e with the highest probability for the given context R and L. We further require paraphrases to have similar syntactic structures to the user-selected phrase in terms of POS tags, Pos .

$$e^* = \operatorname{argmax}_{e \in P} P(e|l_1, l_2, \dots, l_n, r_1, r_2, \dots, r_m, Pos)$$

Translation-based Paraphrase Suggestion (TB-PS)

After the user selects a phrase for paraphrasing, with a L1 phrase F as an additional input, the suggestion problem will be:

$$e^* = \operatorname{argmax}_{e \in T(F)} P(e|l_1, l_2, \dots, l_n, r_1, r_2, \dots, r_m, Pos)$$

The TB-PS model disambiguates paraphrases from the translations of F instead of paraphrases P .

4. Experimental Results

In this section, we describe the experimental setting and the preliminary results. Instead of training a whole machine translation using toolkits such as Moses (Koehn et. al, 2007), we used only bilingual phrase alignment as translations to prevent from the noise produced by the machine translation decoder. Word alignments were produced using Giza++ toolkit (Och and Ney, 2003), over a set of 2,220,570 Chinese-English sentence pairs in Hong Kong Parallel Text (LDC2004T08) with sentences segmented using the CKIP Chinese word segmentation system (Ma and Chen, 2003). In training the phrasal lexical disambiguation model, we used the English part of Hong Kong Parallel Text as our training data.

To assess the effectiveness of FLOW, we selected 10 Chinese sentences and asked two students to translate the Chinese sentences to English sentences using FLOW. We kept track of the sentences the two students entered. Table 1 shows the selected results.

Model	Results
TB-PS	總而言之, the price of rice...
	in short
	all in all
	in a nutshell
	in a word
	to sum up
CS-PS	She looks forward to coming
	look forward to
	looked forward to
	is looking forward to
	forward to
	expect
CS-PS	there is no doubt that ...
	there is no question
	it is beyond doubt
	I have no doubt
	beyond doubt
	it is true
CS-NP	We put forward ...
	the proposal
	additional
	our opinion
	the motion
	the bill
TB-NP	...on ways to identify tackle 洗錢
	money laundering
	money
	his
	forum entitled
	money laundry

Table 1. The preliminary results of FLOW

Both of the paraphrase models CS-PS and TB-PS perform quite well in assisting the user in the writing task. However, there are still some problems such as the redundancy suggestions, e.g., “*look forward to*” and “*looked forward to*”. Besides, although we used the POS tags as features, the syntactic structures of the suggestions are still not consistent to an input or selected phrases. The CS-NP and the TB-NP model also perform a good task. However, the suggested phrases are usually too short to be a semantic unit. The disambiguation model tends to produce shorter phrases because they have more common context features.

5. Conclusion and Future Work

In this paper, we presented FLOW, an interactive writing assistance system, aimed at helping EFL writers compose and revise without interrupting their writing flow. First-language-oriented and context-sensitive features are two main contributions in this work. Based on the studies on second language writing that EFL writers tend to use their native language to produce texts and then translate into English, the first-language-oriented function provides writers with appropriate translation suggestions. On the other hand, due to the fact that selection of words or phrases is sensitive to syntax and context, our system provides suggestions depending on the contexts. Both functions are expected to improve EFL writers' writing performance.

In future work, we will conduct experiments to gain a deeper understanding of EFL writers' writing improvement with the help of FLOW, such as integrating FLOW into the writing courses to observe the quality and quantity of students' writing performance. Many other avenues exist for future research and improvement of our system. For example, we are interested in integrating the error detection and correction functions into FLOW to actively help EFL writers achieve better writing success and further motivate EFL writers to write with confidence.

References

- Valerie Arndt. 1987. Six writers in search of texts: A protocol based study of L1 and L2 writing. *ELT Journal*, 41, 257-267.
- Colin Bannard and Chris Callison-Burch. 2005. Paraphrasing with bilingual parallel corpora. In *Proceedings of ACL*, pp. 597-604.
- Marine Carpuat and Dekai Wu. 2007. Improving Statistical Machine Translation using Word Sense Disambiguation. In *Proceedings of EMNLP-CoNLL*, pp 61-72.
- Alister Cumming. 1989. Writing expertise and second language proficiency. *Language Learning*, 39, 81-141.
- George Foster, Philippe Langlais, and Guy Lapalme. 2002. Transtype: Text prediction for translators. In *Proceedings of ACL Demonstrations*, pp. 93-94.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of ACL Demonstration Session*, pp. 177-180.
- Barbara Kroll. 1990. *Second Language Writing: Research Insights for the Classroom*. Cambridge University Press, Cambridge.
- Aurélien Max. 2009. Sub-sentential Paraphrasing by Contextual Pivot Translation. In *Proceedings of the 2009 Workshop on Applied Textual Inference, ACL-IJCNLP*, pp 18-26.
- Tony Silva. 1993. Toward an Understanding of the Distinct Nature of L2 Writing: The ESL Research and Its Implications. *TESOL Quarterly* 27(4): 657-77.
- Liu, Ting, Mingh Zhou, Jianfeng Gao, Endong Xun, and Changning Huan. 2000. PENS: A Machine-Aided English Writing System for Chinese Users. In *Proceedings of ACL*, pp 529-536.
- Mark Wolfersberger. 2003. L1 to L2 writing process and strategy transfer: a look at lower proficiency writers. *TESL-EJ: Teaching English as a Second or Foreign Language*, 7(2), A6 1-15.