

2

Author Proof

3 **A cost-effective lexical acquisition process**
4 **for large-scale thesaurus translation**

5 **Jimmy Lin · G. Craig Murray · Bonnie J. Dorr · Jan Hajič ·**
6 **Pavel Pecina**

7

8 © Springer Science+Business Media B.V. 2008

9 **Abstract** Thesauri and controlled vocabularies facilitate access to digital
10 collections by explicitly representing the underlying principles of organization.
11 Translation of such resources into multiple languages is an important component for
12 providing multilingual access. However, the specificity of vocabulary terms in most
13 thesauri precludes fully-automatic translation using general-domain lexical resour-
14 ces. In this paper, we present an efficient process for leveraging human translations
15 to construct domain-specific lexical resources. This process is illustrated on a
16 thesaurus of 56,000 concepts used to catalog a large archive of oral histories. We
17 elicited human translations on a small subset of concepts, induced a probabilistic
18 phrase dictionary from these translations, and used the resulting resource to auto-
19 matically translate the rest of the thesaurus. Two separate evaluations demonstrate
20 the acceptability of the automatic translations and the cost-effectiveness of our
21 approach.

22

23 **Keywords** Thesauri · Controlled vocabularies · Manual translation process

24

25 **1 Introduction**

26 Providing multilingual access to digital collections is an important challenge in
27 today's increasingly interconnected world. For the most part, research in multilin-
28 gual access focuses on the content of digital repositories themselves, often
29 neglecting significant knowledge present in the associated metadata. Many

A1 J. Lin (✉) · G. C. Murray · B. J. Dorr
A2 University of Maryland, College Park, MD, USA
A3 e-mail: jimmylin@umd.edu

A4 J. Hajič · P. Pecina
A5 Charles University, Prague, Czech Republic

30 collections employ controlled-vocabulary descriptors, hierarchically arranged in a
31 thesaurus, to characterize the content of items in the collection. Such structures
32 explicitly encode the organizing principles of a collection and facilitate access via
33 searching, browsing, or a combination of both. Multilingual access to such thesauri
34 can enhance content-oriented technologies such as cross-language information
35 retrieval and machine translation in helping users access content in foreign
36 languages.

37 Building on previous work (Murray et al. 2006a, b), this article tackles the question
38 of how one might, given limited resources, efficiently translate a large thesaurus to
39 facilitate multilingual information access. Due to limited vocabulary coverage, off-
40 the-shelf translation technology provides little help for specialized domains. Instead,
41 we propose a process for lexical acquisition that yields high-value reusable resources
42 for automatic translation. The key to a cost-effective process is to model the utility of
43 each descriptor within the thesaurus, taking into account thesaurus structure and the
44 reusability of component phrases. Guided by such a utility function, we elicited
45 manual translations for a small selection of thesaurus terms, and from these induced
46 lexical resources for translating the rest of the terms automatically. Experiments
47 suggest that our approach yields acceptable translations and provides significant cost-
48 savings compared to an unoptimized translation process.

49 2 The problem

50 Our work is situated in the context of MALACH (Multilingual Access to Large
51 Spoken Archives), an effort funded by the U.S. National Science Foundation
52 (Gustman et al. 2002). The USC Shoah Foundation Institute for Visual History and
53 Education manages what is presently the world's largest archive of videotaped oral
54 histories (USC 2006). The archive contains 116,000 h of video testimonies of over
55 52,000 survivors, liberators, rescuers, and witnesses of the Holocaust. The Shoah
56 Foundation uses a hierarchically-arranged thesaurus that contains approximately
57 56,000 domain-specific concepts represented by keyword phrases (descriptors).
58 These descriptors are assigned to time points in the video testimonies as a means for
59 indexing the video content. Although testimonies are available in other languages,
60 the thesaurus is currently available only in English. Translation of this resource into
61 different languages would greatly enhance multilingual access. As a proof-
62 of-concept, this article focuses on translation into Czech.

63 Initial attempts to automatically translate the thesaurus revealed that only 15% of
64 the vocabulary could be found in an available aligned corpus, the Prague Czech-
65 English Dependency Treebank (PCEDT) (Čmejrek et al. 2004). Due to the
66 specificity of the domain, translations for the remaining terms could not be found
67 in general electronic resources, including dictionaries at our disposal. Since reliable
68 access requires high accuracy, we found it necessary to acquire lexical information
69 from humans. However, it would be cost prohibitive to manually translate the entire
70 thesaurus. Our solution involves acquiring human translations for a small selection
71 of phrases from the thesaurus and then leveraging this information to automatically
72 translate the remainder.

73 In this work, we propose a human-assisted translation process that takes into
74 account characteristics of the thesaurus. A relatively small number of keyword
75 phrases provides access to a large portion of the video content. Similarly, a large
76 number of highly specific phrases describe only a small fraction of content.
77 Therefore, not every phrase carries the same utility. The hierarchical arrangement of
78 the keyword phrases presents another challenge: some phrases, while not of great
79 value for directly accessing content, may be important for organizing other concepts
80 and for browsing. These factors must be balanced in developing a cost-effective
81 translation process.

82 3 A proposed solution

83 This article presents a cost-effective, human-in-the-loop approach to translating
84 large thesauri. Using this approach, we collected 3,000 manual translations of
85 keyword phrases from the Shoah Foundation's thesaurus and reused the translated
86 terms to generate a lexicon, which was then employed to automatically translate the
87 rest of the thesaurus. This section describes our process model in detail.

88 3.1 Prioritization

89 Given unlimited financial resources, one could simply elicit manual translations for
90 all concepts in a thesaurus. However, since most projects face resource constraints,
91 one must devise a prioritization scheme for manual translation, placing "more
92 useful" terms before "less useful" ones. We define two measures to quantify the
93 utility of a keyword phrase in our thesaurus: *thesaurus value*, which represents the
94 importance of a particular keyword phrase for providing access to the collection,
95 and *translation value*, which quantifies the usefulness of having the keyword phrase
96 translated. We describe these measures in detail.

97 Keyword phrases in the Shoah Foundation's thesaurus are arranged in a
98 poly-hierarchy where nodes can have multiple parents. Internal (non-leaf) nodes of
99 the hierarchy are primarily used to organize concepts and support browsing,
100 although some of these nodes are also used to index video content. Leaf nodes
101 represent specific concepts and are only used for indexing. Thus, the utility of a
102 keyword phrase for providing access to the collection is directly related to the
103 concept's position in the thesaurus hierarchy.

104 A concrete example will help to make this clear. Consider the fragment of the
105 thesaurus hierarchy shown in Fig. 1. The keyword phrase "Auschwitz II-Birkenau
106 (Poland: Death Camp)", which describes a Nazi death camp, is assigned to 17,555
107 video segments in the collection. It has broader (parent) terms and narrower (child)
108 terms. Some, but not all, of the broader and narrower terms are also assigned to
109 segments. Notably, "German death camps" is not assigned to any video segments,
110 although it is important because it facilitates access to six frequently assigned
111 narrower terms. This example demonstrates the value of internal nodes in providing
112 access to the structure of the thesaurus, even when those concepts are not directly

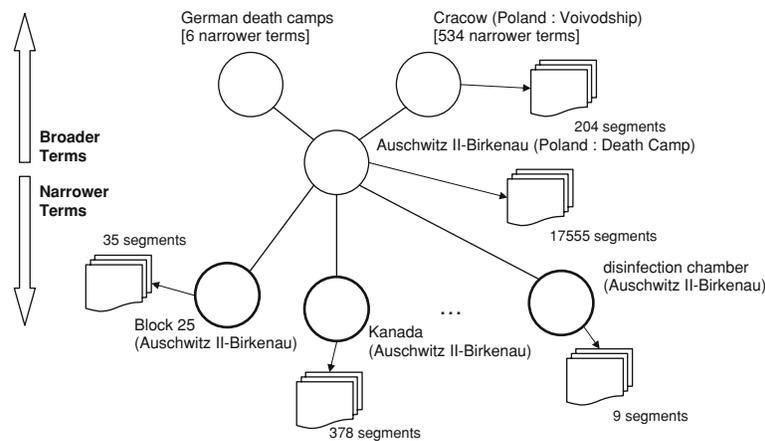


Fig. 1 A fragment from the thesaurus showing broader and narrower terms. Nodes with thick borders indicate leaf nodes

113 assigned to any content segments. As such, translation of these nodes provides
 114 multi-lingual access to the arrangement concepts within the thesaurus.

115 We interpret the number of video segments under any given node in the hierarchy
 116 (directly or via child nodes) as an indication of that node's potential importance for
 117 accessing collection content. We have no principled reason to assume that any
 118 particular video segment is more important than any other. Therefore, we treat each
 119 as equally important. We use these counts to estimate the importance of each
 120 keyword phrase's inclusion in the thesaurus and, by extension, of the utility gained
 121 from translating that keyword phrase.

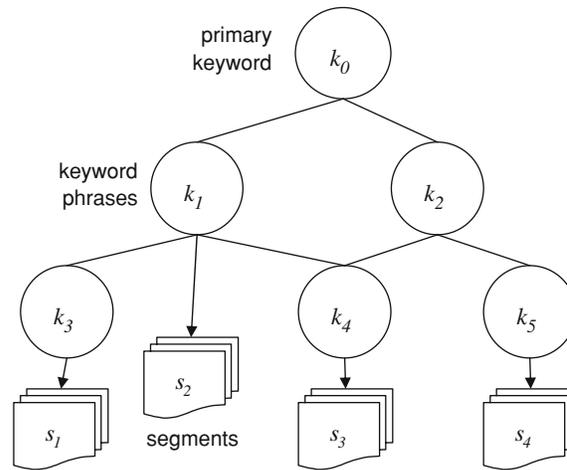
122 Parent nodes high in the hierarchy help users manage concepts, but nodes low in
 123 the hierarchy are closer to the content. The hierarchy organizes content, but
 124 navigating the nodes comes at some cost to the user (e.g., cognitive load, physical
 125 interaction, etc.) Thus, our utility function must balance the value of direct content
 126 access with the value supplied by internal nodes that provide structure and facilitate
 127 browsing. To strike this balance, we introduce *thesaurus value*, which quantifies the
 128 importance of each keyword phrase with respect to the thesaurus:

$$h_k = \text{count}(s_k) + \frac{\sum_{i \in \text{children}(k)} h_i}{|\text{children}(k)|} \quad (1)$$

130 For leaf nodes in our thesaurus, this value is simply the number of video seg-
 131 ments s to which the concept k has been assigned. For non-leaf nodes, the thesaurus
 132 value is the number of segments (if any) to which the concept has been assigned,
 133 plus the average of the thesaurus values of child nodes. This recursive calculation
 134 yields a micro-averaged value that represents the reachability of segments via
 135 downward edge traversals from a given node in the hierarchy. That is, the thesaurus
 136 value captures the number of segments described by a given keyword phrase and the
 137 average number of segments described by its children in the hierarchy (i.e. narrower
 138 terms).

Fig. 2 Thesaurus fragment illustrating the computation of thesaurus values

Author Proof



139 For example, in Fig. 2, each of the leaf nodes (k_3 , k_4 , and k_5) has value only as a
 140 means for directly accessing content (s_1 , s_3 , and s_4). Node k_1 has value both as a
 141 direct access point to segments s_2 and indirectly to segments s_1 and s_3 (via k_3 and
 142 k_4). Other internal nodes, such as k_2 , have value only in providing access to other
 143 keyword phrases (k_4 and k_5). Working our way up from the bottom of the hierarchy,
 144 we can compute the thesaurus value for each node in this simple example as
 145 follows: For nodes k_3 through k_5 , we simply count the number of segments that have
 146 been assigned each keyword phrase. Then we move up to nodes k_1 and k_2 . At k_1 we
 147 count the number of segments s_2 to which k_1 was assigned and add that count to the
 148 average of the thesaurus values for k_3 and k_4 . At k_2 we simply average the thesaurus
 149 values for k_4 and k_5 . And so on up the hierarchy. The final values quantify the utility
 150 of keyword phrases in providing access to video content. Although it would make
 151 some sense to prioritize human translations based simply on these thesaurus values,
 152 we can gain even more efficiency by taking into account the utility of individual
 153 lexical components within the keyword phrase.

154 Our example in Fig. 1 also illustrates the recurrent nature of the individual words
 155 that make up keyword phrases. Note that the term “Auschwitz” appears in four of
 156 the keyword phrases shown. In fact, the term “Auschwitz” occurs in 35 keyword
 157 phrases in the English thesaurus, and these are used as content descriptors for a
 158 significant portion of the archive. Thus, the impact of translating any individual
 159 term (i.e., word) is a function of the cumulative thesaurus value of all the keyword
 160 phrases in which it occurs. As a candidate for translation, “Auschwitz” has high
 161 potential impact, both in the number of keyword phrases that contain this term, and
 162 the value of those keyword phrases (once translated) in providing multi-lingual
 163 access to video segments in the archive.

164 Here, we introduce a measure of the *translation value* for each term (i.e., word)
 165 in the vocabulary. After obtaining the thesaurus values for each keyword phrase, we
 166 compute the translation value as the sum of the thesaurus value for every keyword
 167 phrase in which the term appears:

$$t_w = \sum_{k \in K_w} h_k \quad \text{where } K_w = \{x | \text{phrase } x \text{ contains } w\} \quad (2)$$

169 The end result of computing translation values is a list of terms and the impact
 170 that the correct translation of each term will have on the overall value of the
 171 translated thesaurus.

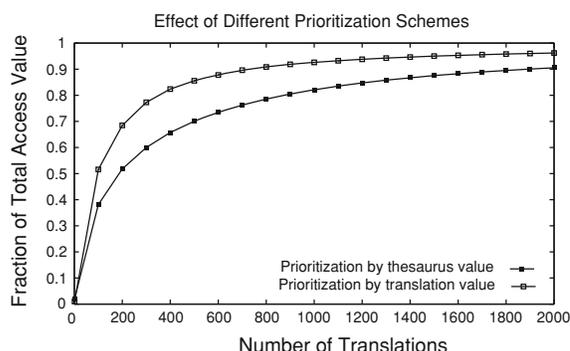
172 Accurate translation of individual terms requires context. Therefore, we elicited
 173 human translations of entire keyword phrases rather than individual terms. But how
 174 best to prioritize these translations? The value that any keyword phrase has for
 175 translation is only indirectly related to its own value as a point of access to the
 176 collection. Some keyword phrases have low thesaurus value but contain terms with
 177 high translation value. The impact of translating those keywords is not directly
 178 reflected by their use in describing the collection (i.e., their thesaurus value). Thus,
 179 the value gained by translating any given keyword phrase is more accurately
 180 estimated by the total value of any untranslated terms it contains. Therefore, we
 181 prioritized keyword phrases based on the translation value of the untranslated terms
 182 in each keyword phrase.

183 This process is implemented as follows: we iterate through the thesaurus
 184 keyword phrases, prioritizing their translation based on the assumption that any
 185 terms contained in a keyword phrase of higher priority would already have been
 186 translated. Starting from the assumption that the entire thesaurus is untranslated, we
 187 choose the keyword phrase that contains the most valuable untranslated terms. This
 188 is done by adding up the translation value of all the untranslated terms in each
 189 keyword phrase and selecting the keyword phrase with the highest sum. We add this
 190 keyword phrase to the prioritized list of items to be manually translated and remove
 191 it from the list of untranslated phrases. We update our vocabulary list, assuming that
 192 all the terms in the keyword phrase are now translated (neglecting issues such as
 193 morphology). Then, we again select the keyword phrase that contains the most
 194 valuable untranslated terms. This process iterates until all terms have been added to
 195 the prioritized list.

196 Note that this prioritization scheme is greedy and biased toward longer keyword
 197 phrases. In addition, some terms may be translated more than once because they
 198 appear in more than one keyword phrase with high (total) translation value.¹ This
 199 side effect is actually desirable. To build an accurate translation dictionary, it is
 200 helpful to have more than one translation of frequently occurring terms, especially
 201 for morphologically rich languages such as Czech. Our approach assumes that
 202 translations of terms gathered in one context can be reused in another context.
 203 Obviously, this is not always true, but contexts of use are relatively stable in
 204 controlled vocabularies. The longer keyword phrases provide richer contextual
 205 support for the translations. Our evaluations examine the validity of this context
 206 assumption and demonstrate that the technique yields acceptable translations.

1FL01 ¹ Even after a term is assumed to be translated, there will be keyword phrases containing that term which
 1FL02 contain other high translation value terms not yet translated. In some cases, the sum of the translation
 1FL03 value of the untranslated terms will be high enough to warrant addition of the keyword phrase to the
 1FL04 prioritized list, despite the already translated term.

Fig. 3 Efficiency of different prioritization schemes



Author Proof

207 Following the process described above, the most important elements of the
 208 thesaurus will be translated first, and the most important vocabulary terms will
 209 quickly become available for automatic translation of those keyword phrases
 210 with high thesaurus value that do not make it onto the prioritized list for
 211 manual translation. To evaluate our prioritization scheme we need to quantify the
 212 accessibility of the collection (via the translated thesaurus) at different levels of
 213 human translation. We measure *access value* as the sum of the thesaurus value of
 214 translated keywords—whether by manual or automatic means. Access value
 215 represents the utility of the thesaurus after machine translation in providing multi-
 216 lingual access to the contents of the archive. Figure 3 plots the rate of gain in access
 217 value after eliciting translations. It can be seen that prioritizing elicited translations
 218 based on translation value yields a more efficient process than prioritization based
 219 on thesaurus value.

220 3.2 Caveats, alternatives, and possible improvements

221 We introduced three measures in this work: *Thesaurus value* is a measure of the
 222 contribution each keyword phrase makes to the overall value of the thesaurus.
 223 *Translation value* is a measure of the contribution each translated vocabulary term
 224 makes to the overall translation of the thesaurus. *Access value* is a measure of the
 225 collection access facilitated by a translated thesaurus. These measures come from a
 226 careful analysis of the problem of prioritized partial translation. Nevertheless, there
 227 are some operational assumptions which deserve further discussion.

228 In our definition of thesaurus value, we did not attempt to quantify the relative
 229 importance of browsing the concept hierarchy vs. accessing collection content. In
 230 some settings it may be more important to facilitate content access over concept
 231 browsing, or vice versa. It would be possible to add a weighting constant to Eq. 1,
 232 giving emphasis to either one or the other.

233 Descriptors in a thesaurus have two functions: indexers use them to index the
 234 collection, and patrons use them to retrieve collection contents. We chose not to
 235 model the frequency of their use to retrieve content. Patrons' use of English keyword
 236 phrases in queries could be used to estimate the expected use of keyword phrases in

237 another language. However, this would reflect only those interests the collection has
238 served in the past. With a large and growing collection of oral histories it is impossible
239 to know what interests it will serve in the future. Smaller collections and narrower
240 domains may profit from an analysis of the patrons' queries.

241 We mention earlier that each video segment was assumed to be of equal
242 importance. In some settings this may not be the case. For example, most patrons
243 will have a preference for video testimonies in their native language. We found that
244 the frequency distribution of keyword assignments to Czech content is similar to
245 that of English, but there may be a biased distribution for other languages. For these
246 it might be advantageous to give higher importance to segments in that language.
247 Purpose of use is also a factor: for example, some of the collection content is cleared
248 for broad use, some only for limited use. If the purpose of thesaurus translation were
249 to provide access only for broad use, limited use segments should be discounted or
250 excluded when calculating thesaurus value.

251 Our measure for the translation value of vocabulary terms is based on the thesaurus
252 value of the keyword phrases, and so it inherits assumptions of that measure.
253 Translation value has the further operational assumption that terms are equally
254 informative about their language. This is obviously false—some terms will carry a
255 great deal of information about the coding scheme of their language (i.e. morphology,
256 syntax, etc.), while others will not. From this view, translations of complex terms may
257 be more valuable. Quantifying that value, however, requires a means of identifying
258 complex terms and of weighing the value of different language features.

259 Access value also inherits the assumptions of translation value and comes with
260 certain caveats. It is an approximation based solely on collection content and is used
261 here only to compare different prioritization schemes. It does not include an
262 assessment of usability and is not intended to measure the translation output in an
263 absolute sense. We report on quality of translations later in this paper.

264 Each of these measures could be expanded in different ways to suit different
265 purposes. Changes to thesaurus value or to translation value will result in changes to
266 the prioritized list of phrases to be translated. Choice of different utility functions
267 for measuring access value will give different views of the success of prioritization.
268 Researchers wishing to expand these measures should be careful to justify any
269 added complexity with clear purpose.

270 3.3 Human translation, alignment, and decomposition

271 Following the prioritization scheme described above, we obtained professional
272 Czech translations for the top 3,000 keyword phrases. We tokenized these
273 translations and presented them to another bilingual Czech speaker for alignment.
274 This second informant linked equivalent Czech and English words using a GUI.
275 Multiple links conveyed the relationship between a single word in one language and
276 a phrase in the other. Details of the alignment step can be found in Murray et al.
277 2006a. Human translation of the keyword phrases took approximately 70 h, and the
278 alignments took 55 h. The overall cost of human input (translation and alignment)
279 was less than 1,000€.

280 From the human output, we constructed a probabilistic English-Czech phrase
281 dictionary based on the distribution of the alignments. For example, in the top 3,000
282 keyword phrases “stills” appeared 29 times. It was aligned with “statické snímky”
283 28 times and only once with “statické záběry”, giving us a translation probability of
284 $28/29 = 0.966$ for “statické snímky”.

285 3.4 Automatic translation

286 To demonstrate the effectiveness of our approach, we show that a probabilistic
287 dictionary, generated using the process we just described, facilitates high quality
288 automatic translation. Our translation system implements a greedy algorithm with a
289 simple back-off strategy. It first scans the English input to find the longest matching
290 substring in our dictionary, and replaces the substring with the most likely Czech
291 translation. For example, given the phrase “monasteries and convents (stills)”, the
292 system first looks for the entire phrase in the dictionary, but finds no translation.
293 Then, the system backs off to “monasteries and convents” and finds the translation
294 “kláštery”. Next, the system tries to find a match for “stills” in the same manner.

295 If the system fails to find a match in our lexical resources, it backs off to a
296 dictionary induced from the PCEDT (Čmejrek et al. 2004). If no match is found in
297 either dictionary for the full token, the process is repeated with the stem. Failing a
298 match on the stem, terms are simply passed through untranslated. A minimal set of
299 heuristic rules is then applied to reorder the Czech tokens, but the output is primarily
300 word-by-word lookup translation.

301 As this work focuses primarily on processes for human translation, we were not
302 concerned about the simplicity of our system (compared to state-of-the-art statistical
303 MT technology). In our case, we interpret measures of translation accuracy as
304 quality measures for lexical resources. The simplicity of our system ensures that
305 improvements to lexical coverage are not conflated with other factors. More
306 sophisticated systems will no doubt also benefit from these resources.

307 4 Evaluation

308 We evaluated our translation process in two different ways. First, we compared
309 automatic translations with human reference translations using BLEU (Papineni et al.
310 2002) and TER (Snover et al. 2005), two commonly-used metrics for automatic
311 evaluation of MT output. Second, we presented automatic translations to Czech
312 speakers and gathered subjective judgments of fluency and accuracy.

313 For evaluation, we selected 418 keyword phrases using a stratified sampling
314 technique so that items with a broad range of thesaurus values would be represented.
315 However, we ensured that there was no overlap between these keyword phrases and
316 the 3,000 manually-translated keyword phrases used to build our lexicon.

317 For the automatic evaluation, we obtained two separate sets of reference
318 translations. First, prior to automatic translation, we gathered at least two independent
319 human translations for each keyword phrase. We refer to this as the “independent

320 reference” set. Second, we asked our informants to correct automatic translations into
 321 fluent Czech, preserving as much of the original machine output as possible. For these,
 322 we automatically translated the test set using a probabilistic dictionary that was
 323 generated using the first 2,500 prioritized translations. The machine output was then
 324 corrected by native Czech speakers, who adjusted word order, word choice,
 325 morphology, etc. We refer to this as the “human corrected” set. These translations
 326 often differed from the independent references, since there are multiple ways to
 327 translate the same phrase.

328 To assess the effectiveness of our translation process, we compared uncorrected
 329 automatic translations to the two different sets of reference translations. These results
 330 are shown in Fig. 4, with BLEU on the top and (1-TER) on the bottom. The x axis shows
 331 the number of aligned human translations used to construct the lexicon. The zero

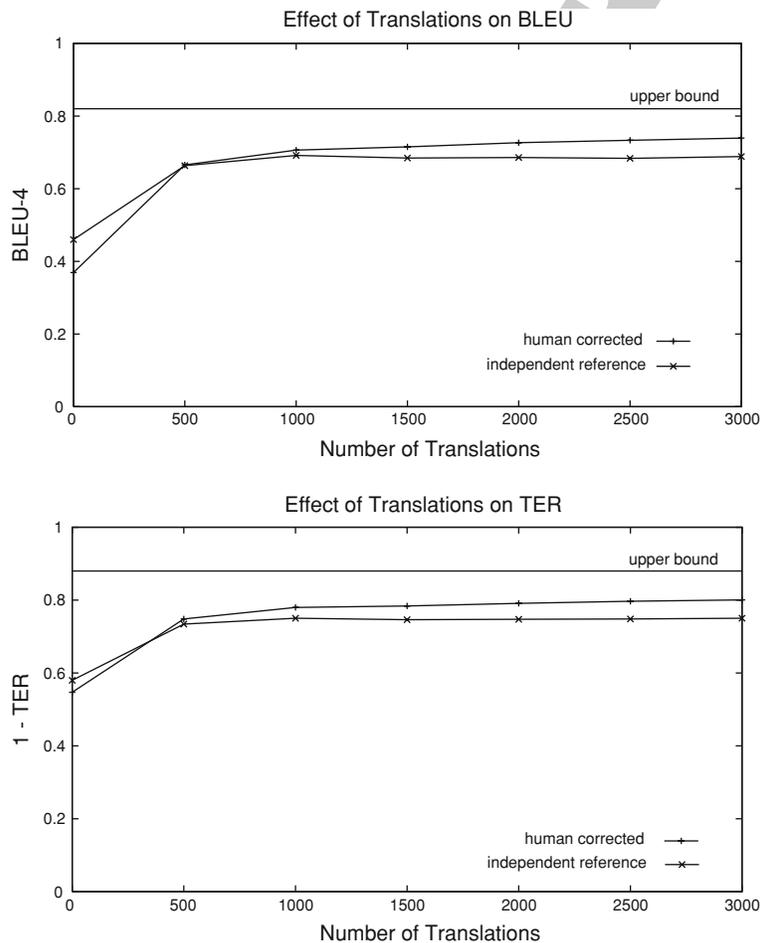


Fig. 4 Automatic evaluation of machine translations with BLEU (top) and TER (bottom) using two reference sets: “independent reference” and “human corrected”

332 condition represents our baseline: translations generated using only the dictionary
333 available in the PCEDT. We take the performance of the human corrected translations
334 with respect to the independent references as the upper bound, shown in both graphs.
335 There is a big jump in both BLEU and (1-TER) scores after the first 500 translations are
336 added to our probabilistic dictionary. Gains thereafter are smaller, but noticeable. In
337 both cases, it appears that performance approaches the upper bound.

338 To determine the impact of external resources, we removed the PCEDT
339 dictionary as a back-off resource and retranslated keyword phrases using only the
340 lexicons induced from our aligned translations. The results of this experiment
341 showed only marginal degradation of the output. Even when as few as 500 aligned
342 translations were used, we still achieved a BLEU score of 0.65 against the
343 independent references. This suggests that even for languages where no resources
344 are available, our process is capable of coping with vocabulary coverage issues.

345 In our subjective evaluation, we presented a random sample of automatic
346 translations and corrected translations (i.e., the “human corrected” set described
347 above) to seven native Czech speakers. They were asked to rate the fluency and
348 accuracy of the phrases on a 5-point Likert scale (1 = good, 5 = bad). Results are
349 shown in Fig. 5. In all cases, the mode is 1 (i.e., “good”). According to our judges,
350 59% of the uncorrected automatic translations were rated 2 or better for fluency;
351 66% were rated 2 or better for accuracy. Disfluencies were primarily caused by
352 errors in morphology and word order; for more details, see (Murray et al. 2006b).
353 We note that lexical accuracy is more important than grammatical fluency for
354 providing information access.

355 5 Related work

356 The notion of human-assisted machine translation is not new, and human input has
357 been used to great effect in the past. The Pangloss project (Frederking et al. 1994)
358 developed an MT system where human assistance was solicited during the
359 translation process. Other approaches to human-in-the-loop translation have
360 involved more sophisticated symbolic representations (Olsen et al. 1998; Sabarís
361 et al. 2001). These detail-oriented approaches tend to be knowledge-intensive and
362 difficult to economize. Our study takes a cost-oriented approach.

363 Several studies have taken a knowledge-acquisition approach to collecting
364 multilingual word pairs. For example, Sadat et al. (2003) automatically extracted
365 bilingual word pairs from comparable corpora. Others have leveraged parallel
366 corpora or bilingual dictionaries for lexical acquisition (Echizen-ya et al. 2006; Kaji
367 and Aizono 1996; Rapp 1999; Tanaka and Iwasaki, 1996). However, our work deals
368 with the fundamentally different task of translating a large thesaurus, where one can
369 leverage the structural properties of the resource.

370 Many recent approaches to dictionary and thesaurus translation are geared toward
371 providing domain-specific thesauri for specialists in a particular field, e.g., medical
372 terminology (Déjean et al. 2005) or agricultural terminology (Chun and Wenlin
373 2002). Researchers on these projects are faced with the choice of either finding
374 human domain experts to manage manual translation or applying automatic

Author Proof

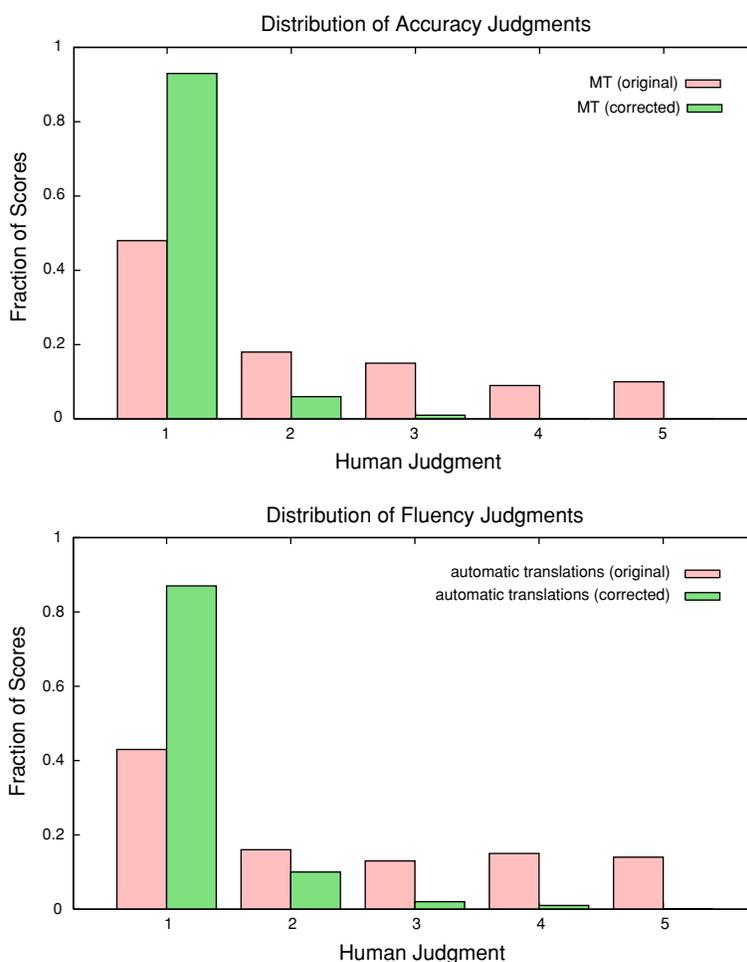


Fig. 5 Results of the subjective evaluation: fluency (left) and accuracy (right)

375 acquisition techniques, where data sparsity poses a problem for low-frequency
 376 terms. we balanced the need for human domain knowledge against the cost of
 377 human input.

378 Our process enriches the field of study with a hybrid alternative to full human
 379 translation or human assisted automated translation. We leverage the structure of the
 380 thesaurus and the recombinant nature of keyword phrases to prioritize the human
 381 input in advance. Then we set a threshold on the cost of translation, effectively
 382 switching from human translation to automated translation when the threshold is
 383 reached. In this way we tie into other research on both sides of the threshold, and
 384 introduce a cost function for managing the tradeoff.

385 Future work in this vein could explore individual influences of factors such as
 386 language complexity, domain specificity, or concept portability. The current work
 387 does not address variability of concepts across different cultures. The keyword

388 phrases we translated are best seen as English labels on distinct concepts. Our
389 translations produced Czech labels for these concepts. However, a thesaurus
390 represents one particular view of the world. A cohesive concept that has clear
391 boundaries in one language (and its culture) may be far less cohesive in another, and
392 the labels for that concept may require more finesse in translation. Any interlingua
393 inherently faces problems of mapping concepts from one culture to concepts in
394 another. In translating thesauri it may be possible to leverage the structural
395 arrangement of concepts to improve translations of other material. It may also be
396 possible to leverage prioritized human input to learn structural mappings between
397 dissimilar concept structures, e.g. competing ontologies. These questions warrant
398 further investigation.

399 6 Conclusion

400 The task of thesauri translation can be recast as the problem of implementing a
401 cost-effective process for acquiring domain-specific lexical resources. We devel-
402 oped a process for eliciting human translations. From 3,000 manually translated
403 keyword phrases, we induced a probabilistic dictionary. Using this resource, we
404 achieved acceptable automatic translation of the complete 56,000-concept thesau-
405 rus. As a rough calculation, the overall cost of human input was less than 1,000€.
406 Had we paid for human translation of the entire thesaurus it would have cost close to
407 20,000€. Naturally, this is a biased comparison since manual translation of the
408 entire thesaurus would have yielded a product much higher in quality. Nevertheless,
409 we are able to implement a solution that *approximates* a gold standard, at a small
410 fraction of the cost.

411 The value of our work lies in the process model we developed for cost-effective
412 acquisition of lexical resources. We have shown that careful prioritization of human
413 translations can efficiently yield reusable lexicons for automatic translation. The
414 development of a utility function that accurately models both the direct and indirect
415 value of a particular concept is the key to a cost-effective prioritization. Our process
416 model aims to address the most critical deficiencies in vocabulary coverage first,
417 such that the value obtained from each additional human translation becomes
418 successively smaller. Under such a framework, choosing the number of human
419 translations to elicit becomes a function of the financial and human resources
420 available for the task.

421 Although this work focuses on thesaurus translation, the process we developed
422 can be extended to other types of structured texts as well. For example, ontologies
423 and knowledge bases have poly-hierarchic structures similar to the typonomic
424 relations in the Shoah thesaurus. Our objective function was based on access to
425 multimedia, but similar objective functions could be developed for different types of
426 structural nodes to guide the translation process. We believe that our process is ideal
427 for languages with scarce resources. Resources tend to be scarce for exactly the
428 same languages and cultures which stand to gain the most from translated structural
429 knowledge representations. The end result of this work will be a step toward a rich
430 multilingual dictionary of Holocaust terms. Similar resources could be developed

431 for legal terms, medical terms, etc. These in turn could serve to educate and
432 empower the peoples of many nations.

433 **Acknowledgements** Our thanks to Doug Oard for helpful discussions; to our Czech informants; and to
434 Soumya Bhat for her programming efforts. This work was supported in part by NSF IIS Award 0122466
435 and NSF CISE RI Award EIA0130422. Additional support also came from grants of the MSMT CR
436 #1P05ME786, #LC536 and #MSM0021620838, and the Grant Agency of the Czech Republic #GA405/
437 06/0589. The first author would like to thank Esther and Kiri for their kind support.

438 References

- 439 Chun, C., & Wenlin, L. (2002). The translation of agricultural multilingual thesaurus. In *Proceedings of*
440 *the Third Asian Conference for Information Technology in Agriculture*.
- 441 Čmejrek, M., Cuřín, J., Havelka, J., Hajič, J., & Kuboň, V. (2004). Prague Czech-English dependency
442 treebank: Syntactically annotated resources for machine translation. In *Proceedings of LREC 2004*.
- 443 Déjean, H., Gaussier, E., Renders, J.-M., & Sadat, F. (2005). Automatic processing of multilingual
444 medical terminology: Applications to thesaurus enrichment and cross-language information
445 retrieval. *Artificial Intelligence in Medicine*, 33(2), 111–124.
- 446 Echizen-ya, H., Araki, K., Momouchi, Y. (2006). Automatic extraction of bilingual word pairs using
447 inductive chain learning in various languages. *Information Processing and Management*, 42(5),
448 1294–1315.
- 449 Frederking, R., Nirenburg, S., Farwell, D., Helmreich, S., Hovy, E., Knight, K., Beale, S., Domashnev, C.,
450 Attardo, D., Grannes, D., & Brown, R. (1994). The Pangloss Mark III machine translation system. In
451 *Proceedings of the 1st AMTA Conference*.
- 452 Gustman, S., Soergel, D., Oard, D. W., Byrne, W. J., Picheny, M., Ramabhadran, B., & Greenberg, D.
453 (2002). Supporting access to large digital oral history archives. In *Proceedings of JCDL 2002*
454 (pp. 18–27).
- 455 Kaji, H., & Aizono, T. (1996). Extracting word correspondences from bilingual corpora based on word
456 co-occurrence information. In *Proceedings of COLING 1996* (pp. 23–28).
- 457 Murray, G. C., Dorr, B., Lin, J., Hajič, J., & Pecina, P. (2006a). Leveraging recurrent phrase structure in
458 large-scale ontology translation. In *Proceedings of EAMT 2006*.
- 459 Murray, G. C., Dorr, B., Lin, J., Hajič, J., & Pecina, P. (2006b). Leveraging reusability: Cost-effective
460 lexical acquisition for large-scale ontology translation'. In *Proceedings of COLING/ACL 2006*
461 (pp. 945–952).
- 462 Olsen, M., Dorr, B., & Thomas, S. (1998). Enhancing automatic acquisition of thematic structure in a
463 large-scale lexicon for Mandarin Chinese. In *Proceedings of the Third Conference of the Association*
464 *for Machine Translation in the Americas (AMTA '98)*.
- 465 Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). BLEU: A method for automatic evaluation of
466 machine translation. In *Proceedings of ACL 2002* (pp. 311–318).
- 467 Rapp, R. (1999). Automatic identification of word translations from unrelated English and German
468 Corpora. In *Proceedings of ACL 1999* (pp. 519–526).
- 469 Sabarís, M., Alonso, J., Dafonte, C., & Arcay, B. (2001). Multilingual authoring through an artificial
470 language. In *Proceedings of MT Summit VIII*.
- 471 Sadat, F., Yoshikawa, M., & Uemura, S. (2003). Enhancing cross-language information retrieval by an
472 automatic acquisition of bilingual terminology from comparable corpora. In *Proceedings of SIGIR*
473 *2003* (pp. 397–398).
- 474 Snover, M., Dorr, B. J., Schwartz, R., Makhoul, J., Micciulla, L., & Weischedel, R. (2005). *A study of*
475 *translation error rate with targeted human annotation*. Technical Report LAMP-TR-126/CS-TR-
476 4755/UMIACS-TR-2005-58, University of Maryland, College Park.
- 477 Tanaka, K., & Iwasaki, H. (1996). Extraction of lexical translations from non-aligned corpora. In
478 *Proceedings of COLING 1996* (pp. 580–585).
- 479 USC. (2006). USC Shoah Foundation Institute for Visual History and Education.
- 480