

Combination of explicit segmentation with Seq2Seq recognition for fine analysis of children handwriting

Omar Krichen, Simon Corbillé, Eric Anquetil, Nathalie Girard, Elisa

Fromont, Pauline Nerdeux

▶ To cite this version:

Omar Krichen, Simon Corbillé, Eric Anquetil, Nathalie Girard, Elisa Fromont, et al.. Combination of explicit segmentation with Seq2Seq recognition for fine analysis of children handwriting. International Journal on Document Analysis and Recognition, 2022, 10.1007/s10032-022-00409-4. hal-03845144

HAL Id: hal-03845144 https://hal.science/hal-03845144

Submitted on 9 Nov 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Combination of explicit segmentation with Seq2Seq

recognition for fine analysis of children handwriting

Omar Krichen^{2*†}, Simon Corbillé^{1†}, Éric Anquetil², Nathalie Girard¹, Élisa Fromont¹ and Pauline Nerdeux²

> ¹Univ Rennes 1, IRISA lab, Rennes, F-35000, France. ²INSA Rennes, IRISA lab, Rennes, F-35000, France.

Contributing authors: firstname.lastname@irisa.fr; [†]These authors contributed equally to this work.

Abstract

We consider the task of analysing children handwriting in the context of a dictation task. The objective is to detect orthographic and phonological errors. To achieve this goal, we extend an existing handwriting analysis engine, based on an explicit segmentation of the handwritten input, originally developed for children copying exercises. We present a new approach, based on the combination of this analysis engine with a deep learning word recognition approach in order to improve both the recognition and segmentation performance. Explicit segmentation needs prior knowledge, and the deep network recognition predictions are a reliable approximation of the ground truth which can guide the analysis process. We propose to combine multiple prior knowledge strategies to further improve the analysis performance. Furthermore, we exploit the deep network approximate implicit segmentation to optimise the existing analysis process in terms of complexity.

Keywords: Online handwriting recognition, Segmentation, Digital learning, Degraded handwriting, Sequence-to-sequence, e-education

1 Introduction

This work aims at designing an educational system targeted towards primary school children, in order to help them master handwriting and spelling skills. More specifically, we deal with online interpretation of children handwritten **French cursive words**. The interpretation task in hand is a word analysis task, which differs from the word recognition task. Fig. 1 illustrates these differences. In a recognition task, the objective of the system is to predict the correct character sequence, whereas the objective of the analysis task is to provide a qualitative evaluation. Consequently, the segmentation quality is instrumental, to enable the system to perform a fine-grained analysis of



Fig. 1 Context: analysis of children handwriting: the dictated instruction is "alors" ("then" in French)

the pupil handwriting, such as highlighting in red the spelling mistakes directly on the ink. (c.f. Fig. 1). Therefore, the educational system needs both an accurate recognition of the child's word but also a good segmentation at character level to precisely locate the spelling mistakes. To achieve this goal, we build on previous works on children handwriting analysis for cursive French words [1]. This approach is based on an explicit segmentation of the input word. A segmentation graph representing all possible segmentations of the word into letters is created. For each node of the graph, letters hypotheses are computed using a letter recognition and analysis system. The analysis result is a set of n best possible pseudo-word hypotheses. In order to be efficient, the explicit segmentation needs to be driven by prior knowledge, especially to deal with degraded children handwriting. Since the instruction to copy was displayed to the child, it served as prior knowledge to guide the letter hypotheses computation phase. This "base system" is discussed in more details in



Fig. 2 Difference between a copy and a dictation task.

section 3. Our new targeted dictation task introduces new challenges, as illustrated in Fig. 2. The instruction is heard, not seen, by the pupil. This may induce a lot more spelling mistakes. In the figure, the written word "mai" is a homophone of the dictated instruction "mes". In this dictation context, the instruction is not directly exploitable to guide the analysis of the handwritten word. To provide a relevant and real-time analysis for this dictation task, new prior knowledge generation strategies are needed. We propose to combine the aforementioned engine, with a deep learning word recognition approach, namely a Seq2Seq architecture. Our contributions consist in exploiting this hybridisation in three different manners: 1) We define the Seq2Seq network recognition process as a new prior knowledge generation strategy, which will drive the analysis process; 2) We combine different prior knowledge strategies to further improve the system's performance; 3) We exploit the Seq2Seq implicit segmentation to prune the explicit segmentation graph and optimise analysis complexity. This paper is organised as follows. Section 2 presents related works about handwriting recognition and segmentation. Section 3 provides a detailed account of the existing engine, while Section 4 describes the deep learning model used for our task. Section 5 presents the approaches combination and our listed contributions. Experiments are presented in Section 6. Conclusion and future works are given in Section 7.

2 Related works

This section presents the latest online and offline methods concerning handwriting recognition and segmentation. Handwriting can be represented offline, through an image, or online through a sequence of points. IAM datasets (offline [2] and online [3] versions) are composed of English adult-written sentences, labelled at line level. They are open and widely used to compare pure recognition methods. To the best of our knowledge, there are no available words datasets with character-level annotation.

2.1 Handwriting text recognition

Deep learning models outperform the previous methods [4][5] on handwriting text recognition (HTR) task. These traditional methods were based on a **bottom-up strategy**, *i.e.* by using expert knowledge to segment input data, then recognising the character in each segmented element. A great advantage of deep learning models lies in the fact that they are end-to-end trainable. There is no need to segment the data, and the feature extraction is learned by the model. The two main deep learning approaches that tackle HTR are Connectionist Temporal Classification (CTC) [6] and Sequence to Sequence (Seq2Seq). The CTC approach divides input into frames for symbol prediction and computes a probability distribution over all possible outputs alignments, while the Seq2Seq approach translates an input sequence represented by an image into a sequence of characters. The CTC-based architectures designed for online recognition use Bidirectional Long-Short Term Memory [7] (BLSTM). The authors of [8] show that this type of architecture outperforms a traditional method based on Hidden Markov Models, whereas the authors of [9] use BLSTM with Bézier curvers encoding of online data to achieve state of the art performances for online recognition on IAM-OnDB [3]. The CTC-based architectures designed for offline recognition are slightly different due to the nature of the input data. Convolutional recurrent neural networks [10] [11] are based on a convolutional neural network coupled with a recurrent network with LSTM cell. The authors of [12] use a Seq2Seq method based on an encoderdecoder model with an **attention module** to do offline recognition. More recently, [13] and [14] use transformers, which need a lot of synthetic data to perform well, for offline recognition. For our work, we use a Seq2seq model since this architecture gets state of the art performances when no synthetic data are used. The next part present methods which focus on handwriting segmentation.

2.2 Handwriting segmentation

The authors of [15] propose regularisation methods on the CTC loss based on entropy and spacing to increase recognition performance and segmentation quality. They present a quantitative analysis on recognition performance and qualitative analysis on segmentation performance. The authors of [16] use a convolutional prototype network and most aligned frame based CTC training for handwriting recognition. They evaluate the recognition performance of their model on IAM [2] dataset whereas the segmentation is evaluated on a synthetic dataset representing a sequence of digits from MNIST [17] dataset. In this work, we choose to combine the Seq2Seq model good recognition performances with the explicit segmentation based existing engine [1] presented in the introduction. The next section 3presents the existing system.

3 Existing analysis engine

In this section, we present the existing analysis engine (for more details, see [1]). Fig. 3 illustrates its global principles. Given the handwritten input and the instruction, the first step of the analysis is the explicit segmentation process. A segmentation graph is constructed based on the extraction of all possible cutting points around descending zones [18], and represents a partition of all possible segmentations given the extracted cutting points. Fig. 4 illustrates the segmentation graph for the French handwritten word "juste". Every node of the graph represents a possible letter hypothesis. The objective is to find the best path in the graph corresponding to the correct segmentation. For each node, confidence-based classifiers [19] compute letters hypotheses. The analysis process is generic and relies on *prior knowledge generation* strategies. Here, prior knowledge is instrumental,



Fig. 3 Existing analysis engine, here the instruction serves as prior knowledge to guide the analysis



Fig. 4 Segmentation graph for the word "juste"

especially in the context of degraded handwriting to avoid recognition confusion at the letter level.

In a copying context, the prior knowledge strategy is straightforward. The instruction drives the letter computation process by filtering the computed hypotheses that belong to the instruction. The best segmentation path is the one which minimises the edit distance with the instruction. This strategy is best suited when the child correctly reproduces the instruction. A first adaptation of this engine to the dictation context was proposed in [20]. Two prior knowledge generation strategies were defined to deal with the fact that driving the analysis by the instruction, in a dictation context, becomes obsolete. The first strategy consisted in asking the child to type was he/she has written on the keyboard. This childtyping drives the analysis, since it is a pretty reliable estimation of the ground truth. However, the objective is to be free from user input and to rely solely on the system capacities. The second prior knowledge generation strategy was to generate, for every instruction, a set of phonetically similar pseudo-words. For example, if the instruction is "alors" (then in French), the generated hypotheses would be "alaur, alor, alord, alort". This generation is based on the *Phonetisaurus engine* [21], a grapheme-to-phoneme WFST (Weighted Finite State Transducer). A Recurrent Neural Network Language Model (RNNLM) is used to extract the best phonetic hypotheses for a given word. This prior knowledge generation strategy enables to cover potential orthographic errors that sound similar to the instruction. The limit of this strategy resides in the fact that it could not cope with written words that were not phonetically similar to the dictated instruction. It is in order to overcome this limitation that we choose to combine the existing analysis engine with the outputs of a Seq2Seq model, namely the predicted word and the correspondent implicit segmentation. The new prior knowledge generation strategy will therefore rely on the Seq2Seq predicted word to drive the generic analysis process. Section 4 describes the Seq2Seq architecture used, whereas section 5 presents the combination of the approaches and its impact.

4 Deep learning model for handwriting recognition

Our Seq2Seq model is derived from [12] for the encoder decoder architecture with hybrid Bahdanau attention mechanism [22] and [10] [11] [23] [14] for the encoder architecture. The encoder's parameters result of an ablation study where the number of convolutional, pooling, blstm layers and dropout are tested.

The authors of [12] demonstrate that using a joint training between encoder and decoder improves recognition performance. The encoder is trained with CTC loss [6] and the decoder with a cross entropy. Thus, the model makes one prediction with the encoder and one prediction with the decoder. The final loss is defined as follows:

Loss = $\lambda * Loss_{ctc} + (1 - \lambda) * Loss_{crossentropy}$, with $\lambda \in [0, 1]$ Fig. 5 illustrates the connection with its three main parts: 1) The Encoder performs the feature extraction of the input image into a feature vector. This vector is used by the encoder to make a word prediction; 2) The Attention module focuses the decoder on a specific area in the feature vector; 3) The Decoder decodes the feature vector and produces a word prediction.

The model takes as input a grayscale image resized proportionally to have a height of 128 pixels. The encoder first extracts spatial features with convolutional layers, then temporal features with recurrent layers, into a feature vector. This feature vector is used by the encoder to make a prediction and by the decoder through the attention module. The table 1 details the encoder's parameters.

Fig. 6 illustrates the attention mechanism. The idea is to focus the decoder on a specific part of



Fig. 5 Global architecture of sequence to sequence model

 Table 1
 Configuration of encoder: k is for kernel size, s

 for stride, p for padding and d for dropout. All

 convolution layers are followed by the Leaky ReLU

 activation function, then a layer normalization.

Type	Configuration
Input	height 128 * width
Convolution	#filters:8, k:3*3, s:1, p:0
Max pooling	k:2*2, s:2, d:0.2
Convolution	#filters:16, k:3*3, s:1, p:0
Max pooling	k:2*2, s:2, d:0.2
Convolution	#filters:32, k:3*3, s:1, p:0
Max pooling	k:2*2, s:2, d:0.2
Convolution	#filters:64, k:3*3, s:1, p:0, d:0.2
Convolution	#filters:128, k:4*2, s:1, p:0, d:0.2
Collapse Convolution	#filters:128, k:9*1, s:1, p:0
Batch normalization	
BLSTM	4 layers, 128 units, d:0.5
Fully connection	size alphabet $+ 1$

the feature vector, and thus ideally use features associated with a sub image representing one letter. The attention module produces at each time a context vector c_t from the feature vector emitted by the encoder and uses the hidden state of the decoder s_t . At each time, the decoder uses an embedding of the precedent prediction and the precedent context vector to update the hidden state s_t of the LSTM layer, then uses the hidden



Fig. 6 Details attention module and decoder: the input is the feature vector produced by the encoder. The decoder produces one character at a time. It starts with the special character $\langle sos \rangle$ (start of sequence) and ends with $\langle eos \rangle$ (end of sequence). FC stands for fully connected layer, Tanh the tangent hyperbolic function and Embed the embedding of one prediction



Fig. 7 Example of segmentation for the encoder/decoder.

state to concatenate with the current context vector to produce the symbol prediction at the time t. The decoder's alphabet uses two extra symbols for the start and the end of the characters sequence (<sos> and <eos>).

For the character segmentation aspect, an approximation can be computed from the encoder or decoder prediction. For the encoder, we compute the receptive fields used to predict a character and extract the associated part of the image to get the segmentation. For the decoder, we re-use the attention map used by the decoder to predict a character and find its position in the associated input image. Fig. 7 illustrate an example of segmentation of the French word "comme" by the Seq2Seq model. The segmentation quality is average due to the fact the network is trained on the recognition task. For the encoder, the letter "o" and "m" are incomplete. The decoder segmentation contains a lot of overlap between the letters "o" and "e". Section 6 details quantitative results for the segmentation and recognition evaluation.

This motivates our choice to combine a deep learning model, which does well recognition-wise, with the existing analysis engine, which does well segmentation-wise. Furthermore, even if the segmentation of deep model is approximate, it can be exploited to prune the explicit segmentation graph. The next section describes the hybridization between the two systems.

5 Combining deep recognition and explicit segmentation

In this section, we present the integration of the Seq2Seq recognition results into the explicit segmentation-based analysis process, the new prior knowledge generation strategies, as well as the pruning of the explicit segmentation graph.

5.1 Seq2Seq prediction as prior knowledge strategy

Fig. 8 illustrates the defined prior knowledge generation strategy, which consists in coupling the explicit segmentation-based analysis approach with the Seq2Seq recognition outputs. The predicted sequence for each written word drives the generic analysis process, especially in the letter hypotheses computation phase, and the word paths search phase. Being a better approximation of the ground truth in a dictation context, this *deep prediction strategy* improves the engine performance, as we will see in section 6.

A valid interrogation would be to question the fact that our system now has two recognition processes. A recognition process for each letter hypothesis (with Evolve classifier [24]), and a Seq2Seq recognition on the whole word. Shouldn't we rely on one or the other? The final goal is to provide feedback to the pupils at the ink level, therefore the segmentation process is as important as the recognition process in our task. The fact that the existing analysis system relies on an explicit segmentation process, with a recognition at letter level, ensures that the predicted result is coherent in terms of letters localisation. However, since we are faced with degraded children handwriting, the system needs some prior knowledge to prioritise the relevant letters hypotheses, hence the guidance of the analysis by the deep predicted sequence.



Fig. 8 Deep prediction as prior knowledge strategy

5.1.1 Deep prediction added value

Fig. 9 illustrates the analysis of the written word "zme", given the dictated instruction "cent" (hundred), with the three strategies: a) Instruction strategy with result="cent"; b) Phonetic strategy with result="cent"; c) Deep recognition strategy with result="zme". The instruction strategy is well suited when there are no errors, but can't cope with the analysis of children mistakes. As for the phonetic strategy, it is not well adapted to this situation either, since the written word "zme" does not sound similar to the dictated instruction "cent". As for the third strategy, since the network was able to predict the correct word, the injection of this prior knowledge enabled the engine to correctly recognize and segment the word.

5.2 Strategies combination

Until now, we have studied the case where the Seq2Seq model is able to predict the correct



Fig. 9 Results of analysis strategies for the word "zme", given the instruction "cent"

sequence, and therefore have a positive impact as prior knowledge on the analysis engine. However, there are cases where it is not able to correctly interpret the input, such as in Fig. 10, which illustrates the analysis results of the written word "biin", given the dictated instruction "bien". We can see that the first two strategies ((a) and (b)) were only able to predict the first written letter "b", which is also the first letter of the dictated instruction, whereas the third strategy (c) was only able to predict the latter part of the word "iin". Intuitively, since every strategy is best suited to a specific scenario, it is fair to assume that they could be complementary. We propose therefore to combine these strategies into a fourth one, named fusion and competition. The latter represents two ways of combining strategies, first a conjunction by merging these prior knowledge, then a dis-junction by introducing a notion of competition between the strategies prediction. We



Fig. 10 Results of analysis strategies for the word "biin", given the instruction "bien"

present now in detail the two steps of this fourth strategy.

5.2.1 Fusion

We propose the fusion of the results of the three mentioned strategies to generate an alternative approximation of the ground truth, which will serve as another prior knowledge source driving the analysis. This fusion is done in two steps: first by aligning the resulting character sequences using dynamic programming techniques, and second by introducing a voting algorithm called Rover [25], which chooses to most occurring character in the alignment. Fig. 11 illustrates the alignment and fusion of the above cited strategies, with the addition of the instruction and the deep model prediction. The fusion result corresponds to the ground truth "biin". Therefore, if used as prior knowledge, it will enable the analysis engine to predict the correct word.

11	V			N		Doop production
	ĸ	•	1	IN		Deep prediction
n /	В	I	E	Ν		Dictated instruction
liim.	В		I	Ν		Instruction driven analysis (a)
	В		I	Ν	S	Phonetic driven analysis (b)
	К	I	I	N		Deep driven analysis (c)
	В	I.	T	Ν		Fusion: alignement and Rover

Fig. 11 Alignment and fusion of multiple prior knowledge for the written word "biin"

5.2.2 Competition

After the fusion step, which adds pertinent prior knowledge information, we introduce the competition step, which enables the system to choose the best strategy, depending on child production. Fig. 12 illustrates this process. To choose the best prediction between instruction strategy, phonetic strategy, deep prediction strategy, and the fusion, we exploit metrics that are already present in the existing analysis engine. As explained in section 3, the result of each analysis process is the segmentation path, which minimises the edition distance with the prior knowledge that guides the instruction. This edition score consists of a Damerau-Leveinshtein [26] distance computed between the word hypothesis and the prior knowledge (e.q. the instruction). In addition, optimised costs are learned by the analyser [1]. Another indication is the handwriting quality, represented by the analysis score. The analysis score S_a of a path



Fig. 12 Fusion and competition strategy

of length n P_n is defined as follows, where $S_a(i)$ is the analysis score of the *ith* element of the path:

$$S_a(P_n) = \sqrt{\prod_{i=0}^n} S_a(i)$$
 [1].

Given these two metrics, we define a phonetic score that combines edition score pertinence and handwriting quality. The phonetic score is defined as follows:

 $PhoneticScore(P) = S_a(P) * 0.7 + \frac{1}{1 + \mathsf{EditScore}(\mathsf{P})} * 0.3$

The strategy chosen is the one where the predicted segmentation path has the best phonetic score. These parameters (0.7, 0.3) are chosen empirically to give more weight to the analysis score of each strategy. We will see in detail the impact of fusion and competition strategy in section 6. In this section, we have presented the integration of the Seq2Seq recognition results in the existing analysis chain and the proposed strategies to optimise the analysis process. Another output of the Seq2Seq model is the result of the implicit segmentation. We choose to use this segmentation result in order to prune the existing analysis process segmentation graph, which would enable to diminish the complexity of the process. Since we are in the context of real-time user interaction, the response time of the system has to be acceptable to the user. However, for long words, the analysis time can be fastidious. Moreover, the fusion and competition strategy increases the analysis complexity. We present in the next section this segmentation graph pruning strategy.



Fig. 13 Segmentation graph of written word "alors"

5.3 Segmentation graph pruning

The word path search step of the analysis (c.f.Fig. 8 in section 5) generates all the possible segmentation paths from the graph. From all the paths generated, the one minimising the edit distance with the prior knowledge is chosen as the prediction of the written word. We exploit the approximate implicit segmentation of the Seq2Seq model to prune the segmentation graph. The implicit segmentation is not directly exploitable to provide feedback, but can help optimise the analysis process. The objective is to have a nice trade-off between the analysis process performance and complexity. Fig. 13 illustrates the word paths search process for the written word "alors". For each node of the first level of the graph (highlighted in blue rectangles), all possible segmentation nodes paths are recursively constructed. Each node having at most four letter hypotheses with their analysis score, all segmentation paths (or word hypotheses) resulting from each segmentation node path are then generated. The Seq2Seq segmentation of the written word "alors" is framed in red in Fig. 14. Each rectangle represents the predicted letters as well as the points used by the attention mechanism to recognise it. This is used to prune the segmentation graph. First, a *deep matching score* (which is in fact an IoU score between the points in a graph segmentation node and the points in a deep segmentation node) is computed for each node of the graph relatively to the deep segmentation, to find the best corresponding deep predicted letter. The deep matching score is defined as follows:

 $DMScore(n_{graph}, n_{deep}) = \frac{\|points_{n\,graph} \cap points_{n\,deep}\|}{\|points_{n\,graph} \cup points_{n\,deep}\|}$ The best deep matching node for a graph segmentation node is defined as follows: $DeepMatch(n_{Graph}) = \max_{n_{Deep} \in Deep} DMScore(n_{Graph}, n_{Deep}).$ Given the computed deep matching scores, the new segmentation paths search process consists in **selecting recursively, at each level, only the nodes whose analysis hypotheses contain the matching deep node predicted letter**, formalised as follows:

 $SelectedNodes(level_i) = n_{Graph} \in level_i$, such as

 $DeepMatch(n_{\mathit{Graph}}) \in AnalysisHypotheses(n_{\mathit{Graph}}).$

Fig. 14 illustrates this pruning process for part of the segmentation graph. Dotted arrows represent the matching process at the first level. Nodes highlighted in red represent the discarded nodes, since their analysis hypotheses do not contain the predicted letter from the matched deep node. We can see that at the first level of the graph, only the relevant nodes have been selected. This is due



Fig. 14 Pruning process for part of the graph

to the fact that the implicit segmentation of the deep network was relatively consistent with the explicit segmentation. In the example in Fig. 14, without the pruning strategy, the number of processed paths is 301, and goes down to only 18 paths when the pruning is activated. In both cases, the correct word and segmentation are predicted. We will see more in detail its impact, as well as the performance of the analysis engine in the next section.

6 Experiments

6.1 Dataset

This work needs **data annotated at character level** to evaluate the system on **recognition** and **segmentation** aspects. To our knowledge, open datasets of children handwriting with character annotation for words do not exist. For our experiments, we use a private dataset, composed of French cursive words written by children. The data were collected in classrooms on pen-based tablets and were recorded as **multivariate time**

puis no surce. ni sur

Fig. 15 Examples of cursive words written by children. series. Each word is a sequence of points represented by their coordinates (x and y), their pressure and their time. Unfortunately, these children data are not publicly available due to RGPD laws¹. Fig. 15 illustrates examples of words in the database (the instruction is in orange). We can see that the handwriting is degraded because children are still learning writing, and naturally they do some mistakes. Another interesting aspect is the diversity of misspelling errors.

Our dataset is split into 6812 words written by more than 500 children for the training set and 1242 words written by more than 300 children for the test set. Train and test datasets come from different data acquisition campaigns (and different classroom). There are no children data present both in train and test set, this enables us to verify the ability of the system to generalise on unseen writing styles.

6.2 Deep learning model evaluation

For each experiment, λ of hybrid loss is set to 0.5 as suggested in [12]. We evaluate our deep learning model on the IAM-OnDB dataset [3] which is

¹https://ec.europa.eu/info/law/law-topic/dataprotection/data-protection-eu_fr

composed of adult handwritten English text. We train the model on a combination of train set and validation set with RMS prop optimiser during 200 epochs, then evaluate it on a test set. We set the learning rate at 0.001 and the batch size at 16. We evaluate the encoder and the decoder of our Seq2Seq model. The table 2 report the error rate on the test set. We can see that the encoder performs better than the decoder and outperforms the state of the art without the use of language model.

The deep learning model performs poorly with only children data. We use the model trained on IAM-OnDB then continue the training on the children handwriting.

Cross-validation with k folds equal to 10 is performed on **the training set** to evaluate the robustness of the system. The training set is split into 10 chunks. A fold is composed of a training part which represent 8 chunks, a validation part of 1 chunk and a test part of 1 chunk. Each fold results in a different splitting of the training set, thus all training set data are used for training and testing. For each fold, the validation set is used to choose the best model. A fold is evaluated on the test fold for the recognition task and

Table 2 Error rates on the IAM-OnDB test set incomparison with the best of state of the art. CER isCharacter Error Rate and WER is Word Error Rate

System	CER (%)	WER (%)
Without model language [9]	5.9	18.6
With model language [9]	4.0	10.6
Our Seq2Seq encoder	5.0	18.3
Our Seq2Seq decoder	5.5	20.2

Table 3Mean and standard deviation for therecognition and segmentation (IoU) evaluation onchildren handwriting. Recognition is evaluate on fold testset and whole test set. Encoder and decoder fromSeq2Seq are evaluated in %.

	Encoder	Decoder
Fold recognition rate	86.65 ± 1.17	86.32 ± 1.25
Test recognition rate	75.08 ± 1.16	69.20 ± 2.17
Segmentation rate (IoU)	51.14 ± 7.06	45.91 ± 3.19

the whole test set for the recognition and segmentation task. The recognition is evaluated with a recognition rate (100 - Word error rate) and the intersection over union to evaluate the segmentation (qualitative results are presented in section 4). The table 3 reports the results. We use the encoder prediction (label and segmentation) for the next experiments because its recognition rate are better on test fold. The recognition rate is better in fold test set because the data in the whole test set are from words written by unseen written styles. The Seq2Seq model has a greater recognition rate than the existing analysis engine (see more details on results in section.6.4) while the segmentation rate is too approximate to make a precise feedback to the children. Combining the Seq2Seq model with the existing analysis engine makes it possible to have a model both efficient in recognition and segmentation. The next section presents the results of the different combination strategies.

6.3 Segmentation evaluation

To study the segmentation from a qualitative viewpoint, Fig. 16 illustrates the analysis results

of the written word "gust". We can see that the raw deep segmentation (e) is approximate, compared to the explicit segmentation driven by the defined strategies (a, b, c, d). In this example, the phonetic strategy performed the best in terms of edition and analysis score, and therefore was chosen within the fusion and competition strategy. Correct segmentation and ground truth detection were performed.

We can observe the same results on the whole dataset, in terms of quality of segmentation. As the ground truth is annotated at the character level, we can therefore study how well the test set was segmented using the IoU metric. Table 4 illustrates the quality of segmentation for each strategy, from a quantitative viewpoint. Deep prediction and fusion/competition strategies are tested on the 10 models generated from the crossvalidation. Mean and standard deviation results are reported.



Fig. 16 Segmentation results for the word "gust", given the instruction "juste"

 Table 4
 Segmentation (IoU) performance of each strategy

Strategy	Segmentation rate (IoU)
Raw Seq2Seq	$51.14 \pm 7.06\%$
Childtying strategy	93.67%
Instruction strategy	88.66 %
Phonetic strategy	88.72%
Deep prediction strategy	$90.4\% \pm 0.54\%$
Fusion and competition	$92.82\% \pm 0.28\%$

Table 5 Recognition performance of each strategy

Strategy	Recognition rate
Childtying strategy	78.98%
Instruction strategy	64.09~%
Phonetic strategy	66.42%
Deep prediction strategy	$72.18\% \pm 0.73\%$
Fusion and competition	83.28% ±0.51%

As we have seen, the raw Seq2Seq segmentation rate is very approximate (51.14%). When we integrate the deep recognition results into the existing analysis engine, the segmentation performance improves with an IOU of 90.4% (better than instruction and phonetic strategies). This demonstrates the merits of combining explicit segmentation with the deep network recognition in the analysis process. Finally, the fusion and competition strategy (92.82%) comes a close second to the Childtyping strategy, which refers to the analysis being guided by the keyboard user input (93.67%). We can consider childtyping analysis performance as a sort of objective to reach for the system, without the aid of the user.

6.4 Recognition evaluation

Table 5 presents the recognition performance of each strategy, without the graph pruning, on the test set. We can see that in a dictation context, the instruction can't guide the analysis effectively, with a recognition rate of 64.09%. The phonetic analysis approach deals well with phonetically coherent misspellings, but fails to reach the ceiling of childtyping recognition performance (66.42%). Even if childtyping is a reliable approximation of the ground truth, the combination of degraded handwriting and in some cases, typing errors, explains the ceiling of 78.98%. The deep prediction strategy achieves better results than the phonetic strategy (72.18%). It is interesting to note that this strategy fails to achieve the recognition performance of the raw Seq2Seq, however this is explained by the explicit segmentation aspect of the analysis engine. While the implicit segmentation is quite approximate, the explicit segmentation driven by the deep prediction is significantly better (c.f. table 4). Finally, the fusion and competition strategy has better performances than the childtyping one (83.28%).

6.5 Impact of the pruning strategy

The deep learning model takes an average of 73 milliseconds per word to make a prediction. This computation is very fast, therefore, it is not included in the following time analysis. Table 6 presents the recognition and segmentation performance of the proposed strategies, as well as their average analysis time per word. In this table, we do not discuss the pruning with childtyping, instruction, or phonetic strategies, since they do Table 6 Impact of pruning strategy.

Strategy	Recognition rate	Segmentation rate (IoU)	Average time (s)
Deep prediction	$72.18\% \pm 0.73\%$	$90.4\% \pm 0.54\%$	1.34
Fusion competition	$83.28\% \pm 0.511\%$	$92.8\% \pm 0.28\%$	4.74
Deep prediction (pruning)	$70.87\% \pm 2.3\%$	$89.3\% \pm 1.12\%$	0.37
Fusion competition (pruning)	$79.87\% \pm 0.94\%$	$91.44\% \pm 0.98\%$	0.67

not exploit the Seq2Seq results, contrary to the other two strategies. As we have seen, the fusion and competition strategy provides the best recognition and segmentation results (barring the childtyping strategy for segmentation), however the analysis time (4.74s per word) is more than 3 times bigger than the deep prediction guidance strategy. This is due to the fact that there are more segmentation paths that are processed for this strategy. Integrating the pruning enables to decrease the analysis time of the fusion strategy to an acceptable 0.67s on average, while loosing about 2%of recognition performance (80.87%), which is still better than the childtyping strategy. The pruning results also in loosing about 1% of segmentation precision. This is due to the approximate nature of the implicit segmentation. The same goes for pruning with the deep prediction guidance strategy. We can therefore conclude that the pruning constitutes an acceptable trade-off between analysis time and performance.

6.6 Feedback typology

This section presents the pedagogical output of our system, providing visual feedbacks on the children mistakes. Since we are in an educational



Fig. 17 Segmentation graph of written word "alors"

context, we have to minimise the analysis system errors. Therefore, the degree of visual feedback precision and detail displayed to the child depends on the analysis confidence. When the analysis confidence is low, we generate more generic feedbacks, *i.e.* a warning on a zone of incertitude, or even no feedback at all. The feedback typology is illustrated in Fig. 17 and decomposed into three different levels: 1) *High confidence*: when the predicted word path corresponds to the prior knowledge strategy (e.g. the deep prediction) \implies precise feedback is given; 2) Medium confidence: when one letter distinguishes between the predicted word and the strategy \implies a warning is generated on an uncertain zone; 3) Reject: when the aforementioned conditions are not met \implies no feedback is given to the child. More details on feedback generation can be seen in [20].

Table 6.6 presents the feedback pertinence results on the fusion competition strategy with pruning. On one hand, the system has a high confidence feedback degree of 88.88% on the test set with an error rate of 15.2% on this type of feedback. On the other hand, the system has a low degree of medium and reject feedback (4.5 and 6.7% respectively). Putting high and medium confidence feedbacks altogether, we can see that the

Table 7 Feedback generation pertinence

Confidence	Ratio (avg)	Errors rate (avg)
High	$1103.9 \pm 11 (88.8\%)$	15.2%
Medium	$55.4 \pm 4.4 (4.5\%)$	0%
Reject	$82.7 \pm 9.48(6.7\%)$	0%
Total feedback	1146 (93.34%)	14.7%

system minimizes its error rate from 21.13% (*c.f.* table 6) to 14.7%, which is positive. However, since we are in an educational context, further work is needed to improve this feedback error ratio.

7 Conclusion

In this paper, we present an approach for the fine analysis, *i.e.* recognition and segmentation, of children handwritten words in a dictation context. This context introduces new challenges, since the handwriting is more degraded than adult handwriting, and the children are prone to misspelling mistakes, which makes the analysis task much harder than in a copying context. An explicit segmentation process is needed to provide precise feedback on the child's mistakes. This explicit segmentation needs to be driven by prior knowledge. We propose to combine an existing explicit segmentation based analysis engine with a Seq2Seq architecture to generate relevant prior knowledge and adapt the system to the dictation context. Using the deep predicted character sequence as prior knowledge compensates for the fact that the dictated instruction cannot drive the analysis, as it has done for the copying context. We then propose to combine multiple strategies, the instruction, phonetically similar pseudo-words, and the deep prediction, in order to further improve analysis performances. Another contribution of this work is to use the implicit segmentation of the Seq2Seq to prune the analysis engine segmentation graph, which resulted in optimising analysis complexity and time, while retaining good analysis performances, in fact outperforming the childtyping strategy, which constituted a "high ceiling baseline" for our task in terms of recognition performances. Our future works consist in further experimenting the system in pilot French schools. Another objective is to improve the Seq2Seq performances, in terms of recognition and segmentation, which will consequently improve the explicit segmentation based analysis engine. We could rely on synthetic data to further improve the network performances. Finally, we could explore the extension of this approach to languages other than French.

References

- D. Simonnet, N. Girard, E. Anquetil, M. Renault, and S. Thomas. Evaluation of Children Cursive Handwritten Words for e-Education. *Pattern Recognition Letters*, 121:133–139, 2019.
- [2] U.-V Marti and H Bunke. A full english sentence database for off-line handwriting

recognition. In Fifth International Conference on Document Analysis and Recognition, ICDAR 1999, 20-22 September, 1999, Bangalore, India, pages 705–708. IEEE Computer Society.

- [3] M. Liwicki and H. Bunke. Iam-ondb an online english sentence database acquired from handwritten text on a whiteboard. In *Eighth International Conference on Document Analysis and Recognition (ICDAR 2005), 29 August - 1 September 2005, Seoul, Korea,* pages 956–961. IEEE Computer Society.
- [4] C. C. Tappert, C. Y. Suen, and T. Wakahara. The state of the art in online handwriting recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 12(8):787–808, 1990.
- [5] R. Plamondon and S.N. Srihari. Online and off-line handwriting recognition: a comprehensive survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(1):63–84, 2000.
- [6] A. Graves, S. Fernández, F. J. Gomez, and J. Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In Machine Learning, Proceedings of the Twenty-Third International Conference (ICML 2006), Pittsburgh, Pennsylvania, USA, June 25-29, 2006, volume 148 of ACM International Conference Proceeding Series, pages 369–376. ACM.

- [7] S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural Comput., 9(8):1735–1780, 1997.
- [8] A. Graves, M. Liwicki, S. Fernández, R. Bertolami, H. Bunke, and J. Schmidhuber. A novel connectionist system for unconstrained handwriting recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(5):855–868, 2009.
- [9] V. Carbune, P. Gonnet, T. Deselaers, H. A. Rowley, A. N. Daryin, M. Calvo, L.-L. Wang, D. Keysers, S. Feuz, and P. Gervais. Fast multi-language lstm-based online handwriting recognition. *Int. J. Document Anal. Recognit.*, 23(2):89–102, 2020.
- [10] B. Shi, X. Bai, and C. Yao. An end-toend trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(11):2298–2304, 2017.
- [11] J. Puigcerver. Are multidimensional recurrent layers really necessary for handwritten text recognition? In 14th IAPR International Conference on Document Analysis and Recognition, ICDAR 2017, Kyoto, Japan, November 9-15, 2017, pages 67–72. IEEE.
- [12] J. Michael, R. Labahn, T. Grüning, and J. Zöllner. Evaluating sequence-to-sequence models for handwritten text recognition. In 2019 International Conference on Document

Analysis and Recognition, ICDAR 2019, Sydney, Australia, September 20-25, 2019, pages 1286–1293. IEEE.

- [13] L. Kang, P. Riba, M. Rusiñol, A. Fornés, and M. Villegas. Pay attention to what you read: Non-recurrent handwritten text-line recognition. *CoRR*, abs/2005.13044, 2020.
- [14] K. Barrere, Y. Soullard, A. Lemaitre, and B. Coüasnon. Transformers for Historical Handwritten Text Recognition. In *Doctoral Consortium - ICDAR 2021*, Lausanne, Switzerland, September. Nibal Nayef and Jean-Christophe Burie.
- [15] H. Liu, S. Jin, and C. Zhang. Connectionist temporal classification with maximum entropy regularization. In Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pages 839–849.
- [16] L. Gao, H. Zhang, and C.-L. Liu. Handwritten text recognition with convolutional prototype network and most aligned frame based CTC training. In 16th International Conference on Document Analysis and Recognition, ICDAR 2021, Lausanne, Switzerland, September 5-10, 2021, Proceedings, Part I, volume 12821 of Lecture Notes in Computer Science, pages 205–220. Springer.
- [17] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied

to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998.

- [18] E. Anquetil and G. Lorette. On-line Handwriting Character Recognition System Based on Hierarchical Qualitative Fuzzy Modelling. In *Progress in Handwriting Recognition*, pages 109–116, 1997.
- [19] D. Simonnet, E. Anquetil, and M. Bouillon. Multi-Criteria Handwriting Quality Analysis with Online Fuzzy Models. *Pattern Recogni*tion, 69:310–324, 2017.
- [20] O. Krichen, S. Corbillé, E. Anquetil, N. Girard, and P. Nerdeux. Online analysis of children handwritten words in dictation context. In 14th International Workshop on Graphics Recognition, Lausanne, Switzerland, September 2021.
- [21] J. R. Novak, N. Minematsu, K. Hirose, C. Hori, H. Kashioka, and P. R. Dixon. Improving wfst-based G2P conversion with alignment constraints and RNNLM n-best rescoring. In *INTERSPEECH 2012, 13th* Annual Conference of the International Speech Communication Association, Portland, Oregon, USA, September 9-13, 2012, pages 2526–2529. ISCA.
- [22] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015,

Conference Track Proceedings.

- [23] T. Bluche and R. O. Messina. Gated convolutional recurrent neural networks for multilingual handwriting recognition. In 14th IAPR International Conference on Document Analysis and Recognition, ICDAR 2017, Kyoto, Japan, November 9-15, 2017, pages 646–651. IEEE.
- [24] A. Almaksour and É. Anquetil. Improving premise structure in evolving takagi-sugeno neuro-fuzzy classifiers. *Evol. Syst.*, 2(1):25– 33, 2011.
- [25] H. Schwenk and J.-L. Gauvain. Improved rover using language model information. 11 2000.
- [26] F. Damerau. A technique for computer detection and correction of spelling errors. *Commun. ACM*, 7(3):171–176, 1964.

Declarations

- Funding: Partial funding was received from the P2IA project (French government) for this study.
- Conflict of interest/Competing interests: The authors have no relevant financial or nonfinancial interests to disclose
- Data availability: The datasets generated during and/or analysed during the current study are not publicly available due to privacy laws (RGDP) in France.

About the authors



Omar Krichen received his PhD degree from INSA Rennes in December 2020. He is currently holding a postDoc position at the Institut National des Sciences Appliquées (INSA) Rennes, within the IntuiDoc team, at the IRISA laboratory. His research interests are online handriwiting recognition, structured documents interpretation, and intelligent tutoring systems.



Simon Corbillé is a French Ph.D student in Computer Science from the University Rennes 1 in Rennes, France. He holds a engineering degree in Compute Science from Polytech, in Tours, France. He is in IntuiDoc research team at the IRISA (Institut de Recherche en Informatique et Systèmes Aléatoires). His works concern handwriting recognition for an education purpose.



Eric Anquetil received his engineering degree from INSA in 1993 and his Ph.D. degree in Computer Science from the University of Rennes in 1997. He received his Accreditation to Supervise Research (HDR) in 2008. Currently he is a full professor at INSA in Rennes. He is the head of IntuiDoc research team at the IRISA laboratory. His research interests include: Pattern Recognition, Fuzzy logic, Evolving Classifiers; Handwriting, Gesture, Symbol and Drawing Recognition; Digital Learning. He directed several scientific projects in collaboration with companies. He is member of the IAPR association and of the GRCE French association on handwriting recognition.



Nathalie Girard is an associate professor in Computer Science department (ISTIC), at "Université de Rennes 1". She works at IRISA in the IntuiDoc research team, after having held lecturer and postdoctoral positions at La Rochelle University, L3I laboratory and University of Tours, RFAI team, LIFAT laboratory. She obtained her Ph.D. in computer science from La Rochelle University in 2013. Her research interests include data classification, image recognition, incremental learning, handwriting recognition, and user interaction, and user modelling. She is member the GRCE French association on handwriting recognition.



Elisa Fromont is professor at "Université de Rennes 1" (UR1) since 2017 and a junior member of the Institut Universitaire de France (IUF) (2019-2024). She works at IRISA in the Inria LACODAM team. From 2008 until 2017, she was associate professor at Université Jean Monnet in Saint-Etienne (UJM), France and worked at the Hubert Curien research institute in the Data Intelligence team. She received her Research Habilitation (HDR) in 2015 from UJM. From 2006 until 2008, she was a postdoctoral researcher in the Machine Learning group of the KUL (Belgium). She received her PhD in 2005 from UR1. She is particularly interested in solving real world problems with data mining techniques especially when the data are heterogeneous, multimodal, temporal, imbalanced, subject to concept drift, ... and when it is important to explain the model's decision to an end-user.



Pauline Nerdeux is a French computer engineer at the Institut National des Sciences Appliquées (INSA) in Rennes, France; she holds an engineering degree from Ecole Centrale Marseille with a specialization in Image Processing. She is a member of the Intuidoc research team at the Institut de Recherche en Informatique et Systèmes Aléatoires (IRISA) laboratory.