



HAL
open science

Improving recognition of proper nouns (in ASR) through generation and filtering of phonetic transcriptions

Antoine Laurent, Sylvain Meignier, Paul Deléglise

► To cite this version:

Antoine Laurent, Sylvain Meignier, Paul Deléglise. Improving recognition of proper nouns (in ASR) through generation and filtering of phonetic transcriptions. *Computer Speech and Language*, 2014, 28 (4), pp.979-996. 10.1016/j.csl.2014.02.006 . hal-01433238

HAL Id: hal-01433238

<https://hal.science/hal-01433238v1>

Submitted on 22 Mar 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Improving recognition of proper nouns in ASR through generating and filtering phonetic transcriptions

Antoine Laurent^{a,b}, Sylvain Meignier^a, Paul Deléglise^a

^a*LIUM (Computer Science Research Department – Université du Maine, Le Mans, France)*

^b*Spécinov (Trélazé, France)*

Abstract

Accurate phonetic transcription of proper nouns can be an important resource for commercial applications that embed speech technologies, such as audio indexing and vocal phone directory lookup. However, accurate phonetic transcription is more difficult to obtain for proper nouns than for regular words. Indeed, phonetic transcription of a proper noun depends on both the origin of the speaker pronouncing it and the origin of the proper noun itself.

This work proposes a method that allows the extraction of phonetic transcriptions of proper nouns using actual utterances of those proper nouns, thus yielding transcriptions based on practical use instead of mere pronunciation rules.

The proposed method consists in a process that first extracts phonetic transcriptions, and then iteratively filters them. In order to initialize the process, an alignment dictionary is used to detect word boundaries. A rule-based grapheme-to-phoneme (G2P) generator (LIA_PHON [1]), a knowledge-based approach (JSM [2]), and a Statistical Machine Translation (SMT)-based system [9] were evaluated for this alignment. As a result, on the ESTER 1 French broadcast news corpus, we were able to obtain a decrease of the Word Error Rate (WER) on segments of speech with proper nouns, without negatively affecting the WER on the rest of the corpus.

Keywords: Speech recognition, Phonetic transcription, Proper nouns, SMT, Moses, G2P

Email addresses: antoine.laurent@lium.univ-lemans.fr (Antoine Laurent),
sylvain.meignier@lium.univ-lemans.fr (Sylvain Meignier),
paul.deleglise@lium.univ-lemans.fr (Paul Deléglise)

URL: <http://www-lium.univ-lemans.fr>, <http://www.specinov.fr>
(Antoine Laurent)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1. Introduction

This work focuses on an approach for enhancing automatic phonetic transcription of proper nouns.

Proper nouns constitute a special case when it comes to phonetic transcription, at least in French, which is the language used for this study. Indeed, there is much less predictability in how proper nouns may be pronounced than for regular words. This is partly due to the fact that, in French, pronunciation rules are much less normalized for proper nouns than for other categories of words: a given sequence of letters is not guaranteed to be pronounced the same way in two different proper nouns.

The lack of predictability also finds its roots in the wide array of origins proper nouns can come from: the more foreign the origin, the less predictable the pronunciation, with variations covering the whole range from correct pronunciation in the original language to a Frenchified interpretation of the spelling.

The high variability induced by this low predictability is a source of difficulty for Automatic Speech Recognition (ASR) systems when dealing with proper nouns. For an ASR system, being confronted with a proper noun pronounced using a phonetic variant very remote from any variant present in its dictionary is a situation similar to encountering an unknown word, if the language model cannot compensate for the acoustic gap. Such errors can have a strong impact on word error rate (WER): according to a comparative study of out-of-vocabulary impact of words in spontaneous and prepared speech [10] the recognition error on an out-of-vocabulary word propagates through the language model to the surrounding words, causing a WER of about 50 % within a window of 5 words to the left and to the right (again, in French). This highlights that the influence of the quality of the phonetic dictionary of proper nouns extends further than just the recognition of proper nouns themselves. It is particularly true in the case of applications where proper nouns are frequently encountered, such as transcription of broadcast news. However, aside from its potential impact on WER, accurate recognition of proper nouns can also be very important—independently of the frequency of their occurrence—in other contexts such as in the case of automatic indexing of multimedia documents, or transcription of meetings.

Setting up a phonetic dictionary of proper nouns (or any other class of words) requires grapheme to phoneme (G2P) conversion, be it manual or automatic. Automatic G2P conversion techniques are widely studied in the literature. The au-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

thors of [6] present an overview of techniques in 1999 and propose to classify the G2P systems into two categories: the knowledge-based approaches, which use existing linguistic knowledge to derive pronunciations, and the data-driven approaches, which derive pronunciation models from acoustic data. Knowledge-based approaches are further divided between formalized (e.g. rule based) and non-formalized (e.g. dictionary lookup). [11] proposes a dictionary look-up strategy (non-formalized knowledge-based). The authors of [1, 12, 3] present rule-based knowledge-based techniques. The authors of [3] propose a rule-based strategy that integrates different type of features (orthographic, syllabic, morphological, ...) to describe the rule context. A large variety of knowledge-based techniques are proposed in the literature: [13, 14, 15, 7] propose local classification strategies and [16, 17, 2] propose some pronunciation-by-analogy approaches. Many data-driven (acoustic-based) strategies can also be found in the literature ([18, 19, 20, 21]).

We propose an acoustic-based method to build a dictionary of phonetic transcriptions of proper nouns by using an iterative filter to retain the most relevant parts of a large set of phonetic variants, the latter being obtained by combining three G2P methods with extraction from actual audio signals [22].

In this work for French, we compare three different G2P systems to initialize the process, and we use a two-level iteration to converge on the best filtered dictionary. The iterative filter is applied in order to reduce noise by invalidating the variants that are deemed irrelevant because never used, and the ones that are found to be too prone to generating confusion with other words.

First, related works will be presented. After proposing an overview of the method, we will focus on the grapheme-to-phoneme systems used to initialize the process. In the next part, the proposed method will be described, and before concluding, experiments and results will be introduced and commented on. The intermediate (before filtering) and final sets of phonetic transcriptions are evaluated in terms of Word Error Rate (WER) and Proper Noun Error Rate (PNER), computed over the corpus of French broadcast news from the ESTER evaluation campaign [23].

2. Related works

Many G2P systems are presented in the literature. Several names are attributed to this task: grapheme-to-phoneme conversion [24, 17], phonetic pronunciation modeling [25], letter-to-sound translation [26], letter-to-phoneme conversion [27,

1
2
3
4
5
6
7
8
9 7], phonetic baseform generation [28, 29], phonetic transcription [30], text-to-
10 phoneme mapping [31], among others.

11 The simplest strategy to get phonetic transcriptions of a word is the dictionary
12 look-up, which consists in searching in a human-made phonetic dictionary. Mak-
13 ing such a dictionary is costly and time-intensive. We have at our disposal the
14 BDLEX dictionary [11]. This dictionary has the advantage of providing a very
15 complete and accurate set of transcriptions for each word it contains. However, it
16 only contains a limited number of entries, and more importantly for our case, it
17 does not contain any proper noun.

18 Rule-based conversion techniques have been developed in order to overcome
19 the kind of issues mentioned above. A rule-based phonetic transcription system
20 generate the possible chains of phones by relying exclusively on the spelling of
21 words. It offers the advantage of providing phonetic variants even for words for
22 which no speech signal is available. In the case of proper nouns, it generates the
23 most “common-sense” variants, *i.e.* the ones which people would use when they
24 have no *a priori* knowledge of the pronunciation of a particular proper noun. It
25 would be prohibitively difficult to establish the complete set of rules needed to
26 automatically find all the possible phonetic transcriptions of every proper noun.

27 In order to do so, an ideal automatic system would have to be able to detect
28 not only the origin of the proper noun, but also the various ways people might
29 pronounce this noun according to their own cultural and linguistic idiosyncrasies.
30 Unfortunately, both tasks are still open problems.

31 In the rest of this section, we will focus on a third approach: *data-driven* G2P
32 conversion systems based on the use of acoustic data. A thorough description of
33 the other methods can be found in [17] and in [2].

34 [32] and [33] use the A* algorithm to find the best phonetic transcriptions
35 from a set of acoustic representations of words. They use a heuristic function to
36 find the phonetic transcription that maximizes the likelihood from a set of acous-
37 tic representations. This method is based on the assumption that one phonetic
38 transcription only is enough to represent a word.

39 [21] improved that approach by computing that heuristic from the best path of
40 every acoustic representation. Unfortunately, the heuristic is too optimistic in con-
41 ditions of high inter-utterance variability. [34] proposes a method to suppress this
42 problem with the introduction of a pre-selection strategy which restricts search to
43 a confusion network built from heuristics. [35] develops a method that consists in
44 searching the k best phonetic transcriptions from a set of extracted pronunciations.
45 Two decision criteria are tested. The first criterion is based on transcription occur-
46 rence frequency, and the second on the maximization of likelihood. The method
47
48
49
50
51
52
53
54
55
56
57
58

1
2
3
4
5
6
7
8
9 that gives the best results is the one based on likelihood maximization. For each
10 acoustic realization, the n -best list (with n set experimentally to 50) is constructed
11 and constrained by the likelihood maximization of the union of those lists. [36]
12 uses the first criterion: the selection of the k most frequently extracted phonetic
13 transcriptions.
14

15 The authors of [30] propose a beam search approach with a two-level (*in-*
16 *tra-arc*, *arc*) pruning criteria. At least 10 samples are needed to get a reason-
17 able (between 5 and 10%) Phoneme Error Rate (PER). The PER is the average
18 edit distance between the found phonetic transcription and the reference phonetic
19 transcription.
20

21 [20, 37] develop a method based on an acoustic phonetic decoding for the
22 addition of words to the personalized vocabulary of their users. To do this, users
23 have to repeat one or two times every word they want to add to their lexicon.
24 [38, 29] describes an almost similar acoustic-phonetic decoding system, which
25 requires the user to repeat the various words to phonetize. Every user has to
26 pronounce twelve different proper nouns and has to call 10 times from different
27 phones (cellular and landline) and in several different acoustic environments (hall,
28 cafeteria, ...). The decoding strategy is based on the combination of speaker-
29 independent acoustic models and a language model that represents the transition
30 probabilities between various phonemes.
31

32 The work presented in [16] is based on the use of a bi-directional n -gram joint
33 sequence model. This model can be used to get a phonetic transcription of a word
34 thanks to its spelling or by using an acoustic representation of it.
35

36 In this part, we only focus on G2P acoustic *data-driven* strategies. Many more
37 G2P methods are compared in [2] on various corpora in English.
38
39
40
41

42 **3. Overview of the proposed method**

43
44 We propose a strategy that allows the extraction of phonetic transcriptions of
45 proper nouns from utterances. It is a multi-step, iterative process. The first step
46 consists in isolating portions of signal corresponding to proper nouns, using the
47 textual transcription of the audio and a forced alignment to get word boundaries.
48 During our study, we noticed that the dictionary used to perform this initialization
49 step had a great influence on our results. The use of bad phonetic transcriptions
50 results in boundary detection errors. Three different G2P techniques were com-
51 pared for this study: a rule-based phonetic transcription generator (LIA_PHON
52 [1]), a Joint-Sequence Model based method (JSM [2]), and a Statistical Machine
53 Translation based grapheme-to-phoneme converter (SMT [9]).
54
55
56
57
58

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Portions of the speech signal assumed to be corresponding to proper nouns are then extracted and fed to an APD (Acoustic Phonetic Decoding) system to obtain their phonetic transcription. Thus, proper nouns which are present several times in the corpus potentially get associated with several distinct phonetic transcriptions. APD yields a high number of phonetic transcriptions per proper noun (specific figures for our experimental corpus can be found in section 8.1). However some of the extracted transcriptions may be flawed: often, some phonemes of neighboring words are added or deleted at the end or at the beginning of the phonetic transcription, and some wrong phonemes are inserted in noisy conditions. Also, the high number of transcriptions increases the risk of generating confusion with other words. Proper nouns could erroneously appear in the ASR output instead of words from other categories. Therefore, it can negatively impact the quality of the decoding for the rest of the corpus. In order to avoid these problems, the result of the extraction is filtered to discard unfit phonetic transcriptions.

The proposed method for filtering is iterative: the filtered dictionary of each iteration is used again to carry out the alignment step, and the process starts again. This process is repeated until two consecutive filtered dictionaries are exactly the same. At least one phonetic transcription of each proper noun is always kept in the proper noun dictionary (*i.e.* there is no out-of-vocabulary word in the ASR lexicon). The method was trained and evaluated using broadcast news in French composed of French, European and world news reports. These data contain a high number of foreign journalist names.

4. Initial dictionary generation

4.1. Rule-based generation of phonetic transcriptions

The rule-based generator we used is LIA_PHON [1]. LIA_PHON is available under the GPL license. It participated in the ARC B3 evaluation campaign of French automatic phonetizers, in which phonetic transcriptions generated by the systems were compared with phonetizations made by human experts. This campaign was held in 1997, and results were published by [39] in 1998. Error rate was calculated according to the same principle as for the classical word error rate used in speech recognition. Compared to human-made phonetic transcriptions, 99.3% of the transcriptions generated by LIA_PHON were correct (for a total of 86938 phonemes) (This measure is computed at the phone level). However, results reveal that transcription errors were not distributed evenly among the various classes of words: erroneous transcription of proper nouns represented 25.6% of the errors generated by LIA_PHON even though proper nouns only represented

5.8% of the test corpus. This reflects poorer performance by LIA_PHON on this class of words.

4.2. Data-driven conversion techniques

In this section, we describe a G2P system based on the use of Joint-Sequence Models (JSM) and a conversion technique based on the use of a Statistical Machine Translation (SMT) system. Both these systems need a bitext corpus for the training step.

4.2.1. Bitext corpus format for data-driven methods

To convert graphemes to phonemes, a bitext associates sequences of letters with sequences of phonemes. Table 1 shows examples of two representations of the bitext corpus, denoted by A and B. In representation A, the sequence of letters corresponds to a word. In representation B, the sequence of letters corresponds to a group of words. A symbol is added to mark the boundary of each word and each phonemic representation of the words. This representation allows to differentiate inter- and intra-word influence. In order to build a bitext corpus for representation B, every sequence of words of the training corpus between two fillers (silence, music, laughter, hesitation, ...) is aligned using the baseline acoustic models and the baseline dictionary. Our baseline dictionary contains variants that take into account the interword coarticulation influence (liaisons in French). Indeed, we hypothesize that the influence of a word on the pronunciation of its neighbors is negligible when they are separated by a filler. Representation B includes word boundaries within each sequence using a dedicated marker symbol.

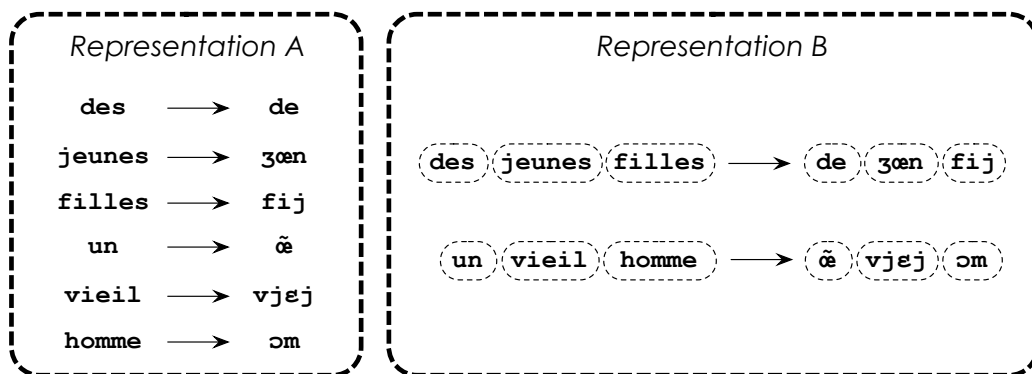


Figure 1: Examples of representations A and B of the bitext corpus

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

4.2.2. Joint-Sequence models (JSM)

This system is a *data-driven* conversion system, available under the GPL license. The system is based on the idea that, given enough examples, it should be possible to predict the pronunciation of unseen words, purely by analogy. The use of joint-sequence models to convert graphemes to phonemes [2] will be denoted by JSM in the rest of this article. JSM being a data-driven conversion system means that we have to give it pronunciation examples in order to train it. Training takes a pronunciation dictionary and creates new model files successively, starting with unigram models and up to 6-gram models. Model files can then be used to transcribe words that were not in the dictionary. The fundamental idea of joint multigram model is that for each word, its orthographic form and its pronunciation are generated by a common sequence of *graphones*. A *graphone*, or grapheme-phoneme joint multigram is a pair $q = (g, \varphi) \in Q \subseteq G^* \times \phi^*$ of a letter sequence g and a phoneme sequence φ of possibly different length. G represents all the letters of the alphabet, ϕ represents the inventory of phonemic symbols, Q represents the inventory of graphones. For example, the pronunciation of “jeunes” may be regarded as a sequence of three graphones:

$$\begin{array}{cccc} \text{“jeunes”} & & \text{j} & \text{eu} & \text{nes} \\ & & = & & \\ & & \text{ʒ} & \text{œ} & \text{n} \end{array}$$

The procedure for having the alignment between graphemes and phonemes is described in [8]. The joint probability distribution $p(\varphi, g)$ is modeled using a standard M -gram:

$$p(q_1^L) = \prod_{i=1}^{L+1} p(q_i | q_{i-1}, \dots, q_{i-M+1}) \quad (1)$$

Phonetic transcriptions are then obtained from words by searching the most likely graph sequence matching the given spelling and projecting it onto the phonemes.

Because computing time on representation B is very expensive using JSM, it is trained only on representation A.

4.2.3. Grapheme to phoneme conversion using Statistical Machine Translation (SMT)

We proposed a method in [9], based on the open source Moses toolkit [40] to convert graphemes to phoneme sequences.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

A Statistical Machine Translation system (SMT) is used to transform text from a source language into a target language. The training step needs a data corpus which is composed of bitext data: source language sentences associated with their translation in the target language.

The SMT system is based on the Moses toolkit. This toolkit is commonly used to translate data in which the elementary unit is the word in both the source and target parts.

The training of a grapheme-to-phoneme translation model is similar to the training of a translation model as described in the Moses documentation.

SMT models. First, the bitext corpus has to be aligned at word level in both directions (source to target and target to source). The phrase pairs are extracted using some heuristics known as *diag-grow-final* which start from the intersection of the two alignments and then adds additional alignment points. After extraction, the phrase pairs are scored. A standard translation model contains 5 different scores, namely direct and inverse phrase translation probabilities, direct and inverse lexical probabilities and a phrase insertion penalty (always set to e^1). Another component of a standard SMT system is the lexicalized reordering model. A distortion model is a model that allows phrase (sequence of words) permutation. As presented in figures 2 and 3, this model takes into account three different features corresponding to three kinds of reordering, namely monotone (phrase pairs are adjacent and in the same order), swap (phrase pairs are adjacent and in the reverse order), and discontinuous (the phrase pairs are not adjacent). For each phrase pair, the relative frequency of each kind of reordering is calculated (a smoothing technique is applied to avoid zero probabilities for unseen orientations). The last main component of a SMT system is the language model which is trained on the target side of the bitexts and all available monolingual data in target language.

Figure 4 shows how the SMT is learned and used for the translation of graphemes to phonemes. We trained a 4-gram language model composed of phonemes learned from a phonemic forced alignment of the ESTER 1 training corpus. The bitext corpus is used to produce a translation model. However two training strategies are proposed: the first one corresponds to the standard Moses training framework based on the maximization of BLEU [42]. The second, based on the Levenshtein metric, minimizes insertion, deletion, and substitution errors of phonemes.

BLEU score. The BLEU score is commonly used for the optimization in order to have the best translation system according to this measure. Training reserves 3% of the corpus for optimization of the parameters according to the BLEU score.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

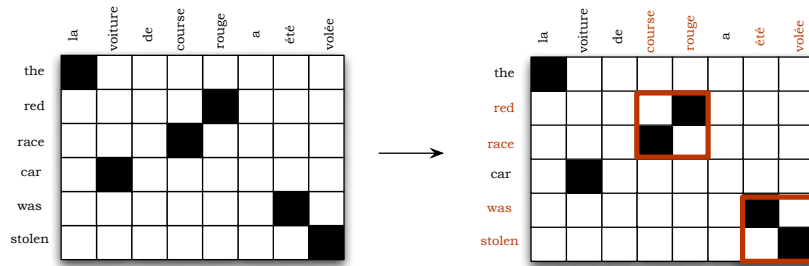


Figure 2: Bitext alignment matrix

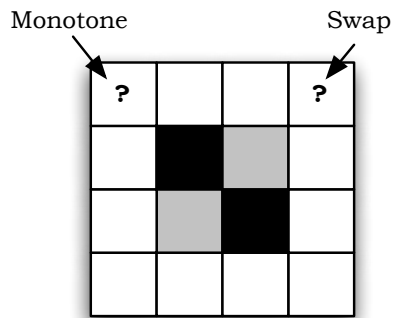


Figure 3: Orientation determination for distortion model

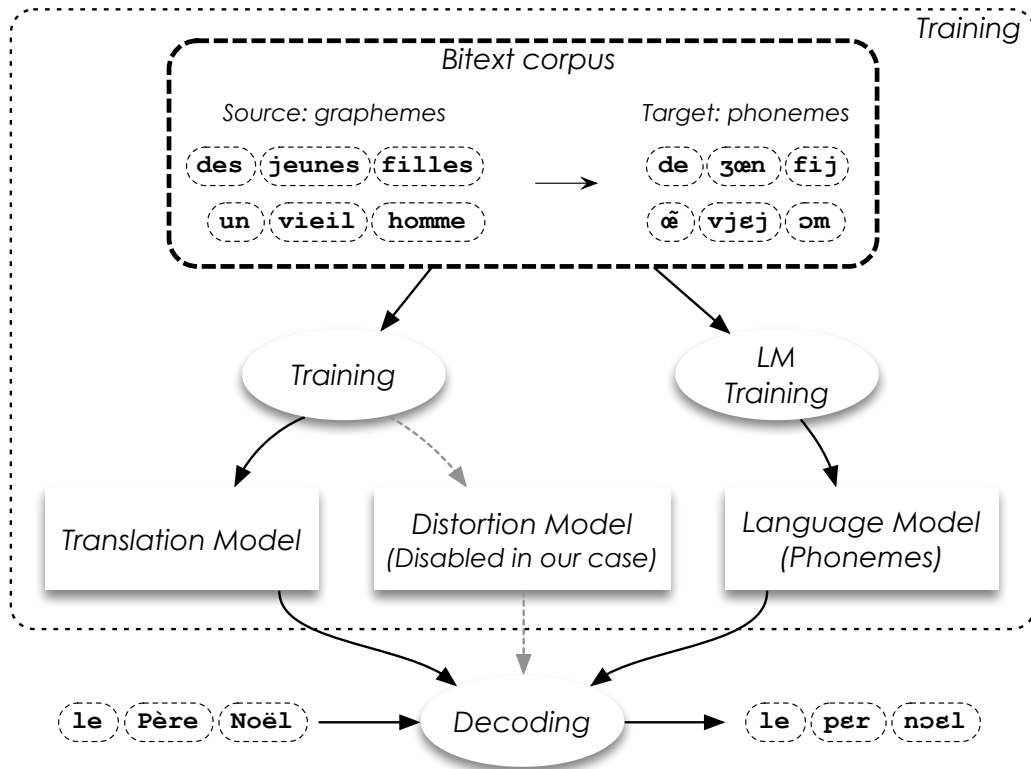


Figure 4: Using SMT for grapheme-to-phoneme translation

Experiments show that the best score (in terms of WER computed over our development corpus) is obtained by using a distortion model for representation A, while for representation B, the best score is obtained without a distortion model. Weights of the different models were trained using the Minimum Error Rate Training ([41]).

Levenshtein score. For our task, we decided to try using a normalized Levenshtein edit distance for parameter optimization.

At the end of a training iteration, 3-best phonetic transcriptions for each training example (sequence of letters) are generated using the current translation model. When using the Levenshtein score optimization, we only optimized the five weights of the translation model (the distortion model was disabled). Assuming that we want to convert a sequence of graphemes \tilde{g} to a sequence of phonemes \tilde{p} , those weights are the phrase translation probability $\varphi(\tilde{p}|\tilde{g})$, the lexical weighting $lex(\tilde{p}|\tilde{g})$, the phrase translation probability $\varphi(\tilde{g}|\tilde{p})$, the lexical weighting $lex(\tilde{g}|\tilde{p})$,

1
2
3
4
5
6
7
8
9 and the phrase penalty.

10 The sum of the normalized Levenshtein measures, S , is computed between the
11 phonetic transcriptions and the references (equation 2).
12

$$13 \quad S = \sum_{e \in E} \log(1 - \min(\forall n \in [1, 3] \frac{d(p_e^n, r_e)}{\max(l_{p_e^n}, l_{r_e})})) \quad (2)$$

14 where p_e^n is the phonetic transcription n of the example e . As stated before, we
15 consider the 3-best phonetic transcriptions, thus n vary from 1 to 3. $d(p_e^n, r_e)$ is
16 the edit distance of Levenshtein of the phonetic transcription p_e^n , with r_e the refer-
17 ence phonetic transcription for example e . E is the set of the generated phonetic
18 transcriptions. $l_{p_e^n}$ is the length of the phonetic transcription p_e^n of the example e
19 and l_{r_e} is the length of the reference phonetic transcription (r_e). Every log argu-
20 ments are floored at 10^{-7} to avoid that just one bad phonetic transcription could
21 impact the measure of the entire database.
22

23 Until getting the lowest S over all the training examples, a simplex framework¹
24 is used to tune the model parameters.
25

26 When using the Levenshtein optimization, the language model weight is set to
27 0.1 and the word penalty weight is set to 0.
28

29 For the task of grapheme-to-phoneme conversion, the best results were ob-
30 tained by using the Levenshtein optimization and representation B. Learning time
31 on our training corpus (ESTER 1 Training corpus, see section 7.1 for details) is
32 about 13 times more for JSM (175.5 hours) than for SMT (13.5 hours).
33

34 5. Extraction of phonetic transcriptions using acoustic-phonetic decoding

35 5.1. Method

36 In order to enrich the set of phonetic transcriptions of proper nouns with some
37 less predictable variants, we gather actual utterances of proper nouns by actual
38 people. This process relies on an acoustic-phonetic decoding system (APD),
39 which generates a phonetic transcription of the speech signal.
40

41 In a corpus consisting of speech with a manual word transcription, portions
42 of the speech signal corresponding to proper nouns are extracted. They are then
43 fed to the APD system to obtain their phonetic transcription. Since the phonetic
44 decoding results of various utterances can be different, proper nouns which are
45
46
47

48 ¹Thanks to the Condor toolkit [43]
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

present several times in the corpus potentially get associated with several different phonetic transcriptions each.

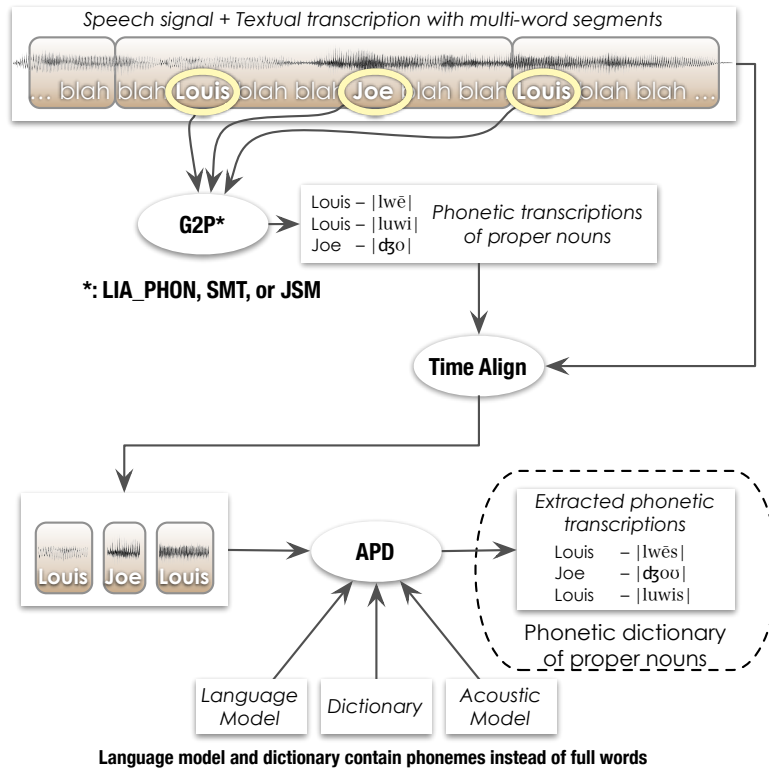


Figure 5: Use of the acoustic-phonetic decoding system

5.2. Proper noun boundaries

As explained above, the first step consists in isolating the portions of signal corresponding to proper nouns using word transcription. However, in the manual transcription we used, words were not aligned with the signal. Start and end times of individual words were not available; only longer segments (composed of several words) had their boundaries annotated. Therefore, the start and end times of each word of the transcription had to be determined by aligning the words with the signal, using a speech recognition system (see figure 5 – “Time Align” step).

The phonetic transcriptions used for proper nouns during this forced alignment were provided by three different G2P systems (see figure 5 – “G2P*” step).

In figure 5, we have two different boxes that contain phonetic transcriptions. The first one represents phonemes that we get directly by using one of our three

1
2
3
4
5
6
7
8
9 different G2P systems (LIA_PHON, SMT, or JSM). The second box represents the
10 phonetic transcriptions that we get from the signal, at the output of our Acoustic
11 Phonetic Decoding system.
12

13 14 *5.3. Effect of inaccurate boundary detection*

15 Because phonetic transcription is not very reliable using these three different
16 G2P systems, boundaries of proper nouns are not very accurate. Portions of signal
17 detected as proper nouns might overlap with neighboring words. As a result,
18 when applied to such portions of signal, the APD system might generate erroneous
19 phonemes at the beginning and/or at the end of the proper nouns, which might in
20 turn introduce errors when the flawed phonetic transcriptions are later used for
21 decoding.
22
23

24 25 *5.4. APD based phonetic transcription*

26 When boundaries of the proper nouns have been determined, APD is applied
27 to the corresponding portions of the signal. The decoding path gives a series of
28 phonemes considered as the phonetic transcription of the proper noun.
29

30 As noted in [30], unconstrained phonetic decoding does not allow the system
31 to obtain reliable phonetic transcriptions. Our own experiments lead us to the
32 same conclusion. The use of a language model allows for some level of guidance
33 for the speech recognition system: it does so by minimizing the risk of having
34 phoneme sequences with a very low probability appear in the transcription results.
35 We set constraints by using tied state triphones and a 3-gram language model as
36 part of the decoding strategy, to generate the best sequence of phonemes. While
37 this setup is close to a speech recognition system, the dictionary and language
38 model contain phonemes instead of full words. The trigram language model was
39 trained using the phonetic dictionary used during the 2005 ESTER evaluation
40 campaign [44]. It contains about 65,000 lexical entries of words, and was gener-
41 ated using BDLEX and LIA_PHON. Only the words which were not part of the
42 BDLEX corpus were phonetised automatically using LIA_PHON. Words which
43 were identified as proper nouns were deleted from this dictionary before learning
44 our 3-gram language model for phonemes.
45
46
47
48
49

50 51 **6. Filtering of phonetic transcriptions**

52 53 *6.1. Motivation*

54 The extraction of phonetic transcriptions for utterances yields an average of 6
55 phonetic transcriptions per proper noun in our experiments (complete results for
56 our experimental corpus can be found in section 8.1).
57

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

However, as stated in the previous section, some of the extracted transcriptions may be flawed. Also, the high number of transcriptions increases the risk of some phonetic transcriptions of proper nouns being erroneously used to decode words of another type. Therefore, it can negatively impact the quality of the decoding for the rest of the corpus. Given that the number of occurrences of the other categories of words is expected to be much higher than the number of occurrences of proper nouns, there is a risk of seeing any gain in performance for proper nouns being outbalanced by a negative impact on the rest of the corpus and on the global WER. The goal of this filtering is to detect and remove the phonetic variants of proper nouns that are the most likely to generate confusion with other words.

6.2. Iterative filtering

In order to minimize the risk of negatively affecting the global WER, it is desirable to filter the set of phonetic transcriptions and keep only the most appropriate. We propose an iterative filtering method to select only those transcriptions deemed to be reliable enough. We have already proposed a different approach to select phonetic transcriptions in previous work [45]; however this early attempt was rendered impractical because of its execution time which was directly proportional to the number of extracted phonetic transcriptions. For a proper noun present in s segments, with v phonetic transcriptions, it was necessary to decode $v \times s$ segments to validate or invalidate the overall set of phonetic variants for this proper noun.

In the present work, we have managed to detect and remove phonetic variants of proper nouns generating confusion with other words by decoding the development corpus using the newly built phonetic dictionary (as well as a separate phonetic dictionary for all the other categories of words, of course). This decoding is unconstrained, with no forced alignment.

Any phonetic transcription that was never used to decode the corresponding proper noun in the right place gets removed from the dictionary, since it either caused an error or was not used at all. However, a heuristic is set in order to keep at least one phonetic transcription for each proper noun.

The process then gets repeated: the corpus is decoded again using the modified dictionary, which then gets filtered according to the results of this decoding. The whole decoding/filtering process is repeated until no more phonetic transcriptions are removed from the dictionary.

This process is illustrated in figure 6, using the same example data as in figure 5.

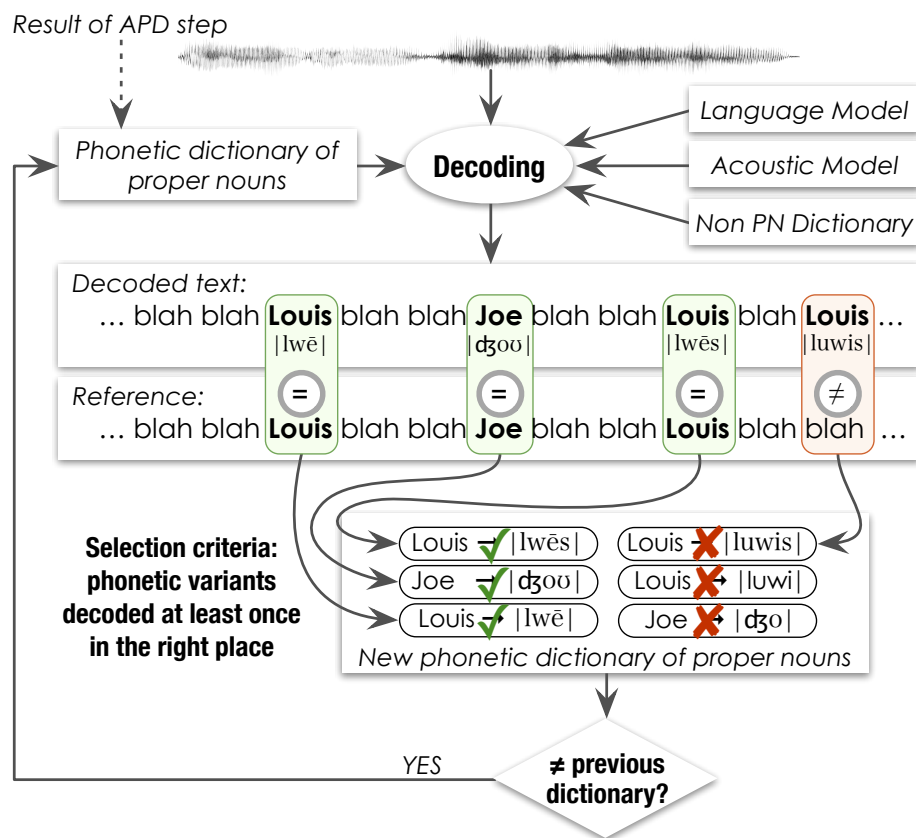


Figure 6: Illustration of iterative filtering of phonetic transcriptions. The initial value of the phonetic dictionary of proper nouns is the union of rule-based and extracted transcriptions.

6.3. Two-level iterative filtering

As stated earlier, the alignment dictionary used to initialize the process has a strong impact on the accuracy of the phonetic transcriptions generated. For this reason, we have decided to rerun the whole process, this time using the iteratively filtered dictionary (the output of the iterative filtering described above) instead of G2P systems to get boundaries of proper nouns inside the audio data during the forced alignment step. This allows the system to call proper noun boundaries into question with the newly built dictionary.

This extraction+filtering cycle, illustrated in figure 7, is repeated until two consecutive iteratively filtered dictionaries are exactly identical.

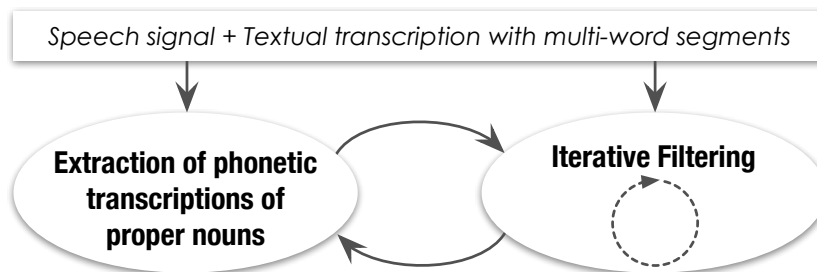


Figure 7: Overview of the double iterative process: the filtered dictionary is used for the initialization of the next extraction+filtering cycle, until the result is stable

7. Experiments

7.1. Corpus

Our experiments were carried out on the ESTER 1 corpus. ESTER is an evaluation campaign of French broadcast news transcription systems which took place in January 2005 [23]. The ESTER corpus was divided into three parts: training, development, and evaluation. The training (81 hours) and the development (12.5 hours) corpora are composed of data recorded from four radio stations in French (France Inter, France Info, Radio France Internationale, and Radio TV Maroc). The test corpus is composed of 10 hours coming from the same four radio stations plus two other stations (France Culture and Radio Classique), all of which were recorded 15 month after the development data. Each corpus is annotated with named entities, allowing easy spotting of proper nouns.

The training corpus was used to learn our automatic speech recognition system. The training corpus and the development corpus are jointly employed to extract phonetic transcriptions and to filter them. The JSM and SMT grapheme-to-phoneme converters were also trained over the ESTER 1 training corpus.

7.2. Metrics

The intermediate and final sets of phonetic transcriptions were evaluated in terms of Word Error Rate (WER) and Proper Noun Error Rate (PNER). PNER is computed the same way as the WER, but it is computed only for proper nouns and not for every word:

$$PNER = \frac{I + S + E}{N} \quad (3)$$

with I being the number of wrong insertions of proper nouns, S the number of substitutions of proper nouns with other words (where the reference word is a

proper noun), E the number of elisions of proper nouns, and N the total number of proper nouns.

The use of PNER as a metric reflects the goal of this work, which is to enhance the recognition of proper nouns, and not merely have an accurate chain of phonemes.

While PNER allows to evaluate the quality of the detection of proper nouns, WER is used to evaluate the impact of the new phonetic transcriptions on the whole test corpus.

7.3. Acoustic and language models

The decoding system is based on CMU Sphinx 3.6 [46].

Our experiments were carried out using a one-pass decoding coming from the LIUM ESTER 1 system [44], using 12 MFCC acoustic features plus the energy, completed with their primary and secondary derivatives. Acoustic models were trained on the ESTER training corpus. These models are composed of 5500 tied states, each state being modeled by a mixture of 22 diagonal Gaussians. Decoding employs tied-state word-position 3-phone acoustic models which are made gender- and bandwidth-dependent through MAP adaptation of means, covariances and weights. The trigram language model was learned on three different data sources :

- On the manual transcriptions of our training and development corpus (81h training + 12.5h development = 93.5 hours recorded from the four radio stations). These transcriptions contain about 1.35M occurrences of 34k distinct words.
- On the articles coming from the French newspaper “Le Monde” from the year 2003 (19M occurrences of 220k distinct words).
- On articles coming from the French newspaper “Le Monde” from 1987 to 2002 (300M word occurrences).

Three 3-gram language models were learned: one using the 81h of our training corpus, and the others on the two other data sources. A linear interpolation was performed to minimize perplexity on the remaining 12.5 hours of data coming from the development corpus. The vocabulary contains all of the 34k distinct word of the manual transcriptions, and words appearing more than ten times in the 2003 articles (about 19k words). The most frequent words in the rest of the articles from

1
2
3
4
5
6
7
8
9 “Le Monde” (from 1987 to 2002) are used to complete the vocabulary, up to 65k
10 words.

11 Using this vocabulary, all the textual data of the training corpus is used to
12 train a trigram language model. To estimate this model, the SRILM toolkit [47]
13 is employed using the modified Kneser-Ney discounting method. Unigrams and
14 bigrams are all kept, but trigrams occurring only once are discarded.

15 The language model includes all the proper nouns present in the development
16 corpus. All the dictionaries contain the same proper nouns, with only their pho-
17 netic transcriptions varying.
18
19
20
21

22 8. Results

23 8.1. Number of phonetic transcriptions per proper noun

24 Table 1 presents the number of phonetic transcriptions generated with the three
25 G2P methods. The ESTER 1 corpus (development plus training) contains 3,348
26 distinct proper nouns, appearing 28,866 times.
27
28
29
30

31 Table 1: Number of phonetic transcriptions generated by each method

32 Method	33 Generated variants (G2P)	34 Extracted variants (APD)	35 After 1 iteration (All process)	36 After 2 iterations (All process)	37 After 3 iterations (All process)
38 LIA_PHON	4,364	20,218	6,776	6,524	6,502
39 SMT	7,031	20,184	7,065	6,813	6,802
40 JSM	3,626	20,008	6,876	6,711	6,708
41 Average	5,007	20,137	6,906	6,683	6,671

42 On average, the number of phonetic transcriptions between G2P generation
43 and APD extraction grows from 5k to 20k. We only consider the best hypothesis
44 generated by the APD.
45

46 One pass of iterative filtering keeps about 7k phonetic transcription variants
47 from the 20k variants generated by the APD. For each of our three grapheme-
48 to-phoneme strategies, filtering is done in 3 iterations. The number of variants
49 contained in the final filtered dictionary slightly decreases compared to the first
50 iteration. The number of variants generated is stable across methods used for the
51 initialization. Figure 8 shows the overlap among the generated dictionaries. As we
52 can see, there are 2467 variants common to the 3 generated dictionaries. Figure
53
54
55
56
57
58
59
60
61
62
63
64
65

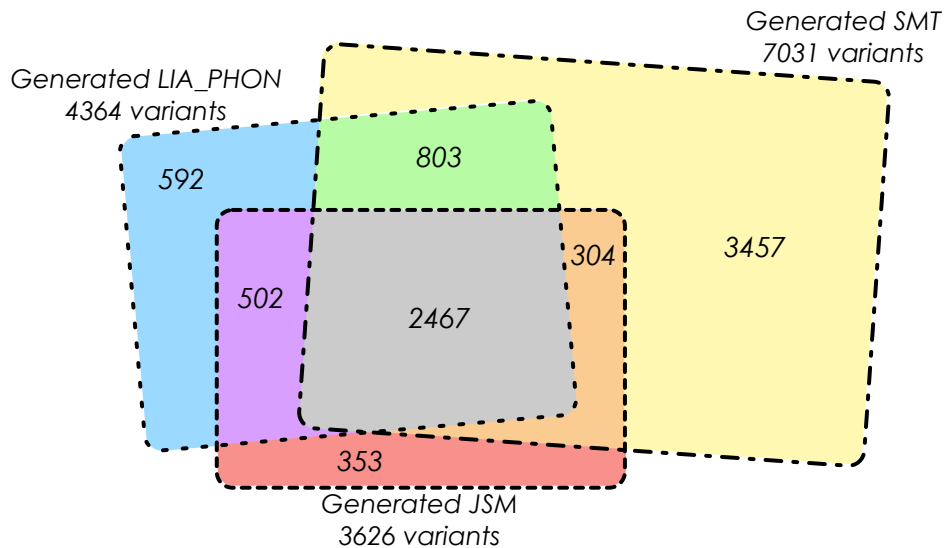


Figure 8: Overlap among the 3 initialization dictionaries

8.2. Results of the first iteration

This section compares the results obtained by directly using the three G2P methods with the use of the extraction and filtering of proper nouns.

Figure 10 shows the PNER obtained using the filtering method after the first iteration for each G2P system on the ESTER test corpus.

These results show that the filtering method produces significant gains in terms of PNER for every G2P system. As we can see, the APD method supplemented by the SMT-based grapheme-to-phoneme conversion system is the one that yields the lowest PNER.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

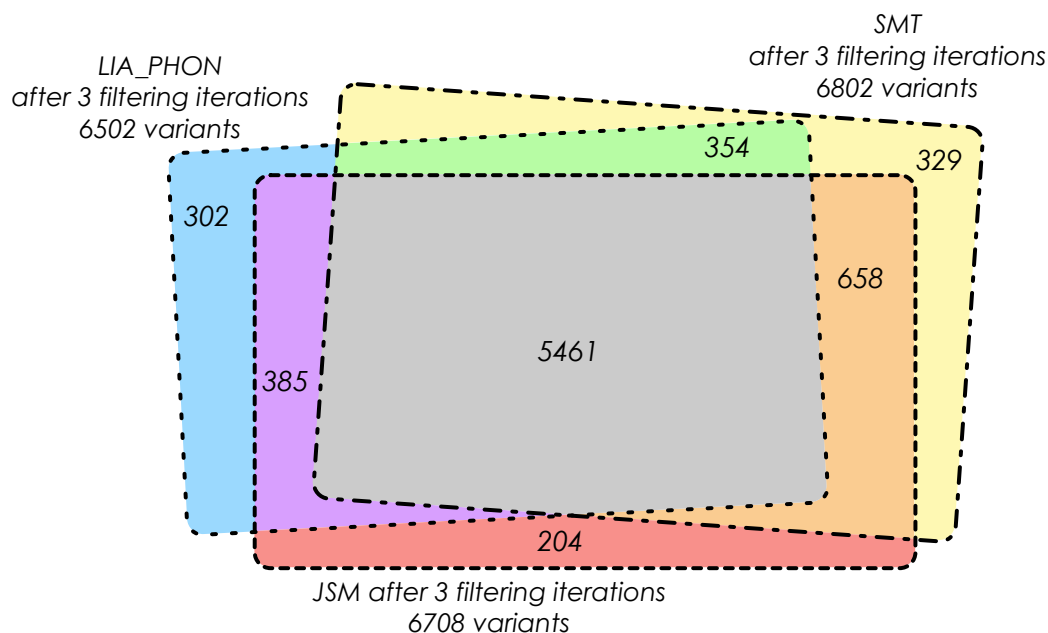


Figure 9: Overlap among the final dictionaries

As explained previously, phonetic transcriptions for non-proper nouns are taken from the BDLEX database, or generated by the rule-based grapheme-to-phoneme tool LIA_PHON for words which are not in the database. The generated dictionaries (SMT, JSM, and LIA_PHON) include the non-proper noun dictionary, supplemented by the phonetic transcriptions of all proper nouns generated using SMT, JSM, or LIA_PHON. Figure 11 compares the results obtained using the three generated dictionaries (SMT, JSM, and LIA_PHON) to initialize the method, in term of WER computed only over segments that contain proper nouns.

Figures 10 and 11 show the interest of filtering: it reduces both the PNER and the WER on segments with proper nouns.

8.3. Using iterative acoustic-based phonetic transcription

Table 2 shows the results obtained with the full iterative process initialized with LIA_PHON, SMT, and JSM G2P systems. The results in bold are those with the best gain in terms of WER and PNER. WER and PNER are computed on segments that contain proper nouns. We can see a small gap between the first filtering iteration and the last one. Using LIA_PHON to initialize our method, the WER decreased from 24.1% (after the first filtering iteration) to 24.0% (at the

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

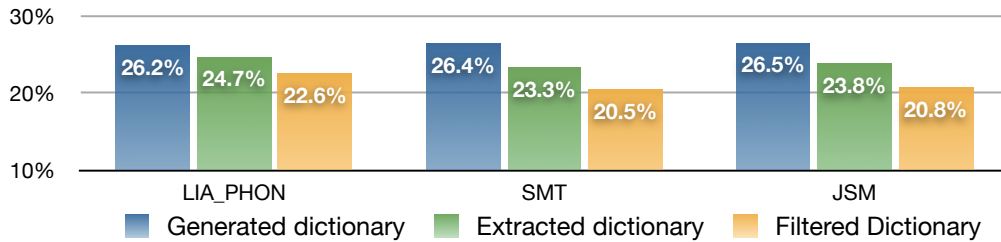


Figure 10: PNER using each G2P method (ESTER test corpus)

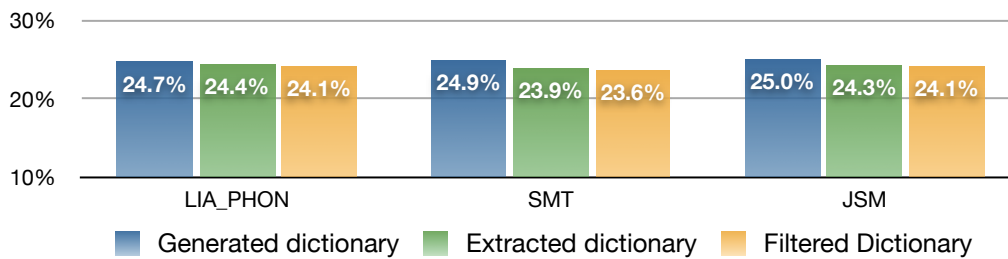


Figure 11: WER on segments with proper nouns in the test corpus

end) and the PNER decreased from 22.6% to 22.5%. With SMT, there is a gain of 0.2 point in terms of WER and a gain of 0.3 point in terms of PNER between the first and the last filtering iterations. Finally, when using JSM, the gains are of 0.2 point in terms of WER and 0.3 point in terms of PNER.

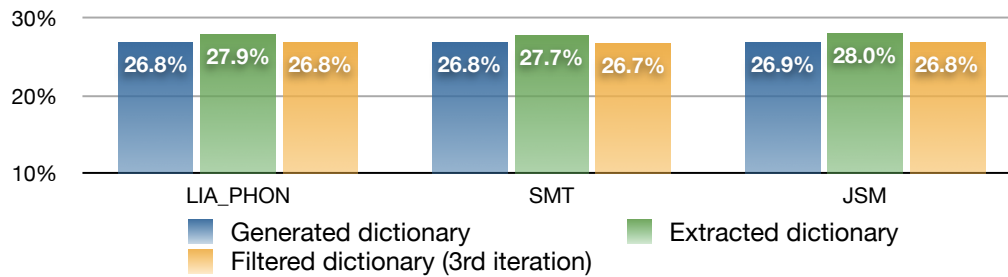


Figure 12: WER on every segment in the test corpus

Figure 12 shows the WER obtained on the whole ESTER 1 test corpus. The test corpus contains 11087 segments. 1412 of them contain proper nouns. With no filtering, extracted dictionaries, while improving the WER on segments that contain proper nouns, also increase the global WER. Errors are introduced: other

Table 2: WER and PNER using the full iterative process

Dictionary	WER (segments with PN)	PNER
LIA_PHON	24.7%	26.2%
SMT Generated	24.9% (+0.2)	26.4% (+0.2)
JSM Generated	25% (+0.3)	26.5% (+0.3)
Two-level filtering iteration 1		
LIA_PHON	24.1% (-0.6)	22.6% (-3.6)
SMT	23.6% (-1.1)	20.5% (-5.7)
JSM	24.1% (-0.6)	20.8% (-5.4)
Two-level filtering iteration 2		
LIA_PHON	24.1% (-0.6)	22.6% (-3.6)
SMT	23.5% (-1.2)	20.3% (-5.9)
JSM	24% (-0.7)	20.5% (-5.7)
Two-level filtering iteration 3		
LIA_PHON	24% (-0.7)	22.5% (-3.7)
SMT	23.4% (-1.3)	20.2% (-6)
JSM	23.9% (-0.8)	20.5% (-5.7)

words are substituted by proper nouns, and some proper nouns are wrongly inserted. The results show that with the filtering step, our method does not generate new errors with other word classes. The WER on segments with no proper nouns remains the same using filtered dictionaries as it is with the generated dictionaries. This highlights the role of filtering, which removes confusable variants from the lexicon.

8.4. Analysis of the results

In our evaluation corpus, 640 different proper nouns are present, with a total of 2080 occurrences. The proposed method decreases the PNER for 152 proper nouns, and increases the PNER for 26 of them. Most of those 152 proper nouns are foreign, therefore they do not follow the usual rules of pronunciation used in French. Examples of those nouns are: Jiantao, Fatima, Rumsfeld, Yahia, Ahmed.

When pronounced in Arabic (Radio TV Maroc station), certain proper nouns contain phonemes that are not present in our French phoneme set. Those phonemes are replaced with French phonemes.

As explained earlier, during filtering, a rule was set in order to avoid eliminating the last phonetic transcription variant of each noun. The average number of phonetic transcriptions per proper noun is about 2. It is only 1.3 for the 26

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 3: Some examples of phonetic transcriptions

Proper Nouns	SMT	Proposed method (initialized with SMT)
Jintao	ʒintao	jintao
Roger	ɾoʒe	ɾoɔʒœɾ
Decaens	dœkɛn	dœkã
Fatima	fatima	fatma
Rumsfeld	ɾyʉmsfɛld	ɾœmsfɛld
Yahia	jaja	jæʒja
Ahmed	amɛd	æʒkmɛd

proper nouns for which the PNER is increased. This actually corresponds to 20 proper nouns with only one variant, which would have been eliminated without this heuristic.

9. Conclusion

In this article, we proposed an iterative, two-step acoustic-based process for phonetic transcription generation, and applied it to the specific case of proper nouns.

The first step adopts a data-driven approach of building a dictionary of phonetic transcriptions, aiming for a closer match to actual usage of proper nouns than knowledge-based approaches can provide. This is accomplished through extraction of phonetic variants from actual audio signals, which is used to filter and enrich an initial set of phonetic transcriptions generated by a knowledge-base grapheme-to-phoneme system—filtering out unused variants and adding variants that the G2P system could not generate.

The second step of our method consists in filtering the resulting dictionary in order to avoid a negative impact on the other classes of words. Indeed, the extraction of phonetic transcriptions for proper nouns in the first step yields a high number of phonetic variants, which generates noise during the decoding. Many of these phonetic variants are too close to the pronunciation of other words of the dictionary. As a result, when used directly, this dictionary has a negative impact on the WER on segments that do not contain any proper noun. The goal of the iterative filtering process is the detection and removal of the phonetic variants that are the most likely to generate confusion with words from other classes.

1
2
3
4
5
6
7
8
9 The method loops, rerunning steps one and two over the resulting dictionary,
10 iterating until stability is reached.

11 The use of the resulting phonetic dictionaries of proper nouns yields a gain in
12 terms of PNER (Proper Noun Error Rate) and WER on the ESTER corpus. The
13 best results are obtained by using an SMT (Statistical Machine Translation [9])
14 system to generate the initial proper noun dictionary for the process. The WER
15 on segments that contain proper nouns decreased by 1.3 points and the PNER de-
16 creased by 6 points compared to the simpler, rule-based system. As was expected,
17 with the filtering step, the WER on segments without proper nouns is unaffected,
18 thus allowing the global WER to improve slightly thanks to better detection of
19 proper nouns.
20
21
22

23 Even though the impact on the global WER is only minor on a corpus such as
24 ESTER, improved detection of proper nouns is crucial for some tasks. An inter-
25 esting field where the proposed method is useful is named speaker identification,
26 which consists in the automatic extraction of speaker identities (first name and last
27 name) from the transcription [4, 5]. The new phonetic transcriptions generated by
28 the proposed method should contribute to render detection easier by improving
29 the decoding of proper nouns.
30
31

32 Finally, one of the advantages of the filtering method described here is that its
33 execution time is not linked to the size of the set of transcriptions to be filtered.
34 This opens up the possibility of applying it to other, larger classes of words.
35
36

37 **References**

- 38
39 [1] F. Béchet, LIA.PHON : un système complet de phonétisation de textes, in:
40 Traitement Automatique des Langues, Vol. 42, 2001, pp. 47–67.
41
42 [2] M. Bisani, H. Ney, Joint-sequence models for grapheme-to-phoneme con-
43 version, in: Speech Communication, Vol. 50, 2008, pp. 434–451.
44
45 [3] B. Réveil, J.-P. Martens, H. van den Heuvel, Improving proper name recog-
46 nition by means of automatically learned pronunciation variants, in: Speech
47 Communication, Vol. 54, 2012, pp. 321–340.
48
49 [4] E. El-Khoury, A. Laurent, S. Meignier, S. Petitrenaud, Combining
50 transcription-based and acoustic-based speaker identifications for broadcast
51 news; in: Proceedings of the International Conference on Acoustics, Speech
52 and Signal Processing (IEEE, ICASSP 2012), 2012, Kyoto, Japan, pp. 4377–
53 4380.
54
55
56
57
58

- 1
2
3
4
5
6
7
8
9 [5] L. Canseco-Rodriguez, L. Lamel J.-L. Gauvain, A comparative study using manual and automatic transcriptions for diarization, in: Proceedings of the Workshop on Automatic Speech Recognition and Understanding (IEEE, ASRU 2005), Puerto Rico, USA, 2005, Vol. 1, pp. 415–419.
- 10
11
12
13
14
15 [6] H. Strik, C. Cucchiarni, Modeling pronunciation variation for ASR: A survey of the literature, in: Speech Communication, Vol. 29, 1999, pp. 224–246.
- 16
17
18 [7] K. Seng, Y. Iribe, T. Nitta, Letter-To-Phoneme Conversion based on Two-Stage Neural Network focusing on Letter and Phoneme Contexts, in: Proceedings of the 12th Annual Conference of the International Speech Communication Association (ISCA, Interspeech 2011), Florence, Italy, 2011, pp. 1885–1888.
- 19
20
21
22
23
24
25 [8] M. Bisani, H. Ney, Investigations on joint-multigram models for grapheme-to-phoneme conversion, in: Proceedings of International Conference on Spoken Language Processing (ISCA, ICSLP 2002), Vol. 1, Denver, CO, USA, 2002, pp. 105–108.
- 26
27
28
29
30
31 [9] A. Laurent, P. Deléglise, S. Meignier, Grapheme-to-phoneme conversion using an SMT system, in: Proceedings of the 10th Annual Conference of the International Speech Communication Association (ISCA, Interspeech 2009), Brighton, England, 2009, pp. 708–711.
- 32
33
34
35
36
37 [10] R. Dufour, From prepared speech to spontaneous speech recognition system: a comparative study applied to French language, in: IEEE/ACM CSTST Student Workshop, Vol. 1, Cergy, France, 2008, pp. 595–599.
- 38
39
40
41
42 [11] M. de Calmès, G. Pérennou, BDLEX: a lexicon for spoken and written French, in: Language Evaluation and Resources Conference (LREC 1998), Grenada, Spain, 1998, pp. 1129–1136.
- 43
44
45
46 [12] J. Tihoni, G. Pérennou, Phonotypical transcription through the GEPH expert system, in: Proceedings of the European Conference on Speech Communication and Technology (ESCA, Eurospeech 1991), Vol. 1, Genoa, Italy, 1991, pp. 767–770.
- 47
48
49
50
51
52 [13] K. Torkkola, An efficient way to learn English grapheme-to-phoneme rules automatically, in: Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE, ICASSP 93), Vol. 2, Minneapolis, MN, USA, 1993, pp. 199–202.
- 53
54
55
56
57
58

- 1
2
3
4
5
6
7
8
9 [14] C. Ma, M. A. Randolph, An approach to automatic phonetic baseform generation based on bayesian networks, in: Proceedings of the International Conference on Speech Communication and Technology (ISCA, Interspeech 2001), Vol. 1, Aalborg, Denmark, 2001, pp. 1453–1456.
- 10
11
12
13
14
15 [15] K. Jensen, S. Riis, Self-organizing letter code-book for text-to-phoneme neural network model, in: Proceedings of the International Conference on Spoken Language Processing (ISCA, ICSLP 2000), Vol. 3, Beijing, China, 2000, pp. 318–321.
- 16
17
18
19
20
21 [16] L. Galescu, J. F. Allen, Bi-directional conversion between graphemes and phonemes using a joint n-gram model, in: Proceedings of 4th ISCA Tutorial and Research Workshop on Speech Synthesis, Perthshire, Scotland, 2001.
- 22
23
24
25 [17] J. R. Bellegarda, Unsupervised, language-independent grapheme-to-phoneme conversion by latent analogy, in: Speech Communication, Vol. 46, 2005, pp. 140–152.
- 26
27
28
29
30 [18] T. Holter, T. Svendsen, Maximum likelihood modelling of pronunciation variation, in: Speech Communication, Vol. 29, 1999, pp. 171–191.
- 31
32
33 [19] W. Byrne, M. Finke, S. Khudanpur, J. McDonough, H. Nock, M. Riley, M. Saraclar, C. Wooters, G. Zavaliagos, Pronunciation modeling using a hand-labelled corpus for conversational speech recognition, in: Proceedings of International Conference on Acoustics Speech and Signal Processing (IEEE, ICASSP 98), Vol. 1, Seattle, WA, USA, 1998, pp. 313–316.
- 34
35
36
37
38
39 [20] S. Deligne, L. Mangu, On the use of lattices for the automatic generation of pronunciations, in: Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE, ICASSP 2003), Vol. 1, Hong-Kong, China, 2003, pp. 204–207.
- 40
41
42
43
44
45 [21] T. Svendsen, F. Soong, H. Purnhagen, Optimizing baseforms for HMM-based speech recognition, in: Proceedings of the European Conference on Speech Communication and Technology (ESCA, Eurospeech 95), Madrid, Spain, 1995, pp. 783–786.
- 46
47
48
49
50
51 [22] A. Laurent, T. Merlin, S. Meignier, Y. Estève, P. Deléglise, Iterative filtering of phonetic transcriptions of proper nouns, in: Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE, ICASSP 2009), Vol. 1, Taipei, Taiwan, 2009, pp. 4265–4268.
- 52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 [23] S. Galliano, É. Geoffrois, D. Mostefa, K. Choukri, J.-F. Bonastre, G. Gravier,
10 The ESTER phase II evaluation campaign for the rich transcription of French
11 broadcast news, in: Proceedings of International Conference on Speech
12 Communication and Technology (ISCA, Interspeech 2005), Vol. 1, Lisbon,
13 Portugal, 2005, pp. 1149–1152.
14
15
16 [24] O. Andersen, R. Kuhn, A. Lazaridès, P. Dalsgaard, J. Haas, E. Nöth, Com-
17 parison of two tree-structured approaches for grapheme-to-phoneme conver-
18 sion, in: Proceedings of the International Conference on Spoken Language
19 Processing (ICSLP 96), Vol. 3, Philadelphia, PA, USA, 1996, pp. 1700–
20 1703.
21
22
23 [25] M. Riley, W. Byrne, M. Finke, S. Khudanpur, A. Ljolje, J. McDonough,
24 H. Nock, M. Saraclar, C. Wooters, G. Zavalagkos, Stochastic pronunciation
25 modelling from hand-labelled phonetic corpora, in: Speech Communication,
26 Vol. 29, 1999, pp. 209–224.
27
28
29 [26] V. Pagel, K. Lenzo, A. W. Black, Letter to sound rules for accented lexicon
30 compression, in: Proceedings of the International Conference on Spoken
31 Language Processing (ICSLP 98), Sydney, Australia, 1998, pp. 2015–2018.
32
33
34 [27] T. Rama, A. K. Singh, S. Kolachina, Modeling letter-to-phoneme conver-
35 sion as a phrase based statistical machina translation problem with mini-
36 mum error rate training, in: Proceedings of North American Chapter of the
37 Association for Computational Linguistics – Human Language Technolo-
38 gies (NAACL HLT) 2009 conference, Vol. 1, Boulder, CO, USA, 2009, pp.
39 90–95.
40
41
42 [28] L. Bahl, S. Das, P. deSouza, M. Epstein, R. Mercer, B. Merialdo, D. Na-
43 hamoo, M. Picheny, J. Powell, Automatic phonetic baseform determination,
44 in: Proceedings of the International Conference on Acoustics, Speech and
45 Signal Processing (IEEE, ICASSP 91), Vol. 1, Toronto, Canada, 1991, pp.
46 173–176.
47
48
49 [29] B. Ramabhadran, L. Bahl, P. deSouza, M. Padmanabhan, Acoustics-only
50 based automatic phonetic baseform generation, in: Proceedings of the In-
51 ternational Conference on Acoustics, Speech and Signal Processing (IEEE,
52 ICASSP 98), Vol. 1, Seattle, WA, USA, 1998, pp. 309–312.
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 [30] M. Bisani, H. Ney, Breadth-first for finding the optimal phonetic transcription from multiple utterances, in: Proceedings of the International Conference on Speech Communication and Technology (ISCA, Interspeech 2001), Vol. 2, Aalborg, Denmark, 2001, pp. 1429–1432.
- 10
11
12
13
14
15 [31] J. Suontausta, J. Häkkinen, Decision tree based text-to-phoneme mapping for speech recognition, in: Proceedings of the International Conference on Spoken Language Processing (ISCA, ICSLP 2000), Vol. 2, Beijin, China, 2000, pp. 831–834.
- 16
17
18
19
20
21 [32] L. Bahl, P. Brown, P. de Souza, R. Mercer, M. Picheny, A method for the construction of acoustic Markov models for words, in: IEEE Transactions on Speech and Audio Processing, Vol. 1, 1993, pp. 443–452.
- 22
23
24
25
26 [33] R. Haeb-Umbach, P. Beyerlein, E. Thelen, Automatic transcription of unknown words in a speech recognition system, in: Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE, ICASSP 95), Vol. 1, Detroit, MI, USA, 1995, pp. 840–843.
- 27
28
29
30
31 [34] J. Wu, V. Gupta, Application of simultaneous decoding algorithms to automatic transcription of known and unknown words, in: Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE, ICASSP 99), Vol. 2, Phoenix, AZ, USA, 1999, pp. 589–592.
- 32
33
34
35
36 [35] H. Mokbel, D. Jouviet, Derivation of the optimal set of phonetic transcriptions for a word from its acoustic realizations, in: Speech Communication, Vol. 29, 1999, pp. 49–64.
- 37
38
39
40
41 [36] T. Sloboda, Dictionary learning: performance through consistency, in: Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE, ICASSP 95), Vol. 1, Detroit, MI, USA, 1995, pp. 453–456.
- 42
43
44
45
46 [37] S. Deligne, B. Maison, R. Gopinath, Automatic generation and selection of multiple pronunciations for dynamic vocabularies, in: Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE, ICASSP 2001), Vol. 1, Salt Lake City, UT, USA, 2001, pp. 565–568.
- 47
48
49
50
51 [38] R. C. Rose, E. Lleida, Speech recognition using automatically derived base-forms, in: Proceedings of the International Conference on Acoustics, Speech
- 52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 and Signal Processing (IEEE, ICASSP 97), Vol. 2, Munich, Germany, 1997,
10 pp. 1271–1274.
11

- 12 [39] F. Yvon, P. Boula De Mareuil, C. D’Alessandro, V. Aubergé, M. Bagin,
13 G. Bailly, F. Béchet, S. Foukia, J.-P. Goldman, E. Keller, D. O’Shaughnessy,
14 V. Pagel, F. Sannier, J. Véronis, B. Zellner, Objective evaluation of grapheme
15 to phoneme conversion for text-to-speech synthesis in French, *Computer*
16 *Speech and Language* 12 (4) (1998) 393–410.
17
18 [40] P. Koehn, H. Hoang, A. Birch, C. Calisson-Burch, M. Federico, N. Bertholdi,
19 B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin,
20 E. Herbst, Moses: Open-source toolkit for statistical machine translation, in:
21 *Proceedings of the Association for Computational Linguistics*, 2007.
22
23 [41] F. J. Och, Minimum Error Rate in Statistical Machine Translation, in: *Pro-*
24 *ceedings of the Association for Computational Linguistics*, 2003, pp. 160–
25 167.
26
27 [42] K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, BLEU: a method for automatic
28 evaluation of machina translation, in: *Proceedings of the Association for*
29 *Computational Linguistics*, 2002.
30
31 [43] F. Vanden Berghen, H. Bersini, CONDOR, a new parallel, constrained exten-
32 sion of Powell’s UOBYQA algorithm: Experimental results and comparison
33 with the DFO algorithm, *Journal of Computational and Applied Mathemat-*
34 *ics* 181 (2005) 157–175.
35
36 [44] P. Deléglise, Y. Estève, S. Meignier, T. Merlin, The LIUM Speech Transcrip-
37 tion System: A CMU Sphinx III-based system for French broadcast news,
38 in: *Proceedings of International Conference on Speech Communication and*
39 *Technology (ISCA, Interspeech 2005)*, Lisbon, Portugal, 2005, pp. 1653–
40 1656.
41
42 [45] A. Laurent, T. Merlin, S. Meignier, Y. Estève, P. Deléglise, Combined sys-
43 tems for automatic phonetic transcription of proper nouns, in: *Language*
44 *Evaluation and Resources Conference (LREC 2008)*, Marrakech, Morocco,
45 2008.
46
47 [46] M. Ravishankar, R. Singh, B. Raj, R. M. Stern, The 1999 CMU 10x real
48 time broadcast news transcription system, in: *Proceedings of the DARPA*
49
50
51
52
53
54
55
56
57
58

1
2
3
4
5
6
7
8
9 workshop on Automatic Transcription of Broadcast News, Washington, DC,
10 USA, 2000.
11

12 [47] A. Stolcke, SRILM—an extensible language modeling toolkit, in: Proceed-
13 ings of International Conference on Spoken Language Processing (ISCA,
14 ICSLP 2002), Vol. 2, Denver, CO, USA, 2002, pp. 901–904.
15
16

17 [48] V. Jousse, S. Petitrenaud, S. Meignier, Y. Estève, C. Jacquin, Automatic
18 named identification of speakers using diarization and ASR systems, in:
19 Proceedings of International Conference on Acoustics Speech and Signal
20 Processing (IEEE, ICASSP 2009), Taipei, Taiwan, 2009, pp. 4557–4560.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Biography of the authors



Antoine Laurent obtained his Ph.D. degree in 2010 from Université du Maine, Le Mans, France, in the field of automatic speech recognition. He is currently R&D project manager at Spécinov (Angers, France), as well as part-time associate professor in the Language and Speech technology team at LIUM (the computer science research department of Université du Maine, Le Mans, France). His research focuses on automatic adaptation of the ASR system.



Sylvain Meignier received his Ph.D. degree in computer science from Université d'Avignon et des Pays de Vaucluse, Avignon, France, in 2002. His work was about speaker recognition. In 2003, he was with LIMSI-CNRS, Orsay, France, in the Spoken Language Processing Group as a Researcher. Since 2004, he has been an associate professor at Université du Maine, where he works on speech processing in the Language and Speech Technology team of the LIUM laboratory.



Paul Deléglise received his Ph.D. in computer science from Pierre & Marie Curie University (Paris, France) in 1983 and his Doctorat d'État in 1991. He worked in the Signal Laboratory of École Nationale Supérieure des Télécommunications (ENST) on automatic speech recognition from 1985 to 1992. Since October 1992, he is full professor at Université du Maine where he works in the LIUM laboratory on data fusion applied to audio-visual speech recognition, and leads the Language and Speech Technology team. Since 2004, he has been working on large vocabulary speech recognition.

Acknowledgements

Special thanks to Dr. Teva Merlin for his help with this work.