

Semantic Annotation for Interlingual Representation of Multilingual Texts

Teruko Mitamura¹, Keith Miller², Bonnie Dorr³, David Farwell⁴,
Nizar Habash³, Stephen Helmreich⁴, Eduard Hovy⁵, Lori Levin¹, Owen Rambow⁶,
Florence Reeder², Advait Siddharthan⁶

¹Carnegie Mellon University {teruko,ls1}@cs.cmu.edu, ²MITRE Corporation {keith,freeder}@mitre.org,
³University of Maryland {bonnie,nizar}@umiacs.umd.edu, ⁴New Mexico State University
{david,shelmrei}@crl.nmsu.edu, ⁵University of Southern California, <hovy@isi.edu>,
⁶Columbia University {rambow, as372}@cs.columbia.edu

Abstract

This paper describes the annotation process being used in a multi-site project to create six sizable bilingual parallel corpora annotated with a consistent interlingua representation. After presenting the background and objectives of the effort, we describe the multilingual corpora and the three stages of interlingual representation being developed. We then focus on the annotation process itself, including an interface environment that supports the annotation task, and the methodology for evaluating the interlingua representation. Finally, we discuss some issues encountered during the annotation tasks. The resulting annotated multilingual corpora will be useful for a wide range of natural language processing research tasks, including machine translation, question answering, text summarization, and information extraction.

1 Introduction

An interlingua is a semantic representation which mediates between source and target languages in interlingua-based machine translation. It is designed to capture the meaning of a sentence that is common to both source and target languages. If a system supports multi-language translation, the design of the interlingua becomes more complex, due to the number of languages represented. Even though the aim of an interlingua is to capture language-independent semantic expressions, it is difficult to design an interlingua that covers all known languages, and there is no universally acceptable interlingua representation currently in existence. In practice, researchers have designed interlingua representations for particular sets of languages, in order to cover the necessary set of semantic expressions for machine translation (Mitamura et al. 1991). More recently, the use of interlingua representations has been extended beyond machine translation to include, for example, applications for question answering (Ogden et al., 1999), representing agent actions (Kipper & Palmer, 2000) and knowledge acquisition from text (Nyberg et al. 2002).

In September 2003, researchers from six sites began a project titled “Interlingual Annotation of Multilingual Corpora” (IAMTC)¹, funded by the National Science Foundation. This project focuses on the creation of a semantic representation system, followed by the development of six semantically-annotated bilingual corpora. The bilingual corpora pair English texts with corresponding text in Japanese, Spanish, Arabic, Hindi, French, and Korean. The

semantically annotated corpora will be useful not only for machine translation development, but also for research in question answering, text summarization and information retrieval. The project participants include the Computing Research Laboratory at NMSU, the Language Technologies Institute at CMU, the Information Sciences Institute at USC, UMIACS at the University of Maryland, the MITRE Corporation, and Columbia University.

In this paper, we first present the objectives of the IAMTC project. We then provide background information on the multilingual corpora and the three stages of interlingual representation being developed. We then focus on the annotation process itself, including a description of an interface environment that supports the annotation task, and a discussion of the evaluation methodology. We conclude with a summary of the current status of the project, and discuss some issues encountered during the annotation tasks.

2 Project Goals

The IAMTC project has the following goals:

- Development of an interlingua representation framework based on a careful study of text corpora in six languages and their translations into English.
- Development of a methodology for accurately and consistently assigning such representations to texts across languages and across annotators.
- Annotation of a corpus of six multilingual parallel subcorpora, using the agreed-upon interlingual representation.
- Development of semantic annotation tools which serve to facilitate more rapid annotation of texts.
- Design of new metrics and evaluations for the interlingual representations, in order to evaluate the

¹ <http://aitc.aitcnet.org/nsf/iamtc/>

degree of annotator agreement and the granularity of meaning representation.

3 Corpus

The data set consists of 6 bilingual parallel corpora. Each corpus is made up of 125 source language news articles along with three independently produced translations into English. (The source news articles for each individual language corpus are different from the source articles in the other language corpora.) The source languages are Japanese, Korean, Hindi, Arabic, French and Spanish. Typically, each article contains between 300 and 400 words (or the equivalent) and thus each corpus has between 150,000 and 200,000 words. Consequently, the size of the entire data set is around 1,000,000 words. The Spanish, French, and Japanese corpora are based on the DARPA MT evaluation data (White and O'Connell 1994). The Arabic corpus is based on LDC's Multiple Translation Arabic, Part 1 (Walker et al., 2003).

For any given subcorpus, the annotation effort is to assign interlingual content to a set of 4 parallel texts (one in the original source language, plus 3 translations to English by different translators), all of which theoretically communicate the same information. A multilingual parallel data set of source language texts and English translations offers a unique perspective and unique problem for annotating texts for meaning.

4 Interlingua

The interlingual representation comprises three levels and incorporates knowledge sources such as the Omega ontology (Philpot et al., 2003) and theta grids (Dorr, 2001). The three levels of representation are referred to as *IL0*, *IL1* and *IL2*. The aim is to perform the annotation process incrementally, with each level of representation incorporating additional semantic features and removing existing syntactic ones. *IL2* is intended as the interlingual level that abstracts away from (most) syntactic idiosyncrasies of the source language. *IL0* and *IL1* are intermediate representations that are useful stepping stones for annotating at the next level.

4.1 IL0

IL0 is a deep syntactic dependency representation. It includes part-of-speech tags for words and a parse tree that makes explicit the syntactic predicate-argument structure of verbs. The parse tree contains labels referring to deep-syntactic grammatical function (normalized for voice alternations). *IL0* does not contain function words (their contribution is represented as features) or semantically void punctuation. While this representation is purely syntactic, many disambiguation decisions, relative clause and PP attachment for example, have been made, and the presentation abstracts as much as possible from surface-syntactic phenomena. (Thus, our *IL0* is intermediate between the analytical and text ogrammatical levels of the Prague School (Hajic et al 2001).) *IL0* is

constructed by hand-correcting the output of a dependency parser (see section 6), and allows annotators to see how textual units relate syntactically when making semantic judgments. Thus, it is a useful starting point for semantic annotation at *IL1*.

4.2 IL1

IL1 is an intermediate semantic representation. It associates semantic concepts with lexical units like nouns, adjectives, adverbs and verbs. It also replaces the syntactic relations in *IL0*, like *subject* and *object*, with thematic roles, like *agent*, *theme* and *goal*. Thus, like PropBank (Kingsbury et al 2002), *IL1* neutralizes different alternations for argument realization. However, *IL1* is not an interlingua; it does not normalize over all linguistic realizations of the same semantics. In particular, it does not address how the meanings of individual lexical units combine to form the meaning of a phrase or clause. It also does not address idioms, metaphors and other non-literal uses of language. Further, *IL1* does not assign semantic features to prepositions; these continue to be encoded as syntactic features of their objects, which may be annotated with thematic roles such as *location* or *time*.

4.3 IL2

IL2 is intended to be an interlingua, a representation of meaning that is (reasonably) independent of language. *IL2* is intended to capture similarities in meaning across languages and across different lexical/syntactic realizations within a language. For example, like FrameNet (Baker et al 1998), *IL2* is expected to normalize over conversives (e.g. X bought a book from Y vs. Y sold a book to X) and also over non-literal language usage (e.g. X started its business vs. X opened its doors to customers). The exact definition of *IL2* is the major research contribution of this project. However, it is important to note that even at the level of *IL2*, it does not include more complex linguistics phenomena, such as speech acts, discourse analysis and pragmatics.

4.4 The Omega Ontology

In progressing from *IL0* to *IL1*, annotators select semantic terms (concepts) to represent the nouns, verbs, adjectives, and adverbs present in each sentence. These terms are represented in the 110,000-node Omega ontology (Philpot et al., 2003), under construction at ISI. Omega has been built semi-automatically from a variety of sources, including Princeton's WordNet (Fellbaum, 1998), New Mexico State University's Mikrokosmos (Mahesh and Nirenburg, 1995), ISI's Upper Model (Bateman et al., 1989) and ISI's SENSUS (Knight and Luk, 1994). The ontology, which has been used in several projects in recent years (Hovy et al., 2001), can be browsed using the DINO browser at <http://blombos.isi.edu:8000/dino>; this browser forms a part of the annotation environment. Omega continues to be developed and extended.

4.5 The Theta Grids

Each verb in Omega is assigned one or more theta grids specifying the theta roles of arguments associated with that verb. Theta roles are abstractions of deep semantic relations that generalize over verb classes. They are by far the most common approach in the field to represent predicate-argument structure. However, there are numerous variant theories with little agreement even on terminology (Fillmore, 1968; Stowell, 1981; Jackendoff, 1972; Levin and Rappaport-Hovav, 1998).

The theta grids used in our project were extracted from the Lexical Conceptual Structure Verb Database (LVD) (Dorr, 2001). The "LCS Database" contains Lexical conceptual Structures built by hand, organized into semantic classes that are a reformulated version of those in Beth Levin (1993) English Verb Classes and Alternations (EVCA), Part 2. The WordNet senses assigned to each entry in the LVD link the theta grids to the verbs in the Omega ontology. In addition to the theta roles, the theta grids specify syntactic realization information, such as Subject, Object or Prepositional Phrase, and the Obligatory/Optional nature of the argument. The set of theta roles used, although based on research in LCS-based MT (Dorr, 1993; Habash et al, 2002), has been simplified for this project.

5 Annotation Tools

We have assembled a suite of tools to be used in the annotation process. Since we are gathering our corpora from disparate sources, we need to standardize the text before presenting it to automated procedures. For English, this involves sentence boundary detection, but for other languages, it may involve segmentation, chunking of text, or other operations. The text is then processed with a dependency parser, the output of which is viewed and corrected in TrED (Hajic, et al., 2001), a graphically-based tree editing program, written in Perl/Tk2. The revised deep dependency structure produced by this process is the IL0 representation for that sentence.

To create IL1 from the IL0 representation, annotators use Tiamat, a tool developed specifically for this project. This tool enables viewing of the IL0 tree with easy reference to all of the IL resources described in section 4 (current IL representation, ontology, and theta grids). Tiamat provides the ability to annotate text via simple point-and-click selections of words, concepts, and theta-roles. The IL0 is displayed in the top left pane, ontological concepts and their associated theta grids, if applicable, are located in the top right, and the sentence itself is located in the bottom right pane. An annotator may select a lexical item (leaf node) to be annotated in the sentence view; this word is highlighted, and the relevant portion of the Omega ontology is displayed in the pane on the left. In addition, if this word has dependents, they are automatically underlined in red in the sentence view. Annotators can view all information pertinent to the

process of deciding on appropriate ontological concepts in this view. Following the procedures described in section 6, selection of concepts, theta grids and roles appropriate to that lexical item can then be made in the appropriate panes.

In order to evaluate the annotators' output, an evaluation tool was also developed to compare the output and to generate the evaluation measures that are described in section 7. The reports generated by the evaluation tool allow the researchers to look at both gross-level phenomena, such as inter-annotator agreement, and at more detailed points of interest, such as lexical items on which agreement was particularly low, possibly indicating gaps or other inconsistencies in the ontology.

6 Annotation Manuals and Process

To describe the annotation task, we first present the annotation manuals and then discuss the annotation process.

6.1 Annotation Manual

We have been developing markup instructions which comprise three manuals: a users' guide for Tiamat (including procedural instructions), a definitional guide to semantic roles, and a manual for creating a dependency structure (IL0). Together these manuals allow the annotator to understand (1) the intention behind aspects of the dependency structure; (2) how to use Tiamat to mark up texts; and (3) how to determine appropriate semantic roles and ontological concepts. In choosing a set of appropriate ontological concepts, annotators were encouraged to look at the name of the concept and its definition, the name and definition of the parent node, example sentences, lexical synonyms attached to the same node, and sub- and super-classes of the node.

6.2 Annotation process

For the initial testing period, only English texts were annotated, and the process described here is for English text. We assume that the process for non-English texts would be the same with a minor modification as needed.

Each sentence of the text is parsed into a dependency tree structure. For English texts, these trees were first provided by the Connexor parser (Tapanainen and Jarvinen, 1997), and then corrected by one of the team PIs. Then the corrected dependency structures (IL0) are provided to annotators.

The annotators were instructed to annotate all nouns, verbs, adjectives, and adverbs. This involves choosing all relevant concepts from Omega – both concepts from Wordnet SYNSETs and those from Mikrokosmos; these sources of information are intertwined in Omega. One of the goals and results of this annotation process will be a simultaneous coding of concepts in both ontologies, facilitating a closer union between them.

In addition, annotators were instructed to provide a semantic case role for each dependent of a

² http://quest.ms.mff.cuni.cz/pdt/Tools/Tree_Editors/TrEd/

verb. LCS verbs were identified with Wordnet classes and the LCS case frames were supplied where possible. The annotator, however, was often required to determine the set of roles or alter them to suit the text. In both cases, the revised or new set of case roles was noted and sent to a reviewer for evaluation and possible permanent inclusion. Thus the set of event concepts in the ontology supplied with roles will grow through the course of the project.

For the initial testing phase of the project, all annotators at all sites worked on the same texts. We have two annotators from each site. Each site, which has different source language texts, provided two texts that were translated into English by two different translators. To test for the effects of coding two texts that are semantically close (since they are both translations of the same source document), the order in which the texts were annotated differed from site to site. Half of the sites marked one translation first, and the other half of the sites marked the second translation first. Another variant tested was to interleave the two translations, so that two similar sentences were coded consecutively.

In the period leading up to the initial test phase, weekly conversations were held at each site by the annotators to review the coded texts. This was followed by a weekly conference call among all the annotators. During the test phase, no discussion was permitted until all the annotation tasks were completed.

7 Evaluation Methodology

We have identified several metrics for evaluation of intercoder agreement on annotations. We are currently measuring intercoder agreement on concept names selected from the Omega ontology and thematic role labels.

Two measures of intercoder agreement are currently used, Kappa (Carletta, 1993) and a Wood Standard similarity (Habash and Dorr, 2002). For expected agreement in the Kappa statistic, $P(E)$ is defined as $1/(N+1)$ where N is the number of choices at a given data point. In the case of Omega nodes, this means the number of matched Omega nodes (by string match) plus one for the possibility of the annotator traversing up or down the hierarchy. The Wood Standard is the category chosen by the most annotators. In cases of no agreement, a random selection is picked from the annotator's selections. Multiple measures were used because it is important to have a mechanism for evaluating inter-coder consistency in the use of the IL representation language which does not depend on the assumption that there is a single correct annotation of a given text.

In addition to intercoder agreement, we are also developing metrics for evaluating the quality of an annotated interlingua. Given the project goal of generating an IL representation which is useful for MT (among other NLP tasks), we measure the ability to generate accurate surface texts from the IL representation as annotated. At this stage, we plan to use an available generator, Halogen (Knight and Langkilde, 2000). A tool to convert the representation

to meet Halogen's requirements is being built. Following the conversion, surface forms will be generated and then compared with the originals through a variety of standard MT metrics (ISLE, 2003). This will serve to determine whether the elements of the representation language are sufficiently well-defined and whether they can serve as a basis for inferring interpretations from semantic representations or (target) semantic representations from interpretations.

8 Annotation Issues

During the test phase, we annotated 144 texts, which come from 2 translations of 6 source texts annotated by 2 annotators in each 6 sites.

A preliminary investigation of intercoder agreement on multiple annotations shows that the more annotators learn the process, the better they become, resulting in an improvement of intercoder agreement. We made two assumptions regarding the training of novice annotators in order to improve intercoder agreement. One assumption is that novice annotators may make inconsistent annotations within the same text. In order to train annotators, we have developed an intra-annotator consistency checking procedure. After the annotators finished an initial annotation pass, they were asked to go over their results to see if there were any inconsistencies within the text. For example, if two nodes in different sentences are co-indexed, then annotators must ensure that the two nodes carry the same meaning in the context of the two different sentences.

Another assumption we made was that if two annotators at the same site discuss their annotation results after their annotation tasks are completed, they can learn more from each other. Under this assumption, we have developed inter-annotator a reconciliation procedure and a voting tool associated with this process. There are three steps to follow. First, we created a combined annotation file, in which disagreements are marked in red. Each annotator votes privately either Yes, Possible, or No for items marked in red. In the second step, annotators get together and discuss the differences. After the open discussion, they vote again privately. We are currently in the process of analyzing the effect of inconsistency checking and inter-annotator reconciliation on overall intercoder agreement.

During the inter-annotator reconciliation process, we have encountered a number of difficult issues. One issue is the granularity of concept selection. The Omega ontology, which is derived from WordNet, contains 110,000 nodes and often provides too many alternatives, whereas Omega-Mikrokosmos, which contains only 6,000 concepts, does not offer all the concepts needed for annotation. For example, the word extremely contains 4 concepts in Omega's WordNet, and each of the senses is hard to distinguish from the others: (1) to a high degree or extent; favorably or with much respect, (2) to an extreme degree, (3) to an extreme degree, super, (4) to an extreme degree or extent, exceedingly. On the other

hand, Omega-Mikrokosmos does not contain a concept for the word extremely.

In the coming months we will be pruning out the extraneous terms from Omega, fleshing out the current procedures for evaluating the accuracy of an annotation and measuring the inter-coder agreement. We will also be working on IL2 design and annotation. Finally, a growing corpus of annotated texts at each stage (IL0, IL1, IL2) will become available.

Additional issues to be addressed include: (1) personal name, temporal and spatial annotation (e.g., Ferro et al., 2001); (2) causality, co-reference, aspectual content, modality, speech acts, etc; (3) reducing vagueness and redundancy in the annotation language; (4) inter-event relations such as entity reference, time reference, place reference, causal relationships, associative relationships, etc; Finally, to incorporate these, cross-sentence phenomena remain a challenge.

From an MT perspective, issues include evaluating consistency in the use of an annotation language, given that any source text can result in multiple, different, legitimate translations (Farwell and Helmreich, 2003). Along these lines, there is the problem of annotating texts for translation without including in the annotations inferences from the source text.

8 Conclusion

The IAMTC project is radically different from those annotation projects that have focused on morphology, syntax or even certain types of semantic content (e.g., for word sense disambiguation). It is most similar to PropBank (Kingsbury et al 2002) and FrameNet (Baker et al 1998). However, our project is novel in its emphasis on: (1) a more abstract level of mark-up (interpretation); (2) the assignment of a well-defined meaning representation to concrete texts; and (3) issues of a community-wide consistent and accurate annotation of meaning.

By providing an essential, and heretofore non-existent, data set for training and evaluating natural language processing systems, the resultant annotated multilingual corpus of translations is expected to lead to significant research and development opportunities for machine translation and a host of other natural language processing technologies, including question answering (e.g., via paraphrase and entailment relations) and information extraction. Because of the unique annotation processes in which the each stage (IL0, IL1, IL2) provides a different level of linguistic and semantic information, a different type of natural language processing can take advantage of the information provided at the different stages. For example, IL1 may be useful for information extraction in question answering, whereas IL2 might be the level that is of most benefit to machine translation. These topics exemplify the research investigations that we can conduct in the future, based on the results of the annotation.

References

- Baker, Collin and J. Fillmore and John B. Lowe, (1998). The Berkeley FrameNet Project. Proceedings of ACL.
- Bateman, J.A., Kasper, R.T., Moore, J.D., and Whitney, R.A. (1989). A General Organization of Knowledge for Natural Language Processing: The Penman Upper Model. Unpublished research report, USC/Information Sciences Institute, Marina del Rey, CA.
- Carletta, J. C. (1996). Assessing agreement on classification tasks: the kappa statistic. *Computational Linguistics*, 22(2), 249-254
- Dorr, Bonnie J. (2001). LCS Verb Database, Online Software Database of Lexical Conceptual Structures and Documentation, University of Maryland. http://www.umiacs.umd.edu/~bonnie/LCS_Database_Documentation.html
- Dorr, Bonnie J., (1993). *Machine Translation: A View from the Lexicon*, MIT Press, Cambridge, MA.
- Farwell, David, and Steve Helmreich. (2003). Pragmatics-based Translation and MT Evaluation. In Proceedings of Towards Systematizing MT Evaluation. MT-Summit Workshop, New Orleans, LA.
- Fellbaum, C. (ed.). (1998). *WordNet: An On-line Lexical Database and Some of its Applications*. MIT Press, Cambridge, MA.
- Ferro, Lisa, Inderjeet Mani, Beth Sundheim and George Wilson. (2001). TIDES Temporal Annotation Guidelines. Version 1.0.2 MITRE Technical Report, MTR 01W0000041
- Fillmore, Charles. (1968). The Case for Case. In E. Bach and R. Harms, editors, *Universals in Linguistic Theory*, pages 1-88. Holt, Rinehart, and Winston.
- Fleischman, M., A. Echihabi, and E.H. Hovy. (2003). Offline Strategies for Online Question Answering: Answering Questions Before They Are Asked. Proceedings of the ACL Conference. Sapporo, Japan.
- Habash, Nizar and Bonnie Dorr. (2002). Interlingua Annotation Experiment Results. AMTA-2002 Interlingua Reliability Workshop. Tiburon, California, USA.
- Habash, Nizar, Bonnie J. Dorr, and David Traum, (2002). "Efficient Language Independent Generation from Lexical Conceptual Structures," *Machine Translation*, 17:4.
- Hajic, Jan; Vidová-Hladká, Barbora; Pajas, Petr. (2001): The Prague Dependency Treebank: Annotation Structure and Support. In Proceeding of the IRCS Workshop on Linguistic Databases, pp. .

- University of Pennsylvania, Philadelphia, USA, pp. 105-114.
- Hovy, E.H., Philpot, A., Ambite, J.L., Arens, Y., Klavans, J., Bourne, W., and Saroz, D. (2001). Data Acquisition and Integration in the DGRC's Energy Data Collection Project, in Proceedings of the NSF's dg.o 2001. Los Angeles, CA.
- Jackendoff, Ray. (1972). Grammatical Relations and Functional Structure. *Semantic Interpretation in Generative Grammar*. The MIT Press, Cambridge, MA.
- Kingsbury, Paul and Martha Palmer and Mitch Marcus, (2002). Adding Semantic Annotation to the Penn TreeBank. Proceedings of the Human Language Technology Conference (HLT 2002).
- Kipper, Karin and Martha Palmer (2000). Representation of Actions as an Interlingua. Proceedings of the Third AMTA SIG-IL Workshop on Interlinguas and Interlingual Approaches, Seattle, WA, April 30.
- Knight, K., and I. Langkilde. (2000). Preserving Ambiguities in Generation via Automata Intersection American Association for Artificial Intelligence conference (AAAI).
- Knight, K, and Luk, S.K. (1994). Building a Large-Scale Knowledge Base for Machine Translation. Proceedings of AAAI. Seattle, WA.
- Levin, Beth. (1993) "English Verb Classes and Alternations: A Preliminary Investigation", University of Chicago Press, Chicago, IL.
- Levin, B. & Rappaport-Hovav, M. (1998). From Lexical Semantics to Argument Realization. Borer, H. (ed.) *Handbook of Morphosyntax and Argument Structure*. Dordrecht: Kluwer Academic Publishers.
- Mahesh, K., and Nirenberg, S. (1995). A Situated Ontology for Practical NLP, in Proceedings of the Workshop on Basic Ontological Issues in Knowledge Sharing at IJCAI-95. Montreal, Canada.
- Mitamura, T., E. Nyberg, J. Carbonell. (1991). An Efficient Interlingua Translation System for Multilingual Document Production, in Proceedings of the Third Machine Translation Summit. Washington, DC.
- Nyberg, E., T. Mitamura, K. Baker, D. Svoboda, B. Peterson, J. Williams. (2002) Deriving Semantic Knowledge from Descriptive Texts using an MT System. Proceeding of the 2002 Conference, Association for Machine Translation in the Americas.
- Ogden, B., J. Cowie, E. Ludovik, H. Molina-Salgado, S. Nirenburg, N. Sharples and S. Sheremtyeva (1999). CRL's TREC-8 Systems: Cross-Lingual IR and Q&A, Proceedings of the Eighth Text Retrieval Conference (TREC -8).
- Philpot, A., M. Fleischman, E.H. Hovy. (2003). Semi-Automatic Construction of a General Purpose Ontology. Proceedings of the International Lisp Conference. New York, NY. Invited.
- Stowell, T. (1981). Origins of Phrase Structure. *PhD thesis*, MIT, Cambridge, MA.
- Tapanainen, P. and T. Jarvinen. (1997). A non-projective dependency parser. In the 5th Conference on Applied Natural Language Processing / Association for Computational Linguistics, Washington, DC.
- White, J., and T. O'Connell. (1994). The ARPA MT evaluation methodologies: evolution, lessons, and future approaches. Proceedings of the 1994 Conference, Association for Machine Translation in the Americas
Walker, Kevin, Moussa Bamba, David Miller, Xiaoyi Ma, Chris Cieri, and George Doddington (2003). Multiple-Translation Arabic Corpus, Part 1. Linguistic Data Consortium (LDC) catalog number LDC2003T18 and ISBN 1-58563-276-7.