# UniMERNet: A Universal Network for Real-World Mathematical Expression Recognition

**Bin Wang**[*1], **Zhuangcheng Gu**[*1], **Guang Liang**[*1],
**Chao Xu** [1], **Bo Zhang** [1], **Botian Shi** [1], **Conghui He**[†]

[1]Shanghai AI Laboratory,

## Abstract

The paper introduces the UniMER dataset, marking the first study on Mathematical Expression Recognition (MER) targeting complex real-world scenarios. The UniMER dataset includes a large-scale training set, UniMER-1M, which offers unprecedented scale and diversity with one million training instances to train high-quality, robust models. Additionally, UniMER features a meticulously designed, diverse test set, UniMER-Test, which covers a variety of formula distributions found in real-world scenarios, providing a more comprehensive and fair evaluation. To better utilize the UniMER dataset, the paper proposes a Universal Mathematical Expression Recognition Network (UniMERNet), tailored to the characteristics of formula recognition. UniMERNet consists of a carefully designed encoder that incorporates detail-aware and local context features, and an optimized decoder for accelerated performance. Extensive experiments conducted using the UniMER-1M dataset and UniMERNet demonstrate that training on the large-scale UniMER-1M dataset can produce a more generalizable formula recognition model, significantly outperforming all previous datasets. Furthermore, the introduction of UniMERNet enhances the model's performance in formula recognition, achieving higher accuracy and speeds. All data, models, and code are available at https://github.com/opendatalab/UniMERNet.

## Introduction

Mathematical Expression Recognition (MER) is a critical task in document analysis, aiming to convert image-based mathematical expressions into corresponding markup languages such as LaTeX or Markdown. MER is essential in applications like scientific document extraction, where a robust MER model helps maintain the logical coherence of documents. Unlike typical Optical Character Recognition (OCR) tasks, MER requires a deeper understanding of complex structures, including superscripts, subscripts, and various special symbols.

Existing research has primarily focused on enhancing the recognition accuracy of relatively simple rendered expressions (Deng et al. 2017) and handwritten data (Mahdavi et al. 2019; Le, Indurkhya, and Nakagawa 2019; Wu et al. 2020; Zhao et al. 2021) through a series of MER algorithms. Some

*Equal contribution.
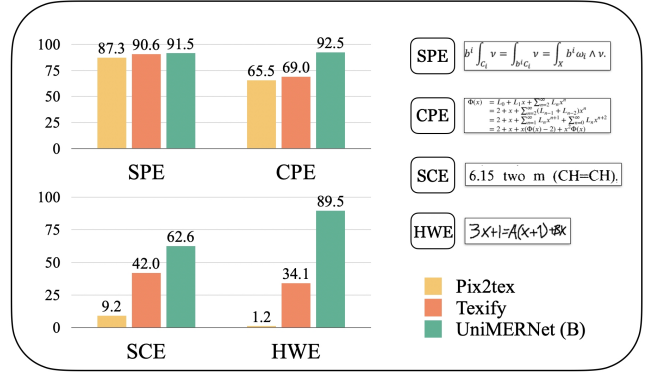†Corresponding author (heconghui@pjlab.org.cn).



Figure 1: Performance comparison (BLEU Score) of mainstream models and UniMERNet in recognizing real-world mathematical expressions: Evaluation across Simple Printed Expressions (SPE), Complex Printed Expressions (CPE), Screen-Captured Expressions (SCE), and Handwritten Expressions (HWE).

researchers have begun to optimize MER algorithms by scaling up the training data and integrating them with transformer models (Vaswani et al. 2017), ensuring their applicability in diverse scenarios (Kim et al. 2022; Blecher 2022; Blecher et al. 2023; Paruchuri 2023). Other researchers have attempted to directly employ Large Vision-Language Models (LVLMs) for document content extraction, including MER (Wei et al. 2023; Blecher et al. 2023). However, existing MER benchmarks (Deng et al. 2017; Mahdavi et al. 2019) primarily focus on simple printed or handwritten expressions. Consequently, these models often struggle with diverse real-world expressions, such as lengthy equations and noisy scanned document screenshots.

In practice, real-world scenarios require the handling of complex, long expressions and noisy, distorted images from scanned documents or webpage screenshots. To fill this gap, we introduce a comprehensive benchmark, UniMER-Test, which extends the existing test set with longer and real-world scenario expressions. Our benchmark aims to stimulate progress in MER by focusing on robustness and practical usage. As depicted in Figure 1, we conduct exhaustive evaluations of state-of-the-art MER methods (Blecher 2022; Paruchuri 2023) using our novel benchmark, UniMER-Test. These methods demonstrate remarkable competence in rec-

ognizing simple printed expressions. However, their performance noticeably declines when tested with more complex printed expressions, particularly long formulas. The performance degradation becomes even more pronounced when these methods are applied to real-world expressions, such as screen-captured expressions embedded in noisy backgrounds and handwritten expressions. Moreover, large vision-language models such as Nougat (Blecher et al. 2023) and Vary (Wei et al. 2023), despite their capacity for convenient end-to-end document content extraction, exhibit only mediocre performance in MER.

To train a high-quality formula recognition model capable of accurately predicting results in diverse scenarios, we have constructed the UniMER-1M dataset. This large-scale dataset is specifically designed for Mathematical Expression Recognition (MER) and includes over one million diverse formula Image-LaTeX pairs. During its construction, we considered various levels of formula complexity, ranging from simple to complex long formulas, as well as different types, including printed and handwritten formulas. This ensures the dataset's suitability for training a model that generalizes well to real-world scenarios. Furthermore, to fully leverage the UniMER dataset, we propose an innovative formula recognition model—UniMERNet. Unlike the mainstream document recognition frameworks that directly use the Swin-Transformer encoder and mBART decoder, we have optimized the model structure specifically for the formula recognition task. In the encoder, we introduce the Fine-Grained Embedding (FGE) module and the Convolutional Enhancement (CE) module for local context awareness. In the decoder, we incorporate the Squeeze Attention (SA) module to accelerate inference. These enhancements result in significant improvements in both inference speed and accuracy.

The main contributions of this paper are as follows:

- We introduce **UniMER**[1] (He et al. 2024), a universal MER dataset, with the training set UniMER-1M and the test set UniMER-Test, which encompasses all types of expressions in practical situations, offering a diverse and comprehensive foundation for MER model development and evaluation.

- We propose a novel network structure, **UniMERNet**, specifically designed for the formula recognition task. By designing a more precise encoder and a faster decoder, we can freely combine models to achieve higher accuracy and faster speed in formula recognition.

- Validation of UniMERNet's superior performance through extensive experiments, establishing it as the new benchmark in open-source MER solutions by outperforming existing models in a variety of scenarios.

## Related Work

### Traditional Machine Learning Methods in MER

Decades ago, researchers recognized the importance of Mathematical Expression Recognition (MER). Anderson (Anderson 1967) pioneered MER in irregular documents

by introducing a parsing algorithm for two-dimensional character configurations. Miller and Viola (Miller and Viola 1998) proposed a system integrating character segmentation with the grammar of mathematical layouts. Chan *et al.* (Chan and Yeung 1999) developed an online MER system featuring error detection and correction mechanisms. INFTY (Suzuki et al. 2003) presented an OCR system for mathematical documents that achieved high character recognition accuracy through novel techniques. However, despite these advancements, MER precision was limited by hand-crafted features in traditional machine learning.

### Deep Learning and Transformer Methods in MER

With the advent of deep learning, various MER algorithms based on Convolutional Neural Networks (CNN)(Krizhevsky, Sutskever, and Hinton 2012; Simonyan and Zisserman 2015) were proposed. Deng *et al.*(Deng et al. 2017) introduced an encoder-decoder model with a coarse-to-fine attention mechanism, demonstrating superior performance over traditional OCR systems using the IM2LATEX-100K dataset. The WAP model (Zhang et al. 2017) autonomously learned mathematical grammar and symbol segmentation, aligning closely with human intuition, while the PAL-v2 model (Wu et al. 2020) used paired adversarial learning to excel in handwritten expression recognition on the CROHME dataset. Zhang *et al.*(Zhang et al. 2020) proposed a tree-structured decoder for complex markups, and Zhao *et al.*(Zhao et al. 2021) and Bian *et al.*(Bian et al. 2022) enhanced MER with bi-directional learning in encoder-decoder models, advancing Handwritten Mathematical Expression Recognition (HMER). The CAN model(Li et al. 2022) improved HMER by incorporating a weakly supervised counting module, while Le *et al.*(Le, Indurkhya, and Nakagawa 2019) and Li *et al.*(Li et al. 2020) employed data augmentation strategies to enhance MER performance.

More recently, the rapid development of Transformer models (Vaswani et al. 2017) and large vision-language models (Zhu et al. 2023; Liu et al. 2024; Dong et al. 2024; Liu et al. 2023; Wang et al. 2024; Zhang et al. 2024; Chen et al. 2024) led researchers to explore document information extraction task based on meticulously constructed evaluation benchmarks, such as DocGenome (Xia et al. 2024) and MM-Sci (Li et al. 2024). For example, Donut (Kim et al. 2022) introduced an end-to-end model that converts document images into structured outputs without relying on OCR, while Nougat (Blecher et al. 2023) utilized auto-generated image-to-markup samples to train a Transformer-based encoder-decoder model. Vary (Wei et al. 2023) offered a fine-grained multimodal model for document parsing. However, these methods often overlook the unique characteristics of mathematical expressions, leading to limitations in their MER capabilities. To address this, Pix2tex (Blecher 2022) and Texify (Paruchuri 2023) trained encoder-decoder models on rendered mathematical expressions, though they struggle with complex or noisy expressions.

In response to these challenges, the UniMERNet model proposed in this paper aims to build a robust and practical MER model that not only achieves state-of-the-art per-

---

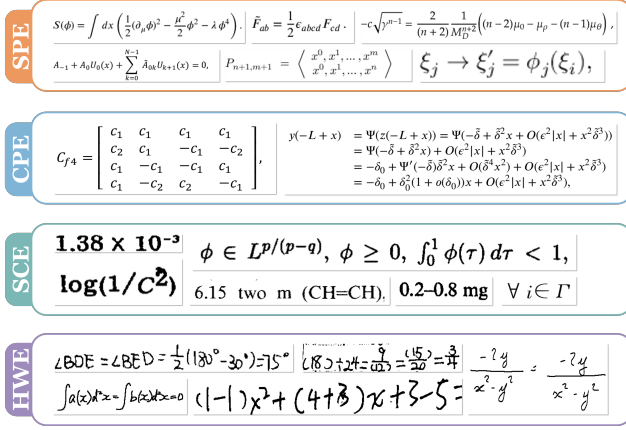[1] https://opendatalab.com/OpenDataLab/UniMER-Dataset

Figure 2: Visualization of the UniMER-Test dataset with four data types: Simple Printed Expressions (SPE), Complex Printed Expressions (CPE), Screen Capture Expressions (SCE), and Handwritten Expressions (HWE).

formance but also optimizes inference speed, enhancing the model's applicability in real-world scenarios.

## UniMER Dataset

The UniMER dataset addresses the diversity of formula recognition challenges in real-world scenarios. It consists of two main components: UniMER-1M, a large-scale training set, and UniMER-Test, a comprehensive evaluation set.

UniMER-1M includes 1,061,791 latex-image pairs, covering both simple and complex printed and handwritten formulas (Table 1). This extensive dataset surpasses existing formula recognition training sets, enabling the development of more robust models.

UniMER-Test, on the other hand, is a test set containing 23,789 samples. Unlike existing evaluation sets that primarily focus on simple printed and handwritten formulas, UniMER-Test comprehensively evaluates formula recognition across varying complexities and types, reflecting real-world scenarios (Figure 2). Specifically, UniMER-Test includes the following types of formulas:

- **SPE**: Formula images rendered from simple LaTeX expressions, characterized by uniform font size, clean background, and relatively short formulas.
- **CPE**: Formula images rendered from complex, long LaTeX expressions, characterized by uniform font size, clean background, and longer, more intricate formulas.
- **SCE**: Screen-captured images of formulas from documents and the web, characterized by inconsistent fonts and sizes, background noise, and image deformation.
- **HWE**: Collected from referenced handwriting recognition datasets (Mouchere et al. 2014; Mouchère et al. 2016; Mahdavi et al. 2019; Yuan et al. 2022), these are complex and diverse, with varying backgrounds, but are relatively short.

This comprehensive approach ensures that UniMER-Test serves as a robust benchmark for evaluating formula recognition systems, setting a new standard for future research and development in the field.

## Data Collection Process

**Printed Rendered Expressions (SPE, CPE)** The assembly of our dataset begins with the Pix2tex (Blecher 2022) public dataset, which serves as the base for our SPE. Due to the limitations in volume and complexity, we expand the dataset by sourcing additional LaTeX expression codes from platforms like Arxiv, Wikipedia, and StackExchange. These codes are regularized (Deng et al. 2017) to resolve LaTeX syntax ambiguities, then compiled into expression PDFs in various fonts using XeLaTeX. Uncompilable expressions are discarded. Subsequently, ImageMagic's conversion function is utilized to transform these images into expressions with multiple DPIs, with data balancing ensuring an even distribution of different lengths.

Following this data expansion pipeline, we sample 725,246 simple formulas from the augmented data and combine them with the Pix2tex training set to form the SPE training data. The Pix2tex test set is designated as the SPE test data. In contrast, the CPE is derived independently of the Pix2tex dataset. We randomly select 110,332 complex formulas from the expanded data for training and test sets.

**Screen-Captured Expressions (SCE)** For SCE, we compile 1,000 diverse PDF pages in both Chinese and English, covering books, papers, textbooks, magazines, and newspapers. This variety ensures a wide range of fonts, sizes, and backgrounds for the formulas. Two annotators identify and label the formula boxes in the documents, capturing the content automatically. This process produces over 6,000 formula boxes, which are processed through Mathpix for formula recognition. After manual corrections and cross-verification by two annotators, redundant formulas are removed, resulting in 4,744 unique mathematical expression for the SCE test set.

**Handwritten Expressions (HWE)** For HWE, we utilize the public datasets CROHME (Mouchere et al. 2014; Mouchère et al. 2016; Mahdavi et al. 2019) and HME100K (Yuan et al. 2022). CROHME, a well-known dataset in HMER, originates from the handwritten digit recognition competition and includes 8,836 training expressions and 3,332 test expressions. HME100K, a real-world handwritten expression dataset, provides 74,502 training and 24,607 test images. Due to the high annotation accuracy of these datasets, we combine them for our HWE data. Specifically, the HWE training set consists of 8,836 formulas from CROHME and 74,502 from HME100K, totaling 83,338 samples. The HWE test set includes 3,332 formulas from CROHME and 3,000 from HME100K, totaling 6,332 test formulas.

## Diversified Training Data Sampling

Existing formula datasets, such as HWE (CHROME & HME100K) (Mouchere et al. 2014; Mouchère et al. 2016; Mahdavi et al. 2019), IM2LATEX (Deng et al. 2017), and Pix2tex (Blecher 2022), primarily consist of rendered and

| Dataset | Type | Train Size | Test Size | Max Len | Avg Len |
|---|---|---|---|---|---|
| HME100K | HWE | 74,502 | 24,607 | 311 | 24.05 |
| CROHME | | 8,836 | 3233 | 147 | 22.27 |
| IM2LATEX-100K | SPE | 83,883 | 10,354 | 440 | 96.01 |
| Pix2tex | | 158,480 | 30,637 | 2949 | 93.35 |
| UniMER | Mixed | 1,061,791 | 23,757 | 7037 | 79.48 |

Table 1: Statistical comparison of the MER dataset. "Max Len" and "Avg Len" mean the maximum length and average string length of the mathematical expression.
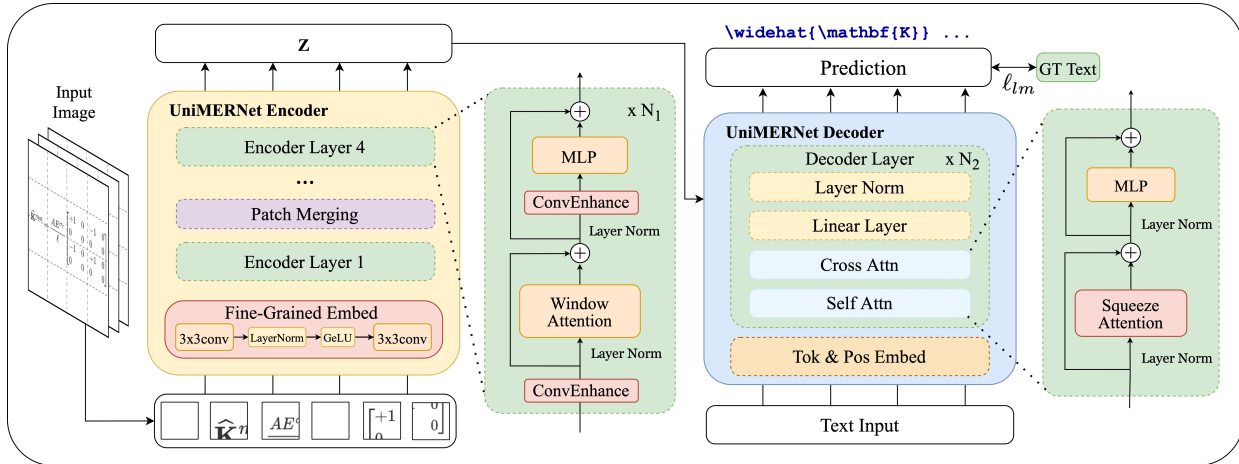


Figure 3: The overall framework of UniMERNet. The UniMER-Encoder incorporates Fine-Grained Embedding (FGE), Convolutional Enhancement (CE), and Removal of Shift Window (RSW) to enhance recognition capabilities. The UniMER-Decoder employs Squeeze Attention (SA) to accelerate inference speed.

handwritten formulas, but they have limitations in formula length and complexity. For example, Pix2tex mostly contains regular formulas, lacking extremely short or complex long formulas, while handwritten formulas are generally short with diverse styles, none exceeding 256 characters.

To address these limitations, we expand our UniMER-1M dataset with a wider range of formulas, sampled from sources like Arxiv and Wikipedia to ensure a balanced distribution of lengths and complexity. This varied sampling strategy enhances the model's ability to recognize formulas across different complexity levels, improving overall performance. The formula length distribution of IM2LATEX, Pix2tex, HWE, and UniMER-1M datasets is shown in Figure 4.

## UniMERNet

In real-world scenarios, mathematical formulas come from diverse sources such as electronic documents, scanned images, screenshots, and photographs. They range from single symbols to complex, lengthy expressions. Unlike general text recognition, formula recognition poses unique challenges in three dimensions. **Visual Similarity:** Many formula symbols look similar, e.g., $\mu$ and $u$, $\beta$ and $B$. This requires the model to have precise recognition capabilities. **Spatial Information:** Formulas often contain superscripts, subscripts, and other spatial arrangements, necessitating model's contextual awareness. **Inference Speed:** For

complex and lengthy formulas, the symbol generation based on an encoder-decoder structure can be time-consuming, slowing down the model's inference speed.

To address these challenges, we design the UniMER-Net formula recognition network. This network enhances recognition capabilities by incorporating fine-grained and context-aware modules in the encoder stage and accelerates inference speed by compressing the attention operation in the decoder stage.

Our architecture is based on the Swin-Transformer Encoder and mBART Decoder, which have been validated in various document processing tasks (Kim et al. 2022; Blecher et al. 2023; Paruchuri 2023). The overall framework of UniMERNet is shown in Figure 3.

During training, each input formula image $\mathbf{I} \in \mathbb{R}^{3 \times H_0 \times W_0}$ undergoes an image augmentation module, transforming a single image representation into a diverse set of images. This effectively handles the varied representations of formulas in real-world scenarios. The UniMERNet encoder processes the image to generate a feature vector $\mathbf{Z}$, which is then fed into the UniMERNet decoder. The decoder interacts with the feature vector $\mathbf{Z}$ and the output text sequence via a cross-attention mechanism to generate the predicted formula. The decoder combines the feature vector $\mathbf{Z}$, token embedding, and position embedding to predict the formula. For language modeling loss, we employ cross-entropy loss to minimize the difference between the predicted prob-
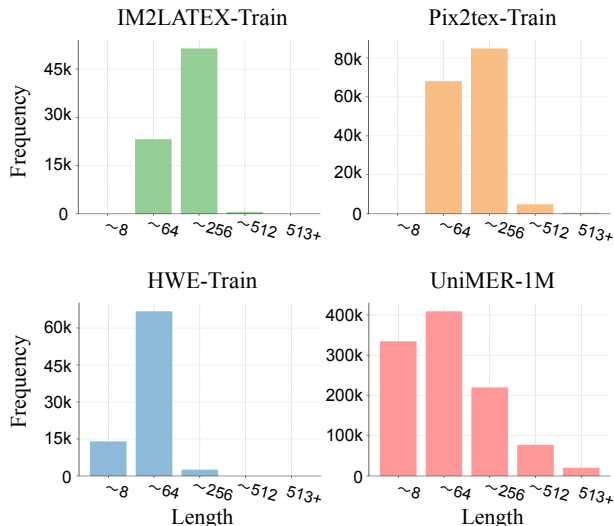
Figure 4: Formula string length distribution across datasets

ability distribution of the next token and the actual distribution observed in the training data. The loss is defined as:

$$\ell_{lm}(\hat{y}, y) = -\sum_{c=1}^{C} y_{o,c} \log(\hat{y}_{o,c}), \qquad (1)$$

Next, we detail the improvements in the UniMERNet network architecture tailored for formula recognition tasks.

### UniMERNet-Encoder

**Fine-Grained Embedding (FGE).** The original Swin-Transformer uses a $4 \times 4$ single-layer convolutional kernel with a stride of 4 for patch embedding, resulting in a fourfold reduction in input resolution. However, these non-overlapping convolutional kernels often extract fragmented features, causing characters within formulas to be separated and adversely affecting recognition accuracy. Conversely, using larger overlapping convolutional kernels may lead to redundancy, as MER typically involves simpler elements like superscripts, subscripts, and various special characters, necessitating precise and streamlined feature extraction.

To address this, we implement an overlapping fine convolution layer composed of two convolutional layers, a Layer Normalization (LN) layer, and a GELU activation layer. Both convolutional layers utilize a kernel size of 3, a stride of 2, and padding of 1. This overlapping fine convolution expands the receptive field while mitigating fragmentation, thereby enhancing the model's overall performance.

**Convolutional Enhancement (CE).** While we retain the Window Attention mechanism from the Swin-Transformer, the receptive field within each window remains relatively large, limiting attention to small details such as labels in formulas and lacking inductive bias. Research in the visual domain (Guo et al. 2022) (Chu et al. 2021) (d'Ascoli et al. 2021) suggests that convolutional and transformer models are complementary. In formula recognition, both global and local information are crucial. Global information enables the model to discern where to pause and where a new line begins, while local information helps the model understand the relationships between adjacent characters, particularly for local relationships like superscripts and subscripts.

Transformer architectures capture global information well, but discerning local details is critical. To address this, we introduce a local perception module called ConvEnhance, which enhances the model's ability to identify superscripts and subscripts. ConvEnhance, placed before each attention and MLP module, consists of 3x3 depthwise convolutions and GELU activation functions. The convolutions provide local perception, while the GELU activation implements gated threshold filtering. This addition allows UniMERNet to alternate between local and global information, significantly improving its performance.

**Removal of Shift Window (RSW).** The Shift Window operation aims to increase the model's receptive field by overlapping windows between layers, similar to convolutional kernels. However, with the introduction of FGE and CE, the receptive field is already significantly enlarged compared to the original Swin architecture. This makes the Shift Window operation redundant. Removing it not only improves the model's performance but also speeds up the model.

### UniMERNet-Decoder

**Squeeze Attention (SA).** Experiments presented in the appendix indicate that the throughput bottleneck of UniMER-Net lies within the language model, mBART. This is primarily because the visual model retrieves all visual features in a single pass, while the language model must iteratively utilize the previous prediction results to generate each successive token. Consequently, the throughput limitation is attributable to the language model. To address this, we introduce the concept of bottleneck from ResNet, implementing Squeeze Attention. Specifically, Squeeze Attention maps the query and key to a lower-dimensional space without excessive loss of information, thereby accelerating the computation of attention.

The enhanced UniMERNet-Encoder captures fine-grained and local contextual information, improving formula recognition accuracy. The optimized UniMER-Decoder, with slightly reduced precision, achieves significant speed improvements, making it highly effective for practical formula recognition systems.

## Experiments

### Datasets and Evaluation Metrics

We utilize the UniMER-1M dataset to train our model and evaluate its formula recognition performance using the UniMER-Test. Our evaluation relies on BLEU, Edit Distance, and ExpRate metrics.

**BLEU:** The BLEU score (Papineni et al. 2002), initially developed for machine translation, quantifies the match of n-grams between candidate and reference sentences. Its application to a similar conversion task of formula recognition provides a robust, quantitative performance measure.

**Edit distance:** The Edit Distance (Levenshtein et al. 1966) measures the minimum character changes needed to convert

| Train Dataset | SPE | | CPE | | SCE | | HWE | |
|---|---|---|---|---|---|---|---|---|
| | BLEU ↑ | EditDis ↓ | BLEU ↑ | EditDis ↓ | BLEU ↑ | EditDis ↓ | BLEU ↑ | EditDis ↓ |
| Pix2tex | 0.911 | 0.063 | 0.773 | 0.194 | 0.527 | 0.371 | 0.067 | 0.800 |
| Pix2tex&HWE | 0.909 | 0.063 | 0.724 | 0.225 | 0.529 | 0.309 | 0.873 | 0.088 |
| UniMER-1M | **0.915** | **0.060** | **0.925** | **0.056** | **0.626** | **0.224** | **0.895** | **0.072** |

Table 2: Ablation results on UniMER-Test for models using different augmentations. Here, "HME" refers to a mixed dataset of CHROME and HME100K.

| FGE | CE | RSW | SA | Params (M) | FPS* (img/s) | SPE BLEU ↑ | CPE BLEU ↑ | SCE BLEU ↑ | HWE BLEU ↑ |
|---|---|---|---|---|---|---|---|---|---|
| | | | | 342 | 4.12 | 0.903 | 0.885 | 0.579 | 0.887 |
| ✔ | | | | 342 | 4.10 | 0.903 | 0.888 | 0.584 | 0.886 |
| ✔ | ✔ | | | 342 | 4.07 | 0.912 | 0.896 | 0.599 | 0.895 |
| ✔ | ✔ | | ✔ | **325** | 5.04 | 0.911 | 0.894 | 0.599 | 0.893 |
| ✔ | ✔ | ✔ | ✔ | **325** | **5.06** | **0.912** | **0.897** | **0.601** | **0.893** |

Table 3: Ablation study of model architecture. **FGE**: replace with Fine-Grained Embedding, **CE**: add ConvEnhance module, **RSW**: Remove Shift Window, **SA**: apply Squeeze Attention in mBART decoder. **FPS***: The model's throughput is tested on an A100 GPU with a batch size of 128, processing each sample up to a maximum sequence length of 1536.

one string to another. Its use in formula recognition offers a precise, character-level accuracy assessment, making it a valuable performance metric.

**ExpRate:** Expression Recognition Rate (Yuan et al. 2022) is a widely used metric for handwritten formula recognition, defined as the percentage of predicted mathematical expressions that perfectly match the actual results.

### Implementation Details

The proposed model, UniMERNet, uses PyTorch with a maximum sequence length set to 1536. Training is conducted on a single GPU equipped with CUDA. Specifically, we utilize an NVIDIA A100 with 80GB of memory. During the training phase, we employ eight such GPUs with a batch size of 64. The learning rate schedule is linear warmup cosine, with an initial learning rate of $1 \times 10^{-4}$, a minimum learning rate of $1 \times 10^{-8}$, and a warmup learning rate of $1 \times 10^{-5}$. Weight decay is set to 0.05. The total iteration is set to 300,000 by our default settings.

The architectural hyperparameters of UniMERNet instances are illustrated in Table 4. Let $N$ denote the depth of the UniMERNet encoder, where [6, 6, 6, 6] indicates that each stage, from the first to the last, consists of six transformer layers. Meanwhile, $M$ represents the depth of the UniMERNet Decoder. $C$ signifies the dimensionality of the vectors after processing through the Encoder, and Params refers to the total number of parameters in the model.

### Ablation Study

**UniMER-1M** The diversity and quantity of training data are crucial for accurate formula recognition. As shown in Table 2, UniMERNet-B, trained solely on Pix2tex, achieves a BLEU score of 0.911 on SPE but performs poorly on CPE

| | $N$ | $M$ | $C$ | Params |
|---|---|---|---|---|
| UniMERNet-T | [6, 6, 6, 6] | 8 | 512 | 100M |
| UniMERNet-S | [6, 6, 6, 6] | 8 | 768 | 202M |
| UniMERNet-B | [6, 6, 6, 6] | 8 | 1024 | 325M |

Table 4: Architectural hyper-parameters of UniMERNet.

and HWE. The simplicity of Pix2tex leads to overfitting on SPE and difficulties with complex and handwritten formulas. Training on both Pix2tex and HWE improves the BLEU score on HWE to 0.873 but slightly declines on CPE. Notably, training with our UniMER-1M dataset, UniMERNet-B excels across all subsets. Compared to Pix2tex&HWE, the CPE BLEU score improves by 0.201, and Edit Distance decreases from 0.225 to 0.056. On SCE, BLEU improves by 0.097, and Edit Distance decreases from 0.309 to 0.224. For HWE, BLEU improves by 0.022, and Edit Distance decreases to 0.072.

**Model Architecture Design** As shown in Table 3, we conduct ablation experiments to validate our proposed optimization modules for the encoder, including Fine-Grained Embedding (FGE) and ConvEnhance (CN), and the decoder optimization module, Squeeze Attention (SA), along with the Remove Shift Window (RSW) operation. Our baseline architecture, **Texify** (Paruchuri 2023), uses the Swin-Transformer Encoder and mBART Decoder, similar to Donut (Kim et al. 2022) and Nougat (Blecher et al. 2023). We train randomly initialized models from scratch using the UniMER-1M dataset.

From the comparison results, it is evident that incorporating the FGE and CE modules leads to a stable performance

| Method | Params (M) | FPS (img/s) | SPE BLEU ↑ | SPE EditDis ↓ | CPE BLEU ↑ | CPE EditDis ↓ | SCE BLEU ↑ | SCE EditDis ↓ | HWE BLEU ↑ | HWE EditDis ↓ |
|---|---|---|---|---|---|---|---|---|---|---|
| Pix2tex (Blecher 2022) | - | - | 0.873 | 0.088 | 0.655 | 0.408 | 0.092 | 0.817 | 0.012 | 0.920 |
| Texify (Paruchuri 2023) | 312 | 4.16 | 0.906 | 0.061 | 0.690 | 0.230 | 0.420 | 0.390 | 0.341 | 0.522 |
| Texify* | 312 | 4.16 | 0.906 | 0.067 | 0.900 | 0.077 | 0.599 | 0.224 | 0.888 | 0.075 |
| UniMERNet-T | 107 | **7.20** | 0.909 | 0.066 | 0.902 | 0.075 | 0.566 | 0.239 | 0.883 | 0.078 |
| UniMERNet-S | 202 | 6.04 | 0.913 | 0.061 | 0.920 | 0.060 | 0.618 | 0.228 | 0.889 | 0.075 |
| UniMERNet-B | 325 | 5.06 | **0.915** | **0.060** | **0.925** | **0.056** | **0.626** | **0.224** | **0.895** | **0.072** |

Table 5: Comparison with SOTA methods on UniMER-Test. *Note:* Texify* is trained using UniMER-1M and the same data augmentation methods described in this paper.

| Model | Pre | SPE BLEU ↑ | CPE BLEU ↑ | SCE BLEU ↑ | HWE BLEU ↑ |
|---|---|---|---|---|---|
| Texify | ✗ | 0.903 | 0.884 | 0.576 | 0.886 |
| Texify | ✓ | 0.906 | 0.903 | 0.599 | 0.888 |
| UniMERNet-B | ✗ | 0.912 | 0.897 | 0.601 | 0.893 |
| UniMERNet-B | ✓ | **0.915** | **0.925** | **0.626** | **0.895** |

Table 6: Ablation of pre-training with text-image pairs. The **Pre** column indicates whether pre-training is used. Note: For Texify and UniMERNet-B, we use the same in-house text-image pre-training data for fair comparison.



Figure 5: Comparative Visualization of Recognition Results Using Different Methods.

improvement across all subsets, with a significant increase observed in the SCE subset. The combination of these two modules improves the BLEU score from 0.579 to 0.599, an increase of nearly 2%. When the SA module is applied to the decoder, the model's inference speed increases from 4.07 to 5.04, a relative improvement of 24%, with minimal loss in accuracy. As mentioned in the methods section, with the CE module, the Shift Window becomes unnecessary. Removing the Shift Window results in slight improvements in both accuracy and speed, achieving the optimal configuration.

At this optimal configuration, the model's accuracy significantly improves compared to the baseline model, with BLEU score increases of 0.9%, 1.2%, 2.2%, and 0.6% on the SPE, CPE, SCE, and HWE subsets, respectively. The speed also improves by 23%, demonstrating the substantial advantages of our UniMERNet in formula recognition.

## Pre-training with Text Image Pairs

Pretraining on large-scale datasets significantly boosts model performance in document recognition. Models like Donut, Texify, and Nougat benefit greatly from this approach. Due to limited available data, we use Arxiv papers, applying text layout detection and OCR to extract text blocks and match them with source code, resulting in 16 million image-text pairs. As shown in Table 6, both the baseline model Texify and our UniMERNet-B exhibit significant improvements in SPE, CPE, and SCE metrics, with a relatively smaller gain in HWE. This is expected, as Arxiv papers predominantly feature printed text, differing from handwritten digit recognition. Incorporating more diverse pretraining
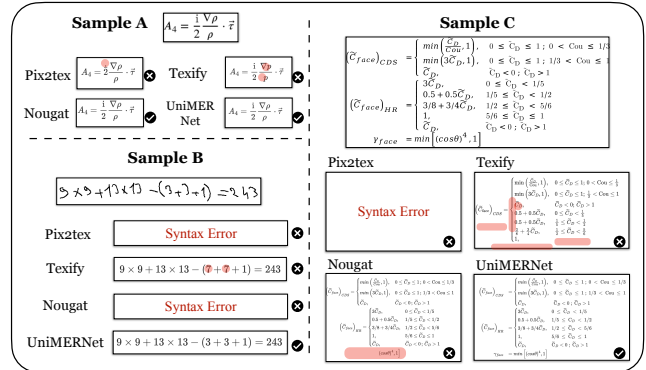
data would further improve overall performance.

**Comparison with SOTA Methods** To more intuitively evaluate the formula recognition performance of UniMER-Net, we compared it with the state-of-the-art methods in the document recognition field. As shown in Table 5, when using the same network architecture as Texify, the baseline model Texify* significantly outperforms original Texify, underscoring the importance of the UniMER-1M dataset. Moreover, our lightweight model UniMERNet-T, with an inference speed 1.73 times faster, already surpasses the previous SOTA model Texify. Using UniMERNet-B, the inference speed is 21.6% faster compared to Texify. The BLEU scores on the SPE, CPE, SCE, and HWE subsets show absolute improvements of 0.009, 0.235, 0.206, and 0.554 respectively, demonstrating the superior recognition accuracy and robustness of our model across various scenarios.

**Qualitative Comparisons** As shown in Figure 5, we selected three representative samples from the UniMER-Test set to thoroughly compare the performance between Pix2tex, Texify, Nougat, and UniMERNet. It's important to highlight that Nougat, being primarily designed for full-page recognition, tends to underperform with isolated formulas; thus, we prepared the test images by integrating random text with the formulas to adapt to Nougat's inference capabilities. Notably, while the other models exhibit certain shortcomings in handling these test samples, our model consistently delivers robust and accurate recognition results.

## Conclusion

This paper introduces the large-scale training dataset UniMER-1M and the diverse evaluation dataset UniMER-Test, contributing significantly to robust formula recognition and fair evaluation. We also designe UniMERNet, a model with superior detail perception and contextual understanding, achieving high accuracy and speed, making it valuable for practical applications. Moving forward, we will explore using larger and more diverse pre-training data to enhance formula recognition. Additionally, we will investigate integrating UniMERNet with large vision-language models to improve the recognition of documents containing text, formulas, and tables, advancing document understanding.

## References

Anderson, R. H. 1967. Syntax-directed recognition of hand-printed two-dimensional mathematics. In *Symposium on interactive systems for experimental applied mathematics: Proceedings of the Association for Computing Machinery Inc. Symposium*, 436–459. 2

Bian, X.; Qin, B.; Xin, X.; Li, J.; Su, X.; and Wang, Y. 2022. Handwritten mathematical expression recognition via attention aggregation based bi-directional mutual learning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 36, 113–121. 2

Blecher, L. 2022. pix2tex - LaTeX OCR. https://github.com/lukas-blecher/LaTeX-OCR. Accessed: 2024-2-29. 1, 2, 3, 7, 10, 11

Blecher, L.; Cucurull, G.; Scialom, T.; and Stojnic, R. 2023. Nougat: Neural optical understanding for academic documents. *arXiv.org*, 2308.13418. 1, 2, 4, 6

Chan, K.; and Yeung, D. 1999. Error detection, error correction and performance evaluation in on-line mathematical expression recognition. 2

Chen, Z.; Wang, W.; Tian, H.; Ye, S.; Gao, Z.; Cui, E.; Tong, W.; Hu, K.; Luo, J.; Ma, Z.; et al. 2024. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*. 2

Chu, X.; Tian, Z.; Zhang, B.; Wang, X.; and Shen, C. 2021. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*. 5

Deng, Y.; Kanervisto, A.; Ling, J.; and Rush, A. M. 2017. Image-to-markup generation with coarse-to-fine attention. In *International Conference on Machine Learning (ICML)*, 980–989. PMLR. 1, 2, 3, 10, 11

Dong, X.; Zhang, P.; Zang, Y.; Cao, Y.; Wang, B.; Ouyang, L.; Wei, X.; Zhang, S.; Duan, H.; Cao, M.; et al. 2024. InternLM-XComposer2: Mastering Free-form Text-Image Composition and Comprehension in Vision-Language Large Model. *arXiv.org*, 2401.16420. 2

d'Ascoli, S.; Touvron, H.; Leavitt, M. L.; Morcos, A. S.; Biroli, G.; and Sagun, L. 2021. Convit: Improving vision transformers with soft convolutional inductive biases. In *International conference on machine learning*, 2286–2296. PMLR. 5

Guo, J.; Han, K.; Wu, H.; Tang, Y.; Chen, X.; Wang, Y.; and Xu, C. 2022. Cmt: Convolutional neural networks meet vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 12175–12185. 5

He, C.; Li, W.; Jin, Z.; Xu, C.; Wang, B.; and Lin, D. 2024. Opendatalab: Empowering general artificial intelligence with open datasets. *arXiv preprint arXiv:2407.13773*. 2

Kim, G.; Hong, T.; Yim, M.; Nam, J.; Park, J.; Yim, J.; Hwang, W.; Yun, S.; Han, D.; and Park, S. 2022. Ocr-free document understanding transformer. In *European Conference on Computer Vision (ECCV)*, 498–517. Springer. 1, 2, 4, 6

Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 25. 2

Le, A. D.; Indurkhya, B.; and Nakagawa, M. 2019. Pattern generation strategies for improving recognition of handwritten mathematical expressions. *Pattern Recognition Letters*, 128: 255–262. 1, 2

Levenshtein, V. I.; et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, 707–710. Soviet Union. 5

Li, B.; Yuan, Y.; Liang, D.; Liu, X.; Ji, Z.; Bai, J.; Liu, W.; and Bai, X. 2022. When counting meets HMER: counting-aware network for handwritten mathematical expression recognition. In *European Conference on Computer Vision (ECCV)*, 197–214. Springer. 2

Li, Z.; Jin, L.; Lai, S.; and Zhu, Y. 2020. Improving attention-based handwritten mathematical expression recognition with scale augmentation and drop attention. In *International Conference on Frontiers in Handwriting Recognition (ICFHR)*, 175–180. IEEE. 2

Li, Z.; Yang, X.; Choi, K.; Zhu, W.; Hsieh, R.; Kim, H.; Lim, J. H.; Ji, S.; Lee, B.; Yan, X.; et al. 2024. MMSci: A Multimodal Multi-Discipline Dataset for PhD-Level Scientific Comprehension. *arXiv preprint arXiv:2407.04903*. 2

Liu, H.; Li, C.; Li, Y.; and Lee, Y. J. 2023. Improved baselines with visual instruction tuning. *arXiv.org*, 2310.03744. 2

Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2024. Visual instruction tuning. volume 36. 2

Mahdavi, M.; Zanibbi, R.; Mouchere, H.; Viard-Gaudin, C.; and Garain, U. 2019. ICDAR 2019 CROHME+ TFD: Competition on recognition of handwritten mathematical expressions and typeset formula detection. In *International Conference on Document Analysis and Recognition (ICDAR)*, 1533–1538. IEEE. 1, 3

Miller, E.; and Viola, P. 1998. Ambiguity and constraint in mathematical expression recognition. *National Conference on Artificial Intelligence (NCAI)*. 2

Mouchere, H.; Viard-Gaudin, C.; Zanibbi, R.; and Garain, U. 2014. ICFHR 2014 competition on recognition of on-line handwritten mathematical expressions (CROHME 2014). In *International Conference on Frontiers in Handwriting Recognition (ICFHR)*, 791–796. IEEE. 3

Mouchère, H.; Viard-Gaudin, C.; Zanibbi, R.; and Garain, U. 2016. ICFHR2016 CROHME: Competition on recognition of online handwritten mathematical expressions. In *International Conference on Frontiers in Handwriting Recognition (ICFHR)*, 607–612. IEEE. 3

Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318. 5

Paruchuri, V. 2023. Texify. https://github.com/VikParuchuri/texify. Accessed: 2024-2-29. 1, 2, 4, 6, 7

Simonyan, K.; and Zisserman, A. 2015. Very deep convolutional networks for large-scale image recognition. In *ICLR*. Computational and Biological Learning Society. 2

Suzuki, M.; Tamari, F.; Fukuda, R.; Uchida, S.; and Kanahori, T. 2003. Infty: an integrated ocr system for mathematical documents. In *Proceedings of the 2003 ACM symposium on Document engineering*, 95–104. 2

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. *NeurIPS*, 30. 1, 2

Wang, B.; Wu, F.; Han, X.; Peng, J.; Zhong, H.; Zhang, P.; Dong, X.; Li, W.; Li, W.; Wang, J.; et al. 2024. Vigc: Visual instruction generation and correction. In *AAAI*. 2

Wei, H.; Kong, L.; Chen, J.; Zhao, L.; Ge, Z.; Yang, J.; Sun, J.; Han, C.; and Zhang, X. 2023. Vary: Scaling up the vision vocabulary for large vision-language models. *arXiv.org*, 2312.06109. 1, 2

Wu, J.-W.; Yin, F.; Zhang, Y.-M.; Zhang, X.-Y.; and Liu, C.-L. 2020. Handwritten mathematical expression recognition via paired adversarial learning. *International Journal of Computer Vision (IJCV)*, 128: 2386–2401. 1, 2

Xia, R.; Mao, S.; Yan, X.; Zhou, H.; Zhang, B.; Peng, H.; Pi, J.; Fu, D.; Wu, W.; Ye, H.; et al. 2024. DocGenome: An Open Large-scale Scientific Document Benchmark for Training and Testing Multi-modal Large Language Models. *arXiv preprint arXiv:2406.11633*. 2

Yuan, Y.; Liu, X.; Dikubab, W.; Liu, H.; Ji, Z.; Wu, Z.; and Bai, X. 2022. Syntax-aware network for handwritten mathematical expression recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 4553–4562. 3, 6

Zhang, J.; Du, J.; Yang, Y.; Song, Y.-Z.; Wei, S.; and Dai, L. 2020. A tree-structured decoder for image-to-markup generation. In *International Conference on Machine Learning (ICML)*, 11076–11085. PMLR. 2

Zhang, J.; Du, J.; Zhang, S.; Liu, D.; Hu, Y.; Hu, J.; Wei, S.; and Dai, L. 2017. Watch, attend and parse: An end-to-end neural network based approach to handwritten mathematical expression recognition. *Pattern Recognition*, 71: 196–206. 2

Zhang, P.; Dong, X.; Zang, Y.; Cao, Y.; Qian, R.; Chen, L.; Guo, Q.; Duan, H.; Wang, B.; Ouyang, L.; et al. 2024. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*. 2

Zhao, W.; Gao, L.; Yan, Z.; Peng, S.; Du, L.; and Zhang, Z. 2021. Handwritten mathematical expression recognition with bidirectionally trained transformer. In *International Conference on Document Analysis and Recognition (ICDAR)*, 570–584. Springer. 1, 2

Zhu, D.; Chen, J.; Shen, X.; Li, X.; and Elhoseiny, M. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv.org*, 2304.10592. 2

# Appendix

## Details of UniMER Dataset

**SPE and CPE Sampling** Existing datasets, such as IM2LATEX-100k and Pix2tex, present two primary challenges. Firstly, the size of these datasets, typically ranging from 100k to 200k formulas, is insufficient for training a precise and robust MER model. Secondly, these datasets contain a limited number of complex formulas, which compromises the model's performance, particularly in handling multi-line complex expressions.
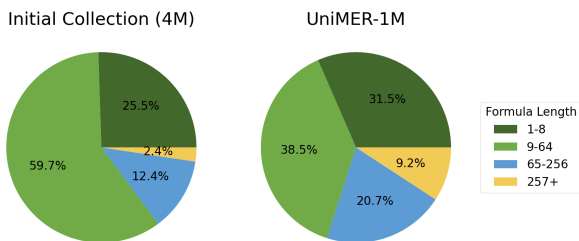
Figure 6: Formula length before and after re-sampling.

To address the limited size of the dataset, we expand it by incorporating an additional 4 million LaTeX expression source codes, building on the previously mentioned open-source datasets. These new entries are predominantly sourced from Arxiv (89%), with supplementary contributions from Wikipedia (9%) and StackExchange (2%). This initial dataset expansion enhances the model's overall capabilities. However, the proportion of long formulas in the initial collection is relatively small (2.3%), which may cause inadequate training for complex expressions. To address this, we extract the longest formulas as CPE and adjust their ratio with randomly sampled SPE. This rearrangement ensures a balanced representation of varying lengths within the dataset, thereby significantly improving the model's ability to recognize complex multi-line mathematical expressions. The distribution after rearrangement is shown in Figure 6.

**SCE Deduplication** When extracting mathematical formulas from PDF pages, we face a unique challenge: formulas originating from the same page often appear identical in content, leading to potential duplicates. Simple deduplication based on textual content alone risks significant data loss, as identical formulas can appear across different pages, each bearing distinct visual characteristics such as font styles, sizes, and backgrounds. To preserve the richness of visual diversity while eliminating true duplicates, we adopted an image-based deduplication strategy, employing Perceptual Hashing to assess image similarity. This method allows us to compare the visual features of the formula images directly, ensuring that only those with high similarity—indicating true duplicates—are removed. Through this meticulous process of image similarity analysis, we effectively reduced the dataset to 4,744 unique Screen-Captured Expressions (SCE), each representing a distinct visual instance of mathematical expressions, thereby constituting our refined SCE test set.
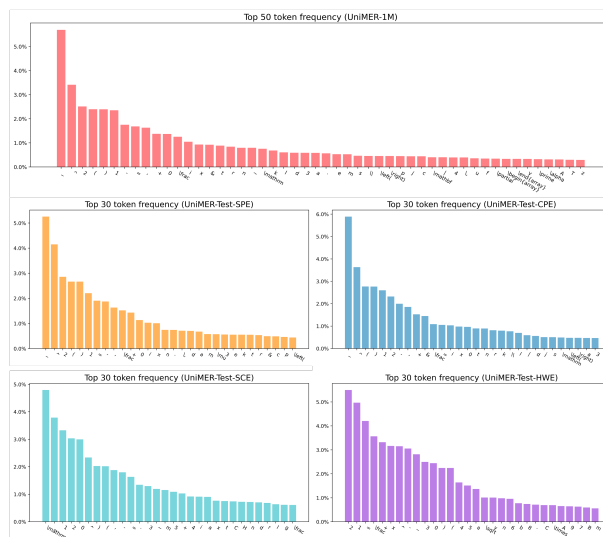
Figure 7: Most frequent occurring latex symbols in UniMER-1M and UniMER-Test subsets

Listing 1: XeLaTeX rendering setting

```
1  \documentclass[varwidth]{standalone}
2  \usepackage{fontspec,unicode-math}
3  \usepackage[active,displaymath,textmath,tightpage]{
       preview}
4  \usepackage[total={16in, 16in}]{geometry}
5  \setmathfont{
6      % MATH_FONT
7  }
8  \begin{document}
9  \thispagestyle{empty}
10 \begin{displaymath}
11     % MATH_FORMULA
12 \end{displaymath}
13 \end{document}
```

**Rendering Settings** For the rendering settings, we follow the similar procedure used in (Deng et al. 2017) and (Blecher 2022). The dataset is rendered using XeLaTeX with a diverse range of math fonts and DPI settings. The chosen fonts included Asana Math, Cambria Math, XITS Math, GFS Neohellenic Math, TeX Gyre Bonum Math, TeX Gyre Dejavu Math, TeX Gyre Pagella Math, and Latin Modern Math, with Latin Modern Math as the default math font being employed in approximately 22% of the cases. To accommodate different levels of clarity and detail, the DPI setting varies between 80 to 350 when converting to PNG format, allowing for adjustments in the resolution and sharpness of the rendered mathematical expressions. The rendering template we use is shown in Listing 1.

**Formula Text Normalization** LaTeX syntax inherently contains ambiguous information, as different source codes can produce the same rendering. This presents significant challenges in the evaluation phase of the math formula recognition task's benchmark because it potentially leads to incorrect assessments of a model's performance despite producing visually identical formula renderings. In handwritten
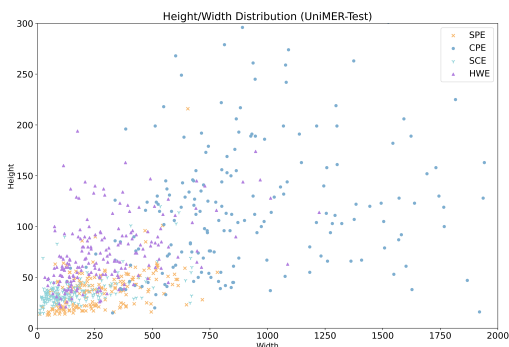
Figure 8: Height width scatter plot in UniMER-Test subsets with sampling

math recognition datasets, such as CROHME, a self-defined label graphs format is used, eliminating ambiguous expressions by employing a character relation-based method. We do not adopt these methods for normalization as they involve format conversions during model training and, more importantly, use only partial LaTeX syntax.

The LaTeX normalization is first introduced in (Deng et al. 2017). This preprocessing operation involves fixing super-script and sub-script order, replacing ambiguity with unified expressions while resulting in no or minimal visual changes in rendering, preserving the integrity of the original mathematical expressions. Subsequent datasets, such as IM2LATEX-100K (Deng et al. 2017) and Pix2tex (Blecher 2022), have adopted similar methods. Building on this foundation of normalization, we have adjusted the normalization rules for certain LaTeX environments, enabling better support for multi-line formula expressions and previously unsupported syntax. All formulas in UniMER-1M and UniMER-Test undergo this normalization process to facilitate better horizontal comparison with previous datasets that employ a similar normalization process.

## Data Statistics

**Most Occurring Symbols** Diving into the dataset's La-TeX symbols offers intriguing insights into the most frequently utilized mathematical notations. The bar chart provided in the Figure 7 illustrates the frequency of specific LaTeX symbols that appear in UniMER. Symbols such as Greek letters, operators, and various mathematical functions are universally prevalent in each dataset, underlining their fundamental role in articulating complex mathematical ideas. A subtle variation is observed in the SCE and HWE datasets, where numbers and letters are noticeably more frequently occurring, as they contain relatively easier and less structured math expressions.

**Image Size Distribution** The scatter plot in Figure 8 provides a visual distribution of image sizes across different subsets within the UniMER-Test. Each point on the plot represents an individual image, with its position determined by the image's width and height. The SPE, CPE, SCE, and HWE subsets each exhibit unique clusters, indicating the variety in dimensionality they encompass. It's evident that

the SPE and SCE subsets tend to have a higher concentration of smaller images, as shown by the dense clustering of points towards the lower end of the spectrum. The distribution of image sizes within the CPE dataset exhibits a considerable spread, highlighting the diversity of dimensions that this particular subset encompasses, indicating its complexity compared to SPE. On the other hand, the HWE subset is characterized by images with generally larger dimensions. This can be attributed to the fact that these images are often photographed and contain noise, necessitating a higher resolution to ensure that the finer details of the handwritten expressions are preserved and recognizable.

## Data Augmentation

While introducing additional training data in UniMER-1M enhances the variety of formulas, it does not account for the diversity of real-world formula images, which can come from scanned documents or photos and can exhibit noise and distortion. We employ various image augmentation techniques during model training to simulate this diversity with extra transformations from Albumentations [2] library and self-defined transformations, which include but are not limited to:

- **Erosion/Dilation** - To simulate the textural imperfections often found in screen-captured formulas, these operations modulate the thickness of characters, mirroring the effects of resolution differences and printer anomalies.

- **Degradation Simulation** (Fog, Frost, Rain, Snow, Shadow) - These augmentations introduce environmental artifacts to mimic the conditions under which documents might be photographed in real-world scenarios, adding layers of complexity such as blurriness and occlusions.

- **Geometric Transformations** (Rotation, Distortion . . . ) - To account for the angle and perspective distortions typical in photographed or scanned documents, these operations adjust the orientation and shape of the mathematical expressions.

Each image undergoes a sequence of these augmentation operations with a given probability. This helps to bridge the gap between the pristine, synthesized training data and the noisy, real-world test images and improves UniMERNet's performance for real piratical use. Figure 9 provides a visualization of selected transformations.

**Impact of Data Augmentation on Performance** Data augmentation proves to be significantly beneficial for real-world formula recognition tasks. As shown in Table 7, incorporating image augmentation into the training process with the Pix2tex dataset leads to varying degrees of improvement across all evaluation subsets. Notably, on the SCE subset, the BLEU score increases by 2.5%. Similar trends are observed when training with the UniMER-1M dataset, where the BLEU score on the SCE subset improves significantly from 0.601 to 0.626, and the edit distance decreases from 0.251 to 0.224,/;.

---

[2]https://albumentations.ai

| Train Dataset | Augment | SPE | | CPE | | SCE | | HWE | |
|---|---|---|---|---|---|---|---|---|---|
| | | BLEU ↑ | EditDis ↓ | BLEU ↑ | EditDis ↓ | BLEU ↑ | EditDis ↓ | BLEU ↑ | EditDis ↓ |
| Pix2tex | ✗ | 0.909 | 0.064 | 0.764 | 0.198 | 0.512 | 0.380 | 0.065 | 0.807 |
| | ✔ | 0.911 | 0.063 | 0.773 | 0.194 | 0.527 | 0.371 | 0.067 | 0.800 |
| UniMER-1M | ✗ | 0.912 | 0.064 | 0.911 | 0.063 | 0.601 | 0.251 | 0.886 | 0.078 |
| | ✔ | **0.915** | **0.060** | **0.925** | **0.056** | **0.626** | **0.224** | **0.895** | **0.072** |

Table 7: Ablation results on UniMER-Test with models using different augmentations.

## Optimal Depth Configuration for UniMERNet

In this section, we investigate the optimal depth configuration for the encoder and decoder of the UniMERNet model. Through a series of comparative experiments, we analyze the impact of varying depths on model performance and throughput.

**Encoder Depth Analysis** Table 8 presents the results of our experiments on different encoder depths while keeping the decoder depth constant. The findings indicate that increasing the encoder depth enhances model performance up to a certain point. Specifically, when the encoder depth reaches six layers per stage, the performance gains begin to plateau. This is evident from the BLEU scores across various evaluation metrics, which show diminishing returns beyond this depth. Additionally, the frames per second (FPS) metric reveals that increasing the encoder depth has a minimal impact on model throughput.

**Decoder Depth Analysis** Table 9 explores the effects of varying decoder depths while keeping the encoder depth fixed at six layers per stage. Similar to the encoder, increasing the decoder depth improves performance up to a depth of eight layers, after which the performance gains diminish significantly. However, unlike the encoder, the decoder depth has a more pronounced impact on throughput, with FPS decreasing notably as the decoder depth increases.

By applying the law of diminishing marginal utility and Pareto optimality principles, we determine that the optimal configuration for UniMERNet is an encoder depth of six layers per stage and a decoder depth of eight layers. This configuration balances model performance and throughput, ensuring efficient and effective processing.
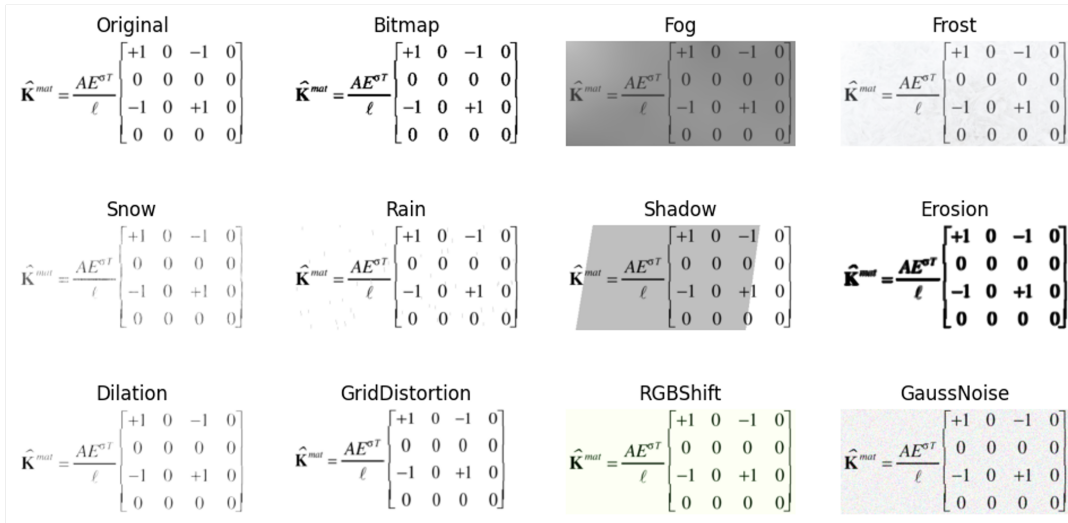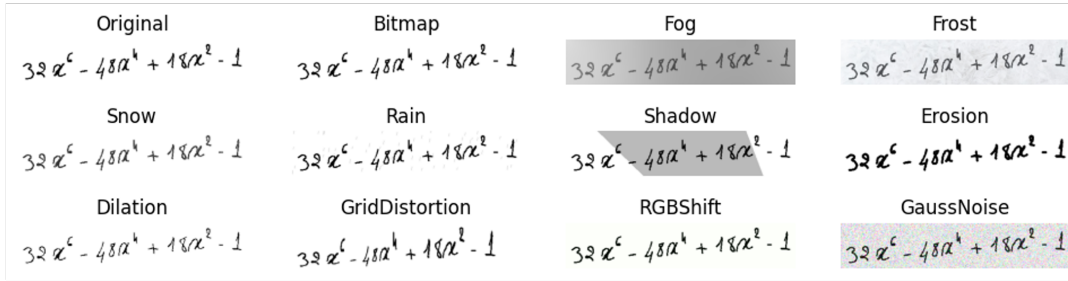
Original  Bitmap  Fog  Frost

$32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$

Snow  Rain  Shadow  Erosion

$32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$

Dilation  GridDistortion  RGBShift  GaussNoise

$32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$   $32x^c - 48x^4 + 18x^2 - 1$

Original  Bitmap  Fog  Frost

$$\hat{K}^{mat} = \frac{AE^{\sigma T}}{\ell}\begin{bmatrix} +1 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & +1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Snow  Rain  Shadow  Erosion

$$\hat{K}^{mat} = \frac{AE^{\sigma T}}{\ell}\begin{bmatrix} +1 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & +1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Dilation  GridDistortion  RGBShift  GaussNoise

$$\hat{K}^{mat} = \frac{AE^{\sigma T}}{\ell}\begin{bmatrix} +1 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & +1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Figure 9: Visualization of selected image augmentations applied during training.

| N | M | Params (M) | FPS (img/s) | SPE BLEU ↑ | CPE BLEU ↑ | SCE BLEU ↑ | HWE BLEU ↑ | AVG BLEU ↑ |
|---|---|---|---|---|---|---|---|---|
| [2, 2, 2, 2] | 6 | 148 | 7.25 | 0.890 | 0.832 | 0.496 | 0.833 | 0.763 |
| [4, 4, 4, 4] | 6 | 167 | 7.23 | 0.897 (+0.7%) | 0.856 (+2.4%) | 0.544 (+4.8%) | 0.870 (+3.7%) | 0.792 (+2.9%) |
| [6, 6, 6, 6] | 6 | 186 | 7.20 | 0.901 (+0.4%) | 0.878 (+1.2%) | 0.564 (+2.0%) | 0.886 (+1.6%) | 0.807 (+1.5%) |
| [8, 8, 8, 8] | 6 | 205 | 7.15 | 0.902 (+0.1%) | 0.880 (+0.2%) | 0.578 (+1.4%) | 0.878 (-0.8%) | 0.809 (+0.2%) |

Table 8: Comparative study of the impact of encoder depth on model performance. Let $N$ denote the depth of the UniMERNet encoder, where [6, 6, 6, 6] indicates that each stage, from the first to the last, consists of six transformer layers. Meanwhile, $M$ represents the depth of the UniMERNet Decoder.

| N | M | Params (M) | FPS (img/s) | SPE BLEU ↑ | CPE BLEU ↑ | SCE BLEU ↑ | HWE BLEU ↑ | AVG BLEU ↑ |
|---|---|---|---|---|---|---|---|---|
| [6, 6, 6, 6] | 4 | 169 | 8.36 | 0.893 | 0.850 | 0.532 | 0.863 | 0.785 |
| [6, 6, 6, 6] | 6 | 186 | 7.20 | 0.901(+0.8%) | 0.878(+1.8%) | 0.564(+3.2%) | 0.886(+2.3%) | 0.807 (+2.2%) |
| [6, 6, 6, 6] | 8 | 202 | 6.04 | 0.905(+0.4%) | 0.892(+1.6%) | 0.587(+2.3%) | 0.888(+0.2%) | 0.818(+1.1%) |
| [6, 6, 6, 6] | 10 | 219 | 4.99 | 0.907(+0.2%) | 0.894(+0.2%) | 0.582(-0.5%) | 0.893(+0.5%) | 0.819(+0.1%) |

Table 9: Comparative study of the impact of decoder depth on model performance. Let $N$ denote the depth of the UniMERNet encoder, where [6, 6, 6, 6] indicates that each stage, from the first to the last, consists of six transformer layers. Meanwhile, $M$ represents the depth of the UniMERNet Decoder.