

Incorporating Syntactic Constraints in Recognizing Handwritten Sentences

Rohini K. Srihari and Charlotte M. Baltus

Center for Document Analysis and Recognition

520 Lee Entrance- Suite 202 (UB Commons)

State University of New York at Buffalo

Amherst, New York 14228-2567

Abstract

The output of handwritten word recognizers (HWR) tends to be very noisy due to various factors. In order to compensate for this behaviour, several choices of the HWR must be initially considered. In the case of handwritten sentence/phrase recognition, linguistic constraints may be applied in order to improve the results of the HWR. This paper discusses two statistical methods of applying syntactic constraints to the output of an HWR on input consisting of sentences/phrases. Both methods are based on syntactic categories (tags) associated with words. The first is a purely statistical method, the second is a hybrid method which combines higher-level syntactic information (hypertags) with statistical information regarding transitions between hypertags.

1 Introduction

Recently, there has been a great deal of interest in handwriting recognition systems. Such systems are referred to as on-line or off-line systems. In the former, a light pen and a special tablet are used by the writer to enter phrases or sentences. In the latter, ordinary handwriting on paper is scanned and digitized. Since handwriting is highly variable, it is advantageous to use linguistic constraints to derive the correct interpretation. It is important to note that the language employed in handwritten text tends to be informal.

This paper describes a method whereby syntactic constraints can be used in the problem of recognizing handwritten sentences. A handwritten word recognizer (HWR) is employed which, for each input word in the sentence, returns the best n words in the lexicon which match the structural features of that input word. In many cases, the HWR does not return the correct word as its top choice.

As an example, consider Fig. 1 which shows the digitized image of the sentence "He will call you when he is back" along with the output of the HWR. If we restrict the recognizer to return only the top choice for each input word, the sentence will be erroneously read as "He with call pen when he us back". Although the

correct words can be found in the top 2 choices, it requires the use of contextual constraints in order to override the (sometimes erroneous) top choice. This example also illustrates the tendency of word recognizers to misread short (less than or equal to 3 characters) words more frequently than longer words. Furthermore, short words tend to be pronouns, prepositions and determiners causing frequent word confusion between very different syntactic categories (e.g., as, an). The HWR problem is further compounded by the possible absence of the correct word in the top n choices (approximately 30% of the time in off-line systems).

Existing methods to apply linguistic knowledge to this problem consist of purely syntactic methods or purely statistical methods. Syntactic methods employ a grammar capturing the syntax of all possible input sentences and reject those sentence possibilities which are not accepted by the grammar. The problems with such a method are (i) the inability of the grammar to cover all possibilities (especially since informal language is frequently ungrammatical) and (ii) the computational expense involved in parsing (increased due to several word choices).

In this paper, we focus on two statistical methods having the same objective: to increase the percentage of correctly recognized (i.e., top choice) words thereby improving sentence recognition. Statistical methods ([Keenan *et al.*, 1991; Hull, 1992]) are formulated as some variation of Markov models and are based on transition frequencies between the syntactic categories represented by pairs (or n -tuples) of entities in the candidate sentence. In the first method, the entities correspond to individual words and first and second-order transitions between syntactic categories (tags) of words are employed. In this paper, we show that although this method does not reliably select the correct sentence candidate, it is effective in reducing the average word/tag neighbourhood size. The second method takes a hybrid approach, where both syntactic and statistical knowledge are employed in arriving at the best sentence choice.

2 Word Recognition

Since the training phase requires the processing of several thousand sentences, the long and computationally expensive procedure of digitizing followed by recognition

He will call you when he is back

he with call pen when he us back
 she will will you were be is bank

Figure 1: Digitized image of sentence “He will call you when he is back”.

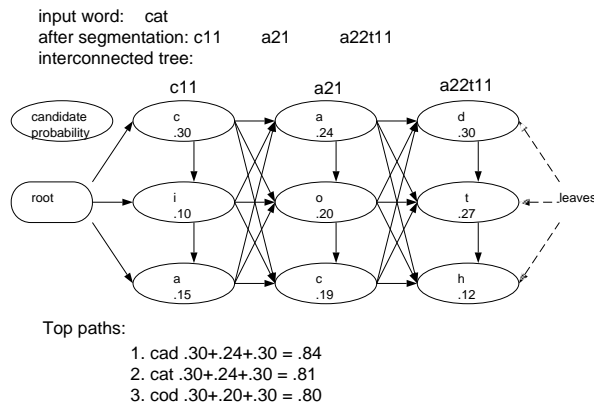


Figure 2: Illustration of HWR simulator on input word ‘cat’. The notation **a22t11** represents character choices for the case where the second part of the (erroneously) split *a* has been merged with the letter *t*.

is avoided by employing a program which simulates the output of an actual HWR. This is implemented as a two-stage procedure (similar to the actual HWR) where the first stage consists of generating candidate words. Based on the intermediate results of the actual word recognizer, we have computed statistics which model the behaviour of the first stage. These include substitution, splitting and merging statistics. Given an input (ASCII) word, and the above statistics, candidate (corrupted) words are generated based on simulating and propagating each of the above three types of errors at each character position. All candidate words which are in the lexicon are returned as the *direct recognition* output of the HWR. In case none of the words are found in the lexicon (62% of the time), the second stage is invoked. Here, the candidate words are matched against the lexicon and the top *n* choices (based on string-edit distances) are output. Figure 2 illustrates the performance of the simulator on the input word ‘cat’. Both substitution and splitting errors are depicted.

Although the HWR simulator does return confidences, we have found the relative ranking between word choices to be more reliable than the absolute confidences. Based on a comparison between the performance of the actual HWR on a given set as well as the simulated performance, the simulator has proven to be a reliable tool in the development stage of our research. It is important to note that the output of the HWR (and the simulator) degrades when the size of the lexicon is increased. Thus, information restricting the size of the lexicon (based on

physical and syntactic features of a word) can be helpful.

3 Model for Handwritten Sentence/Phrase Recognition

Figure 3 illustrates a computational model for a handwritten sentence/phrase recognizer which reflects the use of language models at every stage of the recognition process. The “tokenization” step is where lines in the input image are segmented into words and punctuation thus producing a stream of word/punctuation images for the next stage¹. The word buffer is fed into the wholistic filter whose purpose is to use intrinsic image features (such as word length, presence of ascenders/descenders) in reducing the lexicon. Word recognition consists of “segmentation” of the word image followed by character recognition in order to generate candidate (possibly illegal) words. String matching techniques are employed to produce a neighbourhood of legal word candidates (based on the reduced lexicon).

The next stage, namely semantic filtering, uses collocational information to modify word neighbourhoods. This is beyond the scope of this paper and will not be discussed here. Instead, this paper focuses on syntactic processing. The input to this stage consists of a sentence/phrase buffer containing the word neighbourhoods for each word. The output of this stage is either (i) a final word choice for each word in the sentence/phrase,

¹In the case of online recognition, such a buffer may not be necessary.

(ii) feedback information to the HWR post-processor in terms of eliminating syntactic categories thus enabling restricted lexicons to be utilized, or (iii) feedback to word recognition e.g., if syntactic analysis has determined that a particular token must be alphabetic only (as opposed to mixed alphanumeric), this information could be incorporated in a second “reading” of the word image.

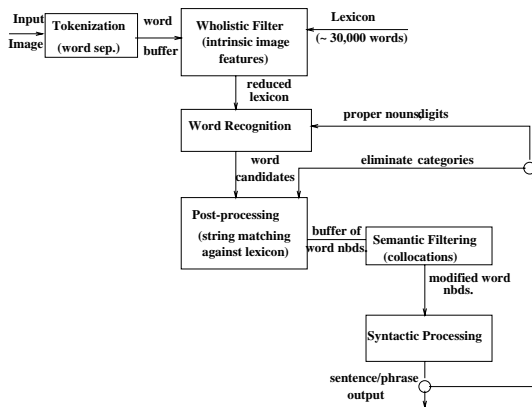


Figure 3: Overall Model of Sentence/Phrase Recognition

4 Method 1 - Incorporating Tag Transition Statistics

This method has been effective in improving the performance of the HWR by bridging the gap between (i) percentage correct in the top choice only, and (ii) percentage correct if the top 3 choices are considered. The performance improvement derives from selection of lower-rank words from the HWR output by considering the context of surrounding words. Ultimately, however, this level of analysis must provide feedback to the HWR postprocess so that errors of omission from the original HWR output can be corrected. A method for detecting typographical errors has been discussed in [Garside *et al.*, 1987]. For example, using a 21,000 word lexicon derived from the tagged Wall Street Journal Corpus, HWR simulation on a portion of the WSJ Corpus produced the correct word as the top choice 58% of the time, while the correct word was among the top three choices 69% of the time. Use of syntax to select among the candidates could produce, at best, a 69% correct output. The system becomes practical only if we can accurately pinpoint the errors, allowing use of feedback to the HWR post-process.

As a first step in achieving this ultimate objective, we have concentrated on improving the HWR post-process output by selecting second and third choice words when indicated by the surrounding context. For this purpose, a purely statistical model of syntax has been implemented. The course of the experiment covers two phases – a training phase, during which the statistics are generated, and a testing phase, during which the statistics are used to select from among HWR output word choices. The resulting selected sentences are then

compared against the actual input sentences to measure the success of the method. This model is patterned after that of [Hull, 1992] and [Keenan *et al.*, 1991], with the notable difference that in our model, the correct word may not be present at all, and is in fact missing in some 30% of the cases. The fact that our system is capable of recovering from these errors is an indication of its strength. It remains to be seen whether we can pinpoint those words which are suspect in order to iteratively improve the HWR performance.

4.1 The Training Set

Our intent is to use an informal language genre, since most handwritten text falls into this category. More specifically, we will be using collections of electronic mail messages as our training and test sets, since they are readily available and reflect the type of language we are modeling. Due to temporary problems with tagging this e-mail corpus, we are currently using the Wall Street Journal Corpus² in which each word has been tagged to indicate its syntactic category (one of 47).

The lexicon used in this method is derived from the entire corpus, including the test set, and includes the following information for each word: (i) all possible tags for the word, and (ii) word probabilities $P(\text{Word} | \text{Tag})$ for each word/tag pair.

First order transition statistics were collected over the entire corpus excluding the test set.

4.2 The Analysis Model

We use the notation $A : B$ to indicate the case where word A has been assigned the tag B . For each sentence analyzed, we form a word:tag lattice representing all possible sentences for the set of word choices output by string matching (see figure 4)³. The problem is to find the best path(s) through this lattice. Computation of the best path requires the following information: (i) tag transition statistics, and (ii) word probabilities.

Transition probabilities describe the likelihood of a tag following some preceding (sequence of) tag(s). These statistics are calculated during training as:

$$P(\text{tag}_B | \text{tag}_A) = \frac{\#(\text{tag}_A \rightarrow \text{tag}_B)}{\#(\text{tag}_A)}$$

Beginning- and end- of-sentence markers are incorporated as tags themselves to obtain necessary sentence-level information.

Word probabilities are defined as:

$$P(\text{Word} | \text{Tag}) = \frac{\#(\text{Word} : \text{Tag})}{\#(\text{AnyWord} : \text{Tag})}$$

The Viterbi algorithm [Forney Jr., 1973] is used to find the best Word:Tag sequence through the lattice, i.e., the

²a collection of 2.5 million words obtained from articles from the Wall Street Journal, as collected and tagged by the UPENN Treebank Project

³the presence of the DT tag in the trellis is explained below

WORD LATTICE

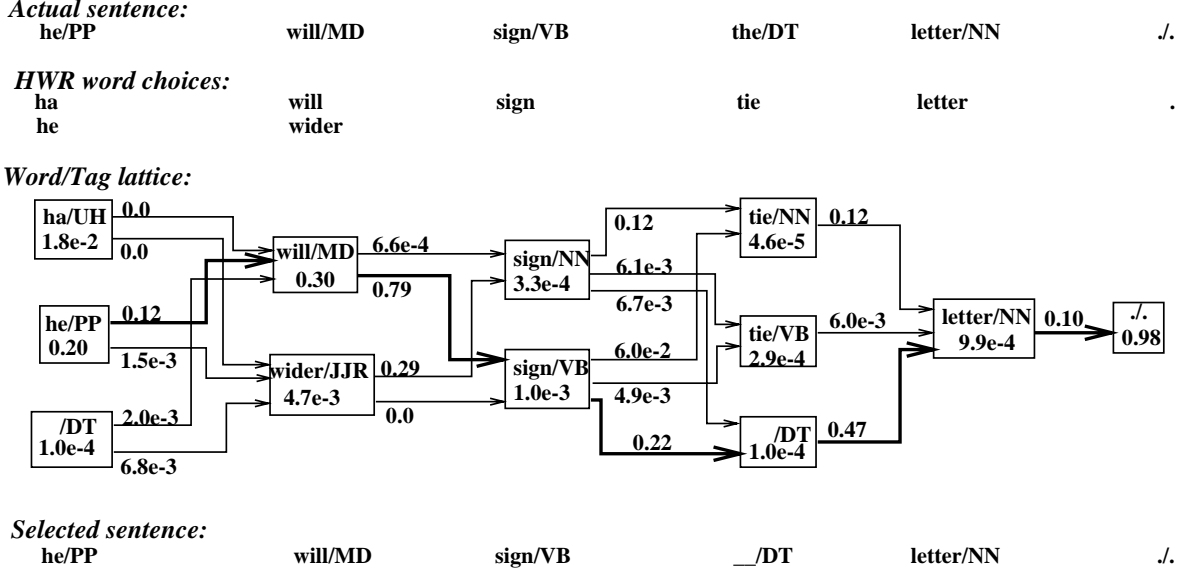


Figure 4: Sample Word:Tag Lattice For Analysis of HWR choices

maximal value of the following quantity:

$$\prod_{i=1}^n P(\text{Word}_i | \text{Tag}_i) P(\text{Tag}_i | \text{Tag}_{i-1})$$

over all possible tag sequences $T = \text{Tag}_0, \text{Tag}_1, \dots, \text{Tag}_{n+1}$ where Tag_0 and Tag_{n+1} are the beginning-of-sentence and end-of-sentence tags respectively. The Viterbi algorithm allows the best path to be selected without explicitly enumerating all possible tag sequences. A modification to this algorithm produces the best n sequences.

The lattice of Figure 4 demonstrates this procedure being used to derive the correct tag sequence even when the correct word (‘the’) was not output by the HWR. The chosen path is illustrated in boldface. The values on the edges represent tag transition probabilities and the node values represent word probabilities. Analysis showed that the correct tag most frequently missing from the lattice was the DT (determiner) tag, which was consistent with the fact that (i) the shorter words are more frequently missed by the HWR, and (ii) words of the category DT tend to be short words. We utilized this information by including the DT tag in the lattice in all cases of short words (< 4 characters) where the DT tag was not otherwise a candidate. It was also necessary to assign a default⁴ word probability to the node.

4.3 Results

The test set contained 93 sentences (approximately 2,000 words) from the WSJ corpus. Figure 5 illustrates the results obtained by running the above described procedure on the test set. For the sake of comparison, we have repeated the experiments with the correct word being

⁴Determined empirically.

added as the fourth choice in cases where the HWR had failed to produce it as one of the top 3 candidates. It is significant that this method is able to reduce the average word neighbourhood size even in cases where the correct word is missing, although the error rate (i.e., correct words erroneously discarded) is slightly higher. If the top 20 sequences are considered, this error rate is approximately 7.9%⁵. While this method is effective in reducing word neighbourhood sizes, it does not seem to be effective in determining the correct/best⁶ sentence or in providing feedback.

5 Hybrid Method

We attempt to overcome the problems associated with the first approach by employing a hybrid method, one which combines both syntactic as well as statistical knowledge. The input to this procedure is the reduced word/tag neighbourhoods obtained from the previous method. The training phase has been completed, the testing phase is in progress.

5.1 Training

The *training stage* consists of partitioning sentences into meaningful syntactic units (*hypertags*⁷) based on the part-of-speech tag that is associated with each word. Most of the 22 hypertag types loosely correspond to noun, verb, prepositional, adjective and adverbial phrases. For example, the sentence “The new

⁵This can be brought down by thresholding on the transition probabilities [Hull, 1992].

⁶The best sentence is one where all the words correctly recognized by the HWR in the top 3 choices are selected; the sentence may not be correct due to failure of the HWR to recognize all words.

⁷This notation was taken from [Garside *et al.*, 1987].

Test set: 93 sentences from WSJ Corpus
HWR results: top choice 57.98% top 3 choices 69.08%

CRITERIA	Correct Word Added				Correct Word Not Added			
	No. of Sequence Choices				No. of Sequence Choices			
	1	5	10	20	1	5	10	20
% correct tags	87.18	92.11	93.47	94.62	63.58	68.04	70.03	71.36
% correct words	87.93	93.45	94.88	96.21	59.68	63.44	64.87	65.75
% correct sentence	32.0	35.48	41.94	46.24	X	X	X	X
% best sentence	32.0	35.48	41.94	46.24	34.0	41.94	51.62	53.76
% reduction in word/tag nbd.	72.0	67.08	62.66	57.87	68.0	63.38	58.11	53.07
% error	8.5	6.56	5.12	3.79	15.8	12.0	9.2	7.9

Average word/tag nbd. size (correct words not added) before Method 1: 3.56

Avg. word/tag nbd. size if top 20 sequence choices are taken: 1.65

10 word sentence: $3.56^{10} \simeq 326,964$ sentences $1.65^{10} \simeq 150$ sentences

Figure 5: Results from Method 1- Tag Transitions

leader of our group will be coming” will be partitioned as follows:

[NPh1 (The DET) (new JJ) (leader NN)]

[Prep (of IN) (our PP\$) (group NN)]

[VPhX (will MD) (be BE) (coming VBG)]

The bracketing indicates the scope of hypertags where the first entry is the hypertag type and the entries following each word represent the tag assigned by the part of speech tagger. There are 4 verb hypertags: (i) **VPh3** denotes hypertags introduced by a 3PS verb form or auxiliary (e.g., ‘sits’, ‘has’), (ii) **VPhP** denotes hypertags introduced by non-3PS verbs (e.g., ‘talk’), (iii) **VPhI** denotes hypertags introduced by the infinitive ‘to’, and (iv) **VPhX** denotes verb hypertags with person indeterminate. Four noun hypertags are used: (i) **NPh1** denoting singular noun phrases, (ii) **NPh2** denoting plural noun phrases, (iii) **NPhF** denoting noun phrases containing foreign words, and (iv) **NPhX** denoting noun phrases where number is uncertain. There are two prepositional phrase hypertags, **Prep** and **VPhG**, where the latter denotes prepositional hypertags containing gerund verb forms (e.g., “in meeting the requirements”). **Adj** and **Adv** denote adjective and adverbial hypertags respectively. In addition, there are 8 single-word hypertags representing various pronoun classes, and some miscellaneous hypertags used to denote conjunctions, punctuation, etc. The complexity of parsing is overcome since a regular expression recognizer is sufficient for grouping the words into hypertags.

This stage is strongly motivated by the idea of “parsing by chunks” [Abney, 1991] and the theory of *phi-phrases* [Gee and Grosjean, 1983]. Due to the problems with tag n-grams described above, a method of naturally partitioning sentences into hypertags was desired. Phi-phrases correspond to prosodic phrases in sentences and thus represent a “natural” partitioning of sentences. Chunks on the other hand represent mid-level sentence constituents and are based on syntactic principles (i.e., a precise syntax is defined). The theory states that the syntax of structures is rigid; it is the relationship between chunks (i.e., the structure of sentences) which is flexible. If so, informal language can be modeled as probabilistic

sequences of transitions between fixed chunk types.

The idea of “hypertags” differs in some key ways from chunks and phi-phrases. Hypertags in our system are a judicious combination of chunks and phi-phrases designed with sentence recognition (rather than sentence understanding) as the goal. For example, possessives within NPh hypertags (e.g., “John’s house”) are allowed, however, constituents such as Wh-NP’s (e.g., that) are not grouped with neighbouring chunks. Secondly, a classification system for hypertags has been developed.

Based on the sentences that have been processed by the above method, we compute the following statistics:

- transition frequencies $P(H_B|H_A)$ where H_A and H_B are hypertags.
- hypertag probabilities

$$P(TS_i|H_j)$$

over all possible hypertags H_j , all possible observed tag sequences TS_i

5.2 Testing

The hypothesis being tested is that the hypertag method produces “better” sentence candidates than the first method where better implies that the percentage of correct words is greater. The procedure to score the sentence candidates is as follows:

- generate candidate sentences based on the reduced word/tag neighbourhoods produced by the previous method.
- Group the candidates into classes based on identical tag sequences.
- Hypertag each of the classes and generate a normalized score (due to varying number of hypertags) based on the (i) transition probabilities, (ii) hypertag probabilities, (iii) ratio of sentence length to number of hypertags.
- Differentiate among identical tag sequences using word probabilities and HWR rankings; output top-scoring sentence.

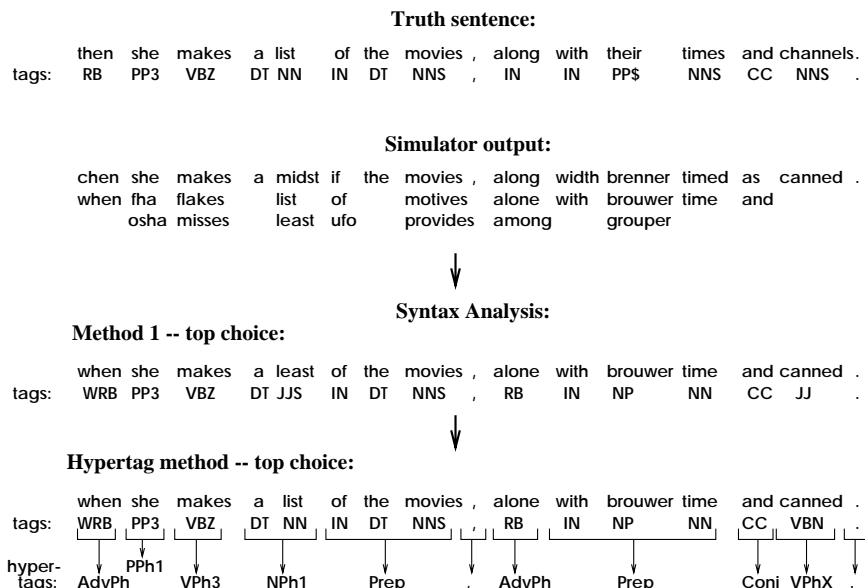


Figure 6: Method 1 (Tag Transitions) vs. Method 2 (Hypertag Transitions)

As an example consider Figure 6 which depicts the improvement of this method over Method 1 alone. In this case, the *best* sentence was not present in the top 20 sentence choices output by Method 1. The hypertag method ranks the *best* sentence as third, but more significantly, produces a better top choice than the first method. The only correct word it has failed to select is the word ‘along’. This error may be overcome if the words ‘along with’ are recognized as a significant collocation.

We are experimenting with a dynamic method (one that avoids enumeration of all possible sentences) which combines both methods. Central to the technique is an adaptive-length Viterbi procedure which is necessitated due to varying number of hypertags/sentence depending on word choices. The hypertag method gives us two more sources for detecting potential errors: (i) the distribution of tags within a hypertag, e.g., an unusual NPh1 pattern may signify an error, and (ii) the ratio of sentence length to the number of hypertags, e.g., a long sequence of single word hypertags may indicate an error. We are also investigating the use of semantic classes of words [Brown *et al.*, 1992] as an additional knowledge source. This would necessitate expanding the tagset to include both syntactic as well as semantic part-of-speech tags. The major advantage this will provide is in choosing between syntactically identical words.

6 Summary

A model for sentence/phrase recognition has been outlined with emphasis on improving the noisy output of the handwritten word recognizer, the central component of the system. We have presented two statistical methods of incorporating syntactic constraints into the system. The first technique is effective in reducing the word/tag neighbourhood; initial experiments indicate that the second method is effective in producing the best sentence.

Tests are currently being conducted in order to validate this claim.

References

- [Abney, 1991] Steven P. Abney. Parsing by Chunks. In Robert C. Berwick, Steven P. Abney, and Carol Tenny, editors, *Principle-Based Parsing: Computation and Psycholinguistics*, pages 257–278. Kluwer Academic Publishers, Boston, 1991.
- [Brown *et al.*, 1992] Peter F. Brown, Peter V. deSouza, Robert L. Mercer, Vincent J. Della Pietra, and Jenifer C. Lai. Class-Based n -gram Models of Natural Language. *Computational Linguistics*, 18(4):467–479, 1992.
- [Forney Jr., 1973] G. E. Forney Jr. The Viterbi Algorithm. *Proceedings of IEEE*, 61:268–278, 1973.
- [Garside *et al.*, 1987] Roger Garside, Geoffrey Leech, and Geoffrey Sampson. *The Computational Analysis of English*. Longman, 1987.
- [Gee and Grosjean, 1983] J. Gee and F. Grosjean. Performance Structures: a Psycholinguistic and Linguistic Appraisal. *Cognitive Psychology*, 15:411–458, 1983.
- [Hull, 1992] Jonathan J. Hull. Combining Syntactic Knowledge and Visual Text Recognition: A Hidden Markov Model for Part of Speech Tagging In a Word Recognition Algorithm. In *Proceedings of AAAI Fall Symposium on Probabilistic Approaches to Natural Language Processing*, pages 77–83, 1992.
- [Keenan *et al.*, 1991] F.G. Keenan, L.J. Evett, and R.J. Whitrow. A large vocabulary stochastic syntax analyser for handwriting recognition. In *Proceedings of the First International Conference on Document Analysis (ICDAR-91)*, pages 794–802, 1991.