

Evaluating Lemmatic Communication

Katherine Everitt, Christopher Lim, Oren Etzioni, Jonathan Pool, Stephen Soderland

Turing Center, University of Washington

Seattle, Washington 98195

{everitt, chrislim, etzioni, pool, soderlan}@cs.washington.edu

ABSTRACT

The need for translation among the world's thousands of natural languages makes information access and communication costly. One possible solution is lemmatic communication: A human sender encodes a message into sequences of lemmata (dictionary words), a massively multilingual lexical translation engine translates them into lemma sequences in a target language, and a human receiver interprets them to infer the sender's intended meanings. Using a 13-million-lemma, 1300-language translation engine, we conducted an experiment in lemmatic communication with Spanish- and Hungarian-speaking subjects. Translingual communication was less successful than intralingual communication, and intralingual communication was less successful when the lemma sequences were artificially randomized before the receiver saw them (simulating word-order differences among languages). In all conditions, however, meanings were transmitted with high or moderate fidelity in at least 40% of the cases. The results suggest interface and translation-algorithm improvements that could increase the efficacy of lemmatic communication.

Keywords

Lexical translation, lemmatic, communication, user studies

INTRODUCTION

The Internet has greatly expanded the ability to share information, enabling communication between physically and culturally distant people. However, there are over 6000 living languages [1], and the need to translate makes communication expensive even when distance is no longer an obstacle. Attempts to make translation inexpensive by automating it have been only partially successful, and they have ignored 99% of the world's languages. For example, the popular Google Translate application covers only 35 languages [3].

If people communicated using only lemmata (words and phrases in their citation, or dictionary, forms), automatic translation would be greatly simplified, permitting translation among thousands of languages. By combining existing resources (bilingual and multilingual dictionaries, thesauri, and glossaries), one could build a system that infers translations of arbitrary lemmata into arbitrary target languages.

In this paper, we evaluate such a system of translingual lemmatic communication. Senders encode message sentences into sequences of lemmata; these are automatically translated; and receivers attempt to decode

the translated sequences of lemmata into sentences that reflect the meanings intended by the senders.

LEMMATIC COMMUNICATION

The lemmatic communication process consists of three steps: encoding, translation, and decoding. Encoding and decoding are done by the sender and receiver respectively, and the translation is done by the system using an automatically constructed translation graph (TransGraph) [2].

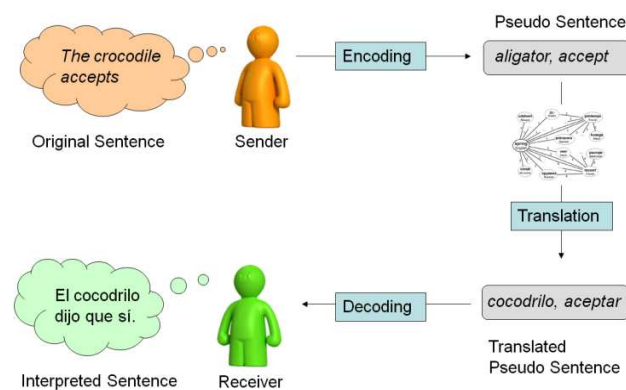


Figure 1. Lemmatic Communication example

Encoding

In the encoding process, the sender selects lemmata to convey a statement or question. For example, the sender might encode "A couple of previous guests recommended your hotel to us." as "two, previous, guest, recommend, hotel". Using autocomplete lists, the system permits the sender to select only lemmata that can be translated automatically into the target language.

Translation

Translation is performed by TransGraph [2], a translation engine based on a graph constructed from about 600 machine-readable lexical resources. The graph currently contains 13 million expressions in 1300 languages; 10 million senses; and 27 million edges, each edge connecting a lemma to a sense. In the translation process, the system translates each lemma into a lemma in the target language and assembles the translated lemmata into sequences corresponding to the original sequences.

When one or more direct translations in the target language exist, the system translates the lemma into one of those. Otherwise, the system infers a translation from paths through intermediate translations. In each case, the system

estimates the probability that each candidate translation is a correct one and selects the candidate with the greatest probability.

Decoding

In the decoding process, the receiver reads sequences of lemmata and attempts to infer the intended meaning of each sequence. For example, the receiver might read “my, home, inside, three, sleep, room, exist” and infer that the intended meaning is “There are three bedrooms in my home.”

USER STUDY

We evaluated a system of translingual lemmatic communication to determine whether such communication can succeed and, if so, what conditions promote success..

Our questions were:

- How satisfying is the process of lemmatic communication to its participants?
- How long do encoding and decoding take?
- How much of the intended meanings is conveyed in lemmatic communication?
- Does the order of lemmata in a sequence convey useful information to the decoder?
- Is lemmatic communication less successful when the lemmata are translated than when they remain in the original language?

We performed the study with two languages, Hungarian and Spanish.

ENCODING PHASE

We created a set of three scenarios, each described with a series of ten sentences. The sentences were written in English, professionally translated from English into Spanish and Hungarian, and checked by bilingual translators.

Subjects in the encoding phase converted each sentence into a sequence of lemmata using our online encoding system. In addition to gathering encoded sequences for a later phase of the study, the purposes of the encoding phase were to get qualitative feedback from the encoders about the process and to get information (length, specificity, etc.) about likely encodings.

The scenarios were:

Visit: Visiting a city and booking a hotel

Fable: The Monkey and the Crocodile

Book Group: Message about a book group

Figure 3 shows the online encoding interface. It was written using .NET aspx pages and the jQuery JavaScript library [4] for dropdown functionality, with data stored in Microsoft SQL Server 2005. When encoders type two letters, a drop down box appears showing the permitted lemmata. There were 18,139 permitted lemmata in Spanish and 24,482 in Hungarian. These were the lemmata in

Hungarian that TransGraph could translate into Spanish and vice versa. If the encoder typed an incorrect string, the box would turn red and disallow it, as seen at the bottom of Figure 3.

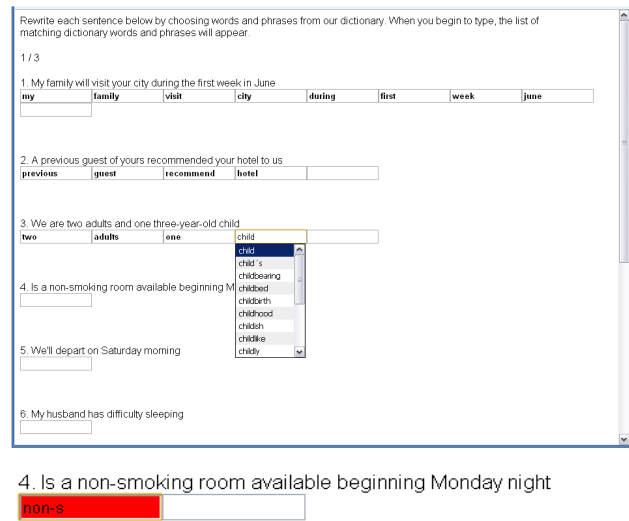


Figure 3: The online encoding interface

There were two Hungarian and two Spanish encoders. Encoders took between 6 and 50 minutes to encode each page, with a mean time of 17 minutes. The mean encoding time for a Spanish page was 9 minutes, versus 24 minutes for a Hungarian page. Because of the small sample size, encoding time may be skewed by one participant. However, the encoding times for Hungarian were always longer than for Spanish, perhaps because of difficulties with characters not in the standard online alphabet. [What does this mean?]

The mean encoded sequence length (in lemmata) was 1.17 more than the mean original sentence length (in words). Of all the sequences created, 68% were longer than, 17 % were equally long as, and 15% were shorter than their source sentences.

Encoding Feedback

The instructions given were deliberately imprecise, in order to explore people’s natural inclinations. The instruction was “Rewrite each sentence below by choosing words and phrases from our dictionary.” We also gave example encodings. From participant comments, we learned that participants often felt they needed to encode every word of the sentence. We also found that they wanted a way to encode information that is not available in the list, such as exclamations, questions and verb tense. In our pilot tests, which were conducted in English, there was excellent coverage of lemmata (40,957), and so participants were surprised when a specific lemma they wanted to use was not in the list. Also, participants were not very aware of the phrasal lemmata and occasionally had to go back and change multiple words into a corresponding phrase.

Encoders expressed some frustration with our list-constrained approach. One criticism was that a space does not move to the next box, so a tab or click is necessary.

This is a necessary feature because phrases require a space to be typed without moving to the next box. Another criticism was that we required people to immediately correct their mistakes.

Encoding Guidelines

We present a series of encoding guidelines based on our encoders' experience.

There is a tradeoff between allowing users to type any words they want in a traditional text-editor format and using a list-constrained approach. While the list setup is limiting, allowing people to type anything may cause a frustrating system response demanding changes to the lemmata that cannot be translated. Potentially, the most appropriate long-term solution is a combination of the two, where people are allowed to type what they wish but receive immediate feedback, such as a colored line under untranslatable lemmata with accompanying suggestions of alternatives. Another useful addition would be to automatically detect and combine phrasal lemmata or give clearer hints about their existence.

Given the requests of participants, the system should allow for the encoding of metadata describing properties of lemmata (e.g., tense) and sequences (e.g., exclamation, question). Observation also suggests that it would be beneficial to encourage shorter encodings and to let encoders know that they do not need to encode particles.

DECODING PHASE

The purposes of the decoding phase were to get qualitative opinions about the clarity of the lemma sequences and to collect sentences produced by decoders for comparison with the originals. Decoding took place under three conditions: Same, Randomized, and Translated. In the Same condition, the decoder was presented with one of the original encodings. In the Randomized condition, the decoder saw an encoding with the lemmata randomly re-ordered. In the Translated condition, decoders worked on an encoding whose lemmata had been translated from the other language by TransGraph, without any change in the order of the lemmata.

There were 49 decoding participants: 30 Hungarian-speaking and 19 Spanish-speaking. We presented the three scenarios in order: Visit, Fable, and Book Group, and counterbalanced the conditions of Same, Randomized, and Translation. All of the ten sentences within a condition were shown in order.

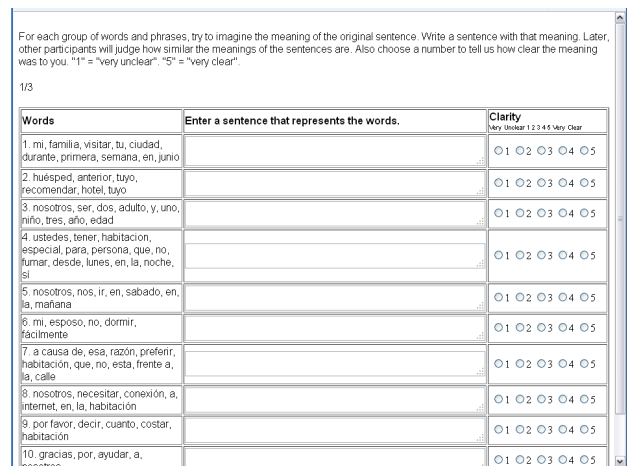


Figure 4: Decoding Interface

Figure 4 shows the decoding interface (instructions have been translated into English here). Decoders expressed their guesses about the sequences' meanings by entering sentences and marked each sequence's subjective clarity on a scale of 1 to 5, where 1 meant very unclear and 5 meant very clear.

Decoding Results

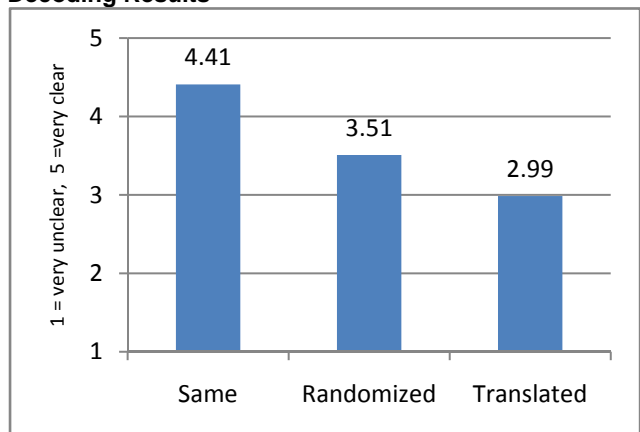


Figure 5: Subjective Clarity by Condition

Figure 5 shows the mean subjective clarity by condition. As one would expect, the Same condition had the highest mean score. We also discovered that translated sequences without randomization were significantly less clear than randomized sequences without translation. Each difference was significant ($p < 0.01$). The mean result for the Translated condition—the one in which lemmatic communication might actually be put to use—was 2.99.

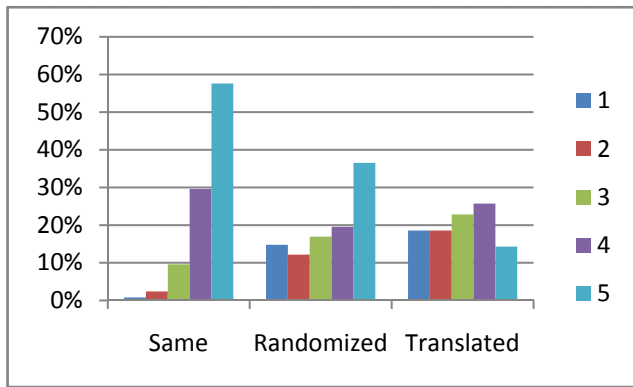


Figure 6: Score distribution by condition

Figure 6 shows the score distributions. For the Same condition, over half of the scores were 4 or 5. In the Randomized condition, over half were 3 or above. The Translate condition scores were fairly uniformly distributed. Almost 90% of sequences in the Same condition received a mean clarity score of 4 or 5, suggesting that lemmatically encoded messages can be understandable under the most favorable conditions. However, these proportions decreased to about 65% in the Randomized condition and 40% in the Translated condition.

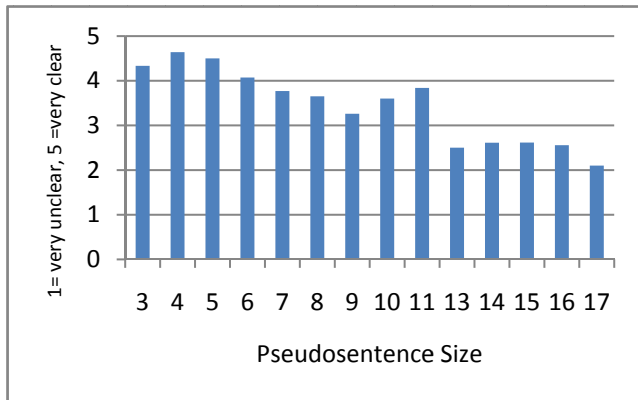


Figure 7. Subjective Clarity by sentence length

Figure 7 shows the subjective clarity by sentence length. Longer sentences, over 11 words in length, had lower subjective clarity. Because longer sentences tend to have more clauses they are more susceptible to reordering effects (in the Randomized and Translate conditions) and mistranslation (in the Translate condition).

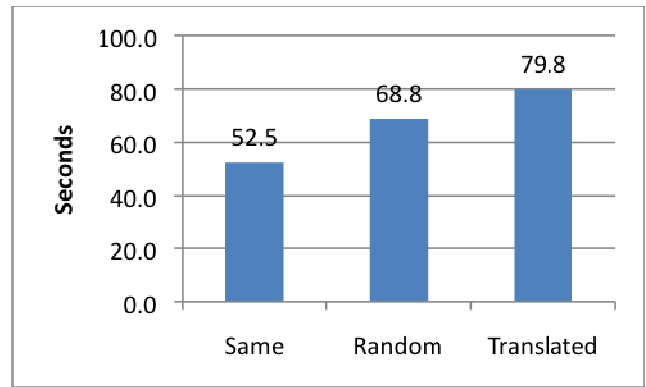


Figure 8: Mean Time to decode a sentence within each condition.

Figure 8 shows the time to decode a sentence within each condition. On average, decoding took about a minute per sentence. The difference between the Same and Randomized conditions is marginally significant ($p=0.067$), and the difference between the Same and Translated conditions is significant ($p < 0.01$).

EVALUATION PHASE

During the evaluation phase, the decoded sentences were compared with the original, professionally translated sentences and evaluated for meaning similarity. There were 10 Spanish-speaking participants and 12 Hungarian-speaking participants in this phase.

Para cada par de oraciones traducidas, seleccione un número para indicarnos qué tan similares son sus significados.
 "1"="Sus significados son iguales", "3"="Sus significados parecen en nada".

1/3 "How similar are the sentences?"
1 = good, 3 = poor

| Oración 1 | Oración 2 | ¿Cómo son similares las sentencias? 1 = más 3 = peor |
|--|--|---|
| 1. Mi familia visitará tu ciudad en la primera semana de junio | Mi familia visitará tu ciudad durante la primera semana de junio. | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| | Voy a visitar a mi familia en tu ciudad durante la primera semana de junio. | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| | Mi ego y mi familia visitaremos luego tu ciudad, pero primero, durante la semana de junio | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| 2. Un anterior huésped suyo anterior nos recomendó su hotel | Tu huésped anterior recomendó tu hotel. | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| | el huésped anterior recomendó que nosotros habitásemos en tu hotel | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| 3. Somos dos adultos y un niño de tres años | Un huésped anterior tuyo me recomendó tu hotel | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| | nosotros somos dos adultos abuelos y el niño es dueño de tres cursos | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| | Nosotros somos dos adultos y un niño con tres años de edad. | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| 4. ¿Tienen una habitación para no fumadores disponible desde la noche del lunes? | Nosotros somos dos adultos y un niño de tres años de edad. | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| | (de verdad no entiendo) | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| 4. ¿Tienen una habitación para no fumadores disponible desde la noche del lunes? | quizá sea que se prohíbe fumar en esta habitación los lunes a la noche, por un ataque de pánico. | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |
| | ¿Tienen ustedes habitación especial para no | <input type="radio"/> 1 <input type="radio"/> 2 <input type="radio"/> 3 |

Figure 9. Evaluation Interface

Figure 9 shows the evaluation interface. Participants were shown an original sentence and one Same, one Randomized, and one Translated version of that sentence, in random order. They were asked to score each output sentence in terms of the similarity of its meaning to the original sentence. A rank of 1 was good and 3 was poor.

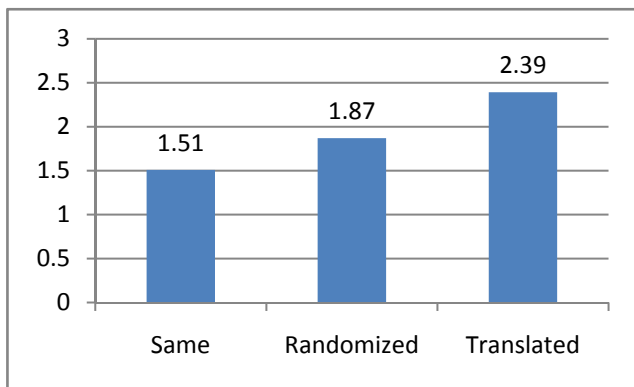


Figure 10: Mean sentence-similarity score by condition.

Figure 10 shows mean score by condition. As before, the best results occurred in the Same condition and the worst in the Translated condition. All differences were statistically significant ($p < 0.01$).

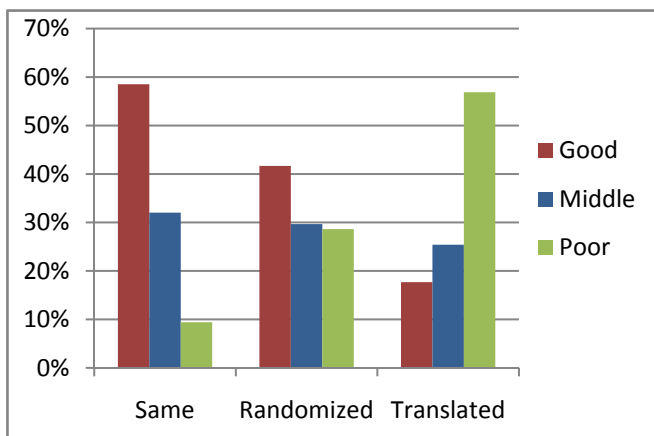


Figure 11. Perceived Sentence Similarity Breakdown by Condition

Figure 11 shows the distribution of perceived sentence similarity by condition. It shows a much higher fraction of good responses in the Same condition than the Translated condition.

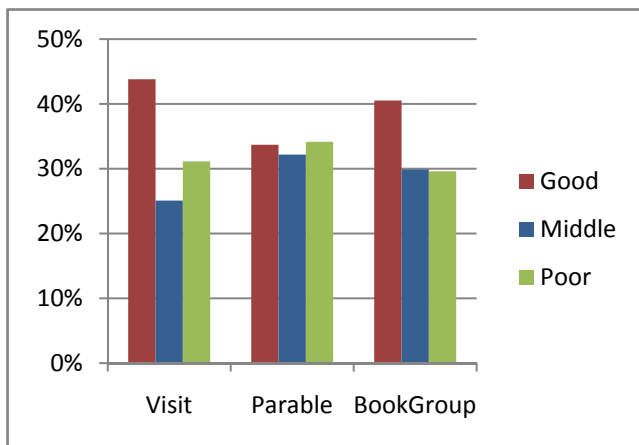


Figure 12. Perceived Sentence Similarity by Scenario

Figure 12 shows the distribution of perceived sentence similarity by scenario. The Visit and Book Group

scenarios had a larger fraction of “good” responses, probably due to the better understanding of context surrounding them. The fable scenario was less familiar and had fewer “good” scores. The pairwise differences in mean score among the Visit, Fable, and Book Group conditions (1.87, 2.00, and 1.89) were all significant ($p < 0.01$).

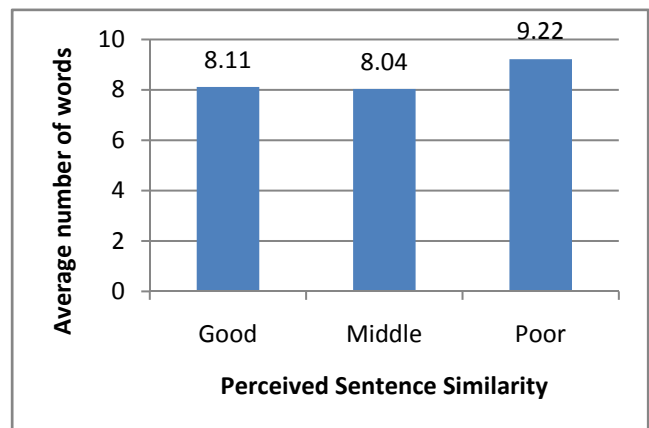


Figure 13. Mean Sentence Length by Perceived Sentence Similarity.

Figure 13 shows the mean sentence length by perceived sentence similarity. Sentences marked Poor were significantly longer, on average, than those marked Good ($p < 0.01$) or Middle ($p < 0.01$).

DISCUSSION

Even with the confusion introduced by the randomization of lemma order and by translation, some successful communication between Hungarian and Spanish speakers occurred in our experiment. The Randomized condition sought to simulate (to an extreme degree) the independent effect of word-order differences among languages, separate from the effect of lemma translation. If this simulation is valid, we have evidence that word-order differences do impair lemmatic communication. Given the diversity of sentential (subject/verb/object) and phrasal (adjective/noun, etc.) word orders among languages and the intuitions that both encoders and decoders might develop to handle lemma ordering, much additional work could be done on the factors that detract from efficacy in lemmatic communication.

The encoding system seems amenable to several improvements: making it faster, with fewer constraints on typing; permitting (and encouraging) encoders to split long sentences into multiple short lemma sequences; and permitting encoders to type freely and then get feedback on the translatability of what they have typed. Further efficacy could result from using more intelligent, context-aware translation; allowing the sender to check tentative translations (e.g., via back-translation feedback); or giving receivers access to multiple translation candidates. Converting our single-pass system to an interactive one, in which receivers can prompt senders for clarifications, might also permit the rapid resolution of linguistic uncertainty and mistranslations.

CONCLUSIONS

We have shown that lemmatic communication can work, but its translation component makes it considerably slower and more error-prone. These results suggest that better interface design, the inclusion of annotation features, more intelligent translation inference, and sender-receiver interactivity could make lemmatic communication effective across thousands of languages.

ACKNOWLEDGMENTS

We thank our participants and our translators. We would also like to thank Susan Colowick for her assistance.

REFERENCES

1. Ethnologue: <http://www.ethnologue.com/>
2. Etzioni, Reiter, Soderland, Sammer, Lexical Translation with Application to Image Search on the Web. *Proceedings of Machine Translation Summit XI, 2007*
3. Google Translate: <http://translate.google.com/>
4. jQuery JavaScript Library: <http://jquery.com/>