

# Improving the Translation of Discourse Markers for Chinese into English

David Steele

Department Of Computer Science

The University of Sheffield

Sheffield, UK

dbsteele1@sheffield.ac.uk

## Abstract

Discourse markers (DMs) are ubiquitous cohesive devices used to connect what is said or written. However, across languages there is divergence in their usage, placement, and frequency, which is considered to be a major problem for machine translation (MT). This paper presents an overview of a proposed thesis, exploring the difficulties around DMs in MT, with a focus on Chinese and English. The thesis will examine two main areas: modelling cohesive devices within sentences and modelling discourse relations (DRs) across sentences. Initial experiments have shown promising results for building a prediction model that uses linguistically inspired features to help improve word alignments with respect to the implicit use of cohesive devices, which in turn leads to improved hierarchical phrase-based MT.

## 1 Introduction

Statistical Machine Translation (SMT) has, in recent years, seen substantial improvements, yet approaches are not able to achieve high quality translations in many cases. The problem is especially prominent with complex composite sentences and distant language pairs, largely due to computational complexity. Rather than considering larger discourse segments as a whole, current SMT approaches focus on the translation of single sentences independently, with clauses and short phrases being treated in isolation. DMs are seen as a vital contextual link between discourse segments and could be used to guide translations in order to improve

accuracy. However, they are often translated into the target language in ways that differ from how they are used in the source language (Hardmeier, 2012a; Meyer and Popescu-Belis, 2012). DMs can also signal numerous DRs and current SMT approaches do not adequately recognise or distinguish between them during the translation process (Hajlaoni and Popescu-Belis, 2013). Recent developments in SMT potentially allow the modelling of wider discourse information, even across sentences (Hardmeier, 2012b), but currently most existing models appear to focus on producing well translated localised sentence fragments, largely ignoring the wider global cohesion.

Five distinct cohesive devices have been identified (Halliday and Hasan, 1976), but for this thesis the pertinent devices that will be examined are conjunction (DMs) and (endophoric) reference. Conjunction is pertinent as it encompasses DMs, whilst reference includes pronouns (amongst other elements), which are often connected with the use of DMs (e.g. ‘Because John ..., therefore he ...’).

The initial focus is on the importance of DMs within sentences, with special attention given to implicit markers (common in Chinese) and a number of related word alignment issues. However, the final thesis will cover two main areas:

- Modelling cohesive devices within sentences
- Modelling discourse relations across sentences and wider discourse segments.

This paper is organized as follows. In Section 2 a survey of related work is conducted. Section 3

outlines the initial motivation and research including a preliminary corpus analysis. It covers examples that highlight various problems with the translation of (implicit) DMs, leading to an initial intuition. Section 4 looks at experiments and word alignment issues following a deeper corpus analysis and discusses how the intuition led towards developing the methodology used to study and improve word alignments. It also includes the results of the experiments that show positive gains in BLEU. Section 5 provides an outline of the future work that needs to be carried out. Finally, Section 6 is the conclusion.

## 2 Literature Review

This section is a brief overview of some of the pertinent important work that has gone into improving SMT with respect to cohesion. Specifically the focus is on the areas of: identifying and annotating DMs, working with lexical and grammatical cohesion, and translating implicit DRs.

### 2.1 Identifying and Annotating Chinese DMs

A study on translating English discourse connectives (DCs) (Hajlaoni and Popescu-Belis, 2013) showed that some of them in English can be ambiguous, signalling a variety of discourse relations. However, other studies have shown that sense labels can be included in corpora and that MT systems can take advantage of such labels to learn better translations (Pitler and Nenkova, 2009; Meyer and Popescu-Belis, 2012). For example, The Penn Discourse Treebank project (PDTB) adds annotation related to structure and discourse semantics with a focus on DRs and can be used to guide the extraction of DR inferences. The Chinese Discourse Treebank (CDTB) adds an extra layer to the annotation in the PDTB (Xue, 2005) focussing on DCs as well as structural and anaphoric relations and follows the lexically grounded approach of the PDTB.

The studies also highlight how anaphoric relations can be difficult to capture as they often have one discourse adverbial linked with a local argument, leaving the other argument to be established from elsewhere in the discourse. Pronouns, for example, are often used to link back to some discourse entity that has already been introduced. This essentially suggests that arguments identified in anaphoric relations

English	Chinese DC
although(1)/but(2)	(1) 虽然, 虽说, 虽 (2) 但, 可是, 却
because(1)/therefore(2)	(1) 因为, 因, 由于 (2) 所以
if(1)/then(2)	(1) 如果, 假如, 若 (2) 就

Table 1: Examples of Interchangeable DMs.

can cover a long distance and Xue (2005) argues that one of the biggest challenges for discourse annotation is establishing the distance of the text span and how to decide on what discourse unit should be included or excluded from the argument.

There are also some additional challenges such as variants or substitutions of DCs. Table 1 (Xue, 2005) shows a range of DCs that can be used interchangeably. The numbers indicate that any marker from (1) can be paired with any marker from (2) to form a compound sentence with the same meaning.

### 2.2 Lexical and Grammatical Cohesion

Previous work has attempted to address lexical and grammatical cohesion in SMT (Gong et al., 2011; Xiao et al., 2011; Wong and Kit, 2012; Xiong et al., 2013b) although their results are still relatively limited (Xiong et al., 2013a). Lexical cohesion is determined by identifying lexical items forming links between sentences in text (also lexical chains). A number of models have been proposed in order to try and capture document-wide lexical cohesion and when implemented they showed significant improvements over the baseline (Xiong et al., 2013a).

Lexical chain information (Morris and Hirst, 1991) can be used to capture lexical cohesion in text and it is already successfully used in a range of fields such as information retrieval and the summarisation of documents (Xiong et al., 2013b). The work of Xiong et al. (2013b) introduces two lexical chain models to incorporate lexical cohesion into document wide SMT and experiments show that, compared to the baseline, implementing these models substantially improves translation quality. Unfortunately with limited grammatical cohesion, propagated by DMs, translations can be difficult to understand, especially if there is no context provided

by local discourse segments.

To achieve improved grammatical cohesion Tu et al. (2014) propose creating a model that generates transitional expressions through using complex sentence structure based translation rules alongside a generative transfer model, which is then incorporated into a hierarchical phrase-based system. The test results show significant improvements leading to smoother and more cohesive translations. One of the key reasons for this is through reserving cohesive information during the training process by converting source sentences into “tagged flattened complex sentence structures”(Tu et al., 2014) and then performing word alignments using the translation rules. It is argued that connecting complex sentence structures with transitional expressions is similar to the human translation process (Tu et al., 2014) and therefore improvements have been made showing the effectiveness of preserving cohesion information.

### 2.3 Translation of Implicit Discourse Relations

It is often assumed that the discourse information captured by the lexical chains is mainly explicit. However, these relations can also be implicitly signalled in text, especially for languages such as Chinese where implicature is used in abundance (Yung, 2014). Yung (2014) explores DM annotation schemes such as the CDTB (2.1) and observes that explicit relations are identified with an accuracy of up to 94%, whereas with implicit relations this can drop as low as 20% (Yung, 2014). To overcome this, Yung proposes implementing a discourse-relation aware SMT system, that can serve as a basis for producing a discourse-structure-aware, document-level MT system. The proposed system will use DC annotated parallel corpora, that enables the integration of discourse knowledge. Yung argues that in Chinese a segment separated by punctuation is considered to be an elementary discourse unit (EDU) and that a running Chinese sentence can contain many such segments. However, the sentence would still be translated into one single English sentence, separated by ungrammatical commas and with a distinct lack of connectives. The connectives are usually explicitly required for the English to make sense, but can remain implicit in the Chinese (Yung, 2014). However, this work is still in the early stages.

## 3 Motivation

This section outlines the initial research, including a preliminary corpus analysis, examining difficulties with automatically translating DMs across distant languages such as Chinese and English. It draws attention to deficiencies caused from under-utilising discourse information and examines divergences in the usage of DMs. The final part of this section outlines the intuition garnered from the given examples and highlights the approach to be undertaken.

For the corpus analysis, research, and experiments three main parallel corpora are used:

- Basic Travel Expression Corpus (BTEC): Primarily made up of short simple phrases that occur in travel conversations. It contains 44,016 sentences in each language with over 250,000 Chinese characters and over 300,000 English words (Takezawa et al., 2012).
- Foreign Broadcast Information Service (FBIS) corpus: This uses a variety of news stories and radio podcasts in Chinese. It contains 302,996 parallel sentences with 215 million Chinese characters and over 237 million English words.
- Ted Talks corpus (TED): Made up of approved translations of the live Ted Talks presentations<sup>1</sup>. It contains over 300,000 Chinese characters and over 2 million English words from 156,805 sentences (Cettolo et al., 2012).

Chinese uses a rich array of DMs including: simple conjunctions, composite conjunctions, and zero connectives where the meaning or context is strongly inferred across clauses with sentences having natural, allowable omissions, which can cause problems for current SMT approaches. Here a few examples<sup>2</sup> are outlined:

Ex (1) 他因为病了，没来上课。

he because ill, not come class.

Because he was sick, he didn't come to class<sup>3</sup>.

He is ill, absent. (Bing)

<sup>1</sup><http://www.ted.com>

<sup>2</sup>These examples (Steele and Specia, 2014) are presented as: Chinese sentence / literal translation / reference translation / automated translation - using either Google or Bing.

<sup>3</sup>(Ross and Sheng, 2006)

Ex (2) 你因为这个在吃什么药吗?  
 you because this (be) eat what medicine?  
 Have you been taking anything for this? (BTEC)  
 What are you eating because of this medicine?  
 (Google)

Both examples show ‘because’ (因为) being used in different ways and in each case the automated translations fall short. In Ex1 the dropped (implied) pronoun in the second clause could be the problem, whilst in Ex2 significant reordering is needed as ‘because’ should be linked to ‘this’ (这个) - the topic - rather than ‘medicine’ (药). The ‘this’ (这个) refers to an ‘ailment’, which is hard to capture from a single sentence. Information preserved from a larger discourse segment may have provided more clues, but as is, the sentence appears somewhat exophoric and the meaning cannot necessarily be gleaned from the text alone.

Ex (3) 一有空位我们就给你打电话。  
 as soon as have space we then give you make phone.  
 We’ll call you as soon as there is an opening.  
 (BTEC)  
 A space that we have to give you a call. (Google)

In Ex3 the characters ‘一’ and ‘就’ are working together as coordinating markers in the form: ...一VP<sup>a</sup> 就 VP<sup>b</sup>. However, individually these characters have significantly different meanings, with ‘一’ meaning ‘a’ or ‘one’ amongst many things. Yet, in the given sentence using the ‘一’ and ‘就’ construct ‘一’ has a meaning akin to ‘as soon as’ or ‘once’, while ‘就’ implies a ‘then’ relation, both of which can be difficult to capture. Figure 1<sup>4</sup> shows an example where word alignment failed to map the ‘as soon as ... then’ structure to ...一... 就... . That is, columns 7, 8, 9, which represent ‘as soon as’ in the English have no alignment points whatsoever. Yet, in this case, all three items should be aligned to the single element ‘一’ which is on row 1 on the Chinese side. Additionally, the word ‘returns’ (column 11), which is currently aligned to ‘一’ (row 1) should in fact be aligned to ‘回来’ (return/come back) in row 2. This misalignment

<sup>4</sup>The boxes with a ‘#’ inside are the alignment points and each coloured block (large or small) is a minimal-biphrase.

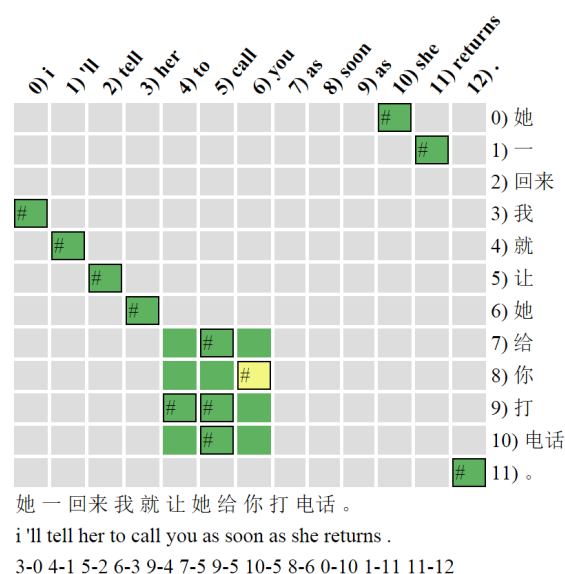


Figure 1: A visualisation of word alignments for the given parallel sentence, showing a non-alignment of ‘as soon as’.

could be a direct side-effect of having no alignment for ‘as soon as’ in the first place. Consequently, the knock-on effect of poor word alignment, especially around markers - as in this case, will lead to the overall generation of poorer translation rules.

Ex (4) 他因为病了, 所以他没来上课。  
 he because ill, so he not come class.  
 Because he was sick, he didn’t come to class.  
 He is ill, so he did not come to class. (Bing)

Ex4 is a modified version of Ex2, with an extra ‘so’(所以) and ‘he’ (他) manually inserted in the second clause of the Chinese sentence. Grammatically these extra characters are not required for the Chinese to make sense, but are still correct. However, the interesting point is that the extra information (namely ‘so’ and ‘he’) has enabled the system to produce a much better final translation.

From the given examples it appears that both implicitation and the use of specific DM structures can cause problems when generating automated translations. The highlighted issues suggest that making markers (and possibly, by extension, pronouns) explicit, due to linguistic clues, more information becomes available, which can support the extraction of word alignments. Although making implicit mark-

ers explicit can seem unnatural and even unnecessary for human readers, it does follow that if the word alignment process is made easier by this explicitation it will lead to better translation rules and ultimately better translation quality.

## 4 Experiments and Word Alignments

This section examines the current ongoing research and experiments that aim to measure the extent of the difficulties caused by DMs. In particular the focus is on automated word alignments and problems around implicit and misaligned DMs. The work discussed in Section 3 highlighted the importance of improving word alignments, and especially how missing alignments around markers can lead to the generation of poorer rules.

Before progressing onto the experiments an initial baseline system was produced according to detailed criteria (Chiang, 2007; Saluja et al., 2014). The initial system was created using the ZH-EN data from the BTE parallel corpus (Paul, 2009) (Section 3). Fast-Align is used to generate the word alignments and the CDEC decoder (Dyer et al., 2010) is used for rule extraction and decoding. The baseline and subsequent systems discussed here are hierarchical phrase-based systems for Chinese to English translation.

Once the alignments were obtained the next step in the methodology was to examine the misalignment information to determine the occurrence of implicit markers. A variance list was created<sup>5</sup> that could be used to cross-reference discourse markers with appropriate substitutable words (as per Table 1). Each DM was then examined in turn (automatically) to look at what it had been aligned to. When the explicit English marker was aligned correctly, according to the variance list, then no change was made. If the marker was aligned to an unsuitable word, then an artificial marker was placed into the Chinese in the nearest free space to that word. Finally if the marker was not aligned at all then an artificial marker was inserted into the nearest free space

<sup>5</sup>The variance list is initially created by filtering good alignments and bad alignments by hand and using both on-line and off-line (bi-lingual) dictionaries/resources.

DM	BTEC	FBIS	TED
if	25.70%	40.75%	23.35%
then	21.00%	50.85 %	40.47%
because	23.95%	32.80%	16.48%
but	29.40%	39.90%	27.08%

Table 2: Misalignment information for the 3 corpora.

System	DEV	TST
BTEC-Dawn (baseline)	34.39	35.02
BTEC-Dawn (if)	34.60	35.03
BTEC-Dawn (then)	34.69	35.04
BTEC-Dawn (but)	34.51	35.21
BTEC-Dawn (because)	34.41	35.02
BTEC-Dawn (all)	34.53	35.46

Table 3: BLEU Scores for the Experimental Systems

by number<sup>6</sup>. A percentage of misalignments<sup>7</sup> across all occurrences of individual markers was also calculated.

Table 2 shows the misalignment percentages for the four given DMs across the three corpora. The average sentence length in the BTE Corpus is eight units, in the FBIS corpus it is 30 units, and in the TED corpus it is 29 units. The scores show that there is a wide variance in the misalignments across the corpora, with FBIS consistently having the highest error rate, but in all cases the percentage is fairly significant.

Initially tokens were inserted for single markers at a time, but then finally with tokens for all markers inserted simultaneously. Table 3 shows the BLEU scores for all the experiments. The first few experiments showed improvements over the baseline of up to +0.30, whereas the final one showed improvements of up to +0.44, which is significant.

After running the experiments the visualisation of a number of word alignments (as per Figures 1,2,3) were examined and a single example of a ‘then’ sentence was chosen at random. Figure 2 shows the word alignments for a sentence from the baseline system, and Figure 3 shows the word alignments for

<sup>6</sup>The inserts are made according to a simple algorithm, and inspired by the examples in Section 3.

<sup>7</sup>A non-alignment is not necessarily a bad alignment. For example: ‘正反’ = ‘positive and negative’, with no ‘and’ in the Chinese. In this case a non-alignment for ‘and’ is acceptable.

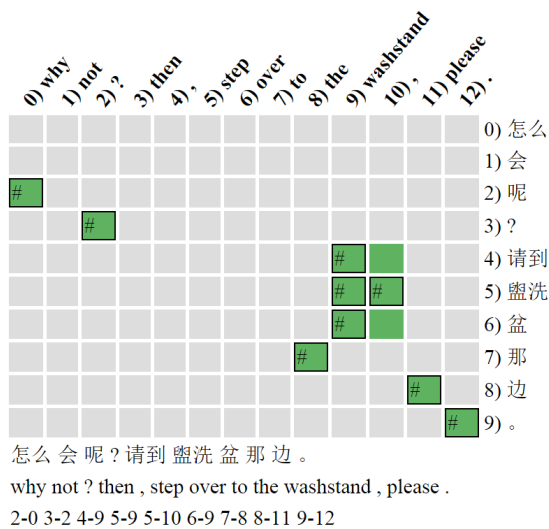


Figure 2: Visualisation of word alignments showing no alignment for ‘then’ in column 3.

the same sentence, but with an artificial marker automatically inserted for the unaligned ‘then’.

The differences between the word alignments in the figures are subtle, but positive. For example, in Figure 3 more of the question to the left of ‘then’ is captured correctly. Moreover, to the right of ‘then’, ‘over’ has now been aligned quite well to ‘那边’ (over there) and ‘to’ has been aligned to ‘请到’ (please - go to). Perhaps most significantly though is the mish-mash of alignments to ‘washstand’ in Figure 2 has now been replaced by a very good alignment to ‘盥洗盆’ (washbasin/washstand) showing an overall smoother alignment. These preliminary findings indicate that there is plenty of scope for further positive investigation and experimentation.

## 5 Ongoing Work

This section outlines the two main research areas (Section 1) that will be tackled in order to feed into the final thesis. Having addressed the limitations of current SMT approaches, the focus has moved on to looking at cohesive devices at the sentential level, but ultimately the overall aim is to better model DRs across wider discourse segments.

### 5.1 Modelling Cohesive Devices Within Sentences

Even at the sentence level there exists a local context, which produces dependencies between certain

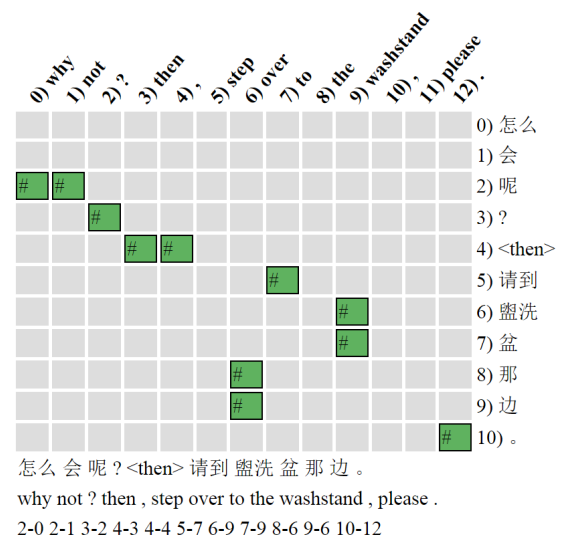


Figure 3: Visualisation of word alignments showing the artificial marker ‘<then>’ and a smoother overall alignment.

words. The cohesion information within the sentence can hold vital clues for tasks such as pronoun resolution, and so it is important to try to capture it.

Simply looking at the analysis in Section 4 provides insight into which other avenues should be explored for this part, including:

- Expanding the number of DMs being explored, including complex markers (e.g. as soon as).
- Improving the variance list to capture more variant translations of marker words. It is also important here to include automated filtering for difficult DMs (e.g. cases where ‘and’ or ‘so’ are not being used as specific markers can perhaps make them more difficult to align). Making significant use of parts of speech tagging and annotated texts could be useful.
- Develop better insertion algorithms to produce an improved range of insertion options, and reduce damage to existing word alignments.
- Looking at using alternative/additional evaluation metrics and tools to either replace or complement BLEU. This could produce more targeted evaluation that is better at picking up on individual linguistic components such as DMs and pronouns.

However, the final aim is to work towards a true prediction model using parallel data as a source of annotation. Creating such a model can be hard monolingually, whereas a bilingual corpus can be used as a source of additional implicit annotation or indeed a source of additional signals for discourse relations. The prediction model should make the word alignment task easier (through either guiding the process or adding constraints), which in turn will generate better translation rules and ultimately should improve MT.

## 5.2 Modelling Discourse Relations Across Sentences

This part will be an extension of the tasks in Section 5.1. The premise is that if the discourse information or local context within a sentence can be captured then it could be applied to wider discourse segments and possibly the whole document. Some inroads into this task have been trialled through using lexical chaining (Xiong et al., 2013b). However, more recently tools are being developed enabling document wide access to the text, which should provide scope for examining the links between larger discourse units - especially sentences and paragraphs.

## 6 Conclusions

The findings in Section 3 highlighted that implicit cohesive information can cause significant problems for MT and that by adding extra information translations can be made smoother. Section 4 extended this idea and outlined the experiments and methodology used to capture some effects of automatically inserting artificial tokens for implicit or misaligned DMs. It showed largely positive results, with some good improvements to the word alignments, indicating that there is scope for further investigation and experimentation. Finally, section 5 highlighted the two main research areas that will guide the thesis, outlining a number of ways in which the current methodology and approach could be developed.

The ultimate aim is to use bilingual data as a source of additional clues for a prediction model of Chinese implicit markers, which can, for instance, guide and improve the word alignment process leading to the generation of better rules and smoother translations.

## References

- Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. *Web Inventory of Transcribed and Translated Talks*. In: EAMT, pages 261-268. Trento, Italy.
- David Chiang. 2007. *Hierarchical phrase-based translation*. Computational Linguistics, 33(2):201-228.
- Chris Dyer, Adam Lopez, Juri Ganitkevitch, Johnathan Weese, Ferhan Ture, Phil Blunsom, Hendra Setiawan, Vladimir Eidelman, and Philip Resnik. 2010. *CDEC: A decoder, Alignment, and Learning Framework for Finite-state and Context-free Translation Models*. In Proceedings of ACL.
- Zhengxian Gong, Min Zhang, and Guodong Zhou. 2011. *Cache-based Document-level Statistical Machine Translation*. In 2011 Conference on Empirical Methods in Natural Language Processing, pages 909-919. Edinburgh, Scotland, UK
- Najeh Hajlaoui and Andre Popescu-Belis. 2013. *Translating English Discourse Connectives into Arabic: a Corpus-based analysis and an Evaluation Metric*. In: CAASL4 Workshop at AMTA (Fourth Workshop on Computational Approaches to Arabic Script-based Languages), San Diego, CA, pages 1-8.
- M.A.K Halliday and Ruqaiya Hasan. 1976. *Cohesion in English (English Language Series)* Longmen, London
- Christian Hardmeier. 2012. *Discourse in Statistical Machine Translation: A Survey and a Case Study* Elanders Sverige, Sweden.
- Christian Hardmeier, Sara Stymne, Jorg Tiedemann, and Joakim Nivre. 2012. *Docent: A Document-Level Decoder for Phrase-Based Statistical Machine Translation*. In: 51st Annual Meeting of the ACL. Sofia, Bulgaria, pages 193-198.
- Christian Hardmeier. 2014. *Discourse in Statistical Machine Translation*. Elanders Sverige, Sweden.
- Thomas Meyer and Andrei Popescu-Belis. 2012. *Using sense-labelled discourse connectives for statistical machine translation*. In: EACL Joint Workshop on Exploiting Synergies between IR and MT, and Hybrid Approaches to MT (ESIRMTHyTra), pages 129-138. Avignon, France.
- Jane Morris and Graeme Hirst. March 1991. *Lexical Cohesion Computed by Thesaural Relations as an Indicator of the Structure of Text*. Computational Linguistics, 17(1):Pages 21-48.
- Joseph Olive, Caitlin Christianson, and John McCary (editors). 2011. *Handbook of Natural Language Processing and Machine Translation: DARPA Global Autonomous Language Exploitation*. Springer Science and Business Media, New York.
- Michael Paul. 2009. *Overview of the IWSLT 2009 evaluation campaign*. In Proceedings of IWSLT.

- Emily Pitler and Ani Nenkova. 2009. *Using Syntax to Disambiguate Explicit Discourse Connectives in Text*. In: ACL-IJCNLP 2009 (47th Annual Meeting of the ACL and 4th International Joint Conference on NLP of the AFNLP), Short Papers, pages 13-16, Singapore.
- Claudia Ross and Jing-heng Sheng Ma. 2006. *Modern Mandarin Chinese Grammar: A Practical Guide*. Routledge, London.
- Avneesh Saluja, Chris Dyer, and Shay B. Cohen. 2014. *Latent-Variable Synchronous CFGs for Hierarchical Translation*. In: Empirical methods in Natural language processing (EMNLP), pages 1953-1964 Doha, Qatar.
- David Steele and Lucia Specia. 2014. *Divergences in the Usage of Discourse Markers in English and Mandarin Chinese*. In: Text, Speech and Dialogue (17th International Conference TSD), pages 189-200, Brno, Czech Republic.
- Toshiyuki Takezawa, Eiichiro Sumita, Fumiaki Sugaya, Hirofumi Yamamoto, and Seiichi Yamamoto. 2002. *Toward a Broad-coverage Bilingual Corpus for Speech Translation of Travel Conversations in the Real World*. In: LREC, pages 147-152. Las Palmas, Spain.
- Mei Tu, Yu Zhou and Chengqing Zong. 2014. *Enhancing Grammatical Cohesion: Generating Transitional Expressions for SMT*. In: 52nd annual meeting of the ACL, June 23-25, Baltimore, USA.
- Billy T.M. Wong and Chunyu Kit. 2012. *Extending Machine Translation Evaluation Metrics with Lexical Cohesion to Document Level*. In: 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, pages 1060-1068. Jeju Island, Korea.
- Tong Xiao, Jingbo Zhu, Shujie Yao, and Hao Zhang. September 2011. *Document-level Consistency Verification in Machine Translation*. In 2011 MT summit XIII, pages 131-138. Xiamen, China:
- Deyi Xiong., Guosheng Ben, Min Zhang, Yajuan Lu, and Qun Liu. August 2013. *Modelling Lexical Cohesion for Document-level Machine Translation*. In: Twenty-Third International Joint Conference on Artificial Intelligence (IJCAI-13) Beijing, China.
- Deyi Xiong, Yang Ding, Min Zhang, and Chew Lim Tan. 2013. *Lexical Chain Based Cohesion Models for Document-Level Statistical Machine Translation*. In: 2013 Conference on Empirical Methods in Natural Language Processing, pages: 1563-1573.
- Jinxi Xu and Roger Bock. 2011. *Combination of Alternative Word Segmentations for Chinese Machine Translation*. DARPA Global Autonomous Language Exploitation. Springer Science and Business Media, New York.
- Nianwen Xue. 2005. *Annotating Discourse Connectives in the Chinese Treebank*. In: ACL Workshop on Frontiers in Corpus Annotation 2: Pie in the Sky.
- Frances Yung. 2014. *Towards a Discourse Relation-aware Approach for Chinese-English Machine Translation*. In: ACL Student Research Workshop, pages 18-25. Baltimore, Maryland USA.