# HTR-VT: Handwritten Text Recognition with Vision Transformer

Yuting Li[†a], Dexiong Chen[b], Tinglong Tang[a], Xi Shen[c]

[a]*China Three Gorges University, China*
[b]*Max Planck Institute of Biochemistry, Germany*
[c]*Intellindust, China*

## Abstract

We explore the application of Vision Transformer (ViT) for handwritten text recognition. The limited availability of labeled data in this domain poses challenges for achieving high performance solely relying on ViT. Previous transformer-based models required external data or extensive pre-training on large datasets to excel. To address this limitation, we introduce a data-efficient ViT method that uses only the encoder of the standard transformer. We find that incorporating a Convolutional Neural Network (CNN) for feature extraction instead of the original patch embedding and employ Sharpness-Aware Minimization (SAM) optimizer to ensure that the model can converge towards flatter minima and yield notable enhancements. Furthermore, our introduction of the span mask technique, which masks interconnected features in the feature map, acts as an effective regularizer. Empirically, our approach competes favorably with traditional CNN-based models on small datasets like IAM and READ2016. Additionally, it establishes a new benchmark on the LAM dataset, currently the largest dataset with 19,830 training text lines. The code is publicly available at: https://github.com/YutingLi0606/HTR-VT.

*Keywords:*
Handwritten Text Recognition, Vision Transformer, Mask Strategy, Sharpness-Aware Minimization, Data-efficient

## 1. Introduction

The task of handwritten text recognition aims at recognizing the text in an image that has been scanned from a document. The standard approach [36, 38, 22, 20] typically involves two steps: first, using a detector to identify the lines of text, and then predicting the sequence of characters that make up each line. This paper focuses on the latter task, which aims to accurately predict the text in a given line image. As described in LAM [34]: "This level of annotation granularity has been chosen as it is a good trade-off between word-level and paragraph-level in terms of required time, cost, and amount of supervision and because it is common in HTR research." This rationale for creating the dataset aligns with our initial decision to focus our research on line-level recognition. The importance of line-level recognition is significant. At the same time, recognizing handwritten text lines is a difficult task due to variations in writing styles between individuals and the presence of cluttered backgrounds (examples are provided in Figure 3 and in the supplementary material). Previous approaches mainly relied on Convolutional Neural Networks (CNNs) [2, 4, 20, 50] or recurrent models [36, 39, 11] to address this challenge .

Recent success of Vision Transformer (ViT) [35] in computer vision tasks has motivated researchers to explore its potential in handwritten text recognition. However, ViT does not introduce any strong inductive bias in its model design and is recognized for its dependency on vast quantities of annotated training data to deliver good performance. Considering the limited number of annotated samples available for handwritten text recognition (as shown in Table 1), earlier transformer-based methods, utilizing the standard transformer architecture (both encoder and decoder), relied on large-scale real-world or synthetic data for pre-training [18, 48] to achieve satisfactory performance.

In this paper, we introduce a simple and data-efficient ViT-based model that solely employs the encoder component of the standard transformer for handwritten text recognition. Our objective is to propose a novel ViT-like model to perform well on this task while making minimal modifications to the standard ViT architecture. Our preliminary findings indicate that ViT can deliver satisfactory results, particularly on the LAM dataset [34], which is the most extensive dataset containing 19,830 training samples. We use this dataset as the basis of our experimental design, which establishes a benchmark for assessing and contrasting our proposed model. Instead of using a patch embedding to generate input tokens, we demonstrate through experimentation that using a widely-used ResNet-18 [23] to extract intermediate visual feature representations as input tokens is much more conducive to stable training and significantly better performance. Additionally, we show that employing Sharpness-Aware Minimization (SAM) [55] as optimizer enforces the model to converge towards flatter minima and randomly replacing span tokens with learnable tokens can alleviate overfitting and achieve consistent improvement across various dataset scales.

Despite its simplicity, our approach achieves promising performance on standard benchmarks. On the largest dataset LAM [34] (containing 19,830 training samples), our approach outperforms

---

[†]Corresponding Author.

both CNN-based and transformer-based approaches by a clear margin. On small-scale datasets such as IAM [33] (containing 6,428 training samples) and READ2016 [11] (containing 8,349 training samples), we achieve better performance than other transformer-based approaches and competitive performance compared with CNN-based approaches.

The main contributions of this paper are summarized as the following:

- We propose a simple and data-efficient approach for handwritten text recognition, with minimal modifications on the ViT.

- We empirically show that without pre-training or any additional data, our ViT-like model can achieve state-of-the-art performance on handwritten text recognition.

## 2. Related Work

*Traditional approaches for Handwritten Text Recognition.* Network architectures for handwritten text recognition today typically use a combination of convolutional layers and recurrent layers. A number of convolutional layers are stacked and placed at the start of the network to extract local features from text-line images, followed by recurrent layers, specifically Bi-directional Long Short-Term Memory (BLSTM) [49] layers. These recurrent layers process the features sequentially to output character probabilities based on contextual dependencies. Such an architecture results in a Convolutional Recurrent Neural Network (CRNN) [36, 22, 38, 37]. The models are typically trained using the Connectionist Temporal Classification (CTC) loss [14], which allows for dealing with label sequences of shorter length than predicted sequences, without knowledge of character segmentation. Encoder-Decoder-based architectures have also been explored for handwritten text recognition [17, 39, 26]. In [17], the CTC loss is replaced with the cross-entropy loss, and the sequence alignment is achieved via an attention-based encoder-decoder architecture. A special end-of-line token is introduced to stop the recurrent process. While these models can obtain lower test error rates, some of them often require complex pre/post-processing steps and suffer from the lack of computation parallelization inherently, which affects both training and inference time. Recently, Fully Convolutional Networks (FCNs) [4, 20, 2] have been proposed as an alternative to traditional CRNNs. FCNs simulate the dependency modeling provided by LSTM by combining them with GateBlocks layers [50], which implement a selection mechanism similar to that of LSTM cells. Each gate in GateBlocks is made up of Depth-wise Separable Convolutions [25], which reduce the number of parameters and speed up the training process. OrigamiNet [4] focuses on learning to unfold the input paragraph image into a single text line. This transformation network enables using the standard CTC loss [14] and processing the image in a single step. In contrast, Coquenet et al. [2] proposed models that incorporate a vertical attention mechanism to recurrently generate line features and perform an implicit line segmentation. While FCNs have obtained state-of-the-art results in recent years, they may still struggle with long-range contextual dependencies.

| Dataset | Training | Validation | Test | Language | Charset |
|---|---|---|---|---|---|
| IAM [33] | 6,482 | 976 | 2,915 | English | 79 |
| READ2016 [11] | 8,349 | 1,040 | 1,138 | German | 89 |
| LAM [34] | 19,830 | 2,470 | 3,523 | Italian | 89 |

Table 1: **Datasets for handwritten text recognition.** Number of training, validation, and testing samples in IAM [33], READ2016 [11] and LAM [34] are presented in the table. We also include the number of characters in their alphabet.

*Transformer-based models for Handwritten Text Recognition.* Transformer-based architectures have not been widely explored in handwritten text recognition, but some recent approaches have used Transformers in place of RNNs. These models often require pre-training on large real or synthetic datasets to achieve comparable performance to mainstream models.

TrOCR [18] is a recent approach to handwritten text recognition that integrates two powerful pre-trained models respectively from computer vision and NLP, BEiT [44] and RoBERTa [45]. BEiT is a vision transformer that functions as an encoder and is pre-trained on ImageNet-1K, a dataset of 1.2 million images, while RoBERTa serves as a decoder that generates texts. To pre-train the TrOCR model, Li et al. [18] synthesize a large-scale dataset consisting of both printed and synthetically generated handwritten text lines in English, totaling approximately 687 million and 18 million in the first stage. In this stage, the dataset is not public. In the second stage, they built two relatively small datasets corresponding to printed and handwritten downstream tasks, containing millions of textline images each. Finally, the model is fine-tuned on real-world data, such as the IAM dataset [33]. Kang et al. [48] use Transformer models with multi-head self-attention layers at the textual and visual stages and trains with a synthetic dataset of 138,000 lines. Another recent approach, Text-DIAE [56], employs a transformer-based architecture that incorporates three pretext tasks as learning objectives to be optimized during pretraining without the usage of labeled data. Some methods [57, 1] explored document-level recognition and also applied transformer architectures. While transformer-based models have shown promising results in line-level handwritten text recognition, they still require large-scale real-world or synthetic data for pre-training.

*Data-efficient Transformer for Handwritten Text Recognition.* The DeiT [30] is the first work to demonstrate that Transformers can be learned on mid-sized datasets (i.e., ImageNet-1k [27])) in relatively shorter training episodes. Besides using augmentation and regularization procedures, the main contribution of DeiT [30] is a novel distillation that relies on a distillation token. Liu et al. [32] propose a dense relative localization loss to improve ViTs' data efficiency. DropKey [51] is a recent data-efficient methodology to effectively improve the dropout technique in ViT by moving dropout operations ahead of attention matrix calculation and setting the Key as the dropout unit, yielding a dropout-before-softmax scheme.

## 3. Method

In this section, we present our approach to handwritten text recognition. Given an input handwritten text line $\mathbf{I} \in \mathbb{R}^{W \times H}$, where $W$ and $H$ are the width and height of the image, our approach encodes the image into a set of spatially-aware features $\{\mathbf{x}_i\}_{i \in [1,2,...,L]}$ using a CNN extractor. The number of features $L = \frac{WH}{S_w S_h}$, determined by the down-sampling ratio of the width and height of the image, denoted as $S_w$ and $S_h$, respectively. We then use a transformer encoder to take these features as input tokens and output character predictions. The entire model is optimized using the Connectionist Temporal Classification [14] (CTC) loss. Our method is summarized in Figure 1.

In Section 3.1, we revisit the architecture of the Vision Transformer (ViT)[35]. In Section 3.2, we describe our data-efficient ViT approach for handwritten text recognition, which involves a CNN feature extractor, Sharpness-Aware Minimization (SAM) [55] and a new masking strategy: span mask strategy. We provide implementation details in Section 3.2.

### 3.1. Preliminary: Vision Transformer (ViT)

Vision Transformer (ViT) [35] decomposes each image into a sequence of tokens with a fixed length, where the tokens represent non-overlapping image patches. Similar to BERT [40], ViT adds an additional class token $\mathbf{x}_{cls}$ to the sequence, which represents the global information of the image. To retain positional information, position embeddings are explicitly added into each patch including the class token. Note that our model removes the additional class token and uses sinusoidal position embeddings by [13] to the encoder's inputs, as used in MAE [12].

Subsequently, all tokens undergo processing via stacked transformer encoders [13], A transformer encoder comprises N blocks, with each block featuring a multi-head self-attention (MSA) layer followed by a feed-forward network (FFN). The FFN, which includes a simple two-layer MLP, is augmented by the GELU activation function [41] after the first linear layer. Furthermore, layer normalization (LN) [42] is applied before every block, and residual shortcuts [23] are used after every block. The processing of the n-th block can be expressed as:

$$\mathbf{y}^n = \mathbf{x}^{n-1} + \mathrm{MSA}\left(\mathrm{LN}\left(\mathbf{x}^{n-1}\right)\right)$$
$$\mathbf{x}^n = \mathbf{y}^n + \mathrm{FFN}\left(\mathrm{LN}\left(\mathbf{y}^n\right)\right) \tag{1}$$

where $\mathbf{x}^{n-1} \in \mathbb{R}^{L \times C}$ is the input of the *n*-th block, $N$ and $C$ denote the number of tokens and the dimension of the embedding, respectively.

### 3.2. ViT for handwritten text recognition

We present a ViT-based model designed for handwritten text recognition with minimal adjustments to the standard ViT [35]. Our proposed network architecture is depicted in Figure 1. ViT alone is not stable for handwritten text recognition (see Section 4.3). Therefore, we suggest three modifications: *i)* a CNN feature extractor to obtain features for each input token, enabling powerful feature extraction, *ii)* a span feature masking strategy
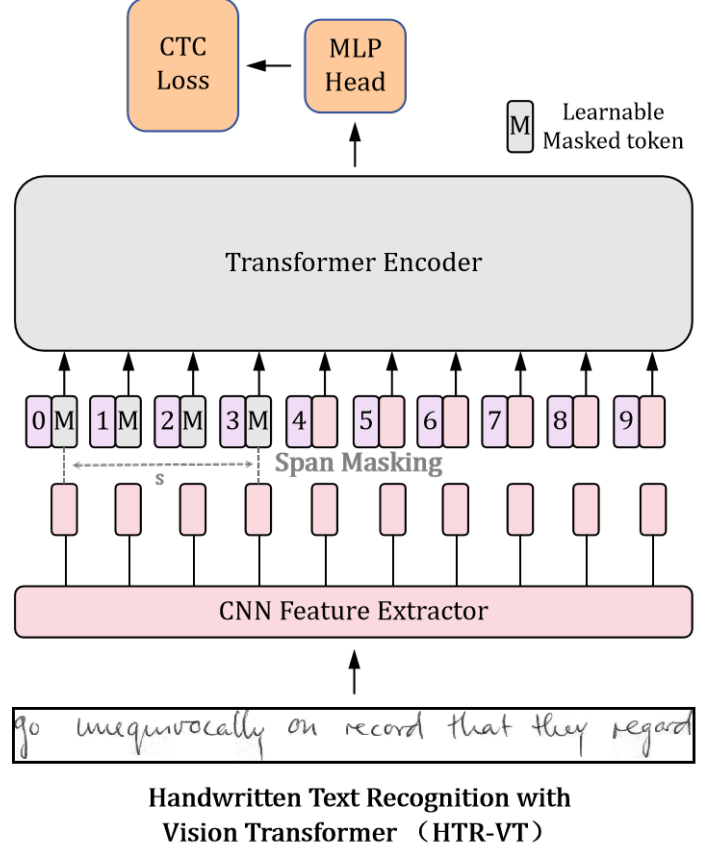


Figure 1: **Architecture overview.** Our approach encodes a text-line image into features using a CNN feature extractor. The transformer encoder takes these features as input tokens output character predictions. During the training, the span input tokens are replaced by learnable mask tokens. The entire model is optimized using CTC [14] loss.

to replace masking tokens with learnable tokens, effectively alleviating the impact of overfitting, and *iii)* employ Sharpness-Aware Minimization (SAM) optimizer to ensure that the model can converge towards flatter minima.

*CNN feature extractor.* To make our pipeline simple, we adopt the widely-used ResNet-18 [23] as our CNN feature extractor, with minor adjustments made to accommodate line-level handwritten text images. Specifically, we remove the final residual block and adjust the stride to produce features with enough information for character recognition while maintaining the two-dimensional nature of the task. More details about the modification as well as experiments of additional CNN feature extractors are provided in the supplementary material.

*Span feature masking strategy.* Our work draws inspiration from BERT [40], SpanBERT [53] and MASS [52], which leverage the prediction of randomly masked words or tokens to learn expressive language representations. We have adapted this methodology to our specific task and observed the benefits of employing random feature masking. Furthermore, our intuition suggests that the feature extractor can capture a board receptive field. To enhance the model's comprehension of contextual information

encompassing neighboring ink pixels, we propose expanding the masking range.

Precisely, the feature map after the CNN feature extractor is flattened to a sequence of tokens with dimensions $L \times C$, where $L$ represents the sequence length and $C$ represents the feature dimension. We randomly mask the span of tokens with a maximum span length $s$ (i.e., the number of interconnected tokens). In total $\tau L$ tokens are masked and replaced with a learnable token, where $\tau$ is a hyperparameter defining the mask ratio. More details of the span mask strategy are provided in the supplementary material.

*Sharpness-Aware Minimization (SAM).* Sharpness-Aware Minimization (SAM), proposed by Foret et al. [55], is an optimization method that enhances the generalization of deep neural networks (DNNs). It aims to find model parameters that reside in flat minima, ensuring a uniformly low loss across the model. Given our objective function $\mathcal{L}_{\text{CTC}}$ and the parameters of the DNN $\theta$, the SAM optimizer is designed to find $\theta$ such that:

$$\min_{\theta} \max_{\|\epsilon\|_2 \leq \rho} \mathcal{L}_{\text{CTC}}(\theta + \epsilon), \quad (2)$$

where $\epsilon$ represents a perturbation vector, and $\rho$ is the size of the neighborhood within which the algorithm minimizes the sharpness of the loss function. The SAM algorithm functions by alternately identifying the worst-case perturbation $\epsilon$ that maximizes the loss within an $\ell_2$-norm ball of radius $\rho$, and then updating the DNN parameters $\theta$ to minimize this perturbed loss.

*Implementation details.* We employ a ViT encoder with 4 layers. Each layer is with a dimension of 768 and 6 heads. The hidden dimension of MLP in the feed-forward network (FFN) is 3,072. Larger ViT models do not bring obvious gain. For our span mask strategy, we set the mask ratio to 0.4 and the span length to 8 in all datasets. An ablation study of the mask ratio and span length is provided in Section 4.3. For all experiments, we use a batch size of 128 and optimize all our models with the AdamW [43] optimizer for 100,000 iterations with a weight decay of 0.5. We perform a warm-up-cosine learning rate schedule with the max learning rate equal to 1e-3 and use 1,000 iterations for warm up. Trainings are performed on a single GPU RTX 4090 (24Gb) and in the following experiments, models are trained for almost 16 hours. Similar to OrigamiNet [4], we use the exponential moving average (EMA) method with a decay rate of 0.9999. For data augmentation, we fix the input image resolution to 512 x 64 and use random transformation, erosion, dilation, color jitter, and elastic distortion. We set the probability of using each data augmentation to 0.5, and they can be combined with each other.

## 4. Experiments

In this section, we evaluate the performance of our model for line-level recognition. Our experimental results demonstrate that our model achieves state-of-the-art results on the LAM [34] and IAM [33] datasets. Moreover, our model competes well with other state-of-the-art models on the READ2016 [11] datasets. It is worth noting that our model achieves good performance without any pre-training or synthetic data and without relying on any pre/post-processing steps.

To further analyze the performance of our model, we conduct an ablation study by modifying the standard ViT [35] architecture. Specifically, we investigate the impact of different mask strategies and hyperparameters on the READ2016 dataset and examine how the SAM optimizer [55] affects our model's performance.

### 4.1. Dataset and evaluation metrics

We evaluated our model's performance on three commonly used datasets for handwritten text recognition: LAM [34], READ2016 [11], and IAM [33]. Among these datasets, READ2016 and IAM are widely recognized as benchmarks for handwritten text recognition, while LAM is currently the largest available line-level handwritten text recognition dataset. The information about the datasets is provided in Table 1. Note that we report the performance on the test set with the model achieving the best performance on the validation sets.

*LAM [34].* The Ludovico Antonio Muratori (LAM) dataset is a massive handwritten text recognition dataset of Italian ancient manuscripts, which was edited by a single author over a span of 60 years. It consists of a total of 25,823 lines and has a lexicon of over 23,000 unique words. The dataset is split into 19,830 lines for training, 2,470 lines for validation, and 3,523 lines for testing, with a charset size of 89. The dataset was annotated at the line level, with each line's bounding box and diplomatic transcription provided. During the transcription process, stroke-out text, illegible words due to stains and scratches, and special symbols not representable in Unicode were replaced with the # symbol. This is currently the largest line-level handwritten text recognition dataset available and could be an ideal choice for demonstrating the potential of our model.

*READ2016 [11].* READ2016 was proposed in the ICFHR 2016 competition on handwritten text recognition. It comprises a subset of the Ratsprotokolle collection used in the READ project, with color images representing Early Modern German handwriting. The dataset provides segmentation at the page, paragraph, and line levels. For line-level tasks, the dataset has a total of 8349 training images, 1040 validation images, and 1138 test images, with a character set size of 89.

*IAM [33].* IAM is a well-known offline handwriting benchmark dataset for modern English. It comprises 1,539 scanned text pages of English texts extracted from the LOB corpus, which were handwritten by 657 different writers. The training set of IAM has 747 documents (6,482 lines), the validation set has 116 documents (976 lines), and the test set has 336 documents (2,915 lines). The IAM dataset consists of grayscale images of English handwriting with a resolution of 300 dpi. In this work, we utilized the line level with the commonly used split, as described in Table 1.

4

| Method | Test CER | Test WER | Param. |
|---|---|---|---|
| CNN + BLSTM★ [36] | 5.8 | 18.4 | 9.3M |
| GFCN★ [20, 34] | 5.2 | 18.5 | 1.4M |
| CRNN★ [22, 34] | 3.8 | 12.9 | 18.2M |
| OrigamiNet-12★ [4, 34] | 3.1 | 11.2 | 39.0M |
| OrigamiNet-18★ [4, 34] | 3.1 | 11.1 | 77.1M |
| OrigamiNet-24★ [4, 34] | 3.0 | 11.0 | 115.3M |
| **Transformer-based models** | | | |
| ViT [35] | 6.1 | 19.1 | 37M |
| ViT + DropKey [35, 51] | 5.7 | 16.5 | 37M |
| DeiT [30] | 5.9 | 18.7 | 6M |
| Transformer§★ [48, 34] | 10.2 | 22.0 | 54.7M |
| TrOCR§★ [18, 34] | 3.6 | 11.6 | 385.0M |
| **HTR-VT** | **2.8** | **7.4** | 53.5M |

§ reports results using extra training data.
★ indicates re-implementations by LAM [34].

Table 2: **Comparison with state-of-the-art approaches on LAM [34] dataset (19,830 training samples).** We outperform all the competitive approaches with a clear margin. The improvement is more important for the transformer-based approaches.

| Method | Test CER | Test WER | Param. |
|---|---|---|---|
| CNN + RNN [11] | 5.1 | 21.1 | - |
| CNN + BLSTM [17] | 4.7 | - | - |
| FCN [24] | 4.6 | 21.1 | 19.2M |
| VAN [2] | 4.1 | 16.3 | 2.7M |
| **Transformer-based models** | | | |
| ViT [35] | 8.5 | 29.6 | 37M |
| ViT + DropKey [35, 51] | 8.1 | 26.4 | 37M |
| DeiT [30] | 8.4 | 28.7 | 6M |
| DAN [1] | 4.1 | 17.6 | 7.6M |
| **HTR-VT** | **3.9** | **16.5** | 53.5M |

Table 3: **Comparison with state-of-the-art approaches on READ2016 [11] dataset (8,349 training samples).** We achieve comparable performance.

| Method | Test CER | Test WER | Param. |
|---|---|---|---|
| GFCN [20] | 8.0 | 28.6 | 1.4M |
| GFCN★ [20, 34] | 8.0 | 28.6 | 1.4M |
| CRNN★ [22, 34] | 7.8 | 27.8 | 18.2M |
| CNN + BLSTM [36] | 8.3 | 24.9 | 9.3M |
| CNN + BLSTM★ [36, 34] | 7.7 | 26.3 | 9.3M |
| OrigamiNet-12 [4] | 5.3 | - | 39.0M |
| OrigamiNet-12★ [4, 34] | 6.0 | 22.3 | 39.0M |
| OrigamiNet-18 [4] | 4.8 | - | 77.1M |
| OrigamiNet-18★ [4, 34] | 6.6 | 24.2 | 77.1M |
| OrigamiNet-24 [4] | 4.8 | - | 115.3M |
| OrigamiNet-24★ [4, 34] | 6.5 | 23.9 | 115.3M |
| VAN [2] | **5.0** | **16.3** | 2.7M |
| **Transformer-based models** | | | |
| ViT [35] | 32.4 | 68.5 | 37.0M |
| ViT + DropKey [35, 51] | 34.2 | 70.1 | 37.0M |
| DeiT [30] | 32.0 | 68.4 | 6.0M |
| Transformer§ [48] | 4.7 | 15.5 | 54.7M |
| Transformer♣ [48, 46] | 7.6 | 24.5 | 54.7M |
| TrOCR§ [18] | 3.4 | - | 385.0M |
| TrOCR★ [18, 34] | 7.3 | 37.5 | 385.0M |
| **HTR-VT** | **4.7** | **14.9** | 53.5M |

§ reports results using extra training data.
★ and ♣ indicate re-implementations by LAM [34] and by [46].

Table 4: **Comparison with state-of-the-art approaches on the test set of IAM [33] dataset (6,482 training samples).** Our approach exceeded the previous state-of-the-art model.

## 4.2. Comparison with state-of-the-art approaches

*Evaluation metrics.* We use Character Error Rate (CER) and Word Error Rate (WER) as performance measures. CER is calculated as the Levenshtein distance between two strings, which is the sum of character substitutions ($SUB_c$), insertions ($INS_c$), and deletions ($DEL_c$) required to transform one string into the other, divided by the total number of characters in the ground truth ($GT_c$). Formally, CER is given by:

$$CER = \frac{SUB_c + INS_c + DEL_c}{GT_c}. \quad (3)$$

Similarly, WER is calculated as the sum of word substitutions ($SUB_w$), insertions ($INS_w$), and deletions ($DEL_w$) needed to transform one string into another, divided by the total number of words in the ground truth ($GT_w$). Mathematically, WER is expressed as:

$$WER = \frac{SUB_w + INS_w + DEL_w}{GT_w} \quad (4)$$

We conducted a comparative study of current state-of-the-art methods on the LAM [34], READ2016 [11], and IAM [33] datasets respectively. Our approach surpassed previous state-of-the-art models on the LAM [34] and IAM [33] datasets and achieved comparable performance on the READ2016 [11] dataset. The results presented in Tables 2, 3 and 4 were achieved without the utilization of any external language models, such as

n-grams or similar techniques. Specifically, on the LAM [34] dataset, our method achieved a CER of 2.8 and a WER of 7.4, outperforming all models tested on this dataset. On the IAM [33] dataset, our approach exceeded the previous state-of-the-art model, VAN [2], with a CER improvement of 0.3 and a WER improvement of 1.4. On the READ2016 [11] dataset, our method reached a CER of 3.9, surpassing the state-of-the-art method VAN [2] and DAN [1] by 0.2, and closely matching its WER.

Furthermore, when compared to all transformer-based methods, our approach consistently led the field, except on the IAM dataset [33] where TrOCR [18] achieved a CER of 3.4. However, it is noteworthy that TrOCR [18] uses pre-trained CV and NLP models and a large-scale synthetic dataset, which is not publicly available, to pre-train their model. Transformer [48] also relies on a large amount of synthetic data for training. Despite this, our method still outperforms it. In addition, we also conduct a fair comparison to two recent works on data-efficient transformers: DeiT [30] and DropKey [51]. We achieve clearly better performance than them on all three datasets. Training details of DeiT [30] and DropKey [51] are provided in the supplementary material. These results demonstrate the data-efficiency of our proposed model.

In summary, our research presents a competitive handwritten text recognition model that stands out against state-of-the-art methods, particularly on the LAM [34] and IAM [33] datasets, and competes well on the READ2016 [11] dataset without resorting to any external language models, pre-training or synthetic data commonly used in the field.

For many years, the CNN + BLSTM paradigm has been the dominant approach in handwritten text recognition. However, our proposed method represents a significant shift in this trend, markedly enhancing the performance of transformer-based models. This breakthrough has the potential to steer the entire field of handwritten text recognition toward new and exciting directions.

| Methods | LAM [34] | | IAM [33] | | READ2016 [11] | |
|---|---|---|---|---|---|---|
| | Val CER | Val WER | Val CER | Val WER | Val CER | Val WER |
| ViT* | 5.7 | 16.7 | 26.6 | 57.1 | 9.4 | 35.2 |
| Ours w/o. CNN extractor | 5.5 | 15.7 | 20.7 | 53.5 | 8.9 | 33.7 |
| Ours w/o. SAM | 2.7 | 7.4 | 3.4 | 11.2 | 4.8 | 20.1 |
| Ours w/o. Span Mask | 2.9 | 7.8 | 3.7 | 12.1 | 5.1 | 21.9 |
| **Ours** | **2.6** | **6.9** | **3.3** | **10.8** | **4.5** | **19.4** |

\* ViT is equivalent to our approach without CNN extractor nor Span Mask.

Table 5: **Ablation study of our approach on LAM [34], IAM [33] and READ2016 [11] datasets.** We reported the performance of the standard ViT and studied the effect of our architecture without CNN feature extractor, SAM, Span Masking, respectively, on the results.

| Layers | Heads | IAM [33] | | | | READ2016 [11] | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Val CER | Val WER | Test CER | Test WER | Val CER | Val WER | Test CER | Test WER |
| 8 | 6 | 3.6 | 11.8 | 5.2 | 16.2 | 4.8 | 20.1 | 4.2 | 17.6 |
| **4** | 6 | 3.3 | 10.8 | 4.7 | 14.9 | 4.5 | 19.4 | 3.9 | 16.5 |
| 2 | 6 | 3.5 | 11.4 | 5.1 | 16.1 | 4.3 | 18.8 | 3.9 | 16.7 |
| 1 | 6 | 4.1 | 13.6 | 6.0 | 18.9 | 5.0 | 21.4 | 4.8 | 20.0 |
| Layers | Heads | Val CER | Val WER | Test CER | Test WER | Val CER | Val WER | Test CER | Test WER |
| 4 | 8 | 3.5 | 11.3 | 4.9 | 15.6 | 4.6 | 20.0 | 4.1 | 17.4 |
| 4 | **6** | 3.3 | 10.8 | 4.7 | 14.9 | 4.5 | 19.4 | 3.9 | 16.5 |
| 4 | 4 | 3.3 | 10.9 | 4.7 | 14.8 | 4.4 | 18.8 | 4.1 | 17.6 |
| 4 | 2 | 3.5 | 11.4 | 5.1 | 16.0 | 4.5 | 19.5 | 4.0 | 17.7 |

Table 6: **Ablation study of more hyperparameters on IAM [33] and READ2016 [11] datasets.** We studied the effect of different transformer encoder layers and attention heads on the results.

### 4.3. Ablation studies and visualization analysis

In this section, we delve into two core areas of our study: ablation studies and visualization analysis. The ablation studies are comprehensive, examining the impact of key building blocks within our model and exploring the influence of decoder and critical hyperparameters. These include the masking ratio and span length, as well as the number of transformer encoder and decoder layers and attention heads. The visualization analysis grants us deeper insights into the effectiveness of our span mask strategy. Additionally, we present several qualitative results that showcase the effectiveness of our model.

We hope that our research can serve as a solid basis that can be readily and swiftly used by future researchers. For this reason, we have intentionally refrained from incorporating intricate and opaque components into our model, which could pose difficulties in explanation.

*Effect of CNN feature extractor.* We achieved relatively good results on LAM [34] and READ2016 [11] datasets using only the standard ViT encoder. This encouraged us to consider the ViT architecture as a promising approach for handwritten text recognition tasks. However, we observed that training with the ViT encoder alone resulted in unstable performance and slow convergence speed on the IAM dataset [33], making it difficult to compete with CNN-based models. To improve performance, we introduced a CNN-based feature extractor before the ViT encoder to combine the transformer's global feature extraction capabilities with the CNN's ability to extract local features via a strong inductive bias. Our experiments showed that this modification significantly improved the model's performance and convergence speed.

*Effect of employing Sharpness-Aware Minimization(SAM) [55] optimizer.* We found that convergence to a flatter minimum
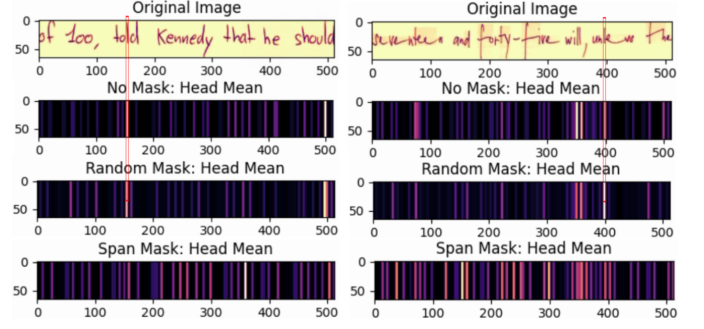


Figure 2: **Visualization of attention maps with different masking strategies on IAM dataset.** In the original image, we highlight the region corresponding to the token of interest with a red bounding box and average the attention across all heads. We observe that when no masking or random masking strategy is employed, each token focuses solely on its own information, as indicated by the illuminated regions in the image. However, when we apply the span masking strategy, a noticeable shift occurs, allowing the token to attend to a broader range of information.

can mitigate overfitting in Handwritten Text Recognition (HTR) models. To facilitate this, we utilized the Sharpness-Aware Minimization (SAM) optimizer, which is straightforward to apply, for locating these flatter minima. Our experimental results show that validation CER and WER on READ2016 increased from 4.8 to 4.5 and from 20.1 to 19.4 with SAM [55]optimizer, indicating that it has a significant impact on HTR tasks.

*Effect of span feature masking.* When labeled data is limited, overfitting can become problematic for transformer-based models. To address this issue, we proposed a new feature masking strategy to reduce overfitting and improve model performance as described in section 3.2. As shown in Table 5, the span feature masking provides consistent and clear improvement across all the datasets.

*Impact of different hyperparameters.* We investigate the impact of different transformer encoder layers and attention heads. The results are illustrated in Table 6. On IAM dataset, taking the number of layers to 4 and attention heads of 6 achieved the best validation CER and WER. To maintain consistency, we employed this set of parameters across all our experiments. We also investigate the impact of different masking strategies. The results are illustrated in Table 7. We can see masking tokens ('*Span Length = 1*') or span feature masking strategy ('*Span Length > 1*') improve the performance for most cases. Span feature masking performs better than random masking tokens and masking none of the tokens. For hyperparameters, taking the masking ratio of 0.4 and a span length of 8 is optimal, which is used for all our experiments. However, larger span lengths (16) reduce the performances, possibly due to the inability to learn context-related information.

*Impact of transformer decoder.* Similar to TrOCR [18], we employ a standard transformer decoder and utilize beam search to produce the final output. We utilized our optimal encoder as the baseline to systematically investigate the impact of the decoder on the overall model performance. The increased number of

| Mask Ratio | Span Length | Val CER | Val WER |
|---|---|---|---|
| 0.0 | None* | 5.1 | 21.9 |
| 0.2 | 1§ | 4.9 | 20.6 |
| 0.4 |  | 4.8 | 20.0 |
| 0.6 |  | 5.1 | 21.8 |
| 0.2 | 4 | 4.6 | 19.8 |
| 0.4 |  | 4.7 | 20.1 |
| 0.6 |  | 4.9 | 20.7 |
| 0.2 | 8 | 4.6 | 19.7 |
| **0.4** | **8** | **4.5** | **19.4** |
| 0.6 |  | 4.9 | 20.9 |
| 0.2 | 16 | 5.0 | 21.5 |
| 0.4 |  | 5.2 | 22.6 |
| 0.6 |  | 5.3 | 23.1 |

\* indicates our approach without any masking strategy.
§ is equivalent to standard random masking.

Table 7: **Ablation study on the masking strategy on READ2016 [11] dataset.** We studied the effect of different mask ratios.

| Decoder Layers | IAM [33] | | | | |
|---|---|---|---|---|---|
|  | Val CER | Val WER | Test CER | Test WER | Param. |
| 0 | 3.3 | 10.8 | 4.7 | 14.9 | 53.5M |
| 8 | - | - | - | - | 129.2M |
| 4 | 5.0 | 15.3 | 7.7 | 21.3 | 91.4M |
| 2 | 5.0 | 15.1 | 7.8 | 21.1 | 72.5M |
| 1 | 5.3 | 15.7 | 8.1 | 21.6 | 63.0M |

Table 8: **Ablation study of adding decoder and training more iterations on IAM [33] dataset.** We studied the effect of different add transformer decoder layers and training iterations on the results. When the number of decoder layers is 8, the model is hard to converge.

| Architecture | 1k iters | Total number of epochs / iters | Training time | Param. | Input size(H x W) |
|---|---|---|---|---|---|
| GFCN [20] | 543 s | 186 (Early stop) / 75.6k iters | 11.4 h | 1.4M | 128 x original W |
| VAN [2] | 420 s | 2100 (Early stop) / 850.4k iters | 99.3 h | 2.7M | original H x W |
| OrigamiNet-24 [4] | 476 s | 100k iters | 13.2 h | 115.3M | 32 x 600 |
| HTR-VT | 586 s | 100k iters | 16.3 h | 53.5M | 64 x 512 |

Note that the reported training times are approximate.

Table 9: **Comparison of the training times across different methods.** We have re-implemented the above methods to compare training times, using the same batch size of 64.

parameters from adding a decoder constrained us to use a batch size of 64 to maintain consistency across all ablations. Our experiments with decoders of varying layer counts, as shown in Table 8, demonstrated that incorporating a transformer decoder did not facilitate better convergence nor prevent overfitting.

*Visualization of attention maps.* In our study, we examine the variations in attention maps when different masking strategies are employed in Figure 2. We averaged the attention across all heads to generate the attention maps displayed. The detailed explanations are as follows:
Firstly, our image size is fixed at 64 x 512, which, after patch embedding, transforms into a shape of 1 x 128, viewed as 128 tokens represented by 128 vertical stripes in the figure. The tokens selected for visualization correspond to the areas enclosed in red boxes in the original image. In the left image, the letter "o" is highlighted, while in the right image, it is the letter "l". According to the principle of self-attention, our selected token should pay more attention to other tokens with higher similarity, which is represented as lighter colors in the attention map. In both no-mask and random-mask scenarios, we can observe that in the left image, the letter "o" in "nvasion" and "bodies" is highlighted, and in the right image, the two "l" letters in "will" are illuminated. This indicates that under no mask and random mask conditions, attention is mainly focused on the token itself. However, a significant change is observed in the span masking scenario. More areas are noticed, indicating that when using span masking, tokens are able to "attend to a broader range of information". This highlights the effectiveness of span masking in enabling tokens to capture more contextual information. The improved contextual awareness provided by span masking facilitates a more comprehensive understanding of the text, which is vital for accurate recognition in handwritten text recognition tasks.
The more examples of the attention maps are provided in the supplementary material.

*Comparison of training time.* Few methods mention the total time required to complete their training, yet this is extremely important for this task. Most approaches that rely on pre-training

or additional data consume significantly expensive computational resources. We compared our method with CNN-based approaches GFCN [20], VAN [2] and OrigamiNet-24 [4] in Table 9. It is important to highlight that the VAN method did not resize images to a fixed resolution but instead used the original image pixels from datasets such as IAM [33]. Similarly, GFCN mentioned that the experiments for the IAM dataset with an image height of 128px, preserving the original width. This resolution is much larger than the fixed resolution we used, which is 512x64. OrigamiNet also used a fixed resolution of 600x32, and our approach of using a fixed resolution follows OrigamiNet. As shown in Table 9, our proposed transformer-based method remains competitive in terms of training time.

*Qualitative results.* We provide visual results in Figure 3 for IAM [33] (First row), READ2016 [11] (Second row) and LAM [34] (Third row). From this, one can recognize the task is challenging, as the visual content present in the text line is not very visible and the background is quite noisy. However, our approach can still produce reasonable predictions on these examples. It is worth noting that in the final image, the ground truth label was annotated incorrectly. Despite this error, our proposed model was still able to accurately recognize the correct handwritten text from the original image, which demonstrates the robustness and effectiveness of the proposed approach.

## 5. Discussion

Although our approach has made notable strides in transformer-based line-level recognition, there is still room for improvement in our current method. The significance of data augmentation for handwriting recognition cannot be overstated, and we have observed that certain data augmentation methods previously utilized in HTR may have adverse effects. Investigating new types of data augmentation specifically tailored for handwriting is a potential direction. Furthermore, delving deeper into mask strategies represents an intriguing avenue for

G. T. : carefully casual. 'The servants are all on
Ours : carefully casual. 'The servants are all on

G. T. : indication pointed to Eve being held.
Ours : indication pointed to - EWe being held.

G. T. : üng geherter Wachten
Ours : üng geherter Wachten

G. T. : Teütsch volckh gebraüche.
Ours : Teütsch volckh gebraüche.

G. T. : Alla bontà di V.P. Rev:ma, che tanto si pren-
Ours : Alla bontà di V.P. Rev:ma, che tanto si pren-

G. T. : superare da chicchessia # stima della persona e del
Ours : superare da chichessia # lla stima della persova e del

Figure 3: Results on example lines from the IAM [33] (First row), READ2016 [11] (Second row) and LAM [34] (Third row) of the best performing model.

exploration; learnable mask strategies adapted for handwriting could prove more beneficial. Lastly, expanding from line-level to paragraph-level or page-level recognition will be the focus of our future research.

## 6. Conclusion

In this work, we have presented a simple and data-efficient approach for handwritten text recognition. With minimal modifications to the ViT architecture, we have successfully developed a ViT-like model that surpasses state-of-the-art performance without requiring pre-training or additional data. Notably, our experiments highlight the remarkable data efficiency of our model compared to ViT and DeiT, while preserving its superior generalizability even in scenarios with vast amounts of available data. These findings provide a promising direction for improving the performance of handwritten text recognition, particularly in limited data scale settings.

## References

[1] D. Coquenet, C. Chatelain, T. Paquet, Dan: a segmentation-free document attention network for handwritten document recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence (2023).

[2] D. Coquenet, C. Chatelain, T. Paquet, End-to-end handwritten paragraph text recognition using a vertical attention network, IEEE Transactions on Pattern Analysis and Machine Intelligence 45 (1) (2022) 508–524.

[3] T. Bluche, Joint line segmentation and transcription for end-to-end handwritten paragraph recognition, Advances in neural information processing systems 29 (2016).

[4] M. Yousef, T. E. Bishop, Origaminet: weakly-supervised, segmentation-free, one-step, full page text recognition by learning to unfold, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 14710–14719.

[5] X. Yang, E. Yumer, P. Asente, M. Kraley, D. Kifer, C. Lee Giles, Learning to extract semantic structure from documents using multimodal fully convolutional neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 5315–5324.

[6] A. K. Bhunia, A. Das, A. K. Bhunia, P. S. R. Kishore, P. P. Roy, Handwriting recognition in low-resource scripts using adversarial learning, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 4767–4776.

[7] A.-L. Bianne-Bernard, F. Menasri, R. A.-H. Mohamad, C. Mokbel, C. Kermorvant, L. Likforman-Sulem, Dynamic and contextual information in hmm modeling for handwritten word recognition, IEEE transactions on pattern analysis and machine intelligence 33 (10) (2011) 2066–2080.

[8] S. Espana-Boquera, M. J. Castro-Bleda, J. Gorbe-Moya, F. Zamora-Martinez, Improving offline handwritten text recognition with hybrid hmm/ann models, IEEE transactions on pattern analysis and machine intelligence 33 (4) (2010) 767–779.

[9] A. Graves, J. Schmidhuber, Offline handwriting recognition with multidimensional recurrent neural networks, Advances in neural information processing systems 21 (2008).

[10] C. Wigington, C. Tensmeyer, B. Davis, W. Barrett, B. Price, S. Cohen, Start, follow, read: End-to-end full-page handwriting recognition, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 367–383.

[11] J. A. Sanchez, V. Romero, A. H. Toselli, E. Vidal, Icfhr2016 competition on handwritten text recognition on the read dataset, in: 2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR), IEEE, 2016, pp. 630–635.

[12] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, R. Girshick, Masked autoencoders are scalable vision learners, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 16000–16009.

[13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, Advances in neural information processing systems 30 (2017).

[14] A. Graves, S. Fernández, F. Gomez, J. Schmidhuber, Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks, in: Proceedings of the 23rd international conference on Machine learning, 2006, pp. 369–376.

[15] P. Voigtlaender, P. Doetsch, H. Ney, Handwriting recognition with large multidimensional long short-term memory recurrent neural networks, in: 2016 15th international conference on frontiers in handwriting recognition (ICFHR), IEEE, 2016, pp. 228–233.

[16] C. Wigington, S. Stewart, B. Davis, B. Barrett, B. Price, S. Cohen, Data augmentation for recognition of handwritten words and lines using a cnn-lstm network, in: 2017 14th IAPR international conference on document analysis and recognition (ICDAR), Vol. 1, IEEE, 2017, pp. 639–645.

[17] J. Michael, R. Labahn, T. Grüning, J. Zöllner, Evaluating sequence-to-sequence models for handwritten text recognition, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), IEEE, 2019, pp. 1286–1293.

[18] M. Li, T. Lv, J. Chen, L. Cui, Y. Lu, D. Florencio, C. Zhang, Z. Li, F. Wei, Trocr: Transformer-based optical character recognition with pre-trained models, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 37, 2023, pp. 13094–13102.

[19] C. Wick, J. Zöllner, T. Grüning, Transformer for handwritten text recognition using bidirectional post-decoding, in: International Conference on Document Analysis and Recognition, Springer, 2021, pp. 112–126.

[20] D. Coquenet, C. Chatelain, T. Paquet, Recurrence-free unconstrained handwritten text recognition using gated fully convolutional network, in: 2020 17th International Conference on Frontiers in Handwriting Recognition (ICFHR), IEEE, 2020, pp. 19–24.

[21] A. El-Yacoubi, M. Gilloux, R. Sabourin, C. Y. Suen, An hmm-based approach for off-line unconstrained handwritten word modeling and recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence 21 (8) (1999) 752–760.

[22] B. Shi, X. Bai, C. Yao, An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition, IEEE transactions on pattern analysis and machine intelligence 39 (11) (2016) 2298–2304.

[23] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[24] D. Coquenet, C. Chatelain, T. Paquet, Span: a simple predict & align network for handwritten paragraph recognition, in: International Conference on Document Analysis and Recognition, Springer, 2021, pp. 70–84.

[25] F. Chollet, Xception: Deep learning with depthwise separable convolutions, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1251–1258.

[26] P. Doetsch, A. Zeyer, H. Ney, Bidirectional decoder networks for attention-based end-to-end offline handwriting recognition, in: 2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR), IEEE, 2016, pp. 361–366.

[27] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge, International journal of computer vision 115 (2015) 211–252.

[28] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, Y. Wei, Deformable convolutional networks, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 764–773.

[29] X. Zhu, H. Hu, S. Lin, J. Dai, Deformable convnets v2: More deformable, better results, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 9308–9316.

[30] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, H. Jégou, Training data-efficient image transformers & distillation through attention, in: International conference on machine learning, PMLR, 2021, pp. 10347–10357.

[31] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, P. Dollár, Designing network design spaces, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 10428–10436.

[32] Y. Liu, E. Sangineto, W. Bi, N. Sebe, B. Lepri, M. Nadai, Efficient training of visual transformers with small datasets, Advances in Neural Information Processing Systems 34 (2021) 23818–23830.

[33] U.-V. Marti, H. Bunke, The iam-database: an english sentence database for offline handwriting recognition, International Journal on Document Analysis and Recognition 5 (2002) 39–46.

[34] S. Cascianelli, V. Pippi, M. Maarand, M. Cornia, L. Baraldi, C. Kermorvant, R. Cucchiara, The lam dataset: A novel benchmark for line-level handwritten text recognition, in: 2022 26th International Conference on Pattern Recognition (ICPR), IEEE, 2022, pp. 1506–1513.

[35] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, arXiv preprint arXiv:2010.11929 (2020).

[36] J. Puigcerver, Are multidimensional recurrent layers really necessary for handwritten text recognition?, in: 2017 14th IAPR international conference on document analysis and recognition (ICDAR), Vol. 1, IEEE, 2017, pp. 67–72.

[37] I. Cojocaru, S. Cascianelli, L. Baraldi, M. Corsini, R. Cucchiara, Watch your strokes: Improving handwritten text recognition with deformable convolutions, in: 2020 25th International Conference on Pattern Recognition (ICPR), IEEE, 2021, pp. 6096–6103.

[38] T. Bluche, R. Messina, Gated convolutional recurrent neural networks for multilingual handwriting recognition, in: 2017 14th IAPR international conference on document analysis and recognition (ICDAR), Vol. 1, IEEE, 2017, pp. 646–651.

[39] T. Bluche, J. Louradour, R. Messina, Scan, attend and read: End-to-end handwritten paragraph recognition with mdlstm attention, in: 2017 14th IAPR international conference on document analysis and recognition (ICDAR), Vol. 1, IEEE, 2017, pp. 1050–1055.

[40] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, arXiv preprint arXiv:1810.04805 (2018).

[41] D. Hendrycks, K. Gimpel, Gaussian error linear units (gelus), arXiv preprint arXiv:1606.08415 (2016).

[42] J. L. Ba, J. R. Kiros, G. E. Hinton, Layer normalization, arXiv preprint arXiv:1607.06450 (2016).

[43] I. Loshchilov, F. Hutter, Decoupled weight decay regularization, arXiv preprint arXiv:1711.05101 (2017).

[44] H. Bao, L. Dong, S. Piao, F. Wei, Beit: Bert pre-training of image transformers, arXiv preprint arXiv:2106.08254 (2021).

[45] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, Roberta: A robustly optimized bert pretraining approach, arXiv preprint arXiv:1907.11692 (2019).

[46] S. Cascianelli, M. Cornia, L. Baraldi, R. Cucchiara, Boosting modern and historical handwritten text recognition with deformable convolutions, International Journal on Document Analysis and Recognition (IJDAR) 25 (3) (2022) 207–217.

[47] M. A. Islam, S. Jia, N. D. Bruce, How much position information do convolutional neural networks encode?, arXiv preprint arXiv:2001.08248 (2020).

[48] L. Kang, P. Riba, M. Rusiñol, A. Fornés, M. Villegas, Pay attention to what you read: Non-recurrent handwritten text-line recognition, Pattern Recognition 129 (2022) 108766.

[49] A. Graves, J. Schmidhuber, Framewise phoneme classification with bidirectional lstm and other neural network architectures, Neural networks 18 (5-6) (2005) 602–610.

[50] M. Yousef, K. F. Hussain, U. S. Mohammed, Accurate, data-efficient, unconstrained text recognition with convolutional neural networks, Pattern Recognition 108 (2020) 107482.

[51] B. Li, Y. Hu, X. Nie, C. Han, X. Jiang, T. Guo, L. Liu, Dropkey for vision transformer, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023, pp. 22700–22709.

[52] K. Song, X. Tan, T. Qin, J. Lu, T.-Y. Liu, Mass: Masked sequence to sequence pre-training for language generation, arXiv preprint arXiv:1905.02450 (2019).

[53] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, O. Levy, Spanbert: Improving pre-training by representing and predicting spans, Transactions of the association for computational linguistics 8 (2020) 64–77.

[54] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556 (2014).

[55] P. Foret, A. Kleiner, H. Mobahi, B. Neyshabur, Sharpness-aware minimization for efficiently improving generalization, in: International Conference on Learning Representations (ICLR), 2020.

[56] M. A. Souibgui, S. Biswas, A. Mafla, A. F. Biten, A. Fornés, Y. Kessentini, J. Lladós, L. Gomez, D. Karatzas, Text-diae: A self-supervised degradation invariant autoencoder for text recognition and document enhancement, in: proceedings of the AAAI conference on artificial intelligence, Vol. 37, 2023, pp. 2330–2338.

[57] M. Dhiaf, A. C. Rouhou, Y. Kessentini, S. B. Salem, Msdoctr-lite: A lite transformer for full page multi-script handwriting recognition, Pattern Recognition Letters 169 (2023) 28–34.

# Appendix

This appendix contains the following sections:

-

-

-

-

-

| Methods | IAM [33] | |
|---|---|---|
| | Test CER | Test WER |
| **ResNet-18** | 4.7 | 14.9 |
| ResNet-50 | 4.9 | 15.6 |
| VGG-16 | 7.2 | 22.1 |

Table .10: **Ablation study on using various CNN backbones on IAM [33] dataset.** We reported the performance of the our approach using different backbones and studied the effect of them. We use **ResNet-18** as our final solution.

## Appendix A. Visual results on the IAM [33], READ2016 [11] and LAM [34] datasets.

We show our handwritten text recognition method's visual results on the IAM [33], READ2016 [11], and LAM [34] datasets.

IAM [33] is a well-known offline handwriting benchmark dataset containing 6 482 images for training, 976 images for validation, and 2 915 images for testing. The image is in grayscale and the font has ligatures and some missing parts. Visual results are provided in Figure B.4.

READ2016 [11] consists of 8 349 train, 1 040 validation, and 1 138 test images. The image contains a noisy background with some blurry fonts. Visual results are provided in Figure B.5.

LAM [34] is currently the largest line-level handwritten text recognition dataset that contains 19 830 lines for training, 2 470 lines for validation, and 3 523 lines for testing. The image contains fonts with stains, some lines of text are skewed and include both upper and lower characters. Visual results are provided in Figure B.6.
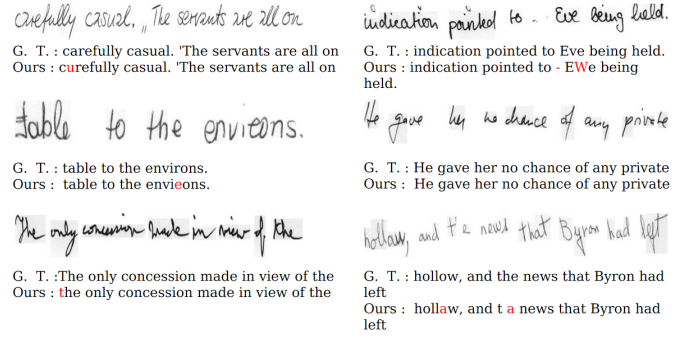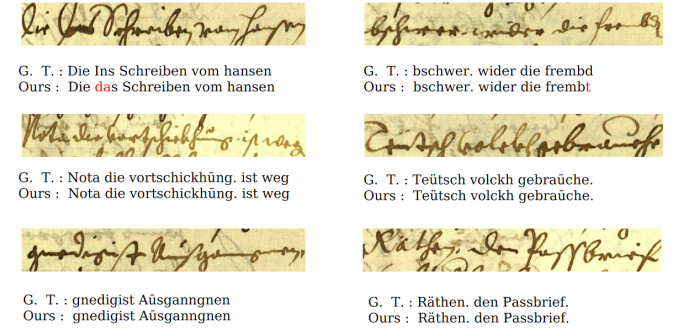


G. T. : carefully casual. 'The servants are all on
Ours : curefully casual. 'The servants are all on

G. T. : indication pointed to Eve being held.
Ours : indication pointed to - EWe being held.

G. T. : table to the environs.
Ours : table to the envieons.

G. T. : He gave her no chance of any private
Ours : He gave her no chance of any private

G. T. :The only concession made in view of the
Ours : the only concession made in view of the

G. T. : hollow, and the news that Byron had left
Ours : hollaw, and t a news that Byron had left

Figure B.4: Visual results on IAM [33]



G. T. : Die Ins Schreiben vom hansen
Ours : Die das Schreiben vom hansen

G. T. : bschwer. wider die frembd
Ours : bschwer. wider die frembt

G. T. : Nota die vortschickhüng. ist weg
Ours : Nota die vortschickhüng. ist weg

G. T. : Teütsch volckh gebräuche.
Ours : Teütsch volckh gebräuche.

G. T. : gnedigist Aüsganngnen
Ours : gnedigist Aüsganngnen

G. T. : Räthen. den Passbrief.
Ours : Räthen. den Passbrief.

Figure B.5: Visual results on READ2016 [11]



G. T. : intanto il lavoro si dee attribuire alla Po-
Ours : intanto il lavoro si dee attribuire alla Po-

G. T. :Romani. Si va intanto avvicinando la Pri-
Ours : Romani. Si va intanto avvicinando la Pri-

G. T. : cercando così d'approfittar-
Ours : cercando così d'approfittar-

G. T. : talm.e sconcertata la mia per altro fievole sa-
Ours : cthelm.e sconcertata la mia per altro fievole sa-

G. T. : Alla bontà di V.P. Rev:ma, che tanto si pren-
Ours : Alla bontà di V.P. Rev:ma, che tanto si pren-

G. T. : superare da chicchessia # stima della persona e del
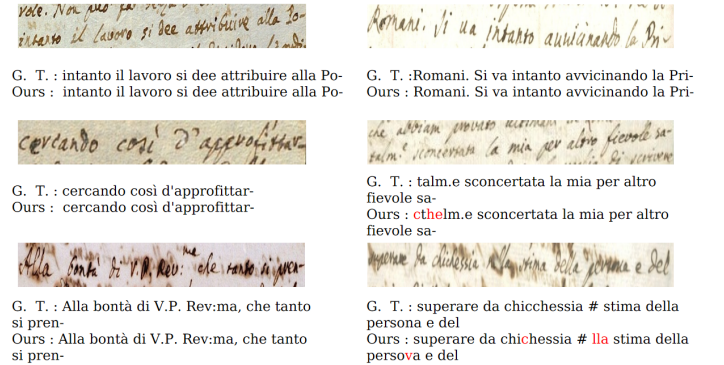Ours : superare da chichessia # lla stima della persova e del

Figure B.6: Visual results on LAM [34]

## Appendix B. CNN Backbones Ablation

In this study, we investigated the impact of different CNN backbones on the overall model performance. We chose the most fundamental ResNet [23]and VGG [54]architectures as our CNN backbones, consistent with the simple and easy-to-implement principles outlined in our paper. The performance of the proposed method is observed to be robust across various backbones. Particularly, ResNet-18 exhibits superior performance compared to other backbones.

## Appendix C. Span mask strategy

The details of our implementation of the span mask strategy are as follows: To achieve the designated mask ratio $R$ (e.g.,
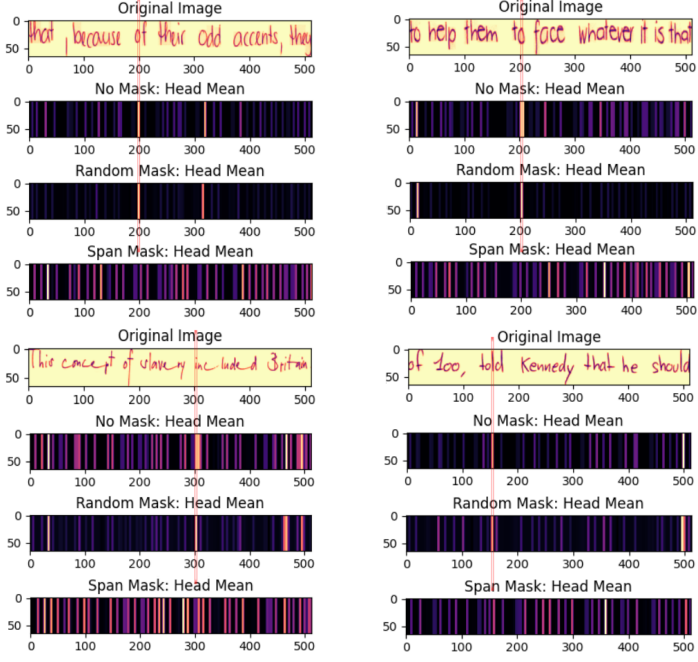
Figure C.7: Visualization of attention maps

there is a conspicuous expansion in the illuminated regions, indicating that the token now engages with a substantially broader contextual landscape.

0.4 of $L$), we adopt an iterative process of sampling spans. In each iteration, we start by defining a maximum span length $l$ (i.e., the number of interconnected tokens), and then randomly select the starting point for each span. Noting that the maximum span length is fixed. This means that the length of the sampled masked segments remains the same for each iteration.

## Appendix D.  Training details about DeiT [30]and DropKey [51]

We implemented it completely following the steps in Drop-Key [51], moving dropout operations ahead of attention matrix calculation and setting the Key as the dropout unit, yielding a dropout-before-softmax scheme. And We set the drop ratio to 0.1. In DeiT [30], each layer has a dimension of 768 and 6 heads as used in our approach. At the same time, we implemented DeiT with no distillation.

## Appendix E.  Visualization results of attention maps

In Figure C.7, we present an extensive set of attention map visualizations that offer valuable insights into the model's behavior. We demarcate the region of interest in the original image corresponding to the token under scrutiny using a red bounding box. It is evident that when employing no mask and random mask strategies, the attention is highly localized, illuminating only the regions that correspond to the annotated characters in the original image. For instance, in the first visualization, the selected token corresponds to the letter 'o' in the word 'of,' and the attention map distinctly highlights this specific region. This suggests that, in these scenarios, each token is predominantly self-attentive. Conversely, when utilizing a span mask strategy,