

Instruction-augmented Multimodal Alignment for Image-Text and Element Matching

Xinli Yue^{*1} JianHui Sun^{*2} Junda Lu² Liangchao Yao² FAN XIA²
Tianyi Wang² Fengyun Rao² JING LYU² Yuetang Deng^{†2}
¹Wuhan University ²WeChat

Abstract

*With the rapid advancement of text-to-image (T2I) generation models, assessing the semantic alignment between generated images and text descriptions has become a significant research challenge. Current methods, including those based on Visual Question Answering (VQA), still struggle with fine-grained assessments and precise quantification of image-text alignment. This paper presents an improved evaluation method named **Instruction-augmented Multimodal Alignment for Image-Text and Element Matching (iMatch)**, which evaluates image-text semantic alignment by fine-tuning multimodal large language models. We introduce four innovative augmentation strategies: First, the *QAlign* strategy creates a precise probabilistic mapping to convert discrete scores from multimodal large language models into continuous matching scores. Second, a validation set augmentation strategy uses pseudo-labels from model predictions to expand training data, boosting the model’s generalization performance. Third, an element augmentation strategy integrates element category labels to refine the model’s understanding of image-text matching. Fourth, an image augmentation strategy employs techniques like random lighting to increase the model’s robustness. Additionally, we propose prompt type augmentation and score perturbation strategies to further enhance the accuracy of element assessments. Our experimental results show that the *iMatch* method significantly surpasses existing methods, confirming its effectiveness and practical value. Furthermore, our *iMatch* won first place in the CVPR NTIRE 2025 Text to Image Generation Model Quality Assessment - Track 1 Image-Text Alignment.*

1. Introduction

In recent years, the rapid development of deep learning technologies has led to significant breakthroughs in text-to-

image (T2I) generation models [2, 5, 9, 10, 16, 21, 23, 30–32], which have demonstrated powerful image generation capabilities. However, objectively and accurately assessing the semantic alignment of these generated images with text descriptions has increasingly become a critical issue and a significant challenge in the research field.

Metrics based on visual language models [22, 29] assess the semantic matching between image and text by measuring cosine similarity in embedding spaces. Several approaches [18, 19, 37, 41] utilize human-annotated data to emulate human judgments of image-text alignment, while others analyze texts broken down into elements evaluated through VQA models. Innovations such as the introduction of question templates improve the alignment of the VQA model with human preferences [12]. Although these methods form a foundation for evaluating image-text alignment, they lack a unified approach for comprehensive and detailed matching assessments, missing a systematic exploration of the complex relationships between text and images.

On the other hand, the emergence of Multimodal Large Language Models (MLLMs) [1, 3, 7, 24–27, 34, 45] in recent years has provided robust technical support for image-text matching tasks. Recently, models such as InternVL2.5 [6], Qwen2.5-VL [4], and the Ovis2 [28] series have shown exceptional performance in multimodal understanding tasks. Nonetheless, further exploration is needed to fully leverage these models for more accurate assessments of image-text matching.

Addressing the challenges highlighted above, this study proposes an enhanced method for comprehensive and fine-grained image-text matching assessment, named **Instruction-augmented Multimodal Alignment for Image-Text and Element Matching (iMatch)**, aimed at precisely measuring the alignment between generated images and textual descriptions. Specifically, our approach encompasses several innovative aspects: Firstly, we employ a fine-tuning strategy based on MLLMs, utilizing fine-grained image-text matching annotations from the EvalMuse-40K dataset [12] to explicitly guide the model in learning the nuanced correspondences between text and images. To fur-

^{*}Equal contribution.

[†]Corresponding author.

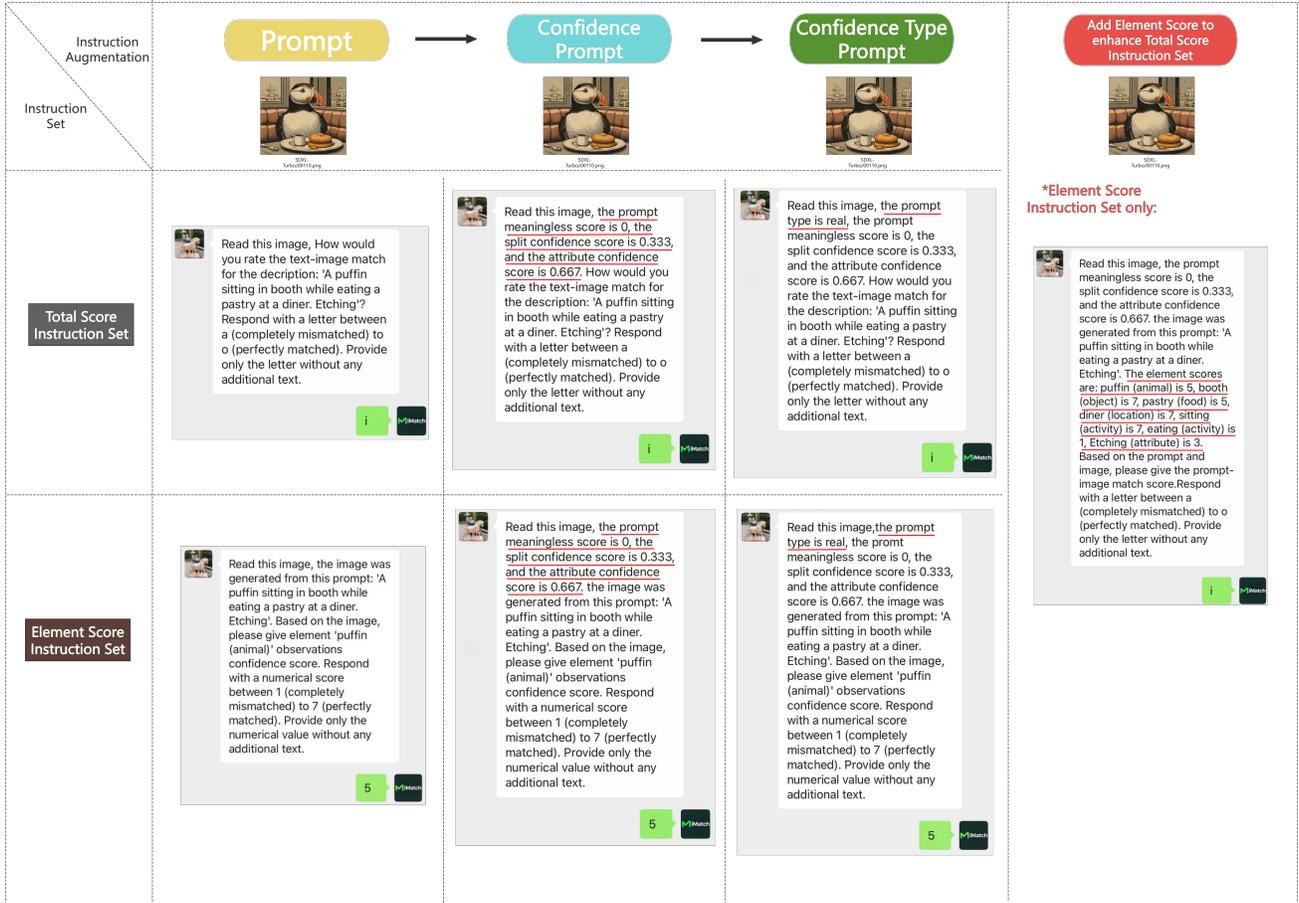


Figure 1. The instruction set augmentation process of the proposed iMatch.

ther enhance model performance, we propose four augmentation strategies: (1) QAlign [39] strategy: This strategy defines a mapping function from textual rating levels to specific numerical scores, coupled with a soft mapping of model prediction probabilities, to more accurately convert overall image-text matching scores. (2) Validation set augmentation strategy: We initially use the model to predict on the validation set during training, generating high-quality pseudo-labels, which are then merged back into the original training set for re-training to enhance the model’s generalization performance. (3) Element augmentation strategy: During training, we explicitly input element labels as additional features into the user query, helping the model deduce the overall matching score from finer-grained information in a Chain-of-Thought-like manner [36]. (4) Image augmentation strategy: We introduce three data augmentation techniques to expand the diversity of training set images and enhance the model’s robustness to image variations.

Additionally, in the element matching task, we propose two additional augmentation techniques: (1) Prompt type augmentation: We explicitly integrate the prompt type (real

or synthetic) into the user query, aiding the model in distinguishing the intrinsic characteristics of different prompt sources. (2) Score perturbation augmentation: By applying slight random perturbations to the target labels of element matching, we reduce the risk of model overfitting to specific training labels, further improving the model’s generalization capabilities. The synergistic effect of these methods has led to outstanding performance on the EvalMuse-40K dataset [12] and in the NTIRE 2025 challenge [13], validating the effectiveness and practicality of our proposed methods in image-text matching tasks.

In summary, our contributions are as follows:

- We present iMatch, an innovative image-text matching method that enhances semantic matching accuracy through fine-tuned multimodal models and strategic augmentations.
- We introduce four augmentation strategies—QAlign, validation set augmentation, element augmentation, and image augmentation—to improve the model’s adaptability and generalization in image-text tasks.
- We develop two techniques—prompt type augmentation

and score perturbation augmentation—to boost model performance and stability.

- Our extensive testing on the EvalMuse-40K dataset and in the NTIRE 2025 challenge shows that iMatch surpasses existing methods across multiple metrics.

2. Related Work

2.1. Image-Text Alignment

Recent advancements in text-to-image (T2I) generation models [2, 5, 9–11, 14, 16, 21, 23, 32, 35] have emphasized the importance of accurately assessing the alignment between images and text descriptions. CLIPScore [15] utilizes a pretrained CLIP [29] model to measure cosine similarity between embeddings, providing an initial automatic evaluation. BLIP2Score [15] follows a similar approach for enhanced evaluation. ImageReward [41] and PickScore [19] refine assessments through fine-tuning with extensive human feedback, aligning more closely with human perceptions. TIFA [18] and VQ2 [42] break down text into element-level questions answered by VQA models, focusing on detailed matching. FGA-BLIP2 [12] improves this with tailored question templates to direct VQA models towards crucial text content, enhancing element-level accuracy. Despite these advances, challenges remain in providing a unified framework for detailed and overall image-text alignment.

2.2. Multimodal Large Language Models

In recent years, advanced MLLMs [1, 3, 7, 24–27, 33, 34, 38, 45] have shown remarkable performance in multimodal tasks. The InternVL2.5 series [6], based on InternVL 2.0, retains the original architecture but introduces improvements in training, testing, and data quality. The InternVL2.5-MPO model employs Mixed Preference Optimization (MPO) to enhance reasoning in multimodal tasks. The Qwen2.5-VL series [4] demonstrates strong visual and linguistic skills, with the Qwen2.5-VL-7B-Instruct model excelling in visual benchmarks and agentic tasks. The Ovis2 models [28], successors to Ovis1.6, feature upgraded dataset organization and training, boosting reasoning in smaller models. These developments in MLLMs advance cross-modal understanding and generation, setting the stage for future AI advancements.

3. Methodology

3.1. Baseline Model

We propose a fine-tuning approach based on MLLMs [4, 6, 28] to address the fine-grained image-text matching task in image generation quality assessment. Utilizing the EvalMuse-40K dataset [12], which provides extensive fine-grained image-text matching annotations, our method ex-

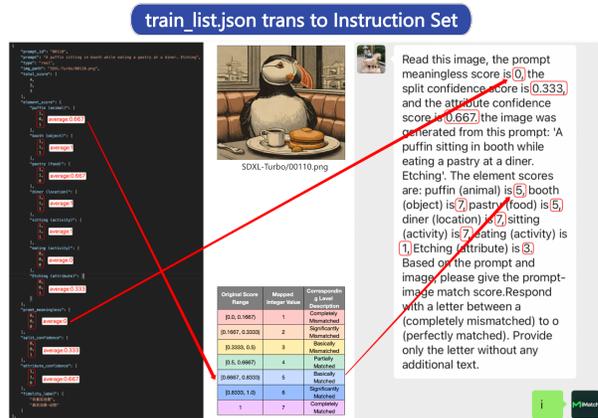


Figure 2. The construction process of the element instruction set augmentation. We map elements to integers between 1 and 7 and incorporate them into the user query. Additionally, confidence scores are also integrated into the user query.

PLICITLY guides the model in learning the detailed correspondences between text and images. Specifically, given a text description and a corresponding generated image, the model is tasked with predicting both an overall image-text matching score and element-level matching hits.

3.1.1. Problem Definition

Let the given text description be P , and its corresponding generated image be I . The EvalMuse-40K dataset [12] provides an overall image-text matching score S_{total} and a set of fine-grained element matching scores $\{S_{e_i}\}_{i=1}^N$ for each image-text pair (I, P) . Here, S_{total} ranges from [1,5], representing the overall matching degree, while each element matching score S_{e_i} ranges from [0,1], indicating whether a specific element is accurately represented in the image. Therefore, the task can be formally defined as:

$$\hat{S}_{\text{total}} = f_{\text{total}}(I, P; \theta), \quad (1)$$

$$\hat{y}_{e_i} = 1 [f_{\text{element}}(I, P, e_i; \phi) > \tau], \quad i = 1, 2, \dots, N, \quad (2)$$

where f_{total} is the model predicting the overall image-text matching score, f_{element} is the model predicting element matching scores, θ and ϕ are model parameters, and τ is the threshold for determining element hits.

3.1.2. Instructional Fine-tuning Strategy

To enhance the model’s capability for learning fine-grained image-text correspondences, we have designed a specific instructional fine-tuning strategy that transforms the task of predicting image-text matching scores into a classification problem. The steps are as follows:

For the overall image-text matching score S_{total} , we first perform a linear scaling:

$$S'_{\text{total}} = \text{round} \left(\frac{S_{\text{total}} - 1}{4} \times 14 + 1 \right). \quad (3)$$

At this point, S'_{total} is transformed into a discrete set of integers $\{1, 2, \dots, 15\}$. Subsequently, we map this integer range to the alphabet set $\{a, b, \dots, o\}$, which serves as the target labels for multimodal model instructional fine-tuning. During the inference phase, the model’s predicted letter labels are mapped back to the 1–15 range and linearly scaled to the original 1–5 range as the final prediction result.

For the element matching task, as illustrated in Figure 2, we discretize the element matching scores S_{e_i} into 7 categories:

$$S'_{e_i} = \text{round}(S_{e_i} \times 6) + 1. \quad (4)$$

During inference, we use a threshold $\tau = 3$, converting the predicted categories into a binary classification task, where scores greater than 3 are considered hits (1), and others are considered misses (0).

Additionally, we incorporate other information provided by the EvalMuse-40K dataset [12] such as prompt meaningfulness, split confidence, and attribute confidence into the problem formulation, as shown in Figure 2. Through this instruction-based fine-tuning strategy, we significantly enhance the multimodal model’s performance on image-text matching tasks, especially in terms of element-level recognition and judgment.

3.2. Image-Text Matching Augmented Model

3.2.1. QAlign Augmentation

To further enhance the accuracy of image-text matching tasks, we introduce a probability distribution-based augmentation strategy to meticulously simulate human scoring behavior, as illustrated in Figure 3. Specifically, we adopt a post-processing strategy similar to the QAlign method [39], which converts the model output from scoring levels to more precise final matching scores.

Firstly, we define a mapping function G from letters to numeric scores:

$$G : l_i \rightarrow i, i \in \{1, 2, \dots, 15\}, \quad (5)$$

where $\{l_i\}_{i=1}^{15} = \{a, b, \dots, o\}$, for example, scoring level a corresponds to numeric score 1, while level o corresponds to numeric score 15.

Next, we calculate the probabilities p_{l_i} predicted by the MLLMs. Specifically, if the logits output by the language model for each scoring level l_i are x_{l_i} , we use the closed-set softmax function to compute the probability distribution for each level:

$$p_{l_i} = \frac{e^{x_{l_i}}}{\sum_{j=1}^{15} e^{x_{l_j}}}, \quad (6)$$

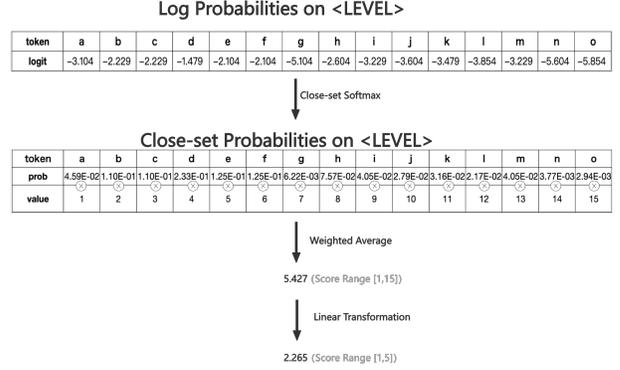
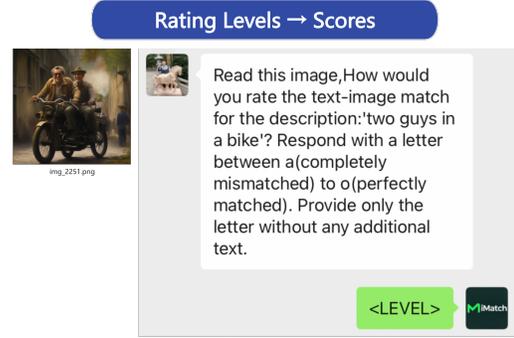


Figure 3. The inference process of the image-text matching augmented model. During inference, we extract closed-set probabilities for rating levels and perform a weighted average to obtain the MLLM-predicted score.

where $\sum_{i=1}^{15} p_{l_i} = 1$.

Consequently, the model’s final continuous predicted score \hat{S}_{total} is:

$$\hat{S}_{\text{total}} = \sum_{i=1}^{15} p_{l_i} \cdot G(l_i). \quad (7)$$

3.2.2. Validation Set Augmentation

The NTIRE 2025 competition has a development phase with a validation set, and a final phase with a test set. To further enhance the model’s generalization ability during the final test phase, we design a pseudo-labeling strategy that fully utilizes the validation set data to augment the training process. In the development phase, we first use a trained model to predict the validation set, generating pseudo labels, denoted as \hat{y}_v . Then, in the final phase, we combine the pseudo-labeled validation set with the original training set to construct an augmented training dataset for subsequent model retraining.

The new training objective can be expressed as:

$$\mathcal{L}_{\text{enhanced}}(\theta) = \mathcal{L}_{\text{train}}(y_t, \hat{y}_t; \theta) + \mathcal{L}_{\text{pseudo}}(\hat{y}_v, \hat{y}'_v; \theta), \quad (8)$$

where $\mathcal{L}_{\text{train}}$ denotes the loss on the original training set, $\mathcal{L}_{\text{pseudo}}$ denotes the loss supervised by pseudo labels, and \hat{y}'_v is the model’s prediction on the validation set samples.

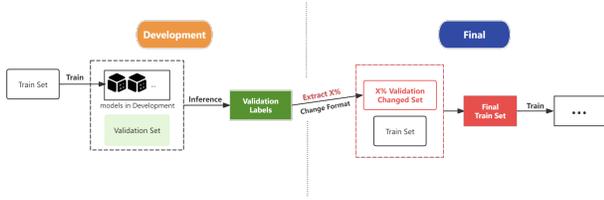


Figure 4. The process of validation set augmentation. We use the model trained on the training set to generate pseudo-labels for the validation set during the development phase, and then use these pseudo-labeled validation data to augment the training dataset.

3.2.3. Element Augmentation

To further enhance the model’s understanding and judgment capabilities for image-text matching, we propose a strategy based on element feature augmentation. Specifically, as shown in Figure 2 during the training phase, element-level scores are explicitly embedded into the user query as additional features to facilitate a reasoning process similar to the Chain-of-Thought (CoT) [36].

However, during the testing phase, due to the lack of real element labels, we cannot directly employ the same augmentation method. Therefore, we adopt a pseudo-label prediction strategy. Initially, we use the model trained in the Section 3.3 to predict pseudo-labels for the element level \hat{y}_{e_i} in the test set. Specifically, the model first predicts element scores:

$$\hat{S}_{e_i} = f_{\text{element}}(I, P, e_i; \phi). \quad (9)$$

Subsequently, these predicted element scores \hat{S}_{e_i} are embedded into the input of the image-text matching task, forming a prompt with pseudo-labeled element scores to predict the final image-text matching score \hat{S}_{total} :

$$\hat{S}_{\text{total}} = f_{\text{total}}(I, P, \hat{S}_{e_i}; \theta). \quad (10)$$

3.2.4. Image Augmentation

To further enhance the model’s generalization performance and improve robustness to image variations, we introduce an image data augmentation strategy, as illustrated in Figure 5. Specifically, we randomly select 10% of the samples from the training set and apply one of three different types of data augmentation methods. These augmented samples are then used in conjunction with the original data for training, effectively enriching the diversity of the training dataset. The specific data augmentation methods include:

Random Lighting Augmentation. We randomly adjust the brightness of images to simulate variations in lighting conditions, thereby improving the model’s robustness to

changes in illumination. Defined as:

$$I_{\text{light}} = T_{\text{brightness}}(I, \alpha), \alpha \sim U(0.1, 0.5), \quad (11)$$

where I represents the original image, and α is a randomly sampled lighting intensity factor.

Random Grid Distortion. We apply random grid deformations to the image to enhance the model’s robustness to spatial transformations. Specifically defined as:

$$I_{\text{grid}} = T_{\text{grid}}(I, \beta), \beta \sim U(0.2, 0.8), \quad (12)$$

where β is a random factor controlling the degree of grid distortion.

Random Crop Augmentation. We randomly crop images to simulate partial occlusions or changes in perspective. Specifically defined as:

$$I_{\text{crop}} = T_{\text{crop}}(I, \gamma), \gamma \sim U(0.1, 0.5), \quad (13)$$

where γ determines the cropping scale.

The above augmentation strategies are implemented as follows: for a randomly selected 10% subset D_{subset} from the training dataset D_{train} :

$$D_{\text{augmented}} = \{T(I) \mid I \in D_{\text{subset}}, T \in \{T_{\text{brightness}}, T_{\text{grid}}, T_{\text{crop}}\}\}. \quad (14)$$

Subsequently, we merge the enhanced dataset with the original dataset to construct the final augmented training set:

$$D_{\text{final}} = D_{\text{train}} \cup D_{\text{augmented}}. \quad (15)$$

3.3. Element Matching Augmented Model

Building upon the baseline model, we further propose an element matching augmented model aimed at improving the model’s fine-grained understanding and predictive accuracy regarding image-text element correspondences. Specifically, we introduce two augmentation strategies, including the incorporation of prompt type information and the introduction of score perturbation.

3.3.1. Prompt Type Augmentation

To more richly characterize the problem features in the element matching task, we propose integrating prompt type information from the EvalMuse-40K dataset—either “real” or “synthetic” (generated using GPT-4)—explicitly into the user query. By explicitly introducing prompt type information $t \in \{\text{real}, \text{synth}\}$, we expand the input formulation, resulting in an enhanced element matching task representation:

$$\hat{y}_{e_i} = f_{\text{element}}(I, P, e_i, t; \phi'). \quad (16)$$



Figure 5. Examples of image augmentation. We employ three types of augmentations: Random Lighting Augmentation, Random Grid Distortion, and Random Crop Augmentation.

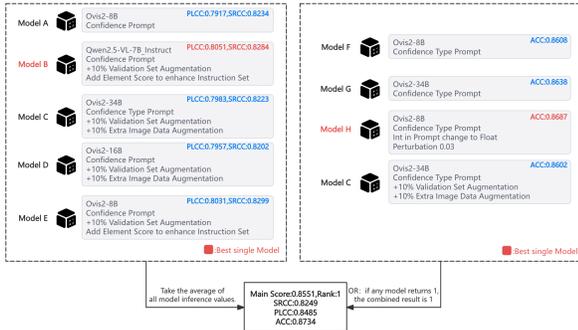


Figure 6. Model Ensemble. We ensemble five image-text matching augmented models and four element matching augmented models to improve the model’s generalization performance.

3.3.2. Score Perturbation

To further enhance the generalization capability of the element matching model and prevent overfitting to specific score annotations in the training data, we propose a perturbation strategy based on element labels. Specifically, during the training phase, we apply slight perturbations to the mapped element labels, which have a discrete numerical range of $\{1, 2, 3, 4, 5, 6, 7\}$, to increase the diversity of the training data.

Let the discretized score of element e_i in the training set be denoted as $S_{e_i}^{(\text{discrete})}$. We introduce a perturbation factor ϵ , and add a random perturbation $\delta \in \{-\epsilon, \epsilon\}$ to each element’s discretized score. The perturbed element score $S_{e_i}^{(\text{perturbed})}$ can then be expressed as:

$$S_{e_i}^{(\text{perturbed})} = S_{e_i}^{(\text{discrete})} + \delta, \quad \delta \sim U\{-\epsilon, \epsilon\}. \quad (17)$$

Table 1. Performance comparison of different methods on the EvalMuse-40K validation set for image-text matching tasks.

Method	SRCC	PLCC
CLIPScore [15]	0.2993	0.2933
BLIPv2Score [15]	0.3583	0.3348
ImageReward [41]	0.4655	0.4585
PickScore [19]	0.4399	0.4328
HPSv2 [40]	0.3745	0.3657
VQAScore [20]	0.4877	0.4841
FGA-BLIP2 [12]	0.7742	0.7722
InternVL2.5-8B-MPO (zero-shot)	0.7262	0.6742
iMatch (ours)	0.8304	0.8294

4. Experiments

4.1. Experimental Setup

Dataset. Our experimental studies employ the EvalMuse-40K dataset [12], which consists of images generated by over 20 different text-to-image (T2I) generation models based on approximately 4,000 prompts. Each image-text pair is annotated with two levels of scoring information: an overall prompt-level image-text matching score and a fine-grained element-level matching score. The dataset is divided into a training set, a validation set, and a test set. The training set contains about 30,000 image-text pairs for model parameter learning; the validation set includes about 10,000 pairs for tuning model parameters and selecting hyperparameters; the test set comprises approximately 5,000 pairs for final performance evaluation.

Models. We fine-tune a variety of advanced MLLMs, including InternVL2.5-8B-MPO [6], Qwen2.5-VL-7B-Instruct [4], Ovis2-8B [28], Ovis2-16B [28], and Ovis2-34B [28], to comprehensively validate the effectiveness and generalization capability of our proposed methods.

Training Parameters. We implement our experiments using the ms-swift framework [44], employing Low-Rank Adaptation (LoRA) [17] for fine-tuning, where the LoRA rank is set to 16. Additionally, we unfreeze the Vision Transformer (ViT) [8] and the cross-modal aligner to fully leverage cross-modal information. Specific training parameters include a maximum image pixel count of 200,704, an initial learning rate of $4e-5$, a weight decay factor of 0.01, a learning rate warm-up proportion of 0.03, and a learning rate decay strategy using cosine annealing. All models are trained for one epoch to ensure fair and consistent experimental conditions.

Evaluation Metrics. To comprehensively assess model performance, we utilize three metrics: the Spearman Rank Correlation Coefficient (SRCC), Pearson Linear Correlation Coefficient (PLCC), and element-level accuracy

Table 2. Performance comparison of different methods on the EvalMuse-40K validation set for element matching tasks. * indicates that the method uses a fixed step search (0.01) to select the optimal binary classification threshold, aiming to maximize overall accuracy.

Method	MLLMs	ACC
TIFA [18]	LLaVA1.6 [26]	0.6210
	mPLUG-Owl3 [43]	0.6450
	Qwen2-VL [4]	0.6450
VQ2*[45]	LLaVA1.6 [26]	0.6750
	mPLUG-Owl3 [43]	0.6640
	Qwen2-VL [4]	0.6790
PN-VQA* [12]	LLaVA1.6 [26]	0.6610
	mPLUG-Owl3 [43]	0.6760
	Qwen2-VL [4]	0.6820
FGA-BLIP2 [12]	BLIP2 [22]	0.7680
zero-shot	InternVL2.5-8B-MPO [6]	0.7681
iMatch (ours)	InternVL2.5-8B-MPO [6]	0.8284
	Qwen2.5-VL-7B-Instruct [4]	0.7948
	Ovis2-8B [28]	0.8317

(ACC). These are combined into a final performance score, calculated as $\text{Final Score} = \text{PLCC}/4 + \text{SRCC}/4 + \text{ACC}/2$.

Comparative Approaches. The comparative schemes include CLIPScore [15], BLIPv2Score [15], ImageReward [41], PickScore [19], HPSv2 [40], VQAScore [20], TIFA [18], VQ2 [42], PN-VQA [12], and FGA-BLIP2 [12].

4.2. Experimental Results

Results on Validation Set. We conduct comparative experiments on the EvalMuse-40K validation set against several mainstream methods to evaluate the effectiveness of our proposed iMatch. The performance results are presented in Tables 1 and 2.

As shown in Table 1, iMatch achieves significant gains in both SRCC and PLCC, outperforming all baselines. Compared to traditional metrics like CLIPScore, it improves SRCC and PLCC by over 0.533 and 0.539, respectively. Even against the strongest fine-tuned baseline, FGA-BLIP2, iMatch achieves absolute improvements of 0.056 (SRCC) and 0.057 (PLCC), demonstrating the effectiveness of our targeted fine-tuning strategy in capturing fine-grained semantic alignment.

We also report zero-shot results using the same backbone (InternVL2.5-8B-MPO). While the zero-shot model performs well, our fine-tuned iMatch surpasses it by a notable margin, highlighting the importance of task-aware fine-tuning.

Table 3. Results of the NTIRE 2025 Text to Image Generation Model Quality Assessment - Track 1 Image-Text Alignment.

Rank	Team	Main Score	SRCC	PLCC	ACC
1	IH-VQA (ours)	0.8551	0.8249	0.8485	0.8734
2	Evalthon	0.8426	0.8002	0.8321	0.8691
3	HCMUS	0.8381	0.8101	0.8306	0.8559
4	MICV	0.8221	0.7864	0.8050	0.8485
5	SJTU-MMLab	0.8158	0.7729	0.8029	0.8438
6	SJTUMM	0.8062	0.7563	0.7993	0.8346
7	WT	0.7913	0.7413	0.7740	0.8249
8	YAG	0.7777	0.7143	0.7456	0.8255
9	SPRank	0.7604	0.6899	0.7280	0.8119
10	AIIG	0.7386	0.6574	0.7073	0.7949
11	Joe1007	0.7359	0.6572	0.7041	0.7912
12	iCOST	0.7350	0.6630	0.7040	0.7865

In Table 2, iMatch also achieves the best accuracy on the element matching task, outperforming both baseline and zero-shot models. This further validates the effectiveness of our approach in improving both global and detailed image-text alignment.

Results on NTIRE 2025 Challenge. We validated the effectiveness of our method through the NTIRE 2025 Challenge on Text-to-Image Generation Model Quality Assessment – Track 1 Image-Text Alignment. As shown in Figure 6, our final solution leveraged a robust model ensemble strategy that played a key role in its success.

Table 3 summarizes the results. Our iMatch method ranked first, outperforming all competitors across all key metrics (SRCC, PLCC, ACC). Compared to the second-place Evalthon team, iMatch improved the Main Score by 0.0125, and achieved gains of 0.0247, 0.0164, and 0.0043 in SRCC, PLCC, and ACC, respectively. The performance margin was even larger relative to other teams, highlighting the strength and consistency of our approach.

4.3. Ablation Study

Results on Image-Text Matching Augmented Model.

To evaluate the effectiveness of the four augmentation strategies proposed in Section 3.2, we conduct ablation experiments on the EvalMuse-40K test set. The results are shown in Table 4. The introduction of the QAlign strategy significantly improved accuracy in image-text matching. The addition of the validation set augmentation strategy further enhanced the model’s generalization performance, as indicated by improvements in both SRCC and PLCC metrics. The element augmentation strategy, which incorporates element-level information, notably increased the model’s understanding of fine-grained aspects of image-text

Table 4. Results of the ablation study for each component of the image-text matching augmented model on the EvalMuse-40K test set.

Base	QAlign Aug.	Validation Aug.	Element Aug.	Image Aug.	InternVL2.5-8B-MPO		Qwen2.5-VL-7B-Instruct		Ovis2-8B	
					SRCC	PLCC	SRCC	PLCC	SRCC	PLCC
✓	×	×	×	×	0.7626	0.7523	0.7696	0.7381	0.7920	0.7757
✓	✓	×	×	×	0.8049	0.7731	0.7922	0.7667	0.8047	0.7939
✓	✓	✓	×	×	0.8121	0.7792	0.7957	0.7656	0.8234	0.7917
✓	✓	✓	✓	×	0.8213	0.7929	0.8284	0.8051	0.8299	0.8031
✓	✓	✓	×	✓	0.8285	0.7846	0.8014	0.7716	0.8124	0.7797
✓	×	×	×	✓	0.7697	0.7646	0.7685	0.7428	0.7967	0.7821

Table 5. Results of the ablation study for each component of the element matching augmented model on the EvalMuse-40K test set.

Base	Type Aug.	Score Perturbation	InternVL2.5-8B-MPO	Qwen2.5-VL-7B-Instruct	Ovis2-8B
			ACC	ACC	ACC
✓	×	×	0.8224	0.8279	0.8550
✓	✓	×	0.8393	0.8342	0.8608
✓	×	✓	0.8500	0.8355	0.8644
✓	✓	✓	0.8505	0.8367	0.8687

Table 6. Inference latency and throughput of iMatch on RTX4090D with images up to 1024×1024 resolution. Batch size is fixed at 4, as larger batches result in OOM.

Model	Latency (s)	Throughput (image/s)
InternVL2.5-8B-MPO	0.51	7.78
Qwen2.5-VL-7B-Instruct	0.99	4.03
Ovis2-8B	0.54	7.46

relations, achieving the highest improvements. Conversely, employing the image augmentation strategy alone showed effectiveness in enhancing the model’s robustness to image variations, although the element augmentation strategy proved more effective for detailed information capture.

Results on Element Matching Augmented Model. To assess the impact of the two augmentation strategies proposed in Section 3.3, we conduct ablation experiments on the EvalMuse-40K test set. The results are presented in Table 5. Introducing the prompt type augmentation strategy significantly improved this accuracy, highlighting its effectiveness in aiding the model’s understanding of image-text element relationships, thereby boosting predictive performance. Further incorporating the score perturbation strategy yielded the best performance, demonstrating its value in enhancing model generalization and robustness by introducing moderate label variations.

4.4. Inference Cost

In Table 6, we report the inference latency and throughput of iMatch on 4 RTX 4090D GPUs. With a fixed batch size of 4 and input resolution up to 1024×1024, the fastest model achieves a latency of 0.51 seconds and a throughput of 7.78 images per second, demonstrating efficient large-scale image processing and potential for real-world deployment.

5. Conclusion

In this paper, we tackle the challenges of image quality and semantic matching in text-to-image generation models with a multimodal approach called iMatch, which quantifies both overall and detailed matching relationships between images and text descriptions. We developed a framework using instructional fine-tuning, alongside several innovative strategies such as QAlign, validation set augmentation, element augmentation, and image augmentation. Additionally, we introduced prompt type augmentation and score perturbation methods for element matching, enhancing the model’s generalization and precision at the element level. Our comprehensive tests on the EvalMuse-40K dataset and the NTIRE 2025 challenge show that iMatch outperforms existing methods in overall and detailed image-text matching, particularly excelling in semantic understanding and fine-grained evaluations. Future work will focus on refining image-text matching strategies and integrating more advanced multimodal models to broaden the applicability of these methods in practical settings.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1, 3
- [2] Vladimir Arhipkin, Andrei Filatov, Viacheslav Vasilev, Anastasia Maltseva, Said Azizov, Igor Pavlov, Julia Agafonova, Andrey Kuznetsov, and Denis Dimitrov. Kandinsky 3.0 technical report. *arXiv preprint arXiv:2312.03511*, 2023. 1, 3
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023. 1, 3
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 3, 6, 7
- [5] blackforestlabs. Flux1.1. <https://blackforestlabs.ai/>, 2024. 1, 3
- [6] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 1, 3, 6, 7
- [7] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *CVPR*, 2024. 1, 3
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 6
- [9] DreaminaAI. Dreamina. <https://dreamina.capcut.com/>, 2023. 1, 3
- [10] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *ICML*, 2024. 1
- [11] Yutong Feng, Biao Gong, Di Chen, Yujun Shen, Yu Liu, and Jingren Zhou. Ranni: Taming text-to-image diffusion for accurate instruction following. *arXiv preprint arXiv:2311.17002*, 2023. 3
- [12] Shuhao Han, Haotian Fan, Jiachen Fu, Liang Li, Tao Li, Junhui Cui, Yunqiu Wang, Yang Tai, Jingwei Sun, Chunle Guo, et al. Evalmuse-40k: A reliable and fine-grained benchmark with comprehensive human annotations for text-to-image generation model evaluation. *arXiv preprint arXiv:2412.18150*, 2024. 1, 2, 3, 4, 6, 7
- [13] Shuhao Han, Haotian Fan, Fangyuan Kong, Wenjie Liao, Chunle Guo, Chongyi Li, Radu Timofte, et al. NTIRE 2025 challenge on text to image generation model quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025. 2
- [14] Wanggui He, Siming Fu, Mushui Liu, Xierui Wang, Wenyi Xiao, Fangxun Shu, Yi Wang, Lei Zhang, Zhelun Yu, Haoyuan Li, et al. Mars: Mixture of auto-regressive models for fine-grained text-to-image synthesis. *arXiv preprint arXiv:2407.07614*, 2024. 3
- [15] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 3, 6, 7
- [16] David Holz. Midjourney. <https://www.midjourney.com>, 2023. 1, 3
- [17] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022. 6
- [18] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *ICCV*, 2023. 1, 3, 7
- [19] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. In *NeurIPS*, 2023. 1, 3, 6, 7
- [20] Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei, Kewen Wu, Xide Xia, Pengchuan Zhang, Graham Neubig, and Deva Ramanan. Evaluating and improving compositional text-to-visual generation. In *CVPR*, 2024. 6, 7
- [21] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linniao Xu, and Suhail Doshi. Playground v2. 5: Three insights towards enhancing aesthetic quality in text-to-image generation. *arXiv preprint arXiv:2402.17245*, 2024. 1, 3
- [22] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 1, 7
- [23] Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xincheng Deng, Yingfang Zhang, Xingchao Liu, Minbin Huang, Zedong Xiao, et al. Hunyuan-dit: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. *arXiv preprint arXiv:2405.08748*, 2024. 1, 3
- [24] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 1, 3
- [25] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *CVPR*, 2024.
- [26] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 7
- [27] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li,

- Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 1, 3
- [28] Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. Ovis: Structural embedding alignment for multimodal large language model. *arXiv preprint arXiv:2405.20797*, 2024. 1, 3, 6, 7
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 1, 3
- [30] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1
- [31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [32] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *ECCV*, 2024. 1, 3
- [33] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. In *NeurIPS*, 2023. 3
- [34] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 1, 3
- [35] Kolos Team. Kolos: Effective training of diffusion model for photorealistic text-to-image synthesis. *arXiv preprint*, 2024. 3
- [36] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022. 2, 5
- [37] Olivia Wiles, Chuhan Zhang, Isabela Albuquerque, Ivana Kajić, Su Wang, Emanuele Bugliarello, Yasumasa Onoe, Chris Knutsen, Cyrus Rashtchian, Jordi Pont-Tuset, et al. Revisiting text-to-image evaluation with gecko: On metrics, prompts, and human ratings. *arXiv preprint arXiv:2404.16820*, 2024. 1
- [38] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023. 3
- [39] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, et al. Q-align: Teaching llms for visual scoring via discrete text-defined levels. In *ICML*, 2024. 2, 4
- [40] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. 6, 7
- [41] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. In *NeurIPS*, 2023. 1, 3, 6, 7
- [42] Michal Yarom, Yonatan Bitton, Soravit Changpinyo, Roei Aharoni, Jonathan Herzig, Oran Lang, Eran Ofek, and Idan Szepeski. What you see is what you read? improving text-image alignment evaluation. In *NeurIPS*, 2023. 3, 7
- [43] Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multimodal large language models. In *ICLR*, 2024. 7
- [44] Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang, Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. Swift: a scalable lightweight infrastructure for fine-tuning, 2024. 6
- [45] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 1, 3