

Lexica and Corpora for Speech-to-Speech Translation: A Trilingual Approach

David Conejero, Jesús Giménez, Victoria Arranz, Antonio Bonafonte,
Neus Pascual, Núria Castell, Asunción Moreno

Talp Research Center
Universitat Politècnica de Catalunya
Jordi Girona, 1-3, 08034 Barcelona
david@talp.upc.es

Abstract

Creation of lexica and corpora for Catalan, Spanish and US-English is described. A lexicon is being created for speech recognition and synthesis including relevant information. The lexicon contains 50K common words selected to achieve a wide coverage on the chosen domains, and 50K additional entries including special application words, and proper nouns.

Furthermore, a large trilingual spontaneous speech corpus has been created. These corpora, together with other available US-English data, have been translated into their counterpart languages. This is being used to investigate the language resources requirements for statistical machine translation.

1. Introduction

The development of a Speech-to-Speech Machine Translation (SST) system involves high performance components for Speech Recognition (SR), Machine Translation (MT) and Text-to-Speech Synthesis (TTS). Therefore, most attempts to develop such systems have been performed in projects including various partners with expertise in these areas. Some well-known references are the C-Star Consortium¹, and the NE-SPOLE², Verbmobil³, and Eutrans⁴ projects.

These projects showed that it is possible to develop robust SST systems for small- to medium-sized domains using sophisticated SR and MT technology. The major problems in this area are: the acquisition of domain-specific training data (either mono- or bilingual), the robust behaviour of the MT component for SR errors and spoken language phenomena, and the development of an efficient recognition and translation component.

Language resources for SR have proved to be very useful to improve SR technology. Oral databases have been created in a large number of different languages. Furthermore, some producing-tools and standards have been created. Many databases are publically available at the European Language Resources Association (ELRA) or the Linguistic Data Consortium (LDC). In the last few years, the performance of speech synthesizers has also improved dramatically due to the usage of oral language resources. However, large lexica specifically designed for speech applications are difficult to find. Speech recognizers and synthesizers need good lexica with a wide coverage of the language, appropriate morphosyntactic and phonetic information.

Moreover, parallel bilingual texts are needed so as to enhance the performance of Statistical Machine Translation (SMT) systems. Electronic texts are only available for a small number of languages and no standard exists so far. Regarding large lexica in Catalan and/or Spanish, EuroWordNet is the only resource available at ELRA that may be useful for this purpose, although it was not designed on the speech technology perspective. In what regards large bilingual corpora (and specially for spoken languages) these are very difficult to find.

This paper shows the creation of large lexica in Catalan and Spanish for SR and TTS. This can be seen in section 2, together with the description of the corpora and word lists used. The paper also describes (section 3) the creation of a large sentence-aligned trilingual corpus (Catalan, Spanish and US-English) from spontaneous speech recordings in a tourist domain. This corpus will be used to investigate the requirements for SST-oriented LR and then to establish the specifications for their development. Section 4 describes the recording platform, the design of the scenarios, as well as the collection, annotation and translation processes involved. Section 5 describes a speech-to-speech translation demonstrator based on the LR generated in the project. This work is being carried out within the framework of the LC-STAR⁵ (*Lexica and Corpora for Speech-to-Speech Translation Components*) project.

2. Word Lists and Lexica for SR and SST

One aim of the project is to create large pronunciation lexica suited for speech recognition and synthesis. Each lexicon will consist of two parts: one covering common words and another one covering proper nouns from a broad range of domains. Both parts will comprise 50K entries each. The proper nouns part will consist of 45K names and 5K application words that were selected by the consortium. The first step taken to create the lexica was the definition of the word lists that would constitute them.

2.1. Corpus and Word Lists Creation

In order to generate the common word list, a corpus with six different domains was defined as listed in Table 1. According to those domains, the text corpora were collected from different sources: on-line newspapers, magazines and other websites.

Once the corpus was collected, it was cleaned and normalized: scripts, frames and html marks were removed. Then, the corpus was lemmatized and POS tagged using MACO+ [1]. Punctuation marks and proper nouns were dropped. Table 1

¹<http://www.c-star.org>

²<http://nespole.itc.it>

³<http://verbmobil.dfki.de/verbmobil>

⁴<http://www.zeres.de/Eutrans/eutrans.html>

⁵<http://www.lc-star.com>

shows the list of domains, corpus size per domain and number of distinct words in the Spanish and Catalan corpora.

Domains	Spanish		Catalan	
	Words per Domain	Diff. Words	Words per Domain	Diff. Words
<i>Sports/Games</i>	4.49 MT	73 KT	1.73 MT	31 KT
<i>News</i>	18.80 MT	161 KT	9.98 MT	99 KT
<i>Finance</i>	3.79 MT	56 KT	1.62 MT	33 KT
<i>Culture/Ent.</i>	4.30 MT	87 KT	5.01 MT	85 KT
<i>Consumer Inf.</i>	1.58 MT	42 KT	1.31 MT	44 KT
<i>Personal Com.</i>	3.68 MT	106 KT	0.55 MT	74 KT
TOTAL	36.65 MT	254 KT	20.20 MT	163 KT

Table 1: Common words per domain distribution

The objective of the word list is to have a wide lexical coverage for every domain. The LC-STAR consortium agreed that the word list should contain at least 50K words reaching a target for each domain of at least 95% self coverage. To reach the 50K different entries, the self coverage target was increased up to 97,5% in Spanish and 98,6% in Catalan. The selection process is explained in [2]. At the end of the process, a manual revision was done to handle problems arisen on the automatic process. A summary of the figures for these word lists is shown in Table 2. Coverage has been calculated with singletons included in the corpus domains, but excluded from the word lists:

Domains	Spanish		Catalan	
	Selected Words	Cover. %	Selected Words	Cover. %
<i>Sports/Games</i>	21,192	98.74	13,672	99.47
<i>News</i>	28,457	98.41	31,463	98.98
<i>Finance</i>	15,301	99.00	15,470	99.45
<i>Culture/Ent.</i>	29,471	98.36	33,794	98.89
<i>Consumer Inf.</i>	19,600	98.83	24,969	98.89
<i>Personal Com.</i>	31,802	97.96	16,293	96.79
TOTAL	55,788	98.48	53,225	99.30

Table 2: Selected word lists and coverage per domain

The proper nouns word list contains more than 45K entries and is divided into 3 different domains: First and Last names, Place names and Organizations. In order to fulfill the specifications, each domain should cover a minimum of 10% and a maximum of 50% of the entries.

2.2. Lexica for Speech Recognition and Synthesis

Once the word lists in the previous section had been created, they were used as input for the creation of monolingual lexica for SR and TTS. This is a complex task that must take into account the requirements and needs imposed by these technologies. Furthermore, the lexica produced must also provide sufficient and appropriate monolingual information so as to be linked with the translation lexica and, thus, contribute to speech-to-speech translation. In order to achieve this, a number of tasks have been considered (for more details see [3]):

- **Creation of language-independent specifications for the content of the lexica as well as addition of some language-dependent grammatical information**

This considers issues such as the required and optional orthographic, phonetic and morpho-syntactic properties of each entry in the lexicon, as well as the size of the lexica.

- **Creation of language-specific specifications for each language**

This deals with the grammatical and morphological representation, as well as with the XML-based exchange format used. The following section will expand on this task.

The creation of the language-specific specifications is as follows: firstly, a detailed POS study and description for the 12 languages of the project have been carried out, considering all possible phenomena. A set of 21 general POS tags has been adopted, where each POS is divided into a detailed set of attributes (e.g., POS for nouns contains attributes such as *Class*, *Number*, *Gender*, *Person*, *Case*, *Appreciative*, etc.). Then, a DTD has been created that implements the basic POS scheme representing all necessary features. Each language has made use of this DTD as a reference point in order to build its own language-specific POS scheme, taking into account only those phenomena relevant for the language. This DTD considers simple entries for unique word forms and entry groups for complex forms (cf. below).

Several attributes require special attention, in particular those handling the complex construction of Catalan verbal forms with *clitic pronouns*⁶. When these pronouns follow verbal forms (up to 3 of them in Catalan), they are assimilated by the verbs building up unique forms. A sample XML entry for one of these complex forms (*dóna-m'ho* - give to me it) can be seen below:

```
<ENTRYGROUP orthography="dóna-m'ho">
  <ENTRY_COMP>
    <PHONETIC>" d o - n @ - m u /PHONETIC>
    <ENTRY_EL orthography="dona">
      <VER person="2" tense="present"
        number="singular" voice="active"
        mood="imperative"/>
      <LEMMA>donar</LEMMA>
    </ENTRY_EL>
    <ENTRY_EL orthography="m' ">
      <PRO person="1" number="singular"
        gender="invariant" type="personal"/>
      <LEMMA>jo</LEMMA>
    </ENTRY_EL>
    <ENTRY_EL orthography="ho">
      <PRO gender="neuter" type="personal"
        person="3" number="invariant"></PRO>
      <LEMMA>ho</LEMMA>
    </ENTRY_EL>
  </ENTRY_COMP>
</ENTRYGROUP>
```

Some other interesting phenomena that take place in Catalan (mostly in Spanish too) and that have required special attention are: 1) *Prepositional expressions*, referring to phrases that function as prepositions and always end up in prepositions; 2) *Politeness attribute*, which takes place in personal pronouns *vostè/vostès*; 3) *Oblique case* for those personal pronouns *mi*, *ti*, *m'*, *l'*... following a preposition or apostrophed; etc. For a detailed description on the specifications generated for Catalan and Spanish languages, please refer to [3].

3. Parallel Corpora and Lexica for SST

The evaluation results obtained in the Verbmobil and the Eutrans projects showed that a very promising approach for the

⁶This phenomenon also takes place in Spanish.

translation component is statistical machine translation. This approach resulted in an error rate which is better by a factor of two in comparison to the other approaches investigated in Verbomobil. This is the approach that LC-STAR has undertaken.

Regarding the type of LRs necessary for speech centered translation, we can distinguish between corpora and lexica. Needless to say that the limitations caused by the present small reference LR are still an objective to be overcome, and this is one of the aims of LC-STAR. In particular, for SST, it is the bi- or multi-lingual collections of on-line data that are highly sought for.

The aim of the LC-STAR project is not only the creation of LRs for SST but also the specification of how lexica and corpora should be created. Taking into account the above mentioned problem of training data, the specification should consider features of LRs in order to be as useful as possible for statistical SST as well as features to reduce the amount of needed data. Always in the tourist domain, the final aligned text corpora (Catalan, Spanish and US-English) will have a size of 750K words, and the monolingual lexica will contain 10K entries per language (Catalan, Finnish, German, Hebrew, Italian, Russian, Spanish and US-English).

One part of the text corpora has been obtained from transcribed speech US-English corpora. The other part of text corpora has been obtained from the transcription of conversations recorded in Catalan and Spanish (see next section). The transcriptions in one language have been manually translated into the other two languages. The golden rule for translation has been to produce target sentences as similar as possible to the source ones, but being at the translator criterion likely to be ever uttered by a competent speaker of the target language. That is, translators tried to be as literal as possible but always preserving the meaning and correctness of the utterances. The translation is not required to be the best one.

This trilingual corpus is the main SMT system training data. The literality criterion was imposed in order to simplify the SMT system training, specially with respect to the alignment process[4]. The trilingual corpus will be adapted to further aligned corpus specification criteria.

From the US-English corpus and other sources, a reference list of 10K words has been created and will be translated into the seven mentioned languages. As the translation must be domain oriented, some 5-grams are provided with each US-English word in order to help human translators. These lists will be completed according to the further specification criteria so as to obtain the eight monolingual lexica for speech centered translation. The Catalan, Spanish and US-English lexica will be used to improve the statistical SST.

4. Oral Database

As already mentioned in the previous section, most of the text corpora come from the transcription of Catalan and Spanish spoken dialogues.

In the LC-STAR project it was decided to focus the research on the tourist domain. We have chosen a subset subset of the tourist domain, namely tourist-employee conversations. Four scenario categories were defined, namely *Hotel*, *Travel Agency*, *Tourism Office* and *Railway/Airline Company*.

4.1. Recording Sessions

In order to avoid non-verbal communication, we decided to record the spoken dialogues through the telephone network: the

Oral DataBase	
Spanish	
<i>speech raw time</i>	31h:7m:32s
<i>#speakers</i>	77
<i>#dialogues</i>	217
<i>#turns</i>	10.998
<i>#sentences</i>	24.372
<i>#words</i>	349.970
<i>#distinct words</i>	11.714
Catalan	
<i>speech raw time</i>	23h:43m:55s
<i>#speakers</i>	56
<i>#dialogues</i>	172
<i>#turns</i>	9.321
<i>#sentences</i>	19.113
<i>#words</i>	277.777
<i>#distinct words</i>	10.057

Table 3: Oral Database: Some figures.

speakers were placed in different rooms and talked on the phone with each other. Dialogues would last ten minutes in average. In order to do so, a recording platform was set up. The platform can be used in two different modes. In the first one, the platform is transparent: the two speakers may overlap each other. But we agreed that this mode was not well matched to the translation system, where a machine is between the speakers. Therefore, we designed and used a second version which imposes a rigid turn strategy, i.e. speakers are not allowed to speak both at the same time. The first turn is given to the speaker receiving the phone call. The speaker in possession of the turn has to indicate a turn exchange by pressing a key. Conversations contain some disfluences such as false starts, corrections, repetitions, filled pauses, and certain ungrammaticalities.

For the conversations to yield the pursued information some specific subscenarios were designed. A series of templates were built so as to assist speakers at conversation time. They were used as a draft or schema containing a description of the information to talk about in every subscenario. In most cases, if not all, ‘*Speaker 0*’ played the role of a tourist/customer, and ‘*Speaker 1*’ played the role of an employee. Below follows an example of one such situation:

Example: Hotel Reservation. ‘*Speaker 0*’ impersonates a tourist trying to book accommodation at a given hotel for a certain date, a given number of people, under certain specific conditions. ‘*Speaker 1*’ is acting as a hotel employee providing the required information.

4.2. Contents

It was agreed that at least 500K words should be recorded considering Catalan and Spanish as a whole. The total figures for both languages can be seen in Table 3. All dialogues have been manually transcribed and validated. During the transcription process, some tags have been added so as to encode acoustic and linguistic information.

Thus are the encoded acoustic phenomena:

- *Lengthening*: Lengthening of a sound within a word denoting hesitation.
- *Filled pause*: Pause filled by a vocalization. It may carry some meaning, i.e. affirmation, negation, hesitation, etcetera.
- *Bad pronunciation*: Speaker mispronunciations possibly

