

ICAL: Implicit Character-Aided Learning for Enhanced Handwritten Mathematical Expression Recognition

Jianhua Zhu¹[0009–0000–3982–2739], Liangcai Gao¹(✉), and Wenqi Zhao¹

Wangxuan Institute of Computer Technology, Peking University, Beijing, China
 zhujianhuapku@pku.edu.cn
 gaoliangcai@pku.edu.cn
 wenqizhao@stu.pku.edu.cn

Abstract. Significant progress has been made in the field of handwritten mathematical expression recognition, while existing encoder-decoder methods are usually difficult to model global information in \LaTeX . Therefore, this paper introduces a novel approach, Implicit Character-Aided Learning (ICAL), to mine the global expression information and enhance handwritten mathematical expression recognition. Specifically, we propose the Implicit Character Construction Module (ICCM) to predict implicit character sequences and use a Fusion Module to merge the outputs of the ICCM and the decoder, thereby producing corrected predictions. By modeling and utilizing implicit character information, ICAL achieves a more accurate and context-aware interpretation of handwritten mathematical expressions. Experimental results demonstrate that ICAL notably surpasses the state-of-the-art(SOTA) models, improving the expression recognition rate (ExpRate) by 2.25%/1.81%/1.39% on the CROHME 2014/2016/2019 datasets respectively, and achieves a remarkable 69.06% on the challenging HME100k test set. We make our code available on the GitHub.¹

Keywords: handwritten mathematical expression recognition · transformer · implicit character-aided learning · encoder-decoder model

1 Introduction

The Handwritten Mathematical Expression Recognition (HMER) task involves taking an image of a handwritten mathematical expression as input and having the model predict the corresponding \LaTeX . The HMER task has a wide range of applications, such as being used for intelligent grading of mathematical assignments, building online grading systems, and improving the efficiency of online education. Therefore, how to improve the recognition accuracy of handwritten mathematical expressions has become a hot topic in previous works.

¹ <https://github.com/qingzhenduyu/ICAL>

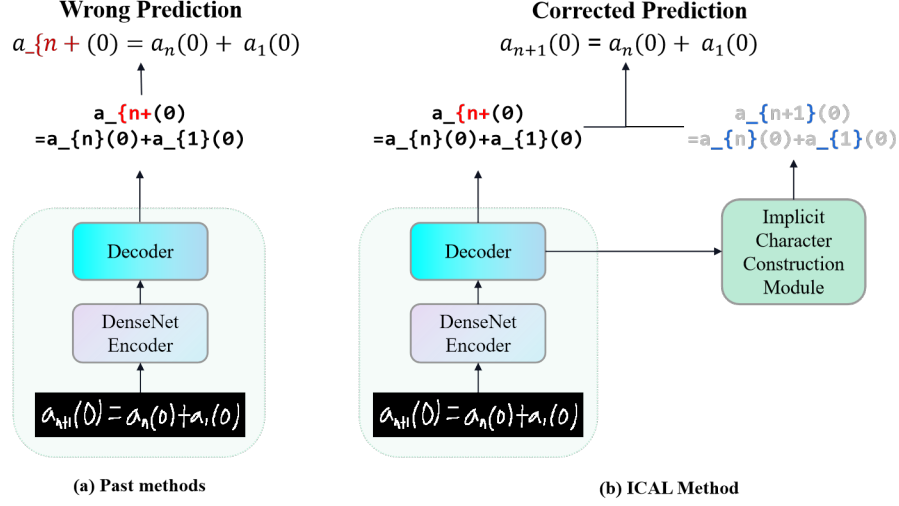


Fig. 1. (a) Illustration of past method which uses DenseNet Encoder and RNN/Transformer Decoder. (b) Our ICAL method aided by implicit character learning. The characters highlighted in red signify inaccuracies in the prediction, whereas the blue highlights denote implicit characters.

Due to the diversity of handwriting and the two-dimensional structure of mathematical expressions, HMER is highly challenging. Compared to the recognition of handwritten text in natural language, the HMER task requires not only the prediction of explicit characters (i.e., characters that are directly represented when written by hand) but also implicit characters, such as “ \sim ”, “ $_$ ”, “ $\{$ ”, and “ $\}$ ”, which is necessary to achieve a complete description of a two-dimensional mathematical expression. Past methods [30, 35, 36] based on encoder-decoder models often extract image features through the encoder, while the decoder aligns visual and textual features and predicts the \LaTeX . However, they often lack modeling of the global information of the expression, which in turn fails to correct prediction errors made by the decoder, as shown in the figure 1 (a):

In this paper, we utilize the task of predicting implicit characters (i.e., “ \sim ”, “ $_$ ”, “ $\{$ ”, and “ $\}$ ”) to assist the decoder in modeling the global information of \LaTeX , which can further correct the output of the decoder and improve recognition performance. To this end, we propose an Implicit Character Construction Module, capable of modeling the sequence of implicit characters from the output of the Transformer Decoder [22]. This global information is then passed to a subsequent Fusion Module which integrates the information with the output of the Transformer decoder to achieve a more accurate prediction of the \LaTeX sequence.

In this work, the main contributions of our work are summarized as follows:

- We introduced the Implicit Character Construction Module (ICCM) to model implicit character information, which can effectively utilize the global information in \LaTeX .
- We proposed the Fusion Module to aggregate the output of the ICCM, thereby correcting the prediction of the Transformer Decoder.
- Experimental results indicate that the ICAL method surpasses previous state-of-the-art methods and achieves expression recognition rate (ExpRate) of 60.63%, 58.79%, and 60.51% on the CROHME 2014 [14]/2016 [15]/2019 [13] test sets, respectively, and an ExpRate of 69.06% on the HME100K test set [28].

2 Related Work

2.1 Traditional Methods

In traditional handwritten mathematical expression recognition, the process mainly involves symbol recognition and structural analysis. Symbol recognition requires segmenting and identifying individual symbols within expressions, utilizing techniques like pixel-based segmentation proposed by OKAMOTO et al [21] and further refined by HA et al [8], into recursive cropping methods. These approaches often depend on predefined thresholds. For symbol classification, methods like Hidden Markov Models (HMM) [1, 11, 24], Elastic Matching [3, 23], and Support Vector Machines (SVM) [10] have been used, where HMMs allow for joint optimization without explicit symbol segmentation, though at a high computational cost.

Structural analysis employs strategies like the two-dimensional Stochastic Context-free Grammar (SCFG) [6] and algorithms such as Cocke-Younger-Kasami (CKY) [16] for parsing, albeit slowly due to the complexity of two-dimensional grammar. Faster parsing has been achieved with Left-to-right Recursive Descent and Tree Transformation methods [29], the latter describing the arrangement and grouping of symbols into a structured tree for parsing. Some approaches bypass two-dimensional grammar altogether, using Define Clause Grammar (DCG) [4] and Formula Description Grammar [18] for one-dimensional parsing, highlighting the challenges in designing comprehensive grammars for the diverse structures of mathematical expressions.

2.2 Deep Learning Methods

In recent years, encoder-decoder-based deep learning models have become the mainstream framework in the field of HMER. Depending on the architecture of the decoder, past deep learning approaches can be categorized into methods based on Recurrent Neural Networks (RNNs) [2, 7, 12, 19, 25–28, 30–33] and those based on the Transformer model [34–36]. Furthermore, based on the decoding strategy, these methods can be divided into those based on sequence decoding [2, 7, 12, 19, 26, 27, 30, 31, 33–36] and those based on tree-structured decoding [25, 28, 32].

RNN-based methods In 2017, Zhang et al. proposed an end-to-end deep learning model, WAP [33], to address the problem of HMER. The encoder part of the model is a fully convolutional neural network similar to VGGnet [17]. The decoder part uses a GRU [5] model to generate the predicted \LaTeX sequence from the extracted visual features. The WAP not only avoids issues caused by inaccurate symbol segmentation but also eliminates the need for manually predefined \LaTeX syntax, thereby becoming a benchmark model for subsequent deep learning methods. Following WAP, Zhang et al. further proposed the DenseWAP [30], which replaces the VGGnet with the DenseNet [9]. Subsequent work has commonly adopted the DenseNet as the backbone network for the encoder. The CAN model [12] introduces a Multi-Scale Counting Module, utilizing a symbol counting task as an auxiliary task to be jointly optimized with the expression recognition task.

Transformer-based methods To alleviate the issue of unbalanced output and fully utilize bidirectional language information, BTTR [36] adopts a bidirectional training strategy on top of the Transformer-based decoder. Following this, CoMER [35] incorporates coverage information [20, 33] into Transformer Decoder, introducing an Attention Refinement Module. This module utilizes the attention weights from the Multi-head Attention mechanism within the Transformer decoder to compute the coverage vector, all while maintaining the characteristic of parallel decoding. Based on CoMER, the GCN [34] incorporates extra symbol categorization information, utilizing the General Category Recognition Task as a supplementary task for joint optimization with the HMER task, resulting in a notable performance. However, the category recognition task introduced by GCN requires manual construction of symbol categories and is limited to specific datasets.

Tree-based decoding methods \LaTeX , as a markup language, can be easily parsed into a tree-like expression due to the influence of delimiters such as brackets. Therefore, by leveraging the inherent two-dimensional structure of mathematical expressions, models can provide certain interpretability for the prediction process. Zhang et al. proposed the DenseWAP-TD [32], which replaces the GRU decoder that directly regresses the \LaTeX sequence with a decoder based on a two-dimensional tree structure. The TDv2 model [25], during training, uses different transformation methods for the same \LaTeX string, weakening the context dependency and endowing the decoder with stronger generalization capabilities. The SAN model [28] converts the \LaTeX sequence into a parsing tree and designs a series of syntactic rules to transform the problem of predicting \LaTeX sequences into a tree traversal process. Additionally, SAN introduces a new Syntax-Aware Attention Module to better utilize the syntactic information in \LaTeX .

3 Methodology

In this section, we will elaborate in detail on the model structure of the ICAL we proposed, as shown in the figure 2. In Section 3.1, we will briefly introduce the DenseNet [9] used in the encoder part. In Section 3.2, we will introduce the

Transformer decoder that adopts coverage attention mechanism [22, 35], which can alleviate the lack of coverage problem. In Sections 3.3 and 3.4, we will introduce the Implicit Character Construction Module (ICCM) proposed in this paper and the Fusion Module that integrates implicit character information. Finally, in Section 3.5, we will discuss how the implicit character loss and fusion loss are introduced while employing a bidirectional training strategy [36].

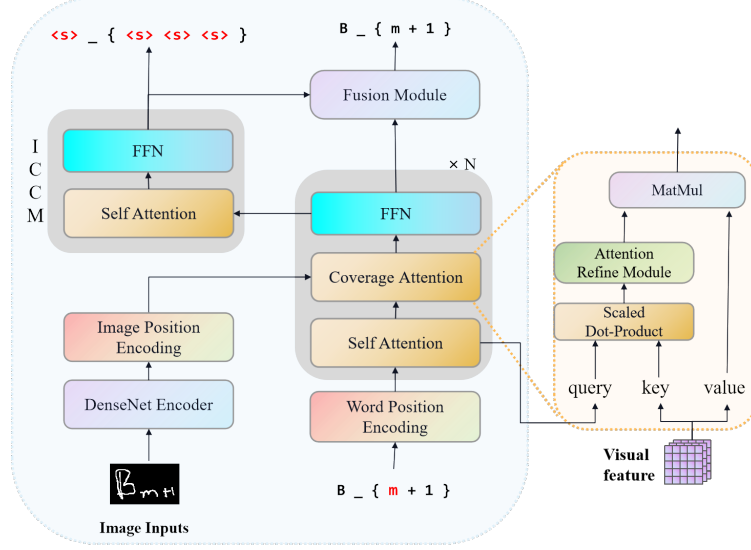


Fig. 2. The architecture of ICAL model (left) and Coverage Attention (right). To simplify the illustration, we have condensed the depiction of bidirectional training in the figure.

3.1 Visual Encoder

Similar to most of the previous work [2, 7, 12, 19, 25–28, 30–32, 34–36], we continue to use DenseNet [9] as the visual encoder to extract features from the input images.

DenseNet consists of multiple dense blocks and transition layers. Within each dense block, the input for each layer is the concatenated outputs from all preceding layers, enhancing the flow of information between layers in this manner. The transition layers reduce the dimensions of the feature maps using 1×1 convolutional kernels, decreasing the number of parameters and controlling the complexity of the model.

For an input grayscale image of size $1 \times H_0 \times W_0$, the output visual feature $\mathbf{V}_{feature} \in \mathbb{R}^{D \times H \times W}$. The ratios of H to H_0 and W to W_0 are both $1/16$. In

this work, a 1×1 convolutional layer is used to adjust the number of channels D to d_{model} , aligning it with the dimension size of the Transformer decoder.

3.2 Transformer Decoder with ARM

In the decoder part, we employ a Transformer Decoder with Attention Refine Module (ARM) as the decoder [22, 35], which mainly consists of three modules: Self-Attention Module, Coverage Attention Module, and Feed-Forward Network.

Self-Attention Module The self-attention module functions as a core component of the standard Transformer decoder, utilizing a masked multi-head attention mechanism. This module processes a set of queries (\mathbf{Q}), keys (\mathbf{K}), and values (\mathbf{V}) to produce an output that leverages both dot-product attention and a multi-head strategy for processing information. For each head within the multi-head attention framework, linear transformations are applied to \mathbf{Q} , \mathbf{K} , and \mathbf{V} using transformation matrices \mathbf{W}_i^q , \mathbf{W}_i^k , and \mathbf{W}_i^v , respectively, where i denotes the index of the head.

The mechanism first calculates a scaled dot-product attention score \mathbf{E}_i for the i -th head by multiplying the transformed queries and keys and then scaling the result by the square root of the dimension of the key vectors \mathbf{K} to avoid large values that could hinder softmax computation:

$$\mathbf{E}_i = \frac{(\mathbf{Q}\mathbf{W}_i^q)(\mathbf{K}\mathbf{W}_i^k)^T}{\sqrt{d_k}}, \quad (1)$$

A softmax function is then applied to these scores to obtain the attention weights \mathbf{A}_i :

$$\mathbf{A}_i = \text{softmax}(\mathbf{E}_i), \quad (2)$$

These weights are used to compute a weighted sum of the values, producing the output \mathbf{H}_i for each head:

$$\mathbf{H}_i = \mathbf{A}_i(\mathbf{V}\mathbf{W}_i^v), \quad (3)$$

Finally, the outputs of all heads are concatenated and linearly transformed to produce the final output of the multi-head attention module:

$$\text{MultiHeadAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\mathbf{H}_0, \mathbf{H}_1, \dots, \mathbf{H}_{\text{head}-1})\mathbf{W}^o. \quad (4)$$

Since the decoder performs decoding in an autoregressive manner, the prediction of the current symbol depends on past predicted symbols and input visual information. To avoid the decoder obtaining information about future symbols when predicting the current symbol and to ensure parallelism, Self-Attention Module uses a masked lower-triangular matrix to constrain the information that the self-attention module can access at the current step.

Coverage Attention Module The CoMER model, without affecting the parallelism of the Transformer, introduces coverage attention commonly used in

RNNs [33] into the Transformer decoder by improving the Cross Attention module. Within the Cross Attention module of the CoMER model, an Attention Refine Module (ARM) is incorporated. By utilizing alignment information from the previous layer and the current layer, it refines the current attention weights A_i , enabling the decoder to faithfully convert the text structure from visual features to corresponding L^AT_EX text. The update formulas for the attention weights in the j -th layer of the decoder, denoted as \hat{A}^j , are as follows in Equations 5 and 6.

$$\hat{\mathbf{E}}^j = \text{ARM}(E^j, \hat{\mathbf{A}}^{j-1}), \quad (5)$$

$$\hat{\mathbf{A}}^j = \text{softmax}(\hat{\mathbf{E}}^j). \quad (6)$$

Feed-Forward Network The feed-forward neural network consists of two linear layers and the ReLU non-linear activation function. For input \mathbf{X} from Coverage Attention Module, it is as follows:

$$\mathbf{E}_{\text{feature}} = \text{FFN}(\mathbf{X}) = \text{ReLU}(\mathbf{X}\mathbf{W}_0 + \mathbf{b}_0)\mathbf{W}_1 + \mathbf{b}_1. \quad (7)$$

3.3 Implicit Character Construction Module

For the target L^AT_EX sequence, we only retain its implicit characters, namely “~”, “_”, “{”, and “}”, replacing other characters with newly constructed ones to form a sequence corresponding to the implicit characters. For example, for the target L^AT_EX sequence `B _ { m + 1 }`, the corresponding implicit character sequence is `<space> _ { <space> <space> <space> }`.

The output of the Transformer Decoder with ARM serves as the input to our Implicit Character Construction Module (ICCM), which consists of a layer of masked self-attention and a Feed-Forward Network (FFN) layer. Consequently, the output of the ICCM, $\mathbf{I}_{\text{feature}}$, is calculated as follows:

$$\mathbf{I}_{\text{feature}} = \text{FFN}(\text{SelfAttention}(\mathbf{E}_{\text{feature}})). \quad (8)$$

3.4 Fusion Module

To integrate the information learned by the ICCM, we introduce a weighted adjustment strategy based on the attention mechanism, capable of aggregating the output of the ICCM, thereby correcting the prediction results of the Transformer Decoder.

The Fusion Module assimilates inputs from both the ICCM’s output, denoted as $\mathbf{I}_{\text{feature}} \in \mathbb{R}^{B \times T \times d_{\text{model}}}$, and the Transformer Decoder’s output, denoted as $\mathbf{E}_{\text{feature}} \in \mathbb{R}^{B \times T \times d_{\text{model}}}$. The attention weights \mathbf{f}_{att} are calculated as follows:

$$\mathbf{f}_{\text{att}} = \sigma(w_{\text{att}}(\text{Concat}(\mathbf{E}_{\text{feature}}, \mathbf{I}_{\text{feature}}))), \quad (9)$$

where a linear layer w_{att} maps the concatenate feature matrix back to the original dimension d_{model} , and attention weights f_{att} are computed through a sigmoid activation function σ .

Finally, f_{att} are used to perform a weighted fusion of the two features, resulting in the final output feature for prediction.

$$\mathbf{F} = \mathbf{f}_{att} \odot \mathbf{E}_{feature} + (1 - \mathbf{f}_{att}) \odot \mathbf{I}_{feature}, \quad (10)$$

where \odot indicates element-wise multiplication.

3.5 Loss Function

To alleviate the issue of unbalanced output, we adhere to the bidirectional training strategy used in BTTR and CoMER, necessitating the computation of losses in both directions within the same batch.

The total loss function is divided into three parts:

$$\mathcal{L} = \mathcal{L}_{Initial} + \mathcal{L}_{Implicit} + \mathcal{L}_{Fusion}, \quad (11)$$

where both $\mathcal{L}_{Initial}$ and \mathcal{L}_{fusion} are standard cross-entropy loss functions.

$\mathcal{L}_{Initial}$ calculates the loss using the Transformer decoder’s predicted probabilities and the ground truth, consistent with the approach used in CoMER. \mathcal{L}_{fusion} uses the predicted probabilities outputted by the Fusion Module and the ground truth.

Since we construct the implicit character sequence directly from the \LaTeX sequence, keeping its length consistent with the original \LaTeX sequence, the constructed target implicit character sequence exhibits an imbalance in the occurrence frequency between $\langle \text{space} \rangle$ and implicit characters. To address this, we use a weighted cross-entropy loss function to calculate the loss for implicit characters, as follows:

$$\mathcal{L}_{implicit} = - \sum_{i=1}^N \sum_{t=1}^{T_i} w_{y_{i,t}} \cdot \log(\hat{y}_{i,t}), \quad (12)$$

where N represents the batch size, T_i is the length of the i -th sequence, $y_{i,t}$ is the ground truth token at position t in the i -th sequence, $\hat{y}_{i,t}$ is the predicted probability of the correct token at position t in the i -th sequence, and $w_{y_{i,t}}$ is the weight associated with the ground truth token $y_{i,t}$.

The weight $w_{y_{i,t}}$ for each token is dynamically adjusted based on the occurrence frequency of the token within the entire batch of sequences. The adjustment is made using a logarithmic function to ensure a smooth transition of weights across different frequencies:

$$w_{y_{i,t}} = 1.0 + \log \left(1 + \frac{1}{f_{y_{i,t}} + \epsilon} \right), \quad (13)$$

where $f_{y_{i,t}}$ is the frequency of the token in the target sequences and ϵ is set to $1e - 6$.

4 Experiments

4.1 Dataset

The CROHME dataset includes data from the Online Handwritten Mathematical Expressions Recognition Competitions (CROHME) [13–15] held over several years and is currently the most widely used dataset for handwritten mathematical expression recognition. The training set of the CROHME dataset consists of 8,836 samples, while the CROHME 2014 [14]/2016 [15]/2019 [13] test sets contain 986, 1147, and 1199 samples respectively. In the CROHME dataset, each handwritten mathematical expression is stored in InkML format, recording the trajectory coordinates of the handwritten strokes. Before training, we convert the handwritten stroke trajectory information in the InkML files into grayscale images, and then carry out model training and testing.

The HME100K dataset [28] is a large-scale collection of real-scene handwritten mathematical expressions. It contains 74,502 training images and 24,607 testing images, making it significantly larger than similar datasets like CROHME. Notably, it features a wide range of real-world challenges such as variations in color, blur, and complex backgrounds, contributed by tens of thousands of writers. With 249 symbol classes, HME100K offers a diverse and realistic dataset for developing advanced handwritten mathematical expression recognition systems.

4.2 Evaluation Metrics

The Expression Recognition Rate (ExpRate) is the most commonly used evaluation metric for handwritten mathematical expression recognition. It is defined as the percentage of expressions that are correctly recognized out of the total number of expressions. Additionally, we use the metrics “ ≤ 1 error” and “ ≤ 2 error” to describe the performance of the model when we tolerate up to 1 or 2 token prediction errors, respectively, in the \LaTeX sequence.

4.3 Implementation Details

We employ DenseNet [9] as the visual encoder to extract visual features from expression images. The visual encoder utilizes 3 layers of DenseNet blocks, with each block containing 16 bottleneck layers. Between every two DenseNet blocks, a transition layer is used to reduce the spatial dimensions and the channel count of the visual features to half of their original sizes. The dropout rate is set at 0.2, and the growth rate of the model is established at 24.

We employ a 3-layer Transformer Decoder with ARM [22, 35] as the backbone of the decoder, with a model dimension (d_{model}) of 256 and head number set to 8. The dimension of the feed-forward layer is set to 1024, and the dropout rate is established at 0.3. The parameter settings for the Attention Rectifying Module (ARM) are consistent with those of CoMER, with the convolutional kernel size set to 5. Additionally, the parameters of masked self-attention and

FFN in Implicit Character Construction Module(ICCM) are consistent with the aforementioned Transformer Decoder.

During the training phase, we utilize Mini-batch Stochastic Gradient Descent(SGD) to learn the model parameters, with weight decay set to $1e-4$ and momentum set at 0.9. The initial learning rate is established at 0.08. We also adopt ReduceOnPlateau as the learning rate scheduler, whereby the learning rate is reduced to 25% of its original value when the ExpRate metric ceases to change. When trained on the CROHME dataset, the CROHME 2014 test set [14] is used as the validation set to select the model with the best performance. During the inference phase, we employ the approximate joint search method previously used in BTTR [36] to predict the output.

4.4 Comparison with State-of-the-art Methods

Table 1 presents the results on the CROHME dataset. To ensure fairness in performance comparison and considering that different methods have used various data augmentation techniques, many of which have not been disclosed, we have limited our comparison to results without the application of data augmentation. Given the relative small size of the CROHME dataset, we conducted experiments with both the baseline CoMER and the proposed ICAL model using five different random seeds (7, 77, 777, 7777, 77777) under the same experimental conditions. The reported results are the averages and standard deviations of these five experiments.

It is noteworthy that while GCN [34] has attained impressive results on CROHME 2016 [15] and 2019 [13], its performance benefits from the additional introduction of category information, while ICAL constructs the implicit character sequence directly from the \LaTeX sequence, eliminating the need for manually constructing additional category information. Consequently, we present the performance of the GCN solely for reference purposes and exclude it from direct comparisons. The CoMER model represents the current state-of-the-art(SOTA) method; however, CoMER [35] did not disclose their results without data augmentation in the original paper. Therefore, we have reproduced the results of CoMER without data augmentation using their open-source code, denoted by † in the table 1.

As shown in table 1, our method achieved the best performance across all metrics. Our method outperforms CoMER by 2.25%/1.81%/1.39% on the CROHME 2014, 2016, and 2019 datasets, respectively. Across all metrics, the ICAL method achieves an average improvement of 1.6% over the CoMER method. The experimental results on CROHME dataset prove the effectiveness of our method.

We also conducted experiments on the challenging and real-world dataset HME100K, as shown in table 3. We have replicated the performance of CoMER on the HME100K dataset, and our method surpasses the state-of-the-art(SOTA) by 0.94%, reaching an impressive 69.06%. The outstanding experimental performance on the HME100K dataset, which is more complex, larger, and more real-

Table 1. Performance comparison on the CROHME dataset. We compare expression recognition rate (ExpRate) between our model and previous state-of-the-art models on the CROHME 2014/2016/2019 test sets. **None of the methods used data augmentation to ensure a fair comparison.** We denote our reproduced results with †. The symbol * signifies the inclusion of supplementary information. All the performance results are reported in percentage (%).

Method	CROHME 2014			CROHME 2016			CROHME 2019		
	ExpRate↑	≤ 1↑	≤ 2↑	ExpRate↑	≤ 1↑	≤ 2↑	ExpRate↑	≤ 1↑	≤ 2↑
WAP	46.55	61.16	65.21	44.55	57.10	61.55	-	-	-
DenseWAP	50.1	-	-	47.5	-	-	-	-	-
DenseWAP-MSA	52.8	68.1	72.0	50.1	63.8	67.4	47.7	59.5	63.3
TAP*	48.47	63.28	67.34	44.81	59.72	62.77	-	-	-
PAL	39.66	56.80	65.11	-	-	-	-	-	-
PAL-v2	48.88	64.50	69.78	49.61	64.08	70.27	-	-	-
WS-WAP	53.65	-	-	51.96	64.34	70.10	-	-	-
ABM	56.85	73.73	81.24	52.92	69.66	78.73	53.96	71.06	78.65
CAN-DWAP	57.00	74.21	80.61	56.06	71.49	79.51	54.88	71.98	79.40
CAN-ABM	57.26	74.52	82.03	56.15	72.71	80.30	55.96	72.73	80.57
DenseWAP-TD	49.1	64.2	67.8	48.5	62.3	65.3	51.4	66.1	69.1
TDv2	53.62	-	-	55.18	-	-	58.72	-	-
SAN	56.2	72.6	79.2	53.6	69.6	76.8	53.5	69.3	70.1
BTTR	53.96	66.02	70.28	52.31	63.90	68.61	52.96	65.97	69.14
GCN*	60.00	-	-	58.94	-	-	61.63	-	-
CoMER†	58.38 ± 0.62	74.48 ± 1.41	81.14 ± 0.91	56.98 ± 1.41	74.44 ± 0.93	81.87 ± 0.73	59.12 ± 0.43	77.45 ± 0.70	83.87 ± 0.80
ICAL	60.63± 0.61	75.99± 0.77	82.80± 0.40	58.79± 0.73	76.06± 0.37	83.38± 0.16	60.51± 0.71	78.00± 0.66	84.63± 0.45

istic compared to CROHME, further proves the superior generalization ability and effectiveness of the ICAL method.

Table 2. Performance comparison on the HME100K dataset. We compare our proposed ICAL with previous models on HME100K. We denote our reproduced results with †. All the performance results are reported in percentage (%).

Method	HME100K		
	ExpRate↑	≤ 1↑	≤ 2↑
DenseWAP	61.85	70.63	77.14
DenseWAP-TD	62.60	79.05	85.67
ABM	65.93	81.16	87.86
SAN	67.1	-	-
CAN-DWAP	67.31	82.93	89.17
CAN-ABM	68.09	83.22	89.91
BTTR	64.1	-	-
CoMER†	68.12	84.20	89.71
ICAL	69.06± 0.16	85.16± 0.13	90.61± 0.09

4.5 Ablation Study

Our research includes a series of ablation studies to corroborate the effectiveness of our proposed method. Presented in Table 4, *Initial Loss* refers to the cross-entropy loss computed from the discrepancy between the Transformer Decoder’s direct output and the ground truth (the intact L^AT_EX code). *Fusion Loss* is the

Table 3. Performance comparison on the HME100K dataset. We compare our proposed ICAL with previous models on HME100K. We denote our reproduced results with †. All the performance results are reported in percentage (%).

Method	HME100K		
	ExpRate	↑ ≤ 1	↑ ≤ 2
DenseWAP	61.85	70.63	77.14
DenseWAP-TD	62.60	79.05	85.67
ABM	65.93	81.16	87.86
SAN	67.1	-	-
CAN-DWAP	67.31	82.93	89.17
CAN-ABM	68.09	83.22	89.91
BTTR	64.1	-	-
CoMER†	68.12	84.20	89.71
ICAL	69.06	85.16	90.61

cross-entropy loss determined by comparing the combined outputs from both the Implicit Character Construction Module (ICCM) and Decoder—synergized through the Fusion Module—with the ground truth (the intact \LaTeX sequence). *Implicit Loss* is calculated using a cross-entropy formula with adaptive weighting, which evaluates the discrepancies between the ICCM’s output and the sequence of implicit characters. When *Implicit Loss* is not applied, the ICCM module is also omitted, serving as a validation of ICCM’s effectiveness.

Additionally, it should be noted that in the ablation experiments mentioned above, when utilizing *Fusion Loss*, we employ the output of the Fusion Module for inference. Conversely, when *Fusion Loss* is not implemented, inference is conducted directly using the output from the Transformer Decoder.

From the 4th and 5th rows of each dataset in the table, it is evident that compared to using only the *Initial Loss* (1st row), which serves as the baseline, CoMER, *Implicit Loss* (4th row) and the *Fusion Loss* (5th row) that we have introduced can both effectively enhance the model’s recognition performance. We also designed experiments that exclusively utilize *Fusion Loss* and *Implicit Loss* (the 3rd row), from which it can be observed that, relative to the baseline approach, our method still manages to achieve a notable improvement in effectiveness.

Due to time limitations, the ablation study conducted on the HME100K dataset only used a single random seed for each experiment, and thus, the results are not averaged over multiple runs. This may introduce some variability in the reported performance, and we plan to address this limitation by conducting additional experiments with multiple seeds in future work for more robust evaluation.

4.6 Inference Speed

As shown in Table 5, we have evaluated the inference speed of our method on a single NVIDIA 2080Ti GPU. Compared to the baseline model, our method has

Table 4. Ablation study on the CROHME 2014/2016/2019 and HME100K test sets(in %). It should be noted that, if *Implicit Loss* is not applied, the ICCM module is also not used, which serves to validate the effectiveness of the ICCM. Similarly, when *Fusion Loss* is not implemented, inference is conducted directly using the output from the Transformer Decoder.

Dataset	Initial Loss	Fusion Loss	Implicit Loss	ExpRate
CROHME 2014	✓			58.38
		✓		58.52
		✓	✓	59.02
	✓	✓	✓	59.25
	✓	✓	✓	60.04
				60.63
CROHME 2016	✓			56.98
		✓		57.04
		✓	✓	58.13
	✓	✓	✓	57.80
	✓	✓	✓	58.44
				58.79
CROHME 2019	✓			59.12
		✓		59.15
		✓	✓	59.70
	✓	✓	✓	60.40
	✓	✓	✓	59.61
				60.51
HME100K	✓			68.12
		✓		68.18
		✓	✓	68.47
	✓	✓	✓	68.46
	✓	✓	✓	69.09
				69.25

a modest increase in the number of parameters with a negligible impact on FPS, and there is also a slight increase in FLOPs.

4.7 Case Study

We provide several typical recognition examples to demonstrate the effectiveness of the proposed method, as shown in Fig. 3. Entries highlighted in red indicate cases where the model made incorrect predictions. 'ICCM' represents the implicit character sequence predicted by the ICCM module, where $\langle s \rangle$ is the abbreviation for the $\langle \text{space} \rangle$ token.

In Case (a), the baseline CoMER model incorrectly identified the first character, π , as $y_{\{0\}}$. In contrast, the ICCM correctly determined that there were no implicit characters in the formula, as indicated on the fourth line of Group a. This accurate detection allowed for the correct \LaTeX sequence to be output by ICAL, as shown on the third line of Group a.

In Case (b), the ICCM's prediction of the implicit character sequence (fourth line of Group b) was crucial. It enabled the ICAL method to correctly place both

Table 5. Comparative Analysis of Parameters (Params), Floating-Point Operations (FLOPs), and Frames Per Second (FPS)

Method	Input Image Size	Params (M)	FLOPs (G)	FPS
CoMER	(1,1,120,800)	6.39	18.81	2.484
ICAL	(1,1,120,800)	7.37	19.81	2.394

characters 3 and 4 in the subscript of q (third line of Group b), unlike CoMER, which misidentified this relationship (second line of Group b).

Case (c) demonstrates that the ICCM’s prediction of implicit characters can also alleviate the lack of coverage issue [20, 33]. Here, ICCM’s prediction indicated that there should be two explicit characters, 9 and 1, within $\{ \}$, and correspondingly, ICAL also successfully predicted these two characters. However, CoMER only predicted the character 9, and missed the character 1.

Moreover, Case (d) also effectively highlights how ICCM and its ability to predict implicit characters can enhance a model’s understanding of the structural relationships within formulas.



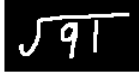

(a)		Ground Truth	$\pi d = 2 \pi r$
		CoMER	$y_{\{0\}} d = 2 \pi r$
		ICAL	$\pi d = 2 \pi r$
		ICCM	$\langle s \rangle \langle s \rangle \langle s \rangle \langle s \rangle \langle s \rangle \langle s \rangle$
(b)		Ground Truth	$q_{\{1\}} + q_{\{2\}} + q_{\{3\}} + q_{\{4\}} = 2$
		CoMER	$q_{\{1\}} + q_{\{2\}} + q_{\{3\}} + q_{\{4\}} = 2$
		ICAL	$q_{\{1\}} + q_{\{2\}} + q_{\{3\}} + q_{\{4\}} = 2$
		ICCM	$\langle s \rangle \{ \langle s \rangle \} \langle s \rangle \langle s \rangle \{ \langle s \rangle \} \langle s \rangle \langle s \rangle \{ \langle s \rangle \} \langle s \rangle \langle s \rangle \{ \langle s \rangle \} \langle s \rangle \langle s \rangle$
(c)		Ground Truth	$\sqrt{91}$
		CoMER	$\sqrt{9}$
		ICAL	$\sqrt{91}$
		ICCM	$\langle s \rangle \{ \langle s \rangle \langle s \rangle \}$
(d)		Ground Truth	$b_{\{n\}} = \lim_{\alpha \rightarrow 0} b_{\{n - \alpha\}}$
		CoMER	$b_{\{n\}} = \lim_{\alpha \rightarrow 0} b_{\{n\}} - \alpha$
		ICAL	$b_{\{n\}} = \lim_{\alpha \rightarrow 0} b_{\{n - \alpha\}}$
		ICCM	$\langle s \rangle \{ \langle s \rangle \} \langle s \rangle \langle s \rangle \{ \langle s \rangle \langle s \rangle \} \langle s \rangle \{ \langle s \rangle \langle s \rangle \} \langle s \rangle \{ \langle s \rangle \langle s \rangle \}$

Fig. 3. Case studies for the Ground Truth and CoMER, ICAL methods. The red symbols represent incorrect predictions. ‘ICCM’ represents the implicit character sequence predicted by the ICCM module, where $\langle s \rangle$ is the abbreviation for the $\langle \text{space} \rangle$ token.

5 Conclusion

In this paper, we propose a novel recognizer framework, ICAL, capable of leveraging global information in \LaTeX to correct the predictions of the decoder. Our main contributions are threefold: (1) We have designed an Implicit Character Construction Module (ICCM) to predict implicit characters in \LaTeX . (2) Additionally, we employ a Fusion Module to aggregate global information from implicit characters, thereby refining the predictions of the Transformer Decoder. We integrate these two modules into the CoMER model to develop our method, ICAL. (3) Experimental results demonstrate that the ICAL method surpasses previous state-of-the-art approaches, achieving expression recognition rate (ExpRate) of 60.63%, 58.79%, and 60.51% on the CROHME 2014, 2016, and 2019 datasets, respectively, and an ExpRate of 69.06% on the HME100K dataset.

Acknowledgements

This work is supported by the projects of National Science and Technology Major Project (2021ZD0113301) and National Natural Science Foundation of China (No. 62376012), which is also a research achievement of Key Laboratory of Science, Technology and Standard in Press Industry (Key Laboratory of Intelligent Press Media Technology).

References

1. Alvaro, F., Sánchez, J.A., Benedí, J.M.: Recognition of on-line handwritten mathematical expressions using 2d stochastic context-free grammars and hidden markov models. *Pattern Recognition Letters* **35**, 58–67 (2014)
2. Bian, X., Qin, B., Xin, X., Li, J., Su, X., Wang, Y.: Handwritten mathematical expression recognition via attention aggregation based bi-directional mutual learning. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 36, pp. 113–121 (2022)
3. Chan, K.F., Yeung, D.Y.: Elastic structural matching for online handwritten alphanumeric character recognition. In: *Proceedings. Fourteenth International Conference on Pattern Recognition (Cat. No. 98EX170)*. vol. 2, pp. 1508–1511. IEEE (1998)
4. Chan, K.F., Yeung, D.Y.: An efficient syntactic approach to structural analysis of on-line handwritten mathematical expressions. *Pattern recognition* **33**(3), 375–384 (2000)
5. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078* (2014)
6. Chou, P.A.: Recognition of equations using a two-dimensional stochastic context-free grammar. In: *Visual Communications and Image Processing IV*. vol. 1199, pp. 852–865. SPIE (1989)
7. Ding, H., Chen, K., Huo, Q.: An encoder-decoder approach to handwritten mathematical expression recognition with multi-head attention and stacked decoder. In: *Document Analysis and Recognition-ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part II* 16. pp. 602–616. Springer (2021)
8. Ha, J., Haralick, R.M., Phillips, I.T.: Understanding mathematical expressions from document images. In: *Proceedings of 3rd International Conference on Document Analysis and Recognition*. vol. 2, pp. 956–959. IEEE (1995)
9. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4700–4708 (2017)
10. Keshari, B., Watt, S.: Hybrid mathematical symbol recognition using support vector machines. In: *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*. vol. 2, pp. 859–863. IEEE (2007)
11. Kosmala, A., Rigoll, G., Lavirotte, S., Pottier, L.: On-line handwritten formula recognition using hidden markov models and context dependent graph grammars. In: *Proceedings of the Fifth International Conference on Document Analysis and Recognition. ICDAR’99 (Cat. No. PR00318)*. pp. 107–110. IEEE (1999)
12. Li, B., Yuan, Y., Liang, D., Liu, X., Ji, Z., Bai, J., Liu, W., Bai, X.: When counting meets hmer: counting-aware network for handwritten mathematical expression recognition. In: *European Conference on Computer Vision*. pp. 197–214. Springer (2022)
13. Mahdavi, M., Zanibbi, R., Mouchere, H., Viard-Gaudin, C., Garain, U.: Icdar 2019 crohme+ tfd: Competition on recognition of handwritten mathematical expressions and typeset formula detection. In: *2019 International Conference on Document Analysis and Recognition (ICDAR)*. pp. 1533–1538 (2019)
14. Mouchere, H., Viard-Gaudin, C., Zanibbi, R., Garain, U.: Icfhr 2014 competition on recognition of on-line handwritten mathematical expressions (crohme 2014). In:

- 2014 14th International Conference on Frontiers in Handwriting Recognition. pp. 791–796 (2014)
15. Mouchère, H., Viard-Gaudin, C., Zanibbi, R., Garain, U.: Icfhr2016 crohme: Competition on recognition of online handwritten mathematical expressions. In: 2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR). pp. 607–612 (2016)
 16. Sakai, I.: Syntax in universal translation. In: Proceedings of the International Conference on Machine Translation and Applied Language Analysis (1961)
 17. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
 18. Toyota, S., Uchida, S., Suzuki, M.: Structural analysis of mathematical formulae with verification based on formula description grammar. In: Document Analysis Systems VII: 7th International Workshop, DAS 2006, Nelson, New Zealand, February 13–15, 2006. Proceedings 7. pp. 153–163. Springer (2006)
 19. Truong, T.N., Nguyen, C.T., Phan, K.M., Nakagawa, M.: Improvement of end-to-end offline handwritten mathematical expression recognition by weakly supervised learning. In: 2020 17th International Conference on Frontiers in Handwriting Recognition (ICFHR). pp. 181–186. IEEE (2020)
 20. Tu, Z., Lu, Z., Liu, Y., Liu, X., Li, H.: Modeling coverage for neural machine translation. arXiv preprint arXiv:1601.04811 (2016)
 21. Twaakyondo, H., Okamoto, M.: Structure analysis and recognition of mathematical expressions. In: Proceedings of 3rd International Conference on Document Analysis and Recognition. vol. 1, pp. 430–437 vol.1 (1995). <https://doi.org/10.1109/ICDAR.1995.599029>
 22. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30**, 5998–6008 (2017)
 23. Vuong, B.Q., He, Y., Hui, S.C.: Towards a web-based progressive handwriting recognition environment for mathematical problem solving. *Expert Systems with Applications* **37**(1), 886–893 (2010)
 24. Winkler, H.J.: Hmm-based handwritten symbol recognition using on-line and off-line features. In: 1996 IEEE International Conference on Acoustics, Speech, and Signal Processing Conference Proceedings. vol. 6, pp. 3438–3441. IEEE (1996)
 25. Wu, C., Du, J., Li, Y., Zhang, J., Yang, C., Ren, B., Hu, Y.: Tdv2: A novel tree-structured decoder for offline mathematical expression recognition. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 2694–2702 (2022)
 26. Wu, J.W., Yin, F., Zhang, Y.M., Zhang, X.Y., Liu, C.L.: Image-to-markup generation via paired adversarial learning. In: Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2018, Dublin, Ireland, September 10–14, 2018, Proceedings, Part I 18. pp. 18–34. Springer (2019)
 27. Wu, J.W., Yin, F., Zhang, Y.M., Zhang, X.Y., Liu, C.L.: Handwritten mathematical expression recognition via paired adversarial learning. *International Journal of Computer Vision* pp. 1–16 (2020)
 28. Yuan, Y., Liu, X., Dikubab, W., Liu, H., Ji, Z., Wu, Z., Bai, X.: Syntax-aware network for handwritten mathematical expression recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4553–4562 (2022)
 29. Zanibbi, R., Blostein, D., Cordy, J.R.: Recognizing mathematical expressions using tree transformation. *IEEE Transactions on pattern analysis and machine intelligence* **24**(11), 1455–1467 (2002)

30. Zhang, J., Du, J., Dai, L.: Multi-scale attention with dense encoder for handwritten mathematical expression recognition. In: 2018 24th international conference on pattern recognition (ICPR). pp. 2245–2250 (2018)
31. Zhang, J., Du, J., Dai, L.: Track, attend, and parse (tap): An end-to-end framework for online handwritten mathematical expression recognition. *IEEE Transactions on Multimedia* **21**(1), 221–233 (2018)
32. Zhang, J., Du, J., Yang, Y., Song, Y.Z., Wei, S., Dai, L.: A tree-structured decoder for image-to-markup generation. In: ICML. p. In Press (2020)
33. Zhang, J., Du, J., Zhang, S., Liu, D., Hu, Y., Hu, J., Wei, S., Dai, L.: Watch, attend and parse: An end-to-end neural network based approach to handwritten mathematical expression recognition. *Pattern Recognition* **71**, 196–206 (2017)
34. Zhang, X., Ying, H., Tao, Y., Xing, Y., Feng, G.: General category network: Handwritten mathematical expression recognition with coarse-grained recognition task. In: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1–5. IEEE (2023)
35. Zhao, W., Gao, L.: Comer: Modeling coverage for transformer-based handwritten mathematical expression recognition. In: Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII. pp. 392–408. Springer (2022)
36. Zhao, W., Gao, L., Yan, Z., Peng, S., Du, L., Zhang, Z.: Handwritten mathematical expression recognition with bidirectionally trained transformer. In: Document Analysis and Recognition–ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part II 16. pp. 570–584. Springer (2021)