

Document downloaded from:

<http://hdl.handle.net/10251/49938>

This paper must be cited as:

Alabau Gonzalvo, V.; Sanchis Navarro, JA.; Casacuberta Nolla, F. (2014). Improving on-line handwritten recognition in interactive machine translation. *Pattern Recognition*. 47(3):1217-1228. doi:10.1016/j.patcog.2013.09.035.



The final publication is available at

<http://dx.doi.org/10.1016/j.patcog.2013.09.035>

Copyright Elsevier

Improving On-line Handwritten Recognition in Interactive Machine Translation

Vicent Alabau*, Alberto Sanchis, Francisco Casacuberta

Institut Tecnològic d'Informàtica, Universitat Politècnica de València, Camino de Vera s/n, 46022 València, Spain

Abstract

On-line handwriting text recognition (HTR) could be used as a more natural way of interaction in many interactive applications. However, current HTR technology is far from developing error-free systems and, consequently, its use in many applications is limited. Despite this, there are many scenarios, as in the correction of the errors of fully-automatic systems using HTR in a post-editing step, in which the information from the specific task allows to constrain the search and therefore to improve the HTR accuracy. For example, in machine translation (MT), the on-line HTR system can also be used to correct translation errors. The HTR can take advantage of information from the translation problem such as the source sentence that is translated, the portion of the translated sentence that has been supervised by the human, or the translation error to be amended. Empirical experimentation suggests that this is a valuable information to improve the robustness of the on-line HTR system achieving remarkable results.

Keywords: interactive pattern recognition, on-line handwritten text recognition, interactive machine translation, human interaction

*Corresponding author: Phone: (34) 96 387 70 69, Fax: (34) 96 387 72 39.

Email addresses: valabau@iti.upv.es (Vicent Alabau), asanchis@iti.upv.es (Alberto Sanchis), fcn@iti.upv.es (Francisco Casacuberta)

1 **1. Introduction**

2 Since the breakout of tactile smartphones, the number of devices featuring a touch-screen has been increasing
3 at a fast pace. The success of tactile smartphones has fostered a new kind of keyboardless technology: the tablet
4 computers. They have been presented as a substitute of paper notebooks although the possibilities this new technology
5 may provide are still to be unveiled.

6 In that context, on-line handwritten text recognition (HTR) can play a crucial role. First, because to input text
7 in such devices using a virtual keyboard is far from the efficiency of regular keyboards. Secondly, handwriting is a
8 natural way to communicate. Withal, an HTR interface can commit recognition errors. Thus, if the HTR system is not
9 robust enough, user experience could be negatively affected hindering its use. In this regard, many works have tried to
10 improve HTR accuracy. Primarily focusing on feature extraction and modeling [1, 2, 3]. Other authors have tackled
11 the problem of automatically correcting errors from the system output in order to provide a more accurate input to
12 higher-level applications. For instance, Quiniou et al. [4] propose a technique to improve the performance of a HTR
13 system by obtaining a consensus hypothesis out of a n -best lists, and then, characterizing the errors and correcting
14 them. Similarly, Farooq et al. [5] use a translation model to conduct an automatic post-editing. Additionally, Devlin
15 et al. [6] used a machine translation system to rerank an OCR n -best list. The idea was that *easily translatable* options
16 would have a better syntax, which in the end resulted in small accuracy improvements. Nevertheless, those works did
17 not leverage any contextual information of the specific task at hand, a topic that, in our opinion, has received little
18 attention. Following this line of research, Toselli et al. [7] explored the use of on-line HTR for interactive transcription
19 of text images. In that work, the user was expected to correct erroneously recognized words by handwriting the
20 correction using a tactile display. The authors took advantage of the erroneously predicted word and the previous one
21 to improve HTR robustness.

22 Inspired by Toselli et al. [7, 8], we address the problem of using an on-line HTR system to correct the errors in a
23 *machine translation* (MT) application. State-of-the-art MT systems usually cannot perform translations to fit quality
24 demands by the translation industry. Hence, it is typical to have the automatically produced output documents revised
25 by a professional translator. In this manual process, known as post-editing (PE), the human expert can spend hours
26 of work to achieve high-quality translations. *Interactive machine translation* (IMT) [9, 10, 11] was developed to deal
27 with this problem. In IMT, a human expert is introduced in the middle of the translation process. This way, she can
28 amend errors from the system output and useful feedback is used by the system to automatically improve the part of
29 the translation to be revised.

30 The usual way to introduce the corrections in IMT is by means of the keyboard where the mouse is used to fix the
31 position [12]. However, other interaction modalities are also possible. For example, speech interaction was studied
32 in [13, 14, 15]. There, several scenarios were proposed, in which the user was expected to utter aloud parts of the
33 current hypothesis along with one or more corrections. Later, we proposed the use of on-line HTR to IMT in [16, 17].
34 To our knowledge, our work has been the first approach to on-line HTR in IMT so far. Nonetheless, those works

35 presented very preliminary results explaining simple contextual models and HTR interaction restricted to isolated
36 words.

37 In this paper we present relevant novelties with respect to previous work that can be summarized in two main
38 improvements. First, we introduce a new HTR model that leverages state-of-the-art phrase-based models, whereas
39 previous work was based only on word-based translation models. Second, we extend the interaction scheme to allow
40 sequences of words (phrases) to be written and not just isolated words. In addition, we propose a method to recover
41 efficiently from HTR errors using contextual menus. Finally, a new and exhaustive experimental study is presented to
42 evaluate all those novel contributions and preliminary ideas.

43 The remainder of this paper is organized as follows. First, the process to produce high-quality translations is
44 introduced in [Sec. 2](#). Second, in [Sec. 3](#) several alternatives to incorporate contextual information from the translation
45 problem into the HTR decoding will be explored. [Section 4](#) is devoted to the evaluation of the proposed models.
46 Finally, conclusions and future work will be discussed in [Sec. 6](#).

47 **2. Producing High-Quality Translations**

48 In the last years, *machine translation* (MT) has become a strategic asset in the translation industry. MT is used
49 to speed up the translation process since it enables the automatic translation of large amounts of documents. In this
50 context, MT is approached under a statistical framework, due to the fact that statistical MT allows companies to
51 build customized, topic-specific MT systems very economically. Here, the problem consists in finding the most likely
52 translation \hat{t} in a target language given a source sentence s in a source language,

$$\hat{t} = \operatorname{argmax}_t Pr(t | s) \quad (1)$$

53 which can be modeled in different ways [18].

54 *2.1. Post-editing a Machine Translation Output*

55 Although leveraging MT can be very convenient, it is usually the case that the translation quality does not meet
56 the user requirements. Thus, the MT output must be revised. The process of revising and amending the system output,
57 known as *post-editing* (PE), consists in deleting, inserting, substituting and swapping text from the MT output to
58 achieve the desired quality in the translation. This is an expensive task, since the users should review the whole output
59 and correct manually the translation errors. In the cases in which the automatically produced translations are of low
60 quality, PE can eventually require more effort than manually translating the source input from the scratch. Moreover,
61 in PE, the system does not take advantage of the human corrections.

62 *2.2. Interactive Machine Translation*

63 The MT paradigm is shifting slowly but steady towards an interactive MT scenario (IMT). In IMT [9, 10, 11] the
64 system goal is not to produce translations in a completely automatic way and then perform a completely unassisted

65 PE. On the contrary, IMT aims at building the translation collaboratively with the user as a professional advisor, so
 66 that the effort to produce a satisfactory output is minimized.

67 A typical approach to IMT is shown in Fig. 1. A source sentence s is given to the IMT system. First, the system
 68 outputs a translation hypothesis \hat{t} in the target language, that would correspond to the output of fully automated MT
 69 system (i.e., based on Eq. (1)). Next, the user analyzes the source sentence and the current hypothesis, and validates
 70 the longest error-free prefix p finding the first error. Then, the user amends the erroneous word by typing the correct
 71 word d . Based on this amendment, the system creates a new validated prefix $p \cdot d$, with \cdot as a concatenation operator.
 72 With that information, the system is able to produce a new, hopefully improved, translation \hat{t} that is coherent with the
 73 information provided, that is, $p \cdot d$ must be a prefix of the new \hat{t} . This process is repeated until the user agrees with
 74 the quality of the resulting translation. In this work we assume that this protocol is performed left-to-right, but other
 75 protocols are also possible.

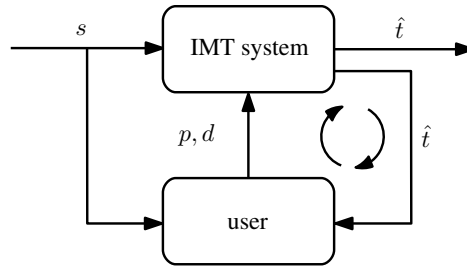


Figure 1: Diagram of a typical approach to IMT

76 The iterative nature of the process is emphasized by the loop in Fig. 1, which indicates that, for a source sentence
 77 to be translated, several interactions between the user and the system could be performed. In each interaction, the
 78 system produces the most probable translation \hat{t} that is coherent with the prefix formed by concatenating the previous
 79 prefix p and the user correction d :

$$\hat{t} = \operatorname{argmax}_{t: \tau(t, p \cdot d)} Pr(t | p, d, s) \quad (2)$$

80 where $\tau(t, p \cdot d)$ is a function that is true if $p \cdot d$ is a prefix of t . It is worthy of note that the main difference between
 81 Eq. (1) and Eq. (2) is that, in the second case, \hat{t} is must be coherent with the validated prefix $p \cdot d$. Since the probabilistic
 82 models in Eq. (1) and Eq. (2) are estimated in the same way, Eq. (2) can be considered as a constrained search problem
 83 of the classical MT problem. In fact, at the beginning, when the user has not validated any prefix, Eq. (1) and Eq. (2)
 84 are equivalent equations. In addition, adaptive approaches can also be assumed, where the system is able to learn from
 85 each user interaction to improve the underlying statistical models [19].

86 For the sake of a better understanding, a typical translation IMT session is exemplified in Fig. 2. First, the system
 87 starts with an empty prefix, so it proposes a full hypothesis. Then, the user corrects the first error, *not*, by typing
 88 ‘is’. Next, the system proposes a new suffix, in which the first word, *not*, has been automatically corrected. The user
 89 amends *at* by typing ‘in’. Finally, as the new proposed suffix is correct, the process ends. Note that 4 operations would

90 have been needed in a PE scenario, whereas only 2 are needed in IMT. In this example, the user types the complete
 91 wrong word. Nevertheless, it is straightforward to extend this operation to the character level instead of word level.

SOURCE (s):		si alguna función no se encuentra disponible en su red
REFERENCE (r):		if any feature is not available in your network
ITER-0	(p)	
ITER-1	(\hat{t})	if any feature not is available on the network
	(p)	<i>if any feature</i>
	(d)	is
ITER-2	(\hat{t})	if any feature is not available at the network
	(p)	if any feature is <i>not available</i>
	(d)	in
FINAL	(\hat{t})	if any feature is not available in your network
	($\hat{t} \equiv r$)	if any feature is not available in your network

Figure 2: Example of an IMT session for translating a Spanish sentence s to an English sentence t . Initially, in iteration 0, the prefix is empty, i.e., the user has not performed any validation. In iteration 1, the system proposes a fully automatic translation \hat{t} . Then, the user finds the first error and amends it by introducing the correct word (d), which is shown in **boldface**. As a result, the user has implicitly validated a prefix (p), shown in *italics*. The concatenation of the prefix and the corrected word constitutes a new prefix for the next iteration (displayed in **blue**). The process continues until the user is satisfied with the solution. Note that 4 operations would have been needed in a PE scenario, whereas only 2 are needed in IMT.

92 3. Using On-Line HTR to Correct MT Output

93 Typically, the correction of MT output is performed using a keyboard and, occasionally, a mouse to position the
 94 cursor [12]. Professional translators agree that this approach has been proved to be efficient. However, the user needs
 95 to be in front of a desktop computer which imposes some restrictions regarding where and how the work is to be
 96 done. Laptop computers can also be used, although arguably performance could be diminished because of the use of
 97 uncomfortable laptop keyboards and track pads. Thus, although e-pen interaction may sound impractical for texts that
 98 need a large amounts of corrections, there is a number of circumstances where e-pen interaction can be more suitable.
 99 For example, it can be well suited for amending sentences with few errors, as the revision of human post-edited
 100 sentences, or translations where the system has a high confidence that the output is of good quality. Furthermore, it
 101 would allow to perform such tasks while commuting, traveling or sitting comfortably on the couch in the living room.

102 Now, imagine an application devised to translate documents. On the one hand, there is a text area with the output
 103 of an automatic machine translation system. As this output may contain errors, the user of the application reads the
 104 output to locate the first error. The reading is performed in a specific order, left-to-right in most western languages, for
 105 instance. Let us assume that when the user finds the first error, all the words before it have already been revised and
 106 validated. Thus, they can be regarded as correct. Once the error has been located, the user introduces the correction
 107 with a stylus. As a result, the system receives a position where the error is located, a word that is incorrect (the word
 108 pointed by the position) and a sequence of pen strokes that represent the correct word in that position. On the other
 109 hand, the source document to be transcribed is shown to the user. There is a strong relationship among the words in
 110 the source sentence and the words in the target sentence.

111 **Figure 3** is a mock-up of a possible application on a tablet device for such scenario. The screen is divided in
 112 two sections. First, the upper part shows the source document, and probably the source sentence being currently
 113 translated, s , is highlighted appropriately. Second, the lower section contains the current state of the translation, t .
 114 Since we assume that post-editing is usually performed from left to right, the text which has already been revised and
 115 validated is highlighted. On the other hand, the text which is to be revised is displayed grayed out. From the sentence
 116 currently being translated we can identify three parts: the revised prefix of the sentence, p , the error committed by the
 117 system, e , and the correction proposed by the user introducing strokes with a stylus, x .

118 In a scenario as described above, the HTR subsystem should make few errors to make the application usable.
 119 The aim of this work is to devise a robust HTR system that allows a potential user to revise and correct the output
 120 of a machine translation system using an electronic pen. To this regard, we assume that the user will introduce the
 121 corrections by writing over the word or sequences of words (phrases) she judges to be incorrect. Thus, the problem of
 122 on-line HTR consists in converting a sequence of strokes, x , into a word or phrase in text format, d . The strokes can
 123 be acquired from a stylus, electronic pen or a touch-screen.

124 3.1. System Baseline

125 The baseline approach to the problem from a statistical point of view is to obtain the most likely decoding d given
 126 the strokes x ,

$$\hat{d} = \underset{d}{\operatorname{argmax}} Pr(d | x) = \underset{d}{\operatorname{argmax}} Pr(d)Pr(x | d) \quad (3)$$

127 where $Pr(d)$ can be represented by a language model and $Pr(x | d)$ by morphological models.

128 The morphological models can be modeled by hidden Markov models [2] or neural networks [1]. On the other
 129 hand, a common and practical approach to model $Pr(d)$ is by means of n -grams [20]. The description of an on-line
 130 HTR system would end here for most applications. However, our purpose is to take advantage of the information
 131 available in the IMT application to make on-line HTR more robust. In the remainder of this section, we will introduce
 132 gradually the different kinds of information sources into the language model. With the addition of each of them, we
 133 aim to make the on-line HTR system more robust.

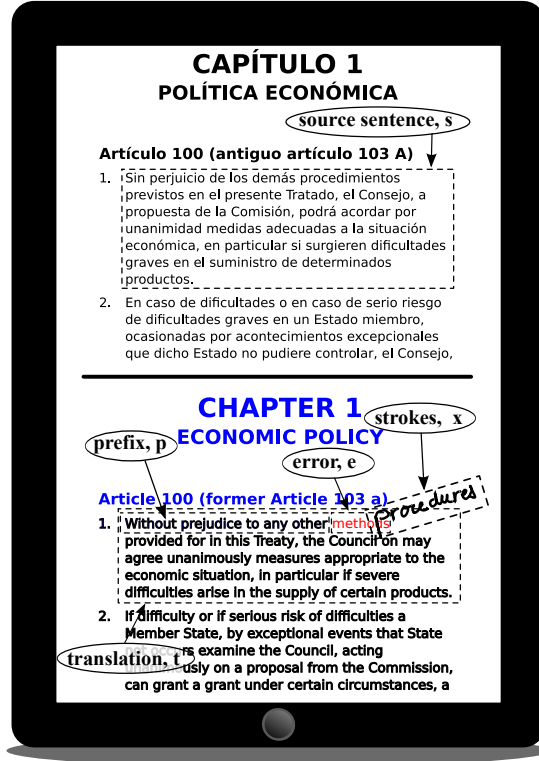


Figure 3: Mock-up of an interactive machine translation application on a tablet device.

134 3.2. Discarding the produced error

135 In the e-pen enabled IMT interface aforementioned, the user is expected to write the strokes over the erroneously
 136 translated word, and thus, the system knows what word the user wants to replace. Therefore, the first and easiest
 137 approach is to remove the erroneous word e from the list of candidate hypotheses. This way, Eq. (3) becomes

$$\hat{d} = \operatorname{argmax}_{d \neq e} Pr(d)Pr(x | d) \quad (4)$$

138 3.3. Exploiting information from the revised translation

139 The second sensible approach to take is to add information regarding the revised translation prefix, p . Again, from
 140 Eq. (3) we can derive an HTR system that takes into account previously validated words:

$$\hat{d} = \operatorname{argmax}_d Pr(d | x, p) \approx \operatorname{argmax}_d Pr(d | p)Pr(x | d) \quad (5)$$

141 under the assumption that $Pr(x | d, p)$ does not depend on p if d is known. In addition, $Pr(d | p)$ is a prefix language
 142 model, i.e., the probability of d depends on the left-context. Of course, we can also discard the erroneous word from
 143 Eq. (5),

$$\hat{d} \approx \operatorname{argmax}_{d \neq e} Pr(d | p)Pr(x | d) \quad (6)$$

144 These techniques can be extrapolated to most post-editing tasks. In fact, Toselli et al. [7] used the erroneous word
 145 and a 2-gram model to improve the HTR performance for interactive transcription of text images. Next, we will show
 146 how the information regarding the translation process can be exploited for further improve HTR decoding.

147 3.4. Leveraging information from the source sentence

148 A specific source of information that can help to improve robustness in the MT scenario is, naturally, the sentence
 149 in the source language. Since the target sentence conveys the meaning of the source sentence, s , user corrections
 150 should be restricted somehow to the possible translations of it. Hence, we can formulate the problem as,

$$\hat{d} = \operatorname{argmax}_d Pr(d | x, p, s) \approx \operatorname{argmax}_d Pr(d | p, s) Pr(x | d) \quad (7)$$

151 assuming that $Pr(x | d, p, s)$ does not depend on p and s if d is known.

152 Nevertheless, the relationship between the target and the source sentence in $Pr(d | p, s)$ is not trivial to establish.
 153 Two possibilities are considered in this work. First, word-based models are the basis for modern statistical MT [21].
 154 Although they cannot provide a good performance when translating complete sentences, they offer a smoothed and
 155 reliable probability distribution for word models. In addition, they serve as initialization for the second kind of models
 156 considered: phrase-based models [18]. These models improve word-based models since they are able to translate
 157 sequences of words (phrases) and constitute the state-of-the-art in MT.

158 3.4.1. Word-based translation models

159 Brown et al. [21] approached the problem of MT in Eq. (1) from a statistical point of view as a search problem
 160 of a translation t . In this approach a hidden variable a is introduced that represents the alignment between the words
 161 in the source and target sentence. Let a be a vector with the length of the target sentence $|t|$ ¹, where each element a_i
 162 represents an index in the source sentence to whom t_i is aligned, i.e., a_i means that t_i is aligned to s_{a_i} . In order to
 163 simplify the notation, from now on we will refer to a_i as j so that j indexes source words. Formally, we can model
 164 the posterior probability of the target sentence t being a translation of the source sentence s by marginalizing over the
 165 set of all possible alignments between the words in t and the words in s ,

$$Pr(t | s) = \sum_a Pr(t, a | s) \quad (8)$$

166 Then, $Pr(t, a | s)$ can be decomposed using the chain rule. After taking some strong assumptions, two distribu-
 167 tions are obtained. First, the alignment model, $Pr(j | i, |s|)$, represents the probability of the target word at position
 168 i to be aligned with the source word at position j for a source sentence of length $|s|$. Second, the word translation
 169 model, $Pr(t_i | s_j)$, models the probability of the target word at position i to be a translation of the source word at

¹We define the length of a sentence as the number of elements in the sentence. The elements are typically words and symbols, but it depends on the tokenization.

170 position j . The above assumptions are necessary to make model estimation tractable and result in the so-called *model*
 171 2 (M2) [21].

172 In M2, the alignment probability, $Pr(j | i, |s|)$, can be approximated by the relative frequency of position j in
 173 the source sentence to be aligned with position i in the target sentence for a source sentence of length $|s|$. On the
 174 other hand, the translation probability, $Pr(t_i | s_j)$, can be approximated by a word-to-word statistical dictionary
 175 which essentially is the relative frequency of t_i being aligned with s_j . Nonetheless, these frequencies cannot be
 176 estimated directly since the real alignments are unknown. Thus, the EM algorithm is needed to reliably estimate
 177 these probabilities [21]. *Model 1* (M1) is a particular case of word-based models where the alignment probability is
 178 approximated by an uniform probability distribution, $Pr(j | i, |s|) \approx (|s| + 1)^{-1}$.

179 Returning to our original problem, we can approach $Pr(d | p, s)$ in Eq. (7) with word-based translation models
 180 with some assumptions. First, from the prefix p we can obtain the position of the erroneous word to be corrected,
 181 $i = |p| + 1$ ignoring the rest of the words in the prefix,

$$Pr(d | p, s) \approx Pr(d | i, s) \quad (9)$$

182 Then, we can introduce the alignment between d and the words from the source sentence by summing for every
 183 possible position j in s ,

$$\begin{aligned} Pr(d | i, s) &= \sum_{j=1}^{|s|} Pr(d, j | i, s) \\ &= \sum_{j=1}^{|s|} Pr(j | i, s) Pr(d | j, i, s) \end{aligned} \quad (10)$$

184 Finally, if we assume, in a similar way to M2, that $Pr(j | i, s)$ does not depend on s but on $|s|$, and that $Pr(d |$
 185 $j, i, s)$ does not depend on the whole s but just the word aligned to d , s_j with j , then we can approximate Eq. (10) as

$$Pr(d | i, s) \approx \sum_{j=1}^{|s|} Pr(j | i, |s|) Pr(d | s_j) \quad (11)$$

186 where $Pr(j | i, |s|)$ is an M1 or M2 alignment model and $Pr(d | s_j)$ is a statistical dictionary.

187 To clarify the role of the alignments and the dictionary, observe Fig. 4. The source sentence is shown in the
 188 middle. Each word has its corresponding position, j , as a subscript. Above each word, there is a list of its most
 189 probable translations using the dictionary. Grey levels are proportional to the probability of the dictionary. On the
 190 other hand, in the bottom, there is a possible translation, which has an error in position $i = 4$. Below that, the user
 191 is trying to correct that mistake by introducing the word υ . Each link between a source word and the target word in
 192 position 4 represents the alignment probability. The stroke boldness is proportional to the M2 alignment probability.
 193 Note that for an M1 model, all alignments would have had the same thickness.

194 If we focus on the possible candidate transcriptions of υ , we realize that there are two possibilities that could
 195 create confusion to the decoder: ‘if’ as translation of ‘si₁’ and ‘in’ as translation of ‘eng’ due to the fact that the

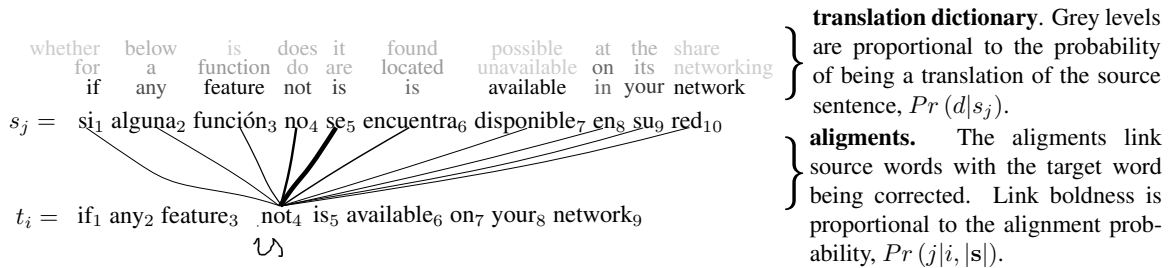


Figure 4: Visualization of alignments and translation dictionary.

196 strokes for ‘is’, ‘if’ and for ‘in’ can be very similar. Both can compete with the correct transcription ‘is’. The first, has
 197 a high probability in the dictionary, $Pr(\text{if} \mid s_1) = 0.88$, whereas $Pr(\text{is} \mid s_5) = 0.46$, $Pr(\text{is} \mid \text{encuentra}_6) = 0.34$
 198 and $Pr(\text{in} \mid \text{en}_8) = 0.40$. Then, since the M1 model has a uniform alignment probability, it would assign a higher
 199 probability to ‘if’ than to ‘is’. However, ‘si₁’ actually has a lower probability of being aligned with ‘not₄’. Therefore,
 200 the M2 model is able to solve this shortcoming thanks to the alignments with high probability to the correct words. In
 201 this case, $Pr(5 \mid 4, 10) = 0.38$ and $Pr(6 \mid 4, 10) = 0.12$, whereas $Pr(1 \mid 4, 10) = 0.04$.

202 3.4.2. Phrase-based translation models

203 Word-based translations provided a basis for MT. However, their performance regarding translation quality was
 204 not sufficient. Their limitation resides in that they cannot model properly context information [22]. Phrase-based
 205 models aim at reducing this problem by translating phrases (fragments of sentences) instead of single words. These
 206 models were popularized by Och and Ney [23], who established the state-of-the-art phrase-based log-linear models.
 207 Phrase-based models offer a great opportunity to estimate $Pr(d \mid p, s)$. However, we cannot use these models directly,
 208 as we did with word-based models. One limitation of phrase-based models is that their probabilities are ‘peaky’ and,
 209 usually, they cannot model all possible translations. As a result, it is possible that $Pr(d \mid p, s)$ is 0 for a user
 210 established prefix like it would be the case in IMT. Then, it is necessary to smooth these probabilities. For instance,
 211 we can generate n -gram-like models from the hypotheses in a word graph (WG) of a MT system [24].

212 Word graphs contain a set of the most likely translations of the source sentence. They can encode a large number
 213 of translations in a more efficient way than n -best lists. Although one may think that the WG could be directly used,
 214 there are some details that must be taken into account. First, WGs do not contain all the possible translations since, in
 215 practice, many pruning techniques must be used to generate the translations efficiently. Second, phrase-based models
 216 are not good dealing with long distance alignments due to the introduction of heuristic length constraints, and thus,
 217 WGs do not present sentences with long distance reorderings. In those cases, a user validating a prefix p that is not
 218 contained in the WG would obtain a zero probability in $Pr(d \mid p, s)$. Hence, it is interesting to smooth the probability
 219 distribution encoded in the WGs. To do so, WGs can be simplified in the way that language modeling is typically
 220 approached: we make each word to depend only on the preceding $n - 1$ words instead of depending on the whole

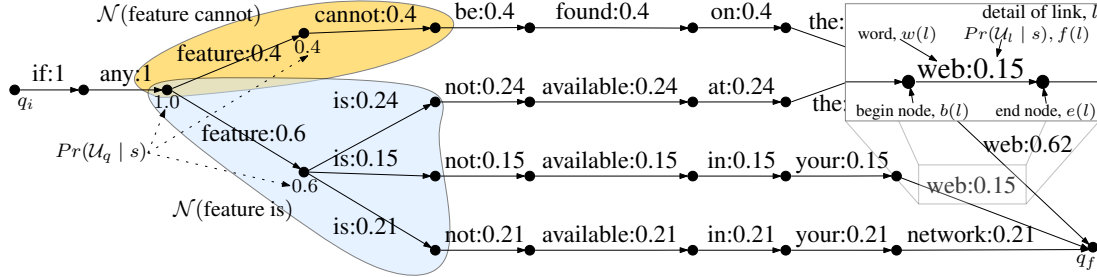


Figure 5: Word graph with posterior probabilities. It represents a subset of hypotheses of the hypothesis space of a state-of-the-art translation model for the source sentence ‘si alguna función no se encuentra disponible en su red’. On the left, in blue the set of links considered when computing the average count of the bi-gram ‘feature is’ whereas in orange the link considered for the bi-gram ‘feature cannot’.

221 prefix. As a result, Eq. (7) can be rewritten as

$$\hat{d} \approx \operatorname{argmax}_d Pr(d | p_{i-n+1}^{i-1}, s) Pr(x | d) \quad (12)$$

222 where p_{i-n+1}^{i-1} are the words in the prefix from position $i - n + 1$ to position $i - 1$, i.e., $Pr(d | p_{i-n+1}^{i-1}, s)$ only
 223 takes into account the latest $n - 1$ words from the prefix. Note that $Pr(d | p_{i-n+1}^{i-1}, s)$ is very similar to a n -gram
 224 language model except for the dependency on s . We used a similar approach for dictation of handwritten historical
 225 documents [25, 26] and speech interaction to IMT [13]. Khadivi and Ney [14] presented a closely related approach for
 226 generating n -gram-like models from n -bests lists instead of WGs. The advantage of the n -gram-like prefix modeling
 227 assumption is that the models only take into account a limited size of the history, and thus, can provide a smoother
 228 probability distribution.

229 Formally speaking, a WG L is a directed, acyclic, weighted graph with an initial node q_i and a final node q_f .
 230 A link l is defined as any edge between two nodes; each link has associated a begin node $b(l)$, an end node $e(l)$, a
 231 hypothesized word $w(l)$, and a score $f(l)$; each link can be considered as a hypothesis $w(l)$ between the nodes $b(l)$
 232 and $e(l)$ with score $f(l)$. Any path from q_i to q_f forms a translation hypothesis t . In MT, $f(l)$ is the score of the
 233 log-linear phrase based model for that particular link. An example of WG is displayed in Fig. 5.

234 Let $Pr(\mathcal{U}_q | s)$ be the posterior probability of all the paths that use the node q and let $Pr(\mathcal{U}_l | s)$ be the poste-
 235 rior probability of all the paths that use the link l . These probabilities can be efficiently computed with a *forward*-
 236 *backward*-based algorithm [27]. Then, the average counts of word sequences for a given source sentence can be
 237 estimated efficiently as in [28]. For a given n -gram length:

$$C^*(d_{i-n+1}^i | s) = \sum_{l_1^n \in \mathcal{N}(d_{i-n+1}^i)} \frac{\prod_{j=1}^n Pr(\mathcal{U}_{l_j} | s)}{\prod_{j=2}^n Pr(\mathcal{U}_{b(l_j)} | s)} \quad (13)$$

238 where $\mathcal{N}(d_{i-n+1}^i)$ is the set of all the sequences of concatenated links in L that produce the sequence of words
 239 d_{i-n+1}^i .

240 An example of such sets on a simplistic WG is shown in Fig. 5 for the 2-grams ‘feature cannot’ and ‘feature is’.

241 Then, $C^*(\text{feature cannot} \mid s)$ and $C^*(\text{feature is} \mid s)$ can be computed as

$$C^*(\text{feature cannot} \mid s) = \frac{0.4 \cdot 0.4}{0.4} = 0.4$$

$$C^*(\text{feature is} \mid s) = \frac{0.6 \cdot 0.24}{0.6} + \frac{0.6 \cdot 0.15}{0.6} + \frac{0.6 \cdot 0.21}{0.6} = 0.6$$

242 That is, ‘feature is’ appears 0.6 times in average in the possible set of translation, whereas ‘feature cannot’ only
 243 appears 0.4 times. Note that if a sequence of words appears more than once in a sentence, the average counts might
 244 exceed 1.

245 Now, n -gram-like probabilities from the WG with posterior probabilities can be calculated after a proper normal-
 246 ization:

$$Pr(d_i \mid d_{i-n+1}^{i-1}, s) = \frac{C^*(d_{i-n+1}^i \mid s)}{C^*(d_{i-n+1}^{i-1} \mid s)} \quad (14)$$

247 Then, Eq. (14) can be used directly in Eq. (7) to approximate $Pr(d \mid p, s)$. In other words, given a sequence
 248 of words d_{i-n+1}^i , $Pr(d_i \mid d_{i-n+1}^{i-1}, s)$ can be estimated by summing up the posterior probabilities of all sentences
 249 containing the sequence d_{i-n+1}^i .

250 The estimation in Eq. (14) presents the problem that many n -grams are not seen in the WG. Then, they will have
 251 zero probability, and the HTR system will fail to recognize them. A common approach is to rely on simpler models to
 252 account for unseen events using back-off models [29]. As the estimated counts are not real counts (they vary from 0
 253 to the number of times the n -gram occurs in a sentence), typical discount methods cannot be applied [30]. However,
 254 absolute discount can be used [31], which consists in subtracting a constant, ϵ , from C^* .

255 Furthermore, only words present in the WG are included into the model (which implies a high number of out-of-
 256 vocabulary words (OOV), since WGs only contain the words of the most likely hypotheses). The OOV problem is
 257 solved by distributing the discounted mass from the unigram among the remaining words of the vocabulary.

258 Finally, to improve the estimation of unseen events, n -grams from the WG can be interpolated linearly with the
 259 standard n -gram model:

$$Pr_\gamma(d \mid p, s) = \gamma Pr(d \mid p, s) + (1 - \gamma) Pr(d \mid p) \quad (15)$$

260 This way, the words that were not used by the MT engine are assigned a meaningful probability.

261 3.5. Integrated HTR and IMT decoding

262 Previous models assume a two-step process, in which the strokes are first decoded into a word or phrase, and then,
 263 the decoded word is used to correct the output of the IMT system. However, this decoding can be performed in an
 264 integrated way by marginalizing over every possible decoding d in Eq. (2):

$$\hat{t} = \operatorname{argmax}_t \sum_d Pr(t, d \mid p, x, s) \quad (16)$$

265 Note that Eq. (16) sums over all possible values of d , but we also are interested in the result of the decoding. Then,
 266 we can decompose Eq. (16) using the chain rule. Approximating the sum by the maximum, and assuming that
 267 $Pr(t | p, x, d, s)$ does not depend on x if d is known,

$$\hat{t} \approx \operatorname{argmax}_t \max_d Pr(d | p, x, s) Pr(t | p, d, s) \quad (17)$$

268 where \hat{d} can be obtained as a byproduct of the decoding of \hat{t} .

269 The first term in Eq. (17) can be approximated as in Eq. (3), Eq. (4), Eq. (5), Eq. (11) or Eq. (12). The second term
 270 is a prefix conditioned translation model as in Eq. (2). This probability forces d not just to be a good translation of s
 271 but to form part of a sentence that is good translation of it. Hence, the decoding of d is benefiting from a new source
 272 of information.

273 4. Experiments

274 In this section, we present a set of experiments to assess the performance of the MT specific HTR systems de-
 275 scribed in the above sections. Two kinds of experiments were conducted. First, the word-based experiments assume
 276 that the user only writes one word at a time. Second, in the phrase-based experiments the user writes a set of consec-
 277 utive erroneous words. Additionally, two corpora were generated from the Xerox corpus, one with Spanish phrases
 278 from translations of English sentences and the other one with English phrases from translations of Spanish sentences.
 279 The details of how the two corpora were generated are given in Sec. 4.3.

280 4.1. IMT corpus: Xerox

281 The Xerox corpus, created in the TT2 project [32], was used for the experiments, since it has been extensively
 282 used in the literature to evaluate IMT systems. It consists of a collection of technical manuals in English, Spanish,
 283 French, and German. The English version is the original document, while the others are professional translations of the
 284 original. The English and Spanish versions were used in the experiments. The training data was used to generate the
 285 translation models. Examples of sentence pairs are shown in Fig. 6. The corpus consists of 56k sentences of training
 286 and a development and test sets of 1.1k sentences. The development set was used to find the tuning parameters that
 287 were used in test. Test perplexities for Spanish and English are 35 and 51, respectively. In addition, the Spanish
 288 test set has 0.7% out-of-vocabulary running words, whereas the English test set has 0.6% out-of-vocabulary running
 289 words.

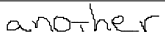
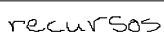

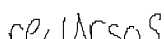
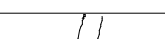

290 4.2. HTR corpus: Unipen

291 For on-line HTR, the UNIPEN corpus [33] was used. The training data was composed of symbols, digits and the
 292 1000 most frequent English and Spanish words. The words were generated by concatenating different instances of
 293 characters from the same writer, with a total of 17 different writers. Overall, 68 character classes and a total of 23.5k
 294 unique character instances were used to generate all the 43.8k training samples. The feature extraction and modeling

Figure 6: Examples of paired sentences in Spanish and English extracted from the Xerox corpus.

Spanish	English
use este botón para ampliar la búsqueda de dispositivos xerox.	use this button to expand the search for xerox devices.
la búsqueda puede ampliarse para incluir otros nombres de comunidades de snmp que se han agregado a la red.	the search may be expanded to include additional snmp community names that have been added to your network.

Figure 7: Examples of pen strokes from the UNIPEN database used for the simulation of HTR. The words were obtained by concatenating random character instances from the corresponding user.

user	<i>another</i>	<i>recursos</i>
User 1		
User 2		
User 3		

295 we used was based on Pastor et al. [2]. Basically, the strokes were preprocessed by eliminating pen-up points and
 296 consecutively repeated points. Then, a low pass filter was applied to reduce noise by replacing each point with the
 297 mean of its neighbors [3]. From the resulting trajectory, 6 features were extracted:

- 298 • the vertical position is normalized by scaling and translating it to [0, 100] keeping aspect ratio.
- 299 • the first and second derivatives for the vertical and horizontal position.
- 300 • the curvature, which is the inverse of the radius of the curve in each point.

301 Next, these feature vectors were used to train the morphological models, which were represented by left-to-right
 302 continuous density Hidden Markov Models (HMM) [34] with Gaussian mixtures and variable number of states per
 303 character. Three users were separated from the training process to produce the words from concatenated characters
 304 for the development sets, which were used to find the optimal tuning parameters, and test sets. Examples of generated
 305 word in Fig. 7.

306 4.3. Procedure

307 For the word-based experiments, the simulation of the user interaction was performed in the following way. First,
308 the publicly available IMT decoder Thot [35] was used to run an off-line simulation for keyboard-based IMT. To do
309 this, we translated each test source sentence. Then, we obtained the longest correct prefix comparing to the reference.
310 Next, we took the word that followed that prefix as the word the user would introduce as a correction. Finally, we
311 used the prefix, and the correct word to obtain a new translation. This was repeated until the reference was obtained.
312 As a result, a list of words that the system failed to predict was obtained. Supposedly, this would be the list of words
313 that the user would correct with handwriting.

314 Then, from UNIPEN corpus, three writers were selected to simulate the user interaction. For each writer and for
315 each of the words in the list of corrections, the handwritten words were generated by concatenating random character
316 instances from the user’s data to form a single stroke. Finally, the generated handwritten words were decoded using
317 the proposed systems with *iAtrós* decoder [36]. The 3-gram perplexities for the generated words are 205 and 226
318 for development and test, respectively, in Spanish, and 242 and 336 for English. It is worthy of note these high
319 perplexities, when for the whole dev and test sets the perplexities are 35 and 51. The word lists were extracted from
320 the erroneous translations that were generated with a decoder using the very same n -grams models used to compute
321 the perplexity. Hence, it is reasonable to assume that if the decoder failed to translate these words it was in part
322 because the language probabilities were low enough, i.e., these probabilities were not well estimated, resulting in a
323 high perplexity. Finally, the number of words in the development sets are 2767 for Spanish and 2398 for English, and
324 in the test sets 2248 and 2102, respectively.

325 For the phrase-based experiments, the development and test sets were constructed in a similar way. In this case,
326 from the word lists aforementioned, we concatenated the strokes of the words that were consecutive in the original
327 text to form strokes of phrases. For instance, if the MT system had translated ‘lista de impresoras’ to ‘list of printers’
328 when the user preferred ‘printer list’, in the word-based scenario we would have generated the word ‘printer’ and the
329 word ‘list’. In the phrase-based scenario, as both errors are consecutive, we would have concatenated them in a single
330 phrase as ‘printer list’. Figure 8 illustrates a box-and-whisker diagram of the phrase lengths in the different sets. We
331 can observe from the whiskers that the majority of the phrases are less than 3 words for Spanish and 6 for English,
332 whereas for the outliers the lengths reach a maximum at 12 and 18, respectively. Note, however, that the interquartile
333 range is between 2 and 3, meaning that half of the phrases are reasonably short. Finally, the number of phrases in the
334 development sets are 941 for Spanish and 896 for English, and in the test sets 1268 and 1130, respectively.

335 4.4. Evaluation measures

336 The performance of the word-based HTR system has been assessed with the *classification error rate* (CER). CER
337 is the ratio between the number of misrecognized words and the total number of words. On the other hand, the
338 phrase-based HTR system has been assessed with the *word error rate* (WER), which can be computed as the number

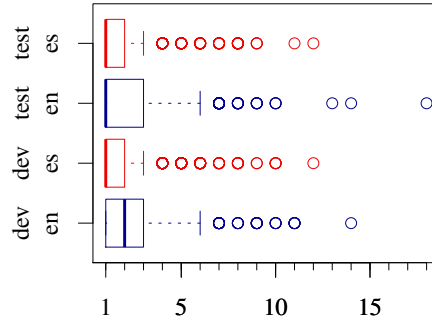


Figure 8: Box-and-whisker plots of the phrase lengths obtained from consecutive errors. In Spanish (es), the system commits typically between 1 and 3 consecutive errors although in rare cases it can commit up to 12 errors. In contrast, up to 6 consecutive errors can be considered normal in English (en). In this case, rare sentences contain at most 18 consecutive errors.

339 of substitutions, deletions and insertions needed to transform the hypothesis into the reference, normalized by the
 340 number of words in the reference. The results present the average error of the three users.

341 4.5. Results

342 In this section, we will compare the performance of the proposed systems. In order to make the references easier,
 343 we will name the different systems as follows:

344 **HTR.** The baseline HTR system as defined in Eq. (3).

345 **ERR.** The baseline HTR system after removing the erroneous word, Eq. (4).

346 **nPREF.** In Eq. (5), the latest n words of validated prefix in the target sentence are taken into account.

347 **M1.** In Eq. (11), information regarding the dictionary is used, but the alignment probabilities are uniform.

348 **M2.** In Eq. (11), the dictionary and the alignment probabilities are used.

349 **nWG.** In Eq. (12), the system uses an n -gram that has been extracted from the translation WG.

350 Furthermore, if the decoding is performed in an integrated way, the system will be marked with **+IMT**. Besides,
 351 several of the proposed systems can be combined by linear interpolation as in Eq. (15). In this case, we will use **+**
 352 symbol to mark which models were interpolated. The interpolation parameters were obtained in the development set
 353 to optimize the accuracy.

354 In addition, the proposed language models were encoded as n -grams. The aim of this is two-folded. First, we
 355 would like to leverage current HTR systems without custom software modifications. Second, since the new sources
 356 of information are added early in the HTR system, we expect to reduce the error cascade produced in post-processing

357 error correcting systems. However, although all the proposed models can be trivially encoded as 1-grams for the case
358 of word-based recognition, some of them cannot be encoded efficiently for n -grams as such and require special search
359 algorithms. As these cases are out of the scope of the current paper, such models will not be evaluated for phrase
360 recognition. Nevertheless, these models could also be applied in a post-processing rescoring stage. For instance, both
361 **M1** and **M2** models can be easily encoded as a 1-gram for word-based recognition. As there is just one possible
362 value for i and s , the 1-gram can be built by computing Eq. (11) for each word of the vocabulary. In contrast,
363 **M2** models cannot be encoded as n -grams for phrase recognition since the probability depends on the position i of
364 the hypothesized word, and then, i should be stored in the search algorithm for every word hypothesis. Luckily,
365 **M1** models assume independence of the position i so they can be encoded as a 1-gram even for the case of phrase
366 recognition.

367 Finally, as it is typical in modern HTR and IMT models, the different probability distributions must be scaled,
368 particularly the language model. Here, the optimum language model scaling factor, λ , was chosen to optimize the
369 average CER or WER in the development set of the three writers with the downhill simplex method [37]. There were
370 not significant differences in the optimum parameters obtained separately for each writer. Therefore, the estimation
371 of these parameters seems rather robust to the variability of writers.

372 Regarding the results for the word-based experiments, Fig. 9 shows the test CER for different values of λ for the
373 most relevant systems. First, it must be pointed out that the optimum λ from the development set approximated quite
374 well the test optimum, i.e., the estimation of λ does not present much overfitting. The only exception was the **2WG**
375 system for which an extra error reduction of 0.5% absolute points could have been achieved.

376 Second, we should note the effect of adding **ERR** to the system on the error rate. A small improvement can be
377 noticed in Spanish. However, the curves in English overlap. The explanation for this is a bit involving. Note that
378 Spanish is a more inflected language than English. For example, ‘both’ (in English) can be translated by ‘ambos’ or
379 ‘ambas’ (in Spanish), depending on the gender, and having very similar writings. In contrast, ‘añade’ (in Spanish)
380 can be translated by ‘adds’ (in English). Thus, we can see how translating from a less inflected language to a more
381 inflected language introduces extra ambiguity. Furthermore, the possible translations of ‘both’ present also a similar
382 spelling. Conversely, the ambiguity is reduced in the opposite direction. Table 1 shows the 5-best list of the HTR
383 scores for the words ‘ambos’ and ‘adds’. In the first case, ‘ambas’ and ‘ambos’ are the two most likely words in the
384 HTR system, which differ in just one character and have similar HTR scores. Now, imagine that the IMT engine
385 mistranslates ‘both’ to ‘ambas’, by changing the gender of the word. Then, by saying that *ambas* is not correct with
386 the **ERR** model, we give the system the opportunity to amend the error himself. However, in the English case, none
387 of the words are synonyms of the word to recognize, and thus is more difficult to find the mistranslated word at the
388 top of the n -best list. As a consequence, it is very unlikely that **ERR** achieves much improvement when translating
389 from Spanish to English.

390 With respect to the n **PREF** models, only **4PREF** has been displayed in the plots. The improvement over the base-
391 line is consistent and significant. The experiments were run on **2PREF**, **3PREF** and **5PREF** as well. However, only

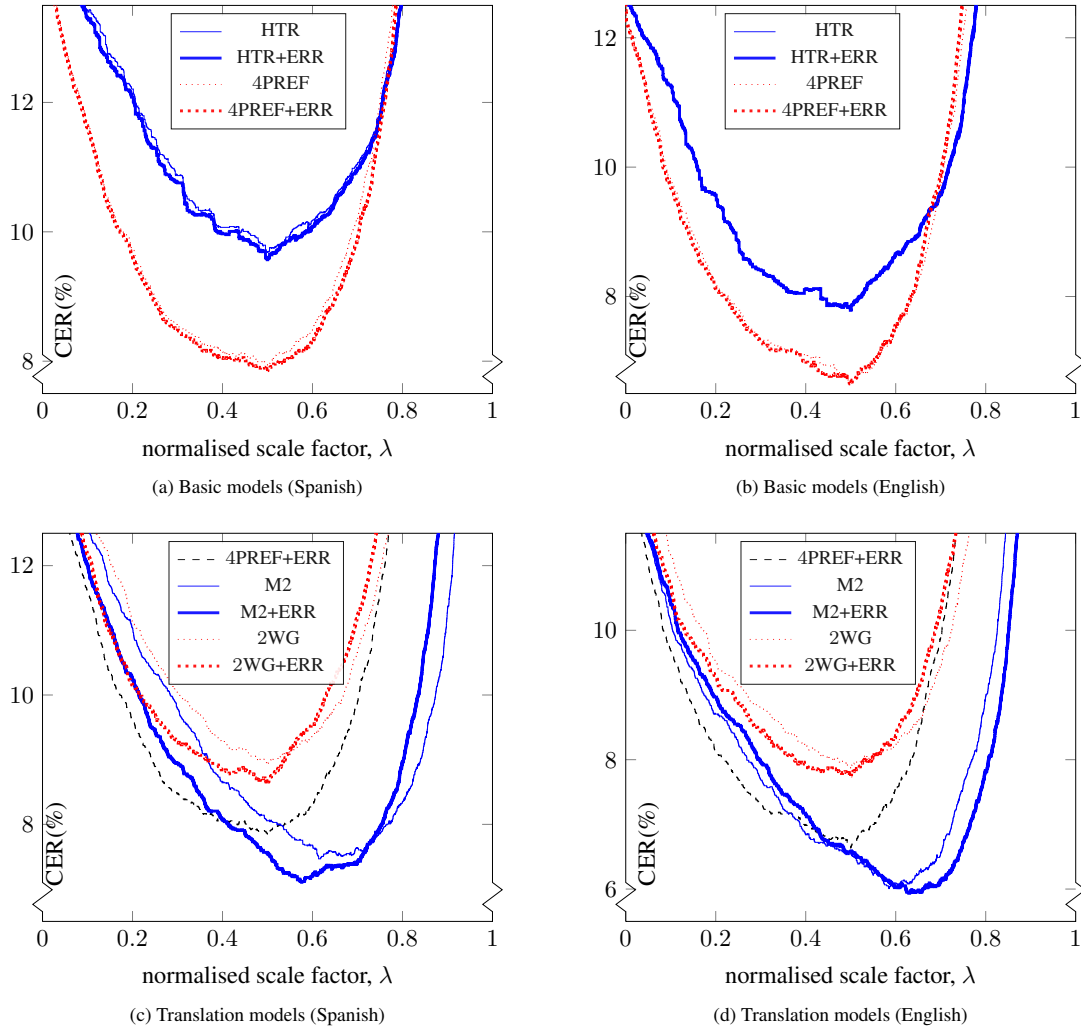


Figure 9: Development CER when modifying the λ scale factor. The x axis represents the variation of the normalized scale factor λ . The y axis shows the classification error rate (CER). At top, the comparison of the basic models described in [Sec. 3.1](#), [Sec. 3.2](#) and [Sec. 3.3](#). At bottom, the most relevant translation models described in [Sec. 3.4](#).

392 **2PREF** for English performed slightly worse than **4PREF**. Longer prefixes achieved almost the same performance.

393 With respect to the systems using the translation models in [Fig. 9c](#) and [Fig. 9d](#), we can see that these systems
 394 usually outperform the best basic system, **4PREF+ERR**. The exception for this is **2WG** for English, which shows a
 395 small performance degradation with respect to **4PREF+ERR**. Still, **2WG** systems do not seem to improve the basic
 396 systems significantly. Although several n WG systems were tested, any of them showed improvements over **2WG**. On
 397 the other hand, **M2** systems achieve good improvements, although they are simpler than **2WG**. A reason for that is

<i>both</i> → <i>ambos</i>		<i>añade</i> → <i>adds</i>	
word	HTR score	word	HTR score
<i>ambas</i>	651.6	<i>aids</i>	137.9
ambos	646.9	cities	105.6
cambios	390.7	cycles	91.6
amplias	384.1	adds	90.7
campos	344.4	circles	85.8

Table 1: 5-best list for the words **ambos** and **adds**, which have been misrecognized. The *cursive* word is the word the IMT system mistranslated and the user is amending.

398 that **M2** models have a smoother distribution probability and *nWG* systems need some sort of hypothesis pruning. In
399 fact, the average number of candidates with probability greater than zero is 292 for **M2** while it is 38 for **4WG**. **IMT**
400 suffer even more from this problem with 2 candidates average.

401 A summary of the different alternatives studied for the word-based experiments is shown in Table 2. First, with
402 only the basic information, **4PREF+ERR** clearly outperforms **HTR**. Second, using translation models we can achieve
403 further improvements. Since **M2** performs much better than **M1** we can deduce that alignment information is crucial
404 for the translation models. On the other hand, *nWG* performance is worse than word-based translation models. As it
405 has been explained before, that might be due to the poorly smoothed probability distribution. Another reason might
406 be that, in the process of obtaining *n*-gram models, information regarding alignments is lost as a result of the *n*-
407 gram assumptions. When interpolating with **4PREF**, **M2** models do not show significant improvements. In fact, for
408 Spanish, the system presents over-fitting, since performance in development improves but in test decreases. However,
409 **4PREF** smooths **2WG** distribution achieving close results to word-based models. Next, by introducing **IMT**, small
410 improvements can be obtained. Not surprisingly, **IMT** suffers from the same problems than *nWG*, but even more
411 prominent. Finally, including all systems we can observe the best results overall, except for the over-fitting in the
412 Spanish test set. Thus, **2WG** seems to contribute slightly to improve the final model accuracy.

413 Table 3 shows the WER for phrase-based recognition. First, it must be noted that the results for **ERR**, **M2**, and
414 **IMT** are not shown, since they would require a different search engine. In addition, it is worth of mention that the
415 baselines for phrase-based HTR have almost the double error rate than the word-based baselines. This is caused
416 primarily because the segmentation for the words in the phrases are unknown. Then, it is the search algorithm that
417 must find the most likely segmentation. As a result, segmentation errors are propagated to word errors. If we look at
418 the results regarding the *nWG* models, they perform unexpectedly bad when used alone. However, when interpolated
419 with **3PREF** they show a good improvement. As in word-based recognition, word-based translation models show the

System	Spanish	English
HTR	11.1	9.9
4PREF+ERR	9.9	9.5
2WG+ERR	9.8	9.4
M1+ERR	9.4	9.0
M2+ERR	8.6	7.7
2WG+4PREF+ERR	9.2	7.9
M2+4PREF+ERR	9.0	7.5
2WG+4PREF+ERR+IMT	9.2	7.9
M2+4PREF+ERR+IMT	8.9	7.5
ALL	8.9	7.4

Table 2: Summary of the CER results for word-based recognition. The results show various language modeling approaches for the test sets. In **boldface** the best systems.

System	Spanish	English
HTR	16.8	18.6
3PREF	16.3	18.0
2WG	18.9	19.7
M1	17.0	17.4
2WG+3PREF	16.2	16.6
M1+3PREF	15.2	15.5
M1+2WG+3PREF	15.2	15.5

Table 3: Summary of the WER results for phrase-based recognition. The results show various language modeling approaches for the test sets. In **boldface** the best systems.

420 best results, especially when interpolated with other models.

421 To sum up, all the proposed systems significantly outperform the baseline recognizer. Basic models obtain a
422 good improvement over the baseline. However, adding information from the translation may achieve remarkable
423 results. Although more complex translation models suffer from smoothing problems, they can also contribute when

424 interpolated with the rest of the models.

425 4.6. Error Analysis

426 An analysis (Table 4) of the results for the best word-based model shows that 49.2% to 54.4% of the recognition
427 errors were produced by punctuation and other symbols. To circumvent this problem, we proposed a contextual
428 menu in [16]. With such menu, errors would have been reduced (best test result) to 4.4% in Spanish and 3.5% in
429 English. Out-of-vocabulary (OOV) words plus zero probability (P0) words (the words for which the decoder assigned
430 zero probability or were pruned out) also summed up a big percentage of the error (40.3% and 28.9%, respectively).
431 Finally, the rest of the errors were mostly due to one-to-three letter words, which can be basically a problem of
432 handwriting morphological modeling.

433 On the other hand, phrase recognition presents a different error distribution. First, note that two new classes of
434 errors have been introduced: deletions and insertions. The former account for the words in the reference that have
435 been omitted, whereas the latter account for words inserted in the output hypothesis but do not correspond to any word
436 in the reference. Both contribute to generate hypotheses with lengths different to their respective references, since the
437 HMM models is not able to perform an accurate segmentation. Then, as a result, the proportion of recognition errors
438 from the 'others' category increases from 3 to 20. In contrast, the proportion of errors regarding punctuation symbols
439 decreases. Finally, it is to be remarked how the errors for short words have increased, probably because of small
440 insertions or deletions.

441 4.7. Reducing Effort Correcting HTR errors

442 In case an HTR error is committed, the user may fall back to the virtual keyboard and type the correct word. The
443 problem with this kind of keyboards is that typing is slow. To minimize this problem, we propose a contextual menu
444 with a list of the n -best candidates (excluding the erroneous word). The aim is to reduce the number of clicks needed
445 to obtain the correct word with respect to a conventional virtual keyboard. As a baseline, for each HTR mistake, we
446 count the number of clicks needed to input the correct word as: one click to pop up the keyboard, plus the number
447 of characters in the word, plus one click to close the keyboard. For the Spanish test set, the average number of clicks
448 per word amounts to 9.3, while for English it is 9.1 for the best word-based models in Table 2. This values can be
449 surprisingly high, since it is known that the average word length is 4.5, i.e. the average number of clicks per word 6.5.
450 However, it must be noticed that longer words are also more difficult to recognize. Thus, the average word length in
451 the erroneous words is higher.

452 If the contextual menu is used, we count: one click for opening the menu plus one for choosing a word. If the
453 correct word cannot be found in the n -best list, then we add: one count for the keyboard, plus the number of characters,
454 plus a closing click. In Fig. 10, we can see, on the left axis, the CER for a given size of the n -best list. Clearly, the
455 error almost reduces to a quarter, around $n = 5$, with respect to the baseline. Between 10 and 15, the error stabilizes.
456 Note that from 5 to 10 is still a reasonable amount of candidates to be shown in a circular menu. For more than 15,

class	words	word-based		phrase-based	
		es (%)	en (%)	es (%)	en (%)
punct.	., ,, :, ;, *, (,), —	49.2	54.4	14.0	18.6
1-char	a, e, y, o, u	4.1	0.9	8.3	2.3
2-char	of, if, la, by, on, is, ...	1.8	7.1	4.4	3.4
3-char	for, off, los, may, ...	0.0	4.3	2.1	4.9
numbers	xxvii, xxvi, xxiii, ...	2.3	0.9	2.1	2.3
OOV + P0	termina, luz, ...	40.3	28.9	20.2	13.6
others	latin, flash, fsma, ...	2.3	3.4	20.3	18.6
substitutions		100	100	71.5	63.8
insertions		—	—	3.0	4.6
deletions		—	—	25.5	31.6

Table 4: Detailed analysis of the word-based and phrase-based recognition errors. Five classes have been identified to produce the most amount of recognition errors. The second column shows samples of misrecognized words for these classes. Columns three and four are the percentage of these classes among the total number of misrecognized words for Spanish (es) and English (en), respectively. Columns five and six are the percentages for the phrase-based experiments. In this case, the percentage of substitutions, insertions and deletions is also shown.

457 the CER almost equals the error for OOV+P0, since they cannot be found in n -best lists. On the right axis, we can
458 observe the average number of clicks per word necessary to correct the mistakes. For $n = 1$ the number of clicks is
459 reduced to 2.0. A trade-off can be found at $n = 7$ with 1.83 (80% relative improvement w.r.t. the baseline) and 1.82
460 (78% relative improvement), for Spanish and English, whereas the lower bounds are 1.75 and 1.73, respectively.

461 5. Final Thoughts and Recommendations

462 While the techniques addressed in this paper have been focused on correcting machine translation output, in re-
463 ality some of them can be generalized to the correction of other automatically generated outputs. In particular, **ERR**
464 and n **PREF** can be used to improve HTR accuracy for any tasks in which n -grams can be used for language mod-
465 eling, e.g., [7]. Obviously, **M1** and **M2** are MT specific, but n **WG** can be used for many other structured prediction
466 problems where a word graph can be generated as an output. In fact in a similar way to this work, n **WG** has been suc-
467 cessfully used for speech-enabled user interfaces for IMT [13] and for dictation of historical documents [25, 26]. In
468 the same way, integrating HTR with interactive systems is possible for other applications as far as n **WG** is available.
469 Nonetheless, using more specific techniques, such as **M2**, although less general, have proven to be more effective.

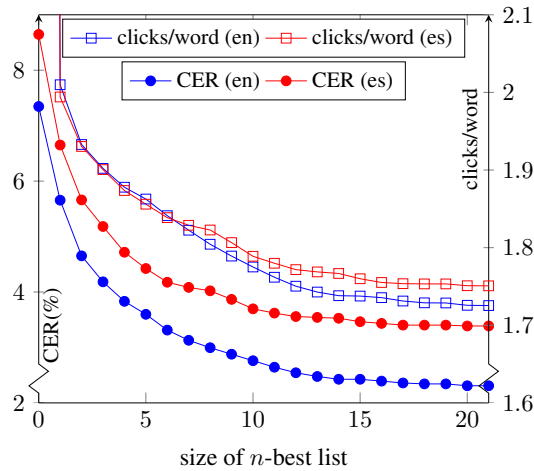


Figure 10: Reduction of CER and number of clicks as a function of the n -best list size.

470 Finally, we recommend that the integration of contextual information in the decoding is performed in an early stage of
 471 the decoding process, avoiding cascade errors. More importantly, if the techniques can be encoded as n -grams, as the
 472 techniques presented here, it will allow practitioners to improve their HTR systems without modifying their preferred
 473 HTR engines.

474 6. Conclusions and Future Work

475 In this paper we have described a task specific on-line HTR system to operate with an IMT application. We have
 476 shown that a tight integration of the HTR and IMT decoding process can produce significant HTR error reductions. It
 477 is worth of note that all the proposed systems significantly outperform the baseline recognizer. Basic models obtain
 478 a good improvement over the baseline. However, translation models achieve remarkable results. Although more
 479 complex translation models suffer from smoothing problems, they also contribute when interpolated with the rest of
 480 the models. We also have introduced a new method for correcting HTR mistakes that consists on a contextual menu
 481 with the n -best candidates. The results show that a list with as few as 7 candidates allows to correct the HTR mistakes
 482 with just 1.83 clicks per word.

483 On the other hand, the analysis of the results has shown two important issues to be tackled. First, the system should
 484 be able to decode unknown words since they are a clear limitation to system performance. A solution for this might
 485 be to use character language models instead of word language models, a technique that has achieved promising results
 486 in other areas. Second, phrase-based models could benefit from better smoothing methods. Alignment information
 487 should be also taken into account more explicitly in these models. Furthermore, other alternatives could also be
 488 explored, as more advanced word-based translation models (such as HMM, M3, M4 or M5) that cannot be used as
 489 n -grams in phrase-based decoding. These models could be used instead in the rescoring of the HTR WGs. Finally, if

490 the rescoring of WGs shows promising results, it would be interesting to directly implement the more advanced MT
491 models into the HTR search algorithm.

492 Acknowledgements

493 The research leading to these results has received funding from the European Union Seventh Framework Pro-
494 gramme (FP7/2007-2013) under grant agreement n° 287576 (CasMaCat), from the EC (FEDER/FSE), and from
495 the Spanish MEC/MICINN under the MIPRCV "Consolider Ingenio 2010" program (CSD2007-00018) and iTrans2
496 (TIN2009-14511) project. It is also supported by the Generalitat Valenciana under grant ALMPR (Prometeo/2009/01)
497 and GV/2010/067.

498 References

- 499 [1] A. Graves, M. Liwicki, S. Fernandez, R. Bertolami, H. Bunke, and J. Schmidhuber. A Novel Connectionist System for Unconstrained
500 Handwriting Recognition. *IEEE Transactions Pattern Analysis and Machine Intelligence*, 31:855–868, 2009.
- 501 [2] M. Pastor, A. Toselli, and E. Vidal. Writing speed normalization for on-line handwritten text recognition. In *Proceedings of the Eighth*
502 *International Conference on Document Analysis and Recognition, 2005*, page 1131–1135, 2005.
- 503 [3] Stefan Jaeger, Stefan Manke, Jrgen Reichert, and Alex Waibel. Online handwriting recognition: the NPen++ recognizer. *International*
504 *Journal on Document Analysis and Recognition*, 3(3):169–180, 2001.
- 505 [4] S. Quiniou, M. Cheriet, and E. Anquetil. Error handling approach using characterization and correction steps for handwritten document
506 analysis. *International Journal on Document Analysis and Recognition*, 15:1–17, 2011.
- 507 [5] F. Farooq, D. Jose, and V. Govindaraju. Phrase-based correction model for improving handwriting recognition accuracies. *Pattern Recogni-*
508 *tion*, 42(12):3271–3277, 2009.
- 509 [6] J. Devlin, M. Kamali, K. Subramanian, R. Prasad, and P. Natarajan. Statistical Machine Translation as a Language Model for Handwriting
510 Recognition. In *Frontiers in Handwriting Recognition (ICFHR), 2012 International Conference on*, page 291–296, 2012.
- 511 [7] A.H. Toselli, V. Romero, M. Pastor, and E. Vidal. Multimodal interactive transcription of text images. *Pattern Recognition*, 43(5):1814–1825,
512 2010.
- 513 [8] A. H. Toselli, E. Vidal, and F. Casacuberta, editors. *Multimodal Interactive Pattern Recognition and Applications*. Springer, 1st edition
514 edition, 2011.
- 515 [9] S. Barrachina, O. Bender, F. Casacuberta, J. Civera, E. Cubel, S. Khadivi, A. L. Lagarda, H. Ney, J. Tomás, E. Vidal, and J. M. Vilar. Statistical
516 Approaches to Computer-Assisted Translation. *Computational Linguistics*, 35(1):3–28, 2009.
- 517 [10] P. Koehn and B. Haddow. Interactive Assistance to Human Translators using Statistical Machine Translation Methods. In *Proceedings of the*
518 *MT Summit XII*, page 73–80, 2009.
- 519 [11] G. Foster, P. Isabelle, and P. Plamondon. Target-Text Mediated Interactive Machine Translation. *Machine Translation*, 12:175–194, 1998.
- 520 [12] G. Sanchis-Trilles, D. Ortiz-Martínez, J. Civera, F. Casacuberta, E. Vidal, and H. Hoang. Improving Interactive Machine Translation via
521 Mouse Actions. In *Proceedings of the conference on Empirical Methods in Natural Language Processing (EMNLP'08)*, page 485–494, 2008.
- 522 [13] V. Alabau, L. Rodríguez-Ruiz, A. Sanchis, P. Martínez-Gómez, and F. Casacuberta. On multimodal interactive machine translation using
523 speech recognition. In *Proceedings of the 13th international conference on multimodal interfaces (ICMI '11)*, page 129–136, 2011.
- 524 [14] S. Khadivi and H. Ney. Integration of speech recognition and machine translation in computer-assisted translation. *IEEE Transactions on*
525 *Audio, Speech, and Language Processing*, 16(8):1551–1564, 2008.
- 526 [15] E. Vidal, F. Casacuberta, L. Rodríguez, J. Civera, and C. Martínez. Computer-Assisted Translation Using Speech Recognition. *IEEE*
527 *Transaction on Audio, Speech and Language Processing*, 14(3):941–951, 2006.

- 528 [16] V. Alabau, A. Sanchis, and F. Casacuberta. Improving On-line Handwritten Recognition using Translation Models in Multimodal Interactive
529 Machine Translation. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language
530 Technologies (ACL'11)*, page 389–394, 2011.
- 531 [17] V. Alabau, D. Ortiz-Martínez, A. Sanchis, and F. Casacuberta. Multimodal interactive machine translation. In *Proceedings of the International
532 Conference on Multimodal Interfaces and the Workshop on Machine Learning for Multimodal Interaction (ICMI-MLMI '10)*, page 46:1–46:4,
533 2010.
- 534 [18] P. Koehn. *Statistical Machine Translation*. Cambridge University Press, 2010.
- 535 [19] Daniel Ortiz-Martínez, Ismael García-Varea, and Francisco Casacuberta. Online Learning for Interactive Statistical Machine Translation.
536 In *Proceedings of the North American Chapter of the Association for Computational Linguistics - Human Language Technologies (NAACL
537 HLT)*, page 546–554, 2010.
- 538 [20] Frederick Jelinek. *Statistical methods for speech recognition*. MIT Press, Cambridge, MA, USA, 1997.
- 539 [21] P. F. Brown, S. A. Della Pietra, V. J. Della Pietra, and R. L. Mercer. The Mathematics of Machine Translation. 19(2):263–311, 1993.
- 540 [22] R. Zens, F. Och, and H. Ney. Phrase-based statistical machine translation. page 35–56. 2002.
- 541 [23] F. J. Och and H. Ney. Discriminative Training and Maximum Entropy Models for Statistical Machine Translation. In *Proceedings of the
542 Association for Computational Linguistics 2002 (ACL'02)*, page 295–302, July 2002.
- 543 [24] N. Ueffing, F.J. Och, and H. Ney. Generation of word graphs in statistical machine translation. In *Proceedings of the ACL-02 conference on
544 Empirical methods in natural language processing (EMNLP'02)*, page 156–163, 2002.
- 545 [25] Vicent Alabau, Carlos D. Martínez-Hinarejos, Verónica Romero, and Antonio L. Lagarda. An iterative multimodal framework for the
546 transcription of handwritten historical documents. *Pattern Recognition Letters*, in press:–, 2012.
- 547 [26] V. Alabau, V. Romero, A. L. Lagarda, and C. D. Martínez-Hinarejos. A Multimodal Approach to Dictation of Handwritten Historical
548 Documents. In *Proceedings of the Interspeech 2011*, page 2245–2248, 2011.
- 549 [27] F. Wessel, R. Schluter, K. Macherey, and H. Ney. Confidence measures for large vocabulary continuous speech recognition. *IEEE Transac-
550 tions on Speech and Audio Processing*, 9(3):288–298, 2001.
- 551 [28] W. Campbell and F. Richardson. Discriminative Keyword Selection Using Support Vector Machines. In *Proceedings of the Neural Information
552 Processing Systems 2007 (NIPS'07)*, page 209–216, 2008.
- 553 [29] S. Katz. Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE Transactions on
554 Acoustics, Speech and Signal Processing*, 35(3):400–401, 1987.
- 555 [30] R. Kneser and H. Ney. Improved backing-off for m-gram language modeling. In *Proceedings of the International Conference on Acoustics,
556 Speech, and Signal Processing (ICASSP'95)*, volume 1, page 181–184. IEEE, 1995.
- 557 [31] H. Ney, U. Essen, and R. Kneser. On the estimation of ‘small’ probabilities by leaving-one-out. *IEEE Transactions Pattern Analysis and
558 Machine Intelligence*, 17(12):1202–1212, 1995.
- 559 [32] SchulmbergerSema S.A., Celer Soluciones, Instituto Técnico de Informática, R.W.T.H. Aachen - Lehrstuhl für Informatik VI, R.A.L.I.
560 Laboratory - University of Montreal, Société Gamma, and Xerox Research Centre Europe. X.R.C.: TT2. TransType2 - Computer assisted
561 translation. Project technical annex, 2001.
- 562 [33] Isabelle Guyon, Lambert Schomaker, Réjean Plamondon, Mark Liberman, and Stan Janet. Unipen project of on-line data exchange and
563 recognizer benchmarks. In *Proceedings of the 12th Int'l Conference on Pattern Recognition (ICPR-94)*, page 29–33, 1994.
- 564 [34] L.R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286,
565 1989.
- 566 [35] D. Ortiz-Martínez, I. García-Varea, and F. Casacuberta. Thot: a Toolkit To Train Phrase-based Statistical Translation Models. In *Proceedings
567 of the MT Summit X*, Phuket, Thailand, 2005.
- 568 [36] M. Luján-Mares, V. Tamarit, V. Alabau, C.D. Martinez-Hinarejos, M.P. i Gadea, A. Sanchis, and A.H. Toselli. iATROS: A speech and
569 handwriting recognition system. In *Proceedings of the V Jornadas en Tecnologías del Habla (VJTH'2008)*, page 75–78, 2008.
- 570 [37] J. A. Nelder and R. Mead. A Simplex Method for Function Minimization. *Computer Journal*, 7:308–313, 1965.

571 **Vicent Alabau** received the Master degree in computer science from the Universidad Politécnica de Valencia, Valen-
572 cia, Spain, in 2003. Since that moment he has been working in several projects for the “Pattern Recognition and
573 Human Language Technology Group” on speech recognition and translation, machine translation, and multi-
574 modal interactive tools for pattern recognition.

575 **Alberto Sanchis** received the Diploma and the Ph.D. degrees in computer science from the Universitat Politècnica
576 de València, Valencia, Spain, in 1997 and 2004, respectively. In 2000, he joined the Departamento de Sistemas
577 Informáticos y Computación, Universitat Politècnica de València, where he is currently serving as an Associate
578 Professor. His main research interests include pattern recognition and its application to speech and handwriting
579 recognition and machine translation. He is now leading a project on handwriting and speech transcription.

580 **Francisco Casacuberta** is a full professor in the Universitat Politècnica de València (UPV), Spain. He is also one
581 of the leaders of the Pattern Recognition and Human Language Technology research group in the Instituto
582 Tecnológico de Informatica. He received his Ph.D. degree from the Universidad de Valencia in 1981. His
583 research interests include the areas of syntactic pattern recognition, statistical pattern recognition, machine
584 translation, speech recognition and machine learning. He is a member of the Spanish AERFAI Society and the
585 IEEE Computer Society.

***Author Biography**

Vicent Alabau received the Master degree in computer science from the Universidad Politécnica de Valencia, Valencia, Spain, in 2003. Since that moment he has been working in several projects for the “Pattern Recognition and Human Language Technology Group” on speech recognition and translation, machine translation, and multimodal interactive tools for pattern recognition.

Alberto Sanchis received the Diploma and the Ph.D. degrees in computer science from the Universitat Politècnica de València, Valencia, Spain, in 1997 and 2004, respectively. In 2000, he joined the Departamento de Sistemas Informáticos y Computación, Universitat Politècnica de València, where he is currently serving as an Associate Professor. His main research interests include pattern recognition and its application to speech and handwriting recognition and machine translation. He is now leading a project on handwriting and speech transcription.

Francisco Casacuberta is a full professor in the Universitat Politècnica de València (UPV), Spain. He is also one of the leaders of the Pattern Recognition and Human Language Technology research group in the Instituto Tecnológico de Informatica. He received his Ph.D. degree from the Universidad de Valencia in 1981. His research interests include the areas of syntactic pattern recognition, statistical pattern recognition, machine translation, speech recognition and machine learning. He is a member of the Spanish AERFAI Society and the IEEE Computer Society.

