

ESTR-CoT: Towards Explainable and Accurate Event Stream based Scene Text Recognition with Chain-of-Thought Reasoning

Xiao Wang, *Member, IEEE*, Jingtao Jiang, Qiang Chen, Lan Chen*, Lin Zhu, Yaowei Wang, *Member, IEEE*, Yonghong Tian, *Fellow, IEEE*, Jin Tang

Abstract—Event stream based scene text recognition is a newly arising research topic in recent years which performs better than the widely used RGB cameras in extremely challenging scenarios, especially the low illumination, fast motion. Existing works either adopt end-to-end encoder-decoder framework or large language models for enhanced recognition, however, they are still limited by the challenges of insufficient interpretability and weak contextual logical reasoning. In this work, we propose a novel chain-of-thought reasoning based event stream scene text recognition framework, termed ESTR-CoT. Specifically, we first adopt the vision encoder EVA-CLIP (ViT-G/14) to transform the input event stream into tokens and utilize a Llama tokenizer to encode the given generation prompt. A Q-former is used to align the vision token to the pre-trained large language model Vicuna-7B and output both the answer and chain-of-thought (CoT) reasoning process simultaneously. Our framework can be optimized using supervised fine-tuning in an end-to-end manner. In addition, we also propose a large-scale CoT dataset to train our framework via a three stage processing (i.e., generation, polish, and expert verification). This dataset provides a solid data foundation for the development of subsequent reasoning-based large models. Extensive experiments on three event stream STR benchmark datasets (i.e., EventSTR, WordArt*, IC15*) fully validated the effectiveness and interpretability of our proposed framework. The source code and pre-trained models will be released on <https://github.com/Event-AHU/ESTR-CoT>.

Index Terms—Scene Text Recognition; Large Language Models; Chain-Of-Thought Reasoning; Event Camera; Explainable Artificial Intelligence

I. INTRODUCTION

SCENE Text Recognition (STR) targets to understand and recognize the words in the given scene using a machine learning model. It has been widely exploited based on the RGB frame cameras and achieves significant improvements with the help of deep neural networks. This task can be applied in

autonomous driving, augmented reality, document digitization, retail and e-commerce, etc. However, the performance of STR in challenging scenarios remains unsatisfactory due to the use of RGB cameras, such as fast motion and extreme illumination. Therefore, there is still a long way to go in the research of robust STR models.

Recently, Large Language Models (LLMs) have achieved great success in the Natural Language Processing (NLP) community, such as ChatGPT [1], GPT-4o [2], DeepSeek [3], Qwen [4], etc. The LLMs are also introduced into the multi-modal scenarios and can be applied to visual question answering [5]–[8], medical report generation [9]–[12], etc. Some researchers also adopt the LLMs for the scene text recognition [13]–[18] and achieve significant improvements over non-LLMs based STR models, as shown in Fig. 1 (a, b). More in detail, TextMonkey [13] proposes a high-resolution, location-aware LLM to unify OCR and VQA tasks; DocPedia [14] directly processes frequency-domain document images without relying on traditional OCR; and Vary [15] extends visual vocabulary and enhances multilingual document understanding, highlighting the potential of LLMs in fine-grained text recognition. Some researchers resort to event cameras for the perception in extremely challenging scenarios (i.e., low illumination, fast motion) to replace or assist RGB cameras, e.g., event-based object detection [19]–[23], tracking [24]–[27], and scene text recognition [18]. EventSTR [18] proposed by Wang et al. introduces a large-scale benchmark dataset for event-based text recognition and proposes SimC-ESTR, a framework that combines event based vision with large language models through vision-text alignment, memory-enhanced reasoning, and glyph-level correction, achieving superior robustness under challenging visual conditions.

Despite these breakthroughs, we can find that these scene text recognition models are still limited by the following issues: 1). Mainstream STR algorithms typically use RGB frames as input, making them susceptible to challenges such as low illumination, high-speed motion, and overexposure. 2). Mainstream STR algorithms lack interpretability and strong reasoning abilities, even when adopting the Large Language Models. It is difficult to explicitly model relationships between different characters, however, the context-based reasoning is key to achieving high-performance recognition in challenging scenarios. 3). The reasoning ability relies on the release of large-scale, high-quality chain-of-thought datasets. However, there is still no dataset in academia specifically for event

• Xiao Wang, Jingtao Jiang, Qiang Chen, Jin Tang are with the School of Computer Science and Technology, Anhui University, Hefei 230601, China. (email: {xiaowang, tangjin}@ahu.edu.cn, e23301220@stu.ahu.edu.cn, jingtaj16@gmail.com)

• Lan Chen is with the School of Electronic and Information Engineering, Anhui University, Hefei 230601, China. (email: chenlan@ahu.edu.cn)

• Lin Zhu is with Beijing Institute of Technology, Beijing, China. (email: lingzhu@pku.edu.cn)

• Yaowei Wang is with Harbin Institute of Technology, Shenzhen, China; Peng Cheng Laboratory, Shenzhen, China. (email: wangyw@pcl.ac.cn)

• Yonghong Tian is with Peng Cheng Laboratory, Shenzhen, China; School of Computer Science, Peking University, China; School of Electronic and Computer Engineering, Shenzhen Graduate School, Peking University, China (email: yhtian@pku.edu.cn)

* Corresponding Author: Lan Chen (chenlan@ahu.edu.cn)

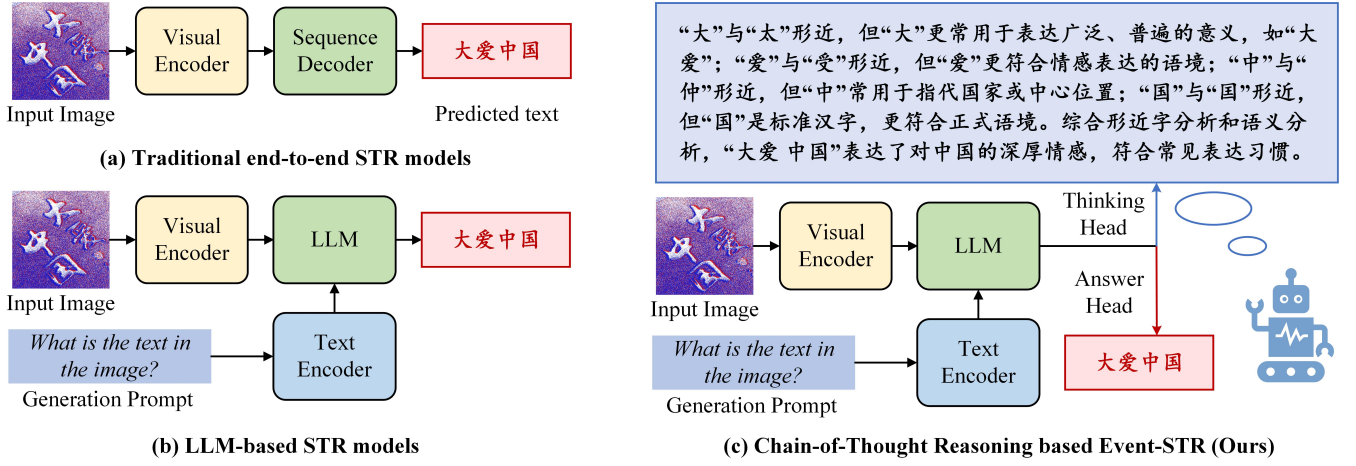


Fig. 1. Comparison between the (a) traditional end-to-end encoder-decoder STR, (b) LLM-based STR, and (c) our newly proposed chain-of-thought reasoning based event stream STR. Traditional STR models rely on task-specific architectures and direct visual-text mappings. LLM-based approaches introduce language understanding but still lack explicit reasoning. Our framework incorporates chain-of-thought reasoning to enable interpretable and logically grounded recognition.

stream based scene text recognition, which further limits the reasoning capabilities of existing models. Therefore, it is natural to raise the following question: “*how to design an event-stream based scene text recognition framework that can effectively integrate contextual information to enable high-quality reasoning, leveraging the capabilities of large language models to achieve more accurate recognition?*”

To address the aforementioned issues, in this paper, we propose a novel chain-of-thought reasoning based event stream scene text recognition framework, termed ESTR-CoT. As shown in Fig. 1 (c), the framework we propose is built upon a large language model, with the core innovation being the introduction of a reasoning task head, which enables high-quality contextual reasoning while simultaneously achieving accurate text recognition. Specifically, given the input event streams, we first adopt a vision encoder (EVA-CLIP (ViT-G/14) [28] adopted in this paper) to transform them into event tokens and use a Llama-Tokenizer [29] to encode the given generation prompt. We concatenate these tokens with the randomly initialized query tokens and feed them into a Q-former network to align the vision tokens for the large language models. Then, we adopt the pre-trained large language model Vicuna-7B [30] to decode the *thinking* and *answer* simultaneously. To bridge the data gaps and achieve the training of our framework, as shown in Fig. 2, we resort to the large language model DeepSeek-V3 to generate the chain-of-thought reasoning data based on the ground truth answers. The additional evaluation on the generated data by the LLMs and human experts is also needed for high-quality CoT data generation. Some representative samples are provided in Fig. 3, and an overview of our proposed framework can be found in Fig. 4.

To sum up, the main contributions of this paper can be summarized as the following three aspects:

1). We propose a novel event stream scene text recognition framework based on chain-of-thought reasoning, termed ESTR-COT. The algorithm we propose can significantly en-

hance the interpretability of scene text recognition, improving recognition accuracy through contextual reasoning.

2). We build a large-scale reasoning dataset for the event stream based scene text recognition. It contains 16,222 image-reasoning pairs, each consisting of an event-based scene text image, the corresponding recognized text (answer), and a detailed CoT rationale explaining how the model arrives at the final answer.

3). Extensive experiments on three widely used event-based STR benchmark datasets (i.e., EventSTR [18], WordArt* [31], IC15* [32]) fully validated the effectiveness and interpretability of our proposed reasoning strategy.

The rest of this paper is organized as follows: We review the related works on the scene text recognition, large language model and reasoning models in section II. In section III, we introduce our proposed reasoning strategy for the large language model based scene text recognition. In section IV, we conduct the experimental analysis and provide both quantitative and qualitative analyses. Finally, in section V, we conclude this paper and propose possible research directions on this work.

II. RELATED WORKS

In this section, we introduce the related works from Scene Text Recognition, LLM and Reasoning Models, Event-based Vision. More details can be found in the following surveys [33]–[35] and paper list ¹.

A. Scene Text Recognition

Scene text recognition [36]–[39] naturally involves both vision and language processing. Traditional research efforts often emphasized either visual feature extraction or language modeling, but recent advances have sought a more balanced integration of both modalities to enhance robustness across diverse scenarios. E2STR [40] enhances adaptability by introducing context-rich text sequences and a context training

¹https://github.com/Event-AHU/OCR_Paper_List

strategy, offering improved flexibility across environments. CCD [41] leverages a self-supervised segmentation module and character-to-character distillation to improve text representation learning, while SIGA [42] further refines segmentation via implicit attention alignment. CDistNet [43] incorporates both visual and semantic positional embeddings into its transformer-based design to address irregular text layouts and complex backgrounds.

In parallel, several works introduce iterative error correction strategies using language models. VOLTER [44], BUS-Net [45], MATRNet [46], LevOCR [47], and ABINet [48] exemplify this trend, integrating language models to refine recognition through feedback loops, improving both robustness and interoperability. Building on this foundation, recent research has moved toward large language model (LLM)-based scene text recognition. These models exploit the generative and contextual reasoning capabilities of LLMs to unify visual and linguistic understanding. For instance, TextMonkey [13] is a multimodal LLM optimized for text-centric tasks, utilizing high-resolution inputs and location-aware responses to support enhanced interaction and interpretability. DocPedia [14] eliminates the need for traditional OCR by processing high-resolution document images in the frequency domain, efficiently capturing both visual and textual cues. Vary [15] enriches the visual vocabulary of vision-language models, enabling fine-grained document OCR and chart understanding, especially in multilingual contexts. mPLUG-DocOwl 1.5 [16] introduces Unified Structure Learning for improved document layout comprehension, while OCR2.0 [17] presents a powerful 580M-parameter model capable of handling OCR tasks ranging from text recognition to formula parsing. Despite these advances, LLM-based models still struggle under extreme conditions such as low lighting, blur, or complex noise. These limitations highlight the ongoing challenge of achieving robust, generalizable scene text recognition across varied and adverse environments. To address these challenges, EventSTR [18] pioneers the use of event cameras for scene text recognition.

B. LLM and Reasoning Models

Large Language Models (LLMs) are pretrained models exhibiting strong capabilities in natural language understanding and generation. They serve as the foundational backbone for many advanced reasoning techniques by enabling coherent and contextually relevant text generation. Notable examples include the GPT series from OpenAI [49], [50], Google's PaLM [51], Meta's LLaMA [29], as well as instruction-tuned variants such as InstructGPT [52] and Vicuna [30]. Although these models differ in scale and training data, they share a common versatility, being adaptable to a broad spectrum of downstream tasks, including reasoning, dialogue systems, and code generation.

Building upon LLMs, reasoning models [53] explicitly structure the inference process to enhance interpretability and accuracy. One effective technique is Chain of Thought (CoT) prompting, which encourages the model to generate intermediate, step-by-step rationales instead of only final answers [54]. This method substantially improves performance

on complex reasoning tasks by rendering the model's decision-making process more transparent. Further extending CoT, frameworks such as Tree-of-Thoughts and Graph-of-Thoughts facilitate the parallel exploration and integration of multiple reasoning paths, thereby increasing robustness, especially in scenarios involving ambiguity or multiple valid solutions [55], [56]. To reduce noise and hallucinations inherent in single-chain reasoning, self-consistency approaches aggregate outputs from multiple independent reasoning chains [57]. Additionally, reinforcement learning techniques optimize reasoning processes by refining reward functions that emphasize coherence, verifiability, and task-specific criteria, exemplified by models like LLaVA-Reasoner [58]. Complementing these advances, the integration of vision-language foundation models such as BLIP-2 and MiniGPT-4 provides powerful multimodal encoders, enabling reasoning pipelines to extend beyond text to encompass images, audio, and other modalities [7], [59]. Collectively, these developments constitute a versatile and comprehensive toolkit for eliciting interpretable, high-quality reasoning across diverse application domains.

C. Event-based Vision

Event cameras [60]–[62] are bio-inspired sensors that asynchronously record per-pixel brightness changes with microsecond-level latency, offering high temporal resolution, low power consumption, and wide dynamic range. These properties make event cameras well-suited for tasks in fast-moving or poorly lit environments, including autonomous driving, robotic perception, and medical imaging.

In human activity recognition, ESTF [63] projects raw event streams into learned spatial-temporal embeddings, enabling robust classification under rapid motion and low illumination. For visual tracking, EventVOT [27] provides a large-scale (1280×720) event-only dataset of 1,141 sequences spanning pedestrians, vehicles, drones, and sports objects, alongside a hierarchical distillation framework that yields high-speed, low-latency tracking. Recurrent Vision Transformers (RVTs) [19] exploit event cameras' temporal fidelity to achieve accurate object detection in dynamic scenes. SAFE [64] fuses event streams with RGB frames and semantic labels via a pretrained vision-language backbone, bridging the modality gap and overcoming the limitations of small-scale networks. Despite these advances, event-based methods have seen little application in scene text recognition. To fill this gap, EventSTR [18] introduces a novel task and dataset for event-stream-based STR, comprising 9,928 high-definition (1280×720) event sequences annotated with both Chinese and English text under diverse lighting, motion, and occlusion conditions.

III. OUR PROPOSED APPROACH

A. Overview

In this paper, we propose ESTR-CoT, a reasoning-enhanced framework designed to improve both the accuracy and interpretability of scene text recognition in visually challenging environments. The core idea is to enable the model to produce

not only the correct answer but also a coherent reasoning process, thereby improving transparency and interoperability. As shown in Fig. 4, ESTR-CoT consists of the following four core components, i.e., visual encoder, Q-Former, pre-trained LLM, and prompt-based control architecture. Specifically, the visual encoder is a pre-trained EVA-CLIP (ViT-G/14) [28] which is used to extract discriminative visual features from event-based scene images. Q-Former aligns visual features with a tokenized textual prompt to generate query embeddings. These embeddings act as a bridge between the visual encoder and the language model, capturing the multi-modal context. A large language model is the foundation of our framework, which receives the projected visual features and prompt embeddings to auto-regressively generate textual outputs. It produces either the final recognition result or the corresponding reasoning chain. For the prompt-based control architecture, different prompt suffixes (`<answer>` and `<thinking>`) are used to steer the generation target. During training, the model is supervised with both answer annotations and reasoning chains from high-quality CoT data. A multi-task loss is applied to jointly optimize both objectives. At inference time, ESTR-CoT can generate only the final answer for efficiency or include the reasoning chain for better transparency.

B. Chain-of-Thought Data Generation

To enhance the reasoning capability of large vision-language models in the context of scene text recognition, we construct a high-quality Chain-of-Thought dataset via a structured multi-stage pipeline, as illustrated in Fig. 2. Unlike conventional text recognition annotations that only provide flat textual labels, our CoT corpus incorporates fine-grained visual-semantic reasoning, enabling models to explicitly consider ambiguities such as lookalike characters, semantic context, and domain-specific constraints.

• **Stage 1: Initial CoT Generation.** Given a raw dataset of OCR labels or question-answer pairs (e.g., the text string “LOVEL” predicted from an image), we first utilize a powerful LLM (e.g., Deepseek-V3) to generate an initial CoT explanation. The LLM is prompted with task-specific instructions that require it to analyze both visual similarity (e.g., “LOVE” vs. “NOVEL”) and semantic coherence. The resulting outputs are structured in the format:

```
<answer>LOVEL</answer><thinking>"LOVEL"
could be a stylized version of
"LOVE" or a misspelling of "NOVEL".
The letters "L", "O", "V", "E", and
"L" are clearly present. Lookalike
words such as "LEVEL" or "LOVELY"
are considered but ruled out due
to differences in letter count and
semantic context.</thinking>
```

This produces the first-stage dataset D_1 , which consists of (input text, reasoning) pairs.

• **Stage 2: Automatic Evaluation and Rewriting.** To guarantee logical validity, informativeness, and structural clarity of the generated reasoning, we employ an automatic evaluation module that examines each CoT explanation against the following criteria:

- (1) *Length Constraint*: reasoning chains exceeding a preset maximum token length (e.g., 100 tokens) are flagged for rewriting, preventing overly verbose or unfocused explanations;
- (2) *Visual-Semantic Completeness*: the explanation must explicitly include both visual form analysis (e.g., character shape comparisons) and semantic context reasoning (e.g., word plausibility), ensuring comprehensive reasoning coverage;
- (3) *Logical Consistency*: the chain must maintain coherent argumentative flow without contradictions or irrelevant content.

Samples passing all checks are directly added to the curated dataset D_2 . Failed cases are sent to an LLM-powered rewriter module for improved reformulation. The rewritten explanations undergo re-evaluation and are included in D_2 only upon satisfying the quality standards. This iterative process effectively filters out noisy or superficial chains of thought while enhancing reasoning quality.

• **Stage 3: Expert Review and Validation.** While the automatic modules are highly effective, some complex or ambiguous cases require human judgment. Therefore, samples in D_2 are reviewed by human experts with background in scene text understanding. The experts assess correctness, clarity, and linguistic fluency. Only the most reliable samples are retained to construct the final dataset D_3 , which serves as the training source for CoT-aware STR models.

Consider the OCR output “LOVEL”, which can be ambiguously interpreted. Initially, the CoT reasoning explores multiple candidate words: it could be a stylized version of “LOVE” or a misspelling/truncation related to “NOVEL”. The letters “L”, “O”, “V”, “E”, and “L” are clearly present, while lookalike words such as “LEVEL” or “LOVELY” are also considered but ruled out due to inconsistencies in letter count or semantic context.

Further analysis focuses on the visual similarity of individual letters and their semantic alignment. Although “LOVEL” shares visual traits with “LEVEL” and “NOVEL”, the latter candidates are discarded based on their differing meanings—“LEVEL” denotes a stage or flat surface, and “NOVEL” refers to a type of book, which do not align with the romantic or affectionate context suggested by the imagery. Consequently, the reasoning favors “LOVE” as the most plausible interpretation, viewing “LOVEL” as a stylized or truncated variant emphasizing romantic connotations.

As shown in Fig. 2, this case exemplifies how the Stage 2 automatic evaluation and rewriting module refines ambiguous initial explanations by emphasizing letter shape details and semantic context, resulting in a more precise, semantically coherent, and visually robust Chain-of-Thought.

Benefits. This pipeline ensures that the final CoT dataset not only captures domain-specific reasoning patterns but also minimizes annotation noise. It supports multiple downstream tasks, including OCR error correction, few-shot generalization, and instruction-tuned STR modeling. By incorporating both machine-filtered and human-verified samples, we achieve a balance between scalability and quality, which is essential for training reliable reasoning-augmented vision models. For-

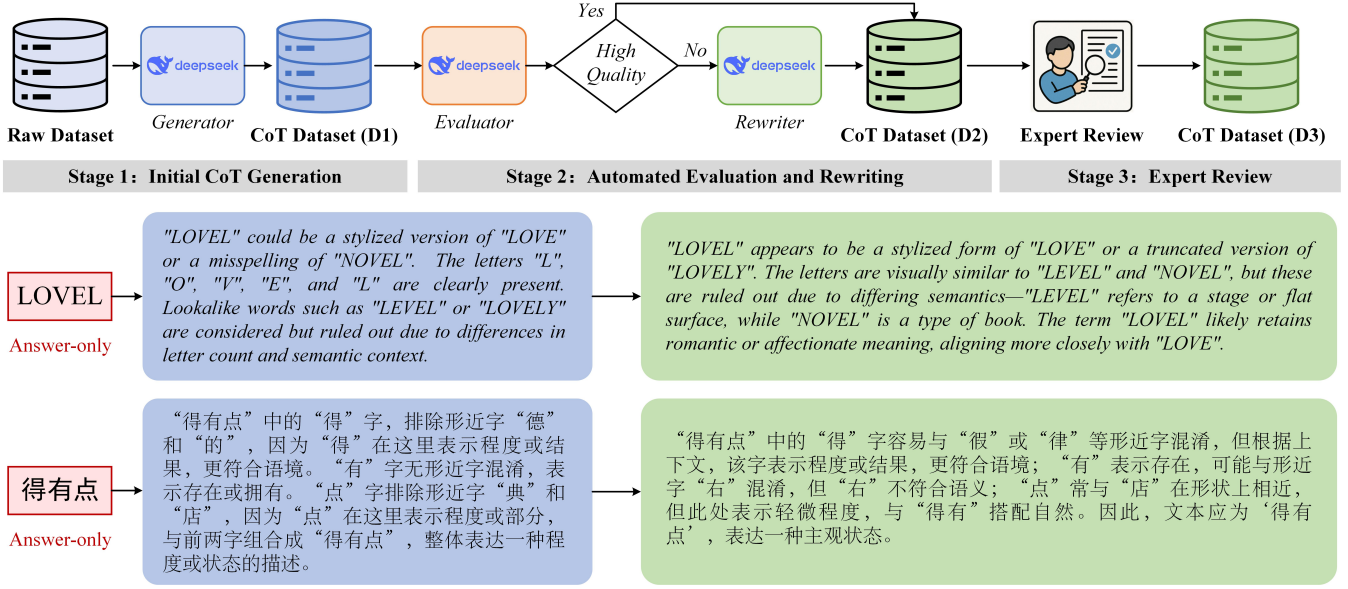


Fig. 2. **Data Generation Pipeline of ESTR-CoT.** The figure illustrates a three-stage pipeline for constructing high-quality CoT data in ESTR-CoT. Beginning with a raw dataset, a generator produces initial CoT responses (D1). An evaluator filters these outputs: high-quality samples proceed to D2, while suboptimal ones are revised by a rewriter based on evaluator feedback. All refined samples are added to D2. Finally, human experts review D2 to curate the final expert-approved dataset D3. The bottom example shows how reasoning is refined from a basic explanation (Stage 1) to a more precise and semantically coherent version (Stage 2).

mally, our pipeline constructs a reasoning-augmented dataset through the following process:

$$D_1 = \{(a_i, c_i^{(0)})\}_{i=1}^N \quad (1)$$

$$D_2 = \{(a_i, c_i^{(1)}) \mid \text{Eval}(c_i^{(0)}) = \text{Pass or Revised}\} \quad (2)$$

$$D_3 = \{(a_i, c_i^*) \mid (a_i, c_i^{(1)}) \in D_2, \text{HumanPass}(c_i^{(1)}) = \text{True}\} \quad (3)$$

Here, a_i denotes the original OCR output or answer, $c_i^{(0)}$ is the initial CoT generated by the LLM, $c_i^{(1)}$ is the reasoning after automatic evaluation and possible rewriting, and c_i^* is the final expert-approved reasoning.

The automatic evaluation function $\text{Eval}(\cdot)$ checks each reasoning chain c as:

$$\text{Eval}(c) = \begin{cases} \text{Pass,} & \text{if } \text{Len}(c) < L_{\max} \\ & \wedge \text{HasVisual}(c) \\ & \wedge \text{HasSemantic}(c) \\ \text{Fail,} & \text{otherwise} \end{cases} \quad (4)$$

where L_{\max} is the maximum allowed length, and the predicates $\text{HasVisual}(\cdot)$ and $\text{HasSemantic}(\cdot)$ ensure the reasoning contains both visual and semantic analysis.

Each final sample in D_3 is formatted as:

$$x_i = \langle \text{answer} \rangle a_i \langle / \text{answer} \rangle \langle \text{thinking} \rangle c_i^* \langle / \text{thinking} \rangle \quad (5)$$

which allows straightforward downstream usage in training CoT-aware models. An illustration of the dataset structure and examples is shown in Fig. 3.

C. Model Design

After acquiring high-quality Chain-of-Thought data D_3 via the previously proposed data generation pipeline, we leverage this reasoning-enriched dataset to improve the model's reasoning ability. As illustrated in Fig. 4, we first explain the motivation for producing two separate outputs answer and think, and then outline two architectural variants developed to effectively model the interaction between visual and textual modalities.

• **Input Construction.** Given an event-based scene image $\mathcal{I} \in \mathbb{R}^{C \times H \times W}$, we use a pre-trained visual encoder, EVA-CLIP [28] (ViT-G/14), to extract discriminative visual features. The encoder divides the image \mathcal{I} into fixed-size patches (14×14 px), flattens them into tokens, and applies multi-head self-attention to obtain the visual feature map F_v . A global [CLS] token provides a holistic representation of the image. To combine the visual features with textual information, we introduce a Q-Former module that aligns F_v with tokenized prompt embeddings F_ℓ to produce query embeddings F_q . The textual prompt is set to:

$$\mathcal{P} = \text{"What is the text in the image?"},$$

which is then tokenized and encoded into embeddings F_ℓ . The Q-Former processes the visual features F_v and prompt embeddings F_ℓ to generate the query embeddings F_q .

The final model input is constructed by combining the visual and textual components:

$$\text{Input}_{\text{LLM}} = [\text{Proj}(F_q), \text{Proj}(F_v), F_\ell] \quad (6)$$

where F_ℓ is the tokenized prompt embedding, $\text{Proj}(F_q)$ is the projection of query embeddings, and $\text{Proj}(F_v)$ is the projection of the visual feature map. This input is used for

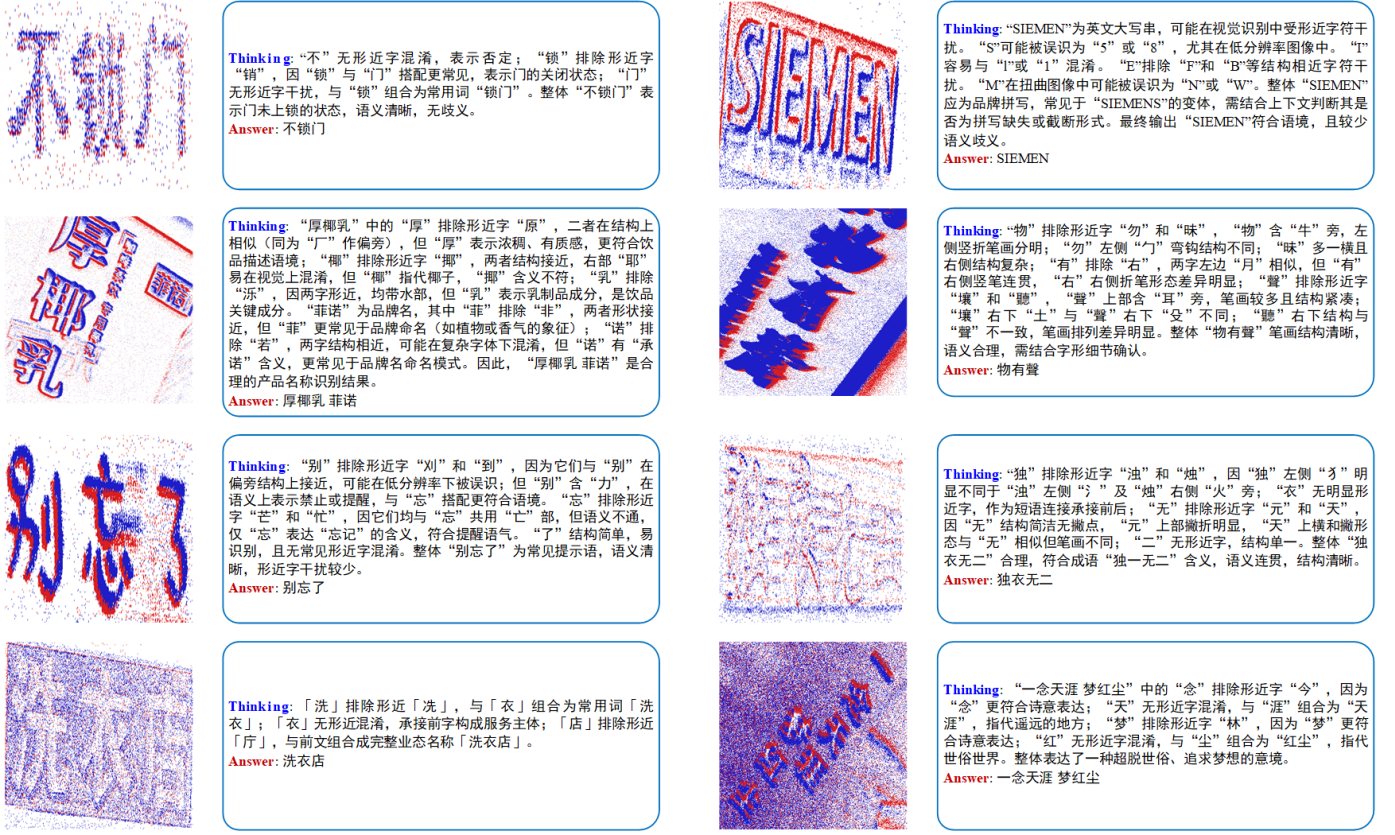


Fig. 3. Representative samples from our newly proposed CoT ESTR Dataset. Each example consists of an event stream scene image, a recognition result **<Answer>**, and an accompanying reasoning chain **<Thinking>**.

both generating the answer and the reasoning chain in the subsequent architecture.

• Why Separate Outputs for answer and thinking?

The decision to output both answer and thinking arises from the need to balance high recognition accuracy with interpretability in complex scene text recognition tasks.

1). *Improved Interpretability*: By outputting thinking, the model provides not only the correct answer but also articulates the reasoning behind its decision. This transparency helps users understand the model’s decision-making process, which is especially valuable in domains like document analysis or autonomous systems where reasoning is crucial.

2). *Enhanced Model Robustness*: Separate outputs allow the model to treat answer generation and reasoning as distinct tasks. This enables better management of the interplay between visual features and textual context. The thinking output ensures that the answer is contextually grounded and logically sound.

3). *Multi-task Learning Synergy*: By training the model to predict both the answer and its reasoning, the model learns to optimize the reasoning and answering process simultaneously. This dual-head setup helps refine the model’s understanding, improving both the accuracy and coherence of the output.

• **Prompt-based Control Architecture**. We adopt a prompt-based control approach to guide the model’s generation behavior. Specifically, we use a shared prompt prefix for both the final text output (answer) and the reasoning chain

(thinking), distinguishing the two by appending different suffixes:

$$\mathcal{P}_{\text{answer}} = \mathcal{P} \text{ <answer>}, \quad \mathcal{P}_{\text{thinking}} = \mathcal{P} \text{ <thinking>} \quad (7)$$

These prompts are tokenized and encoded into embeddings:

$$F_{\ell_{\text{answer}}} = \text{Encode}(\mathcal{P}_{\text{answer}}), \quad F_{\ell_{\text{thinking}}} = \text{Encode}(\mathcal{P}_{\text{thinking}}) \quad (8)$$

The model input for both tasks is formed by concatenating the projected query embeddings $\text{Proj}(F_q)$, projected visual features $\text{Proj}(F_v)$, and the corresponding prompt embeddings:

$$\text{Input}_{\text{LLM,answer}} = [\text{Proj}(F_q), \text{Proj}(F_v), F_{\ell_{\text{answer}}}] \quad (9)$$

$$\text{Input}_{\text{LLM,thinking}} = [\text{Proj}(F_q), \text{Proj}(F_v), F_{\ell_{\text{thinking}}}] \quad (10)$$

The model then autoregressively generates the answer and reasoning outputs respectively:

$$\hat{y}_{\text{answer}} = \text{LLM}(\text{Input}_{\text{LLM,answer}}), \quad (11)$$

$$\hat{y}_{\text{thinking}} = \text{LLM}(\text{Input}_{\text{LLM,thinking}}). \quad (12)$$

• **Projection-separated Control Architecture**. We also explore an alternative architecture that employs separate projection layers for visual and query embeddings for the answer and reasoning tasks, respectively. This variant aims to spe-

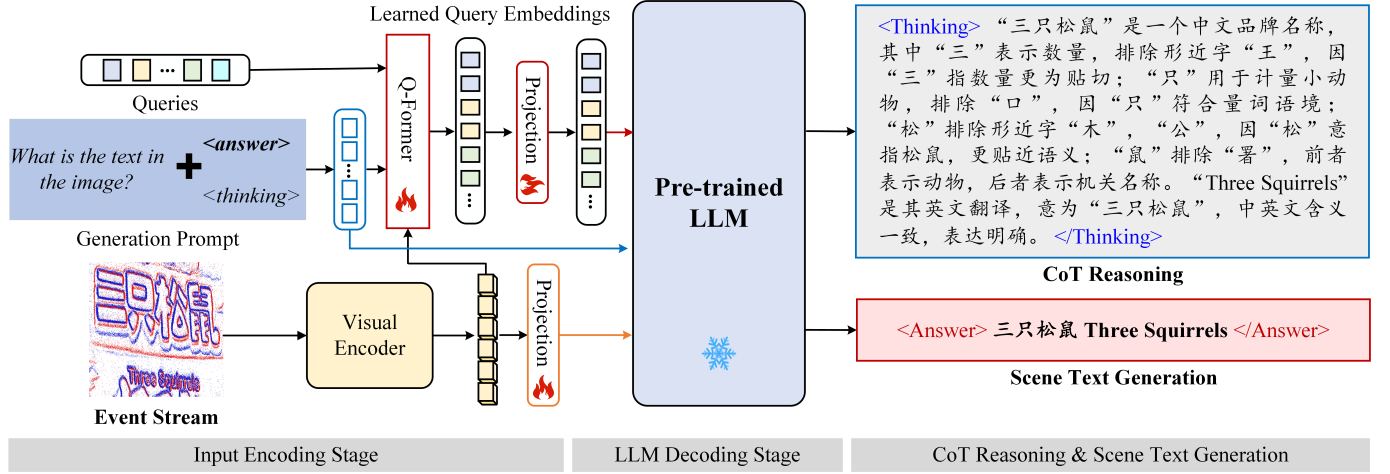


Fig. 4. An overview of the proposed ESTR-CoT framework for the reasoning-based event stream scene text recognition. This figure illustrates the prompt-based control architecture for generating answers and reasoning chains from event-based scene text inputs. A Q-Former extracts learned query embeddings from visual features captured by the event stream. These embeddings are fused with a fixed textual prompt and passed to a pre-trained LLM. By appending different suffixes (<answer> or <thinking>) to the prompt, the model is guided to generate either the final recognition result or the corresponding reasoning chain.

cialize feature processing for each output by using distinct projections:

$$\text{Input}_{\text{LLM,answer}} = [\text{Proj}_{\text{ans}}(F_q), \text{Proj}_{\text{ans}}(F_v), F_\ell], \quad (13)$$

$$\text{Input}_{\text{LLM,thinking}} = [\text{Proj}_{\text{think}}(F_q), \text{Proj}_{\text{think}}(F_v), F_\ell]. \quad (14)$$

The resulting outputs \hat{y}_{answer} and $\hat{y}_{\text{thinking}}$ are generated in the same autoregressive manner as above. In this work, we primarily focus on the prompt-based control architecture due to its simplicity and efficiency. The projection-separated variant is briefly introduced here and further analyzed in the experiments section IV-D.

D. Head and Loss Function

To train the model to generate both the final answer and its corresponding reasoning process, we define the total loss as the sum of two components:

$$\mathcal{L} = \mathcal{L}_{\text{answer}} + \mathcal{L}_{\text{thinking}}, \quad (15)$$

where $\mathcal{L}_{\text{answer}}$ is the cross-entropy loss for answer prediction, and $\mathcal{L}_{\text{thinking}}$ is the loss for generating the reasoning chain. Specifically,

$$\mathcal{L}_{\text{answer}} = \text{CE}(\hat{y}_j^{\text{answer}}, y_j), \quad (16)$$

$$\mathcal{L}_{\text{thinking}} = \text{CE}(\hat{y}_j^{\text{thinking}}, c_j), \quad (17)$$

where $\hat{y}_j^{\text{answer}}$ denotes the predicted answer and y_j is the corresponding ground-truth label. Similarly, $\hat{y}_j^{\text{thinking}}$ is the generated reasoning chain, and c_j is the ground-truth chain-of-thought. Although we explored weighted combinations of the two loss components in our experiments (see Section IV-D), we found that the simple summation strategy performs comparably and is more stable. Therefore, we adopt the unweighted sum as the default objective for training.

IV. EXPERIMENTS

A. Dataset and Evaluation Metric

We evaluate our method on three event-based scene text recognition benchmark datasets, i.e., **WordArt***², **IC15***³, and **EventSTR** [18]. These datasets are consistent with the evaluation protocol in our previous work [18], allowing for direct comparison and fair benchmarking.

- **WordArt***: Converted from the original WordArt [31] dataset using the event simulator ESIM [69], this dataset includes artistic text samples from posters, greeting cards, covers, and handwritten notes. It contains 4,805 training images and 1,511 validation images.

- **IC15***: Derived from the ICDAR2015 [32] dataset and converted to event-based images. It includes 4,468 training samples and 2,077 test samples from natural scenes.

- **EventSTR**: Proposed by Wang et al., is a large-scale benchmark dataset for event-based scene text recognition, encompassing diverse character layouts, motion patterns, and illumination conditions. It provides a comprehensive testbed for evaluating both recognition accuracy and temporal robustness, and serves as the foundation for exploring LLM-based reasoning under challenging visual dynamics.

We follow the same evaluation setting as in [18]. For EventSTR, we report BLEU-1 to BLEU-4 scores to assess both linguistic accuracy and reasoning quality. BLEU is computed at the character level for Chinese and word level (case-insensitive) for English. For WordArt* and IC15*, we report word-level recognition accuracy.

B. Implementation Details

We build our framework upon the BLIVA architecture [8], which serves as our vision-language baseline. The model is initialized with pre-trained weights from BLIVA and then

²<https://opendatalab.com/OpenDataLab/WordArt>

³<https://aistudio.baidu.com/datasetdetail/96799>

TABLE I
COMPARISON OF BLEU SCORES WITH SOTA METHODS ON THE EVENTSTR DATASET.

Algorithm	Publish	Backbone	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Params(M)	Code
CCD [41]	ICCV 2023	ViT	0.365	0.254	0.172	0.145	52.0	URL
SIGA [42]	CVPR 2023	ResNet	0.434	0.393	0.346	0.307	40.4	URL
CDistNet [43]	IJCV 2023	ResNet+Transformer	0.333	0.242	0.157	0.135	65.5	URL
PARSeq [65]	ECCV 2022	ViT	0.450	0.357	0.281	0.224	23.4	URL
MGP-STR [66]	ECCV 2022	Transformer	0.427	0.339	0.278	0.232	148.0	URL
GOT-OCR2.0 [17]	arXiv 2024	ViT	0.426	0.390	0.358	0.332	580.0	URL
BLIVA [8]	AAAI 2024	ViT	0.584	0.528	0.450	0.386	7531.3	URL
SimC-ESTR [18]	arXiv 2025	ViT	0.638	0.583	0.500	0.430	7531.3	URL
ESTR-CoT (Ours)	-	ViT	0.648	0.586	0.500	0.430	7531.3	URL

TABLE II
THE ACCURACY COMPARISONS WITH SOTA METHODS ON WORDART* AND IC15*.

Algorithm	Publish	Backbone	Accuracy		Params(M)	Code
			WordArt*	IC15*		
LISTER [67]	ICCV 2023	CNN	55.3	69.0	49.9	URL
CCD [41]	ICCV 2023	ViT	62.1	55.4	52.0	URL
SIGA [42]	CVPR 2023	ResNet	69.0	66.2	40.4	URL
CDistNet [43]	IJCV 2023	ResNet+Transformer	66.6	62.3	65.5	URL
DiG [68]	ACM MM 2022	ViT	62.7	53.2	52.0	URL
PARSeq [65]	ECCV 2022	ViT	75.0	72.7	23.4	URL
MGP-STR [66]	ECCV 2022	Transformer	69.6	67.5	148.0	URL
BLIVA [8]	AAAI 2024	ViT	56.7	51.3	7531.3	URL
SimC-ESTR [18]	arXiv 2025	ViT	65.1	56.8	7531.3	URL
ESTR-CoT (Ours)	-	ViT	65.6	57.1	7531.3	URL

fine-tuned on each target dataset individually, allowing it to adapt to domain-specific characteristics in *WordArt**, *IC15**, and *EventSTR*. For optimization, we adopt the AdamW optimizer [70] with hyper-parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, and weight decay of 0.05. The learning rate is warmed up linearly from 10^{-8} to 10^{-5} over the first 1,000 steps, followed by a cosine decay schedule that reduces the learning rate to a minimum of 0. All experiments are conducted on an Nvidia A800 GPU. More implementation details, including prompt construction and loss design, can be found in our source code on GitHub.

C. Comparison on Public Benchmark Datasets

• **Results on EventSTR.** Table I reports BLEU score comparisons on the EventSTR dataset with recent state-of-the-art scene text recognition methods. Traditional approaches such as CCD [41], SIGA [42], and CDistNet [43] obtain relatively low BLEU-4 scores (<0.31), indicating weak reasoning ability. Transformer-based methods like PARSeq [65] and MGP-STR [66] perform better but are not tailored for event-based inputs. GOT-OCR2.0 [17] achieves a BLEU-4 of 0.332, while SimC-ESTR [18] boosts this to 0.430 by incorporating a vision-language backbone. ESTR-CoT consistently improves performance across all BLEU metrics. In particular, it achieves the highest BLEU-1 score of 0.648, matching or surpassing all prior baselines. These results demonstrate the effectiveness of explicit reasoning supervision in event-based scene text recognition.

• **Results on WordArt* and IC15* Datasets.** As shown in Table II, our method, which leverages pre-training on Visual Question Answering (VQA) data, results in a significant

improvement over the baseline. However, the performance on the WordArt* and IC15* datasets is still suboptimal when compared to methods specifically trained on large-scale text recognition datasets like MJ [71] and ST [72]. While VQA pre-training helps the model learn visual-textual relationships, it is not fully optimized for complex text recognition tasks, particularly those involving noisy or highly variable backgrounds. These are areas better addressed by models specifically trained on OCR data.

Additionally, the synthetic datasets used for fine-tuning our model (WordArt* and IC15*) have relatively low resolution and lack the complexity and diversity found in larger OCR datasets. This limits the model’s ability to reach the performance levels of established methods such as LISTER, CCD, and PARSeq, which benefit from extensive, high-quality OCR training data. Despite this, our approach demonstrates a clear improvement over the baseline, highlighting the potential benefits of VQA pre-training for enhancing visual-textual understanding in OCR tasks. In future work, we plan to explore how large-scale datasets like MJ and ST can be leveraged more efficiently without excessive resource consumption, to further boost performance.

D. Ablation Study

• **Analysis on the Performance of Different Architectures.** In this experiment, we compare the Prompt-based Control architecture and the Projection-separated Control architecture, corresponding to the *diff-prompt* and *diff-projection* variants, respectively. The BLEU scores across different n-gram levels (BLEU-1 to BLEU-4) are reported in Table III.

TABLE III
PERFORMANCE OF DIFFERENT ARCHITECTURES.

Architectures	BLEU-1	BLEU-2	BLEU-3	BLEU-4
diff-projection	0.638	0.581	0.498	0.430
diff-prompt	0.648	0.586	0.500	0.430

The Prompt-based Control architecture appends different suffixes (`<answer>` and `<thinking>`) to the common prompt, guiding the model to generate specific outputs related to the final answer or the reasoning process. This results in a more controlled generation process, where the model is explicitly instructed to focus on different aspects of the task (i.e., the answer or the reasoning). From the table, we can see that this architecture yields slightly better BLEU scores compared to the Projection-separated Control architecture, especially in BLEU-1 and BLEU-2. This suggests that the model benefits from having more explicit guidance through the prompt suffix, improving its ability to generate high-quality answers and reasoning chains.

On the other hand, the Projection-separated Control architecture uses distinct projections for both visual and textual components, allowing each modality to be treated separately before combining them. This separation can help the model handle visual and textual information more effectively but may also lead to challenges in generating a coherent final output when the modalities interact. The slightly lower BLEU scores for this architecture suggest that while it might have advantages in certain visual-textual tasks, the overall output quality could be improved with a more explicit prompt-based control mechanism, as seen in the diff-prompt configuration.

In summary, the slight differences in BLEU scores between these two architectures demonstrate the impact of controlling the generation process. The Prompt-based Control method, by directly influencing the output with suffixes, appears to yield better results for text generation tasks, particularly when reasoning is involved. Future work could explore further refinements in combining both strategies for even higher performance.

• **Analysis on CoT filtering.** In this experiment, we evaluate the effectiveness of our CoT filtering pipeline by comparing the performance of models trained with unfiltered CoT data versus those trained with the filtered CoT data. The models were trained using identical configurations except for the inclusion of the CoT filtering process. Specifically, for the unfiltered case, we removed the automatic evaluation and expert validation steps, directly using raw CoT samples for supervision.

The results, shown in Table IV, indicate a noticeable improvement in the model's performance after applying CoT filtering. Specifically, the filtered CoT data leads to a higher BLEU score across all n-gram levels (BLEU-1 to BLEU-4). The unfiltered model achieved a BLEU-1 score of 0.632, while the filtered model showed an improvement to 0.648. Similarly, for BLEU-2, BLEU-3, and BLEU-4, the filtered CoT model outperforms the unfiltered model, achieving scores of 0.586, 0.500, and 0.430, respectively, compared to the unfiltered model's 0.574, 0.492, and 0.423.

TABLE IV
PERFORMANCE COMPARISON WITH AND WITHOUT CoT FILTERING.

Filtering	BLEU-1	BLEU-2	BLEU-3	BLEU-4
✗	0.632	0.574	0.492	0.423
✓	0.648	0.586	0.500	0.430

This demonstrates that the CoT filtering pipeline significantly enhances the quality of the training data, leading to better generalization and more accurate text generation. By removing noisy or irrelevant samples and focusing on high-quality CoT data, the model is able to generate more coherent and contextually accurate reasoning chains, which translates into improved overall performance.

• **Analysis on Loss Weighting.** We investigate the impact of different loss weighting strategies when jointly training the model to generate both the final answer and the reasoning chain. Specifically, we compare two schemes:

1). *Weighted Loss Combination:* The total loss is a weighted sum of the answer and reasoning losses:

$$\mathcal{L} = \lambda \cdot \mathcal{L}_{\text{answer}} + (1 - \lambda) \cdot \mathcal{L}_{\text{thinking}},$$

where $\lambda \in [0, 1]$ controls the relative importance of each objective.

2). *Direct Summation:* The total loss is a simple summation of both components without reweighting:

$$\mathcal{L} = \mathcal{L}_{\text{answer}} + \mathcal{L}_{\text{thinking}}.$$

We experiment with different values of λ (0.3, 0.5, and 0.7) to examine how emphasizing either the answer or reasoning affects overall performance. As shown in the Table V, all weighted schemes perform comparably, with only minor variations across BLEU scores. Notably, the direct summation strategy achieves the highest performance overall. This suggests that giving equal emphasis to both the answer and reasoning generation may encourage better joint optimization and lead to more coherent outputs.

TABLE V
PERFORMANCE COMPARISON OF DIFFERENT LOSS WEIGHTING SCHEMES.

Loss Scheme	BLEU-1	BLEU-2	BLEU-3	BLEU-4
Weighted ($\lambda = 0.5$)	0.638	0.581	0.498	0.430
Weighted ($\lambda = 0.7$)	0.636	0.580	0.497	0.431
Weighted ($\lambda = 0.3$)	0.637	0.581	0.498	0.430
Direct Summation	0.648	0.586	0.500	0.430

E. Visualization

Fig. 5 presents a series of example images processed by our model, accompanied by the corresponding reasoning chains generated for each case. These images showcase various challenging visual-textual scenarios, and the reasoning chains highlight our model's ability to identify and resolve ambiguities in the text, even in noisy or complex backgrounds. Each image is paired with the question prompt ("What is the text in the image?"), followed by the model's reasoning process

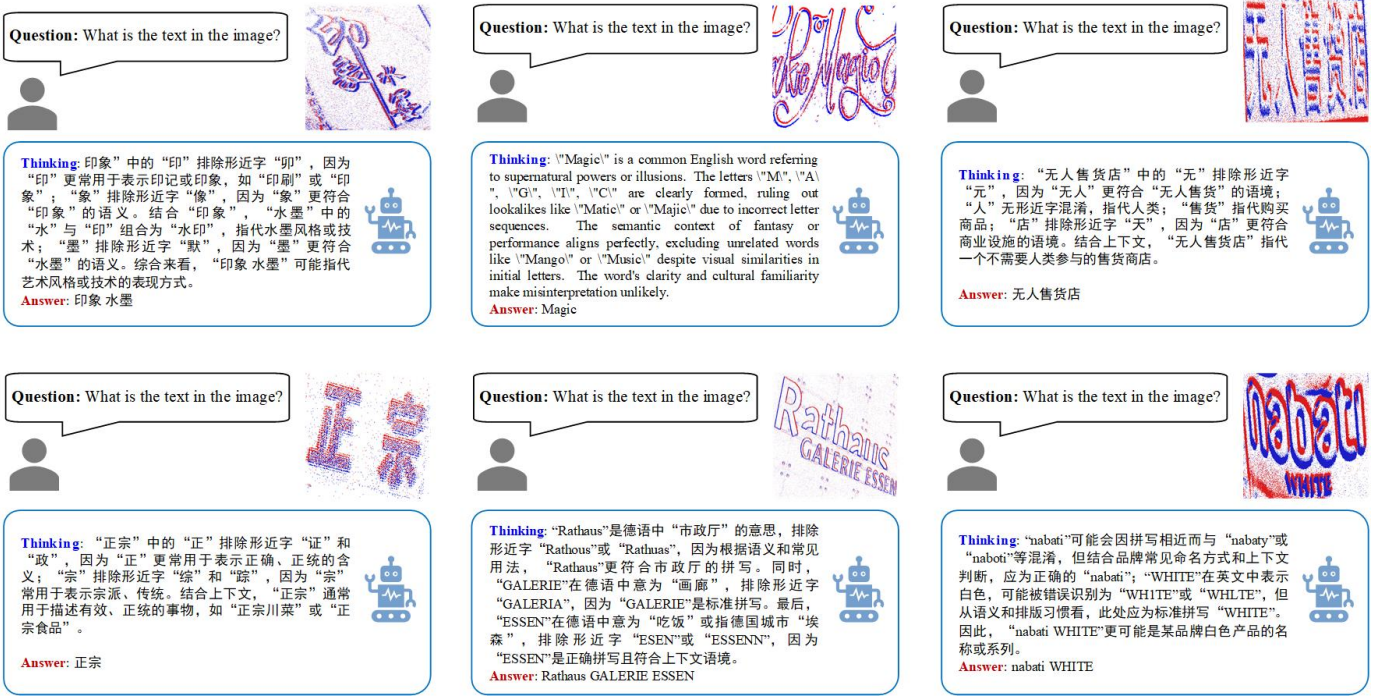


Fig. 5. **Visualization of rationales.** Our model can generate rationales that not only recognize the correct textual content from event-based images but also provide step-by-step reasoning to justify the prediction. These rationales typically include visual disambiguation cues (e.g., excluding visually similar distractors) and semantic justifications (e.g., matching word meaning with context), enabling interpretable and trustworthy scene text recognition under complex or ambiguous visual conditions.



Fig. 6. **Comparison of text recognition results: Baseline, Our Proposed Approach, and Ground Truth (GT).**

that explains how it arrives at the final output. These examples demonstrate the effectiveness of our model in handling complex text recognition tasks, generating coherent reasoning chains that align with the textual content in the images.

For instance, in the first image, the model analyzes the text "印象" (impression) and identifies that certain characters can be excluded to clarify the meaning, ultimately recognizing "印象水墨" (impression ink). Similarly, in the second image, the model distinguishes the word "Magic" from visually similar words like "Matic" or "Majic" based on the context, leading to the correct identification. The reasoning behind these decisions helps the model to generate more accurate final answers, improving both interpretability and performance.

Fig. 6 presents a comparison of text recognition results

between the baseline method, our proposed approach, and the ground truth (GT) for a sample image. The input image, along with the recognition results from the baseline, our method, and the GT text, are displayed for comparison. As shown in the figure, the baseline method struggles with accurately identifying the text, especially in challenging conditions such as distorted or noisy text. In contrast, our approach shows a significant improvement in recognition accuracy, demonstrating the effectiveness of our model in handling complex and visually degraded text. The GT text serves as the ideal reference, highlighting the target recognition outcome. This comparison underscores the strengths of our method in achieving better performance on text recognition tasks, particularly in scenarios where the baseline model fails to handle distorted or noisy text.

F. Limitation Analysis

While our proposed method significantly enhances the accuracy and interpretability of scene text recognition, there are still some limitations to address: First, even when configured to output only the final answer without the reasoning chain, our model still suffers from slower inference speed compared to traditional scene text recognition models. This is primarily due to the inherently autoregressive nature and LLM used in our architecture. Such latency can be a drawback in real-time or resource-constrained scenarios. Second, our model is pre-trained on VQA datasets instead of large-scale OCR-specific corpora (e.g., MJ [71] and ST [72]). This limits the model's performance upper bound on standard OCR benchmarks, as domain-specific pre-training has been shown to provide strong inductive biases for text recognition tasks.

V. CONCLUSION

To address the limitations of current event stream based scene text recognition (STR) on the lack interpretability and struggle with contextual logical reasoning, in this paper, we proposed a novel chain-of-thought reasoning-based framework for event stream STR, termed ESTR-CoT. Our approach leverages a vision encoder (EVA-CLIP) to extract visual features from the event stream and aligns them with a pre-trained LLM (Vicuna-7B) via a Q-former module. This allows the model to generate not only accurate text predictions but also interpretable reasoning processes in the form of chain-of-thought. The framework is trained end-to-end through supervised fine-tuning. Furthermore, we introduced a large-scale CoT dataset constructed through a three-stage process (i.e., generation, polishing, and expert verification) to support the training of reasoning-capable models. This dataset lays a solid foundation for future development in this area. Extensive experiments on three benchmark datasets, i.e., EventSTR, WordArt*, and IC15*, fully demonstrate that ESTR-CoT achieves strong performance while significantly improving model interpretability and transparency.

In our future works, we will further attempt to generate more high-quality reasoning datasets using reinforcement learning technique, such as GRPO (Group Relative Policy Optimization). Also, we will explore how to leverage large-scale OCR datasets (e.g., MJ and ST) in a resource-efficient manner, enabling better recognition performance without significantly increasing computational cost.

REFERENCES

- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.
- [2] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford *et al.*, "Gpt-4o system card," *arXiv preprint arXiv:2410.21276*, 2024.
- [3] X. Bi, D. Chen, G. Chen, S. Chen, D. Dai, C. Deng, H. Ding, K. Dong, Q. Du, Z. Fu *et al.*, "Deepseek llm: Scaling open-source language models with longtermism," *arXiv preprint arXiv:2401.02954*, 2024.
- [4] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang *et al.*, "Qwen technical report," *arXiv preprint arXiv:2309.16609*, 2023.
- [5] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2425–2433.
- [6] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, "Flamingo: a visual language model for few-shot learning," *Advances in neural information processing systems*, vol. 35, pp. 23 716–23 736, 2022.
- [7] J. Li, D. Li, S. Savarese, and S. Hoi, "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," in *International conference on machine learning*. PMLR, 2023, pp. 19 730–19 742.
- [8] W. Hu, Y. Xu, Y. Li, W. Li, Z. Chen, and Z. Tu, "Bliva: A simple multimodal llm for better handling of text-rich visual questions," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 3, 2024, pp. 2256–2264.
- [9] B. Jing, P. Xie, and E. Xing, "On the automatic generation of medical imaging reports," *arXiv preprint arXiv:1711.08195*, 2017.
- [10] Y. Cao, L. Cui, L. Zhang, F. Yu, Z. Li, and Y. Xu, "Mmtm: multi-modal memory transformer network for image-report consistent medical report generation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 37, no. 1, 2023, pp. 277–285.
- [11] Q. Chen and Y. Hong, "Medblip: Bootstrapping language-image pre-training from 3d medical images and texts," in *Proceedings of the Asian Conference on Computer Vision*, 2024, pp. 2404–2420.
- [12] X. Wang, F. Wang, Y. Li, Q. Ma, S. Wang, B. Jiang, and J. Tang, "Cxpimg-bench: Pre-training and benchmarking for x-ray medical report generation on chexpert plus dataset," in *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, June 2025, pp. 5123–5133.
- [13] Y. Liu, B. Yang, Q. Liu, Z. Li, Z. Ma, S. Zhang, and X. Bai, "Textmonkey: An ocr-free large multimodal model for understanding document," *arXiv preprint arXiv:2403.04473*, 2024.
- [14] H. Feng, Q. Liu, H. Liu, W. Zhou, H. Li, and C. Huang, "Docpedia: Unleashing the power of large multimodal model in the frequency domain for versatile document understanding," *arXiv preprint arXiv:2311.11810*, 2023.
- [15] H. Wei, L. Kong, J. Chen, L. Zhao, Z. Ge, J. Yang, J. Sun, C. Han, and X. Zhang, "Vary: Scaling up the vision vocabulary for large vision-language model," in *European Conference on Computer Vision*. Springer, 2025, pp. 408–424.
- [16] A. Hu, H. Xu, J. Ye, M. Yan, L. Zhang, B. Zhang, C. Li, J. Zhang, Q. Jin, F. Huang *et al.*, "mplug-docowl 1.5: Unified structure learning for ocr-free document understanding," *arXiv preprint arXiv:2403.12895*, 2024.
- [17] H. Wei, C. Liu, J. Chen, J. Wang, L. Kong, Y. Xu, Z. Ge, L. Zhao, J. Sun, Y. Peng *et al.*, "General ocr theory: Towards ocr-2.0 via a unified end-to-end model," *arXiv preprint arXiv:2409.01704*, 2024.
- [18] X. Wang, J. Jiang, D. Li, F. Wang, L. Zhu, Y. Wang, Y. Tian, and J. Tang, "Eventstr: A benchmark dataset and baselines for event stream based scene text recognition," *arXiv preprint arXiv:2502.09020*, 2025.
- [19] M. Gehrig and D. Scaramuzza, "Recurrent vision transformers for object detection with event cameras," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 13 884–13 893.
- [20] Z. Guo, J. Gao, G. Ma, and J. Xu, "Spatiotemporal aggregation transformer for object detection with neuromorphic vision sensors," *IEEE Sensors Journal*, vol. 24, no. 12, pp. 19 397–19 406, 2024.
- [21] Y. Peng, H. Li, Y. Zhang, X. Sun, and F. Wu, "Scene adaptive sparse transformer for event-based object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 16 794–16 804.
- [22] X. Wang, Y. Jin, W. Wu, W. Zhang, L. Zhu, B. Jiang, and Y. Tian, "Object detection using event camera: A moe heat conduction based detector and a new benchmark dataset," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 29 321–29 330.
- [23] X. Wang, Y. Jin, L. Chen, B. Jiang, L. Zhu, Y. Tian, J. Tang, and B. Luo, "Dynamic graph induced contour-aware heat conduction network for event-based object detection," *arXiv preprint arXiv:2505.12908*, 2025.
- [24] J. Zhang, Y. Wang, W. Liu, M. Li, J. Bai, B. Yin, and X. Yang, "Frame-event alignment and fusion network for high frame rate tracking," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 9781–9790.
- [25] H. Chen, Q. Wu, Y. Liang, X. Gao, and H. Wang, "Asynchronous tracking-by-detection on adaptive time surfaces for event-based object tracking," in *Proceedings of the 27th ACM International Conference on Multimedia*, 2019, pp. 473–481.
- [26] J. Zhang, X. Yang, Y. Fu, X. Wei, B. Yin, and B. Dong, "Object tracking by jointly exploiting frame and event domain," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13 043–13 052.

- [27] X. Wang, S. Wang, C. Tang, L. Zhu, B. Jiang, Y. Tian, and J. Tang, "Event stream-based visual object tracking: A high-resolution benchmark dataset and a novel baseline," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 248–19 257.
- [28] Q. Sun, Y. Fang, L. Wu, X. Wang, and Y. Cao, "Eva-clip: Improved training techniques for clip at scale," *arXiv preprint arXiv:2303.15389*, 2023.
- [29] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [30] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez *et al.*, "Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality," See <https://vicuna.lmsys.org> (accessed 14 April 2023), vol. 2, no. 3, p. 6, 2023.
- [31] X. Xie, L. Fu, Z. Zhang, Z. Wang, and X. Bai, "Toward understanding wordart: Corner-guided transformer for scene text recognition," in *European conference on computer vision*. Springer, 2022, pp. 303–321.
- [32] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu *et al.*, "Icdar 2015 competition on robust reading," in *2015 13th international conference on document analysis and recognition (ICDAR)*. IEEE, 2015, pp. 1156–1160.
- [33] X. Wang, G. Chen, G. Qian, P. Gao, X.-Y. Wei, Y. Wang, Y. Tian, and W. Gao, "Large-scale multi-modal pre-trained models: A comprehensive survey," *Machine Intelligence Research*, vol. 20, no. 4, pp. 447–482, 2023.
- [34] J. Huang and K. C.-C. Chang, "Towards reasoning in large language models: A survey," in *Findings of the Association for Computational Linguistics: ACL 2023*, 2023, pp. 1049–1065.
- [35] S. Long, X. He, and C. Yao, "Scene text detection and recognition: The deep learning era," *International Journal of Computer Vision*, vol. 129, no. 1, pp. 161–184, 2021.
- [36] K. Wang, B. Babenko, and S. Belongie, "End-to-end scene text recognition," in *2011 International conference on computer vision*. IEEE, 2011, pp. 1457–1464.
- [37] B. Shi, X. Bai, and C. Yao, "An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 11, pp. 2298–2304, 2016.
- [38] M. Liao, J. Zhang, Z. Wan, F. Xie, J. Liang, P. Lyu, C. Yao, and X. Bai, "Scene text recognition from two-dimensional perspective," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 8714–8721.
- [39] X. Han, J. Gao, C. Yang, Y. Yuan, and Q. Wang, "Spotlight text detector: Spotlight on candidate regions like a camera," *IEEE Transactions on Multimedia*, 2024.
- [40] Z. Zhao, J. Tang, C. Lin, B. Wu, C. Huang, H. Liu, X. Tan, Z. Zhang, and Y. Xie, "Multi-modal in-context learning makes an ego-evolving scene text recognizer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 15 567–15 576.
- [41] T. Guan, W. Shen, X. Yang, Q. Feng, Z. Jiang, and X. Yang, "Self-supervised character-to-character distillation for text recognition," in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, pp. 19 416–19 427.
- [42] T. Guan, C. Gu, J. Tu, X. Yang, Q. Feng, Y. Zhao, and W. Shen, "Self-supervised implicit glyph attention for text recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 15 285–15 294.
- [43] T. Zheng, Z. Chen, S. Fang, H. Xie, and Y.-G. Jiang, "Cdistnet: Perceiving multi-domain character distance for robust text recognition," *International Journal of Computer Vision*, vol. 132, no. 2, pp. 300–318, 2024.
- [44] J.-N. Li, X.-Q. Liu, X. Luo, and X.-S. Xu, "Volter: Visual collaboration and dual-stream fusion for scene text recognition," *IEEE Transactions on Multimedia*, 2024.
- [45] J. Wei, H. Zhan, Y. Lu, X. Tu, B. Yin, C. Liu, and U. Pal, "Image as a language: Revisiting scene text recognition via balanced, unified and synchronized vision-language reasoning network," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 6, 2024, pp. 5885–5893.
- [46] B. Na, Y. Kim, and S. Park, "Multi-modal text recognition networks: Interactive enhancements between visual and semantic features," in *European Conference on Computer Vision*. Springer, 2022, pp. 446–463.
- [47] C. Da, P. Wang, and C. Yao, "Levenshtein ocr," in *European Conference on Computer Vision*. Springer, 2022, pp. 322–338.
- [48] S. Fang, H. Xie, Y. Wang, Z. Mao, and Y. Zhang, "Read like humans: Autonomous, bidirectional and iterative language modeling for scene text recognition," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 7098–7107.
- [49] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [50] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [51] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann *et al.*, "Palm: Scaling language modeling with pathways," *Journal of Machine Learning Research*, vol. 24, no. 240, pp. 1–113, 2023.
- [52] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, "Training language models to follow instructions with human feedback," *Advances in neural information processing systems*, vol. 35, pp. 27 730–27 744, 2022.
- [53] D. Xue, S. Qian, Q. Fang, and C. Xu, "Linin: Logic integrated neural inference network for explanatory visual question answering," *IEEE Transactions on Multimedia*, vol. 27, pp. 16–27, 2025.
- [54] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are zero-shot reasoners," *Advances in neural information processing systems*, vol. 35, pp. 22 199–22 213, 2022.
- [55] S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, and K. Narasimhan, "Tree of thoughts: Deliberate problem solving with large language models," *Advances in neural information processing systems*, vol. 36, pp. 11 809–11 822, 2023.
- [56] M. Besta, F. Memedi, Z. Zhang, R. Gerstenberger, N. Blach, P. Nyczzyk, M. Copik, G. Kwasniewski, J. Müller, L. Gianinazzi *et al.*, "Topologies of reasoning: Demystifying chains, trees, and graphs of thoughts," *CoRR*, 2024.
- [57] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," *arXiv preprint arXiv:2203.11171*, 2022.
- [58] R. Zhang, B. Zhang, Y. Li, H. Zhang, Z. Sun, Z. Gan, Y. Yang, R. Pang, and Y. Yang, "Improve vision language model chain-of-thought reasoning," *arXiv preprint arXiv:2410.16198*, 2024.
- [59] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "Minigtpt-4: Enhancing vision-language understanding with advanced large language models," *arXiv preprint arXiv:2304.10592*, 2023.
- [60] G. Gallego, T. Delbrück, G. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. J. Davison, J. Conradt, K. Daniilidis *et al.*, "Event-based vision: A survey," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 1, pp. 154–180, 2020.
- [61] Y. Jiang, Y. Wang, S. Li, Y. Zhang, Q. Guo, Q. Chu, and Y. Gao, "Evcslr: Event-guided continuous sign language recognition and benchmark," *IEEE Transactions on Multimedia*, 2024.
- [62] J. Jiang, X. Lu, L. Zhao, R. Dazaley, and M. Wang, "Masked autoencoders in 3d point cloud representation learning," *IEEE Transactions on Multimedia*, 2023.
- [63] X. Wang, Z. Wu, B. Jiang, Z. Bao, L. Zhu, G. Li, Y. Wang, and Y. Tian, "Hardvs: Revisiting human activity recognition with dynamic vision sensors," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 6, 2024, pp. 5615–5623.
- [64] D. Li, J. Jin, Y. Zhang, Y. Zhong, Y. Wu, L. Chen, X. Wang, and B. Luo, "Semantic-aware frame-event fusion based pattern recognition via large vision-language models," *Pattern Recognition*, vol. 158, p. 111080, 2025.
- [65] D. Bautista and R. Atienza, "Scene text recognition with permuted autoregressive sequence models," in *European conference on computer vision*. Springer, 2022, pp. 178–196.
- [66] P. Wang, C. Da, and C. Yao, "Multi-granularity prediction for scene text recognition," in *European Conference on Computer Vision*. Springer, 2022, pp. 339–355.
- [67] C. Cheng, P. Wang, C. Da, Q. Zheng, and C. Yao, "Lister: Neighbor decoding for length-insensitive scene text recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 19 541–19 551.
- [68] M. Yang, M. Liao, P. Lu, J. Wang, S. Zhu, H. Luo, Q. Tian, and X. Bai, "Reading and writing: Discriminative and generative modeling for self-supervised text recognition," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 4214–4223.

- [69] H. Rebecq, D. Gehrig, and D. Scaramuzza, “Esim: an open event camera simulator,” in *Conference on robot learning*. PMLR, 2018, pp. 969–982.
- [70] I. Loshchilov, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [71] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman, “Synthetic data and artificial neural networks for natural scene text recognition,” *arXiv preprint arXiv:1406.2227*, 2014.
- [72] A. Gupta, A. Vedaldi, and A. Zisserman, “Synthetic data for text localisation in natural images,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2315–2324.