

# LUMOS : Empowering Multimodal LLMs with Scene Text Recognition

Ashish Shenoy\*  
ashishvs@meta.com  
Meta Reality Labs

Yichao Lu\*  
yichaol@meta.com  
Meta Reality Labs

Srihari Jayakumar\*  
srihari2761@meta.com  
Meta Reality Labs

Debojeet Chatterjee\*  
debo@meta.com  
Meta Reality Labs

Mohsen Moslehpour\*  
moslehpour@meta.com  
Meta Reality Labs

Pierce Chuang\*  
pichuang@meta.com  
Meta Reality Labs

Abhay Harpale\*  
harpale@meta.com  
Meta Reality Labs

Vikas Bhardwaj\*  
vikasb@meta.com  
Meta Reality Labs

Di Xu  
dixu@meta.com  
Meta Reality Labs

Shicong Zhao  
zsc@meta.com  
Meta Reality Labs

Longfang Zhao  
longfangzhao@meta.com  
Meta Reality Labs

Ankit Ramchandani  
ankit61@meta.com  
Meta

Xin Luna Dong  
lunadong@meta.com  
Meta Reality Labs

Anuj Kumar  
anujk@meta.com  
Meta Reality Labs

## ABSTRACT

We introduce LUMOS, the first end-to-end multimodal question-answering system with text understanding capabilities. At the core of LUMOS is a Scene Text Recognition (STR) component that extracts text from first person point-of-view images, the output of which is used to augment input to a Multimodal Large Language Model (MM-LLM). While building LUMOS, we encountered numerous challenges related to STR quality, overall latency, and model inference. In this paper, we delve into those challenges, and discuss the system architecture, design choices, and modeling techniques employed to overcome these obstacles. We also provide a comprehensive evaluation for each component, showcasing high quality and efficiency.

## CCS CONCEPTS

• **Applied computing**; • **Computing methodologies** → **Computer vision tasks**; *Discourse, dialogue and pragmatics*;

## KEYWORDS

OCR, Scene Text Recognition, On-device, NLP, Multimodal LLMs, Hand-Object Interaction, Salient Region of Interest Detection

\*Joint First Authors

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*KDD '24, August 25–29, 2024, Barcelona, Spain*

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## ACM Reference Format:

Ashish Shenoy, Yichao Lu, Srihari Jayakumar, Debojeet Chatterjee, Mohsen Moslehpour, Pierce Chuang, Abhay Harpale, Vikas Bhardwaj, Di Xu, Shicong Zhao, Longfang Zhao, Ankit Ramchandani, Xin Luna Dong, and Anuj Kumar. 2024. LUMOS : Empowering Multimodal LLMs with Scene Text Recognition. In *Proceedings of Knowledge Discovery and Data Mining (KDD '24)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Visual question answering has been a research area for nearly a decade [5] and has gained increased attention with recent progress on LLMs and vision language pre-training (Multi-Modal LLMs). Industry anticipates that very soon, we will have smart assistants that understand scenes/images just as well as humans [3, 29]. In this paper, we focus on one key abilities needed for scene understanding, visual understanding and question-answering related to text in the scene. We describe LUMOS<sup>1</sup>, an end-to-end system for multimodal text-understanding. Lumos is well-suited for real-world application and readily leverages on-device processing to enable smooth user experiences.

Figure 1 shows example user interactions for some of LUMOS's use-cases. At the first glance, one may think this problem is already solved by Multimodal Large Language Models (MM-LLMs). In [3, 29], MM-LLMs demonstrated capabilities understanding texts from images without a standalone STR component. It would seem straight-forward to enable this feature for edge devices by taking a picture using the on-device camera, transferring to a cloud-based service, along with any voice recognition of user questions, and then having an MM-LLM answer the question using the image. If text recognition is sub-optimal when relying purely on the MM-LLM,

<sup>1</sup>Derived from the Latin word 'lumen', meaning 'light'.

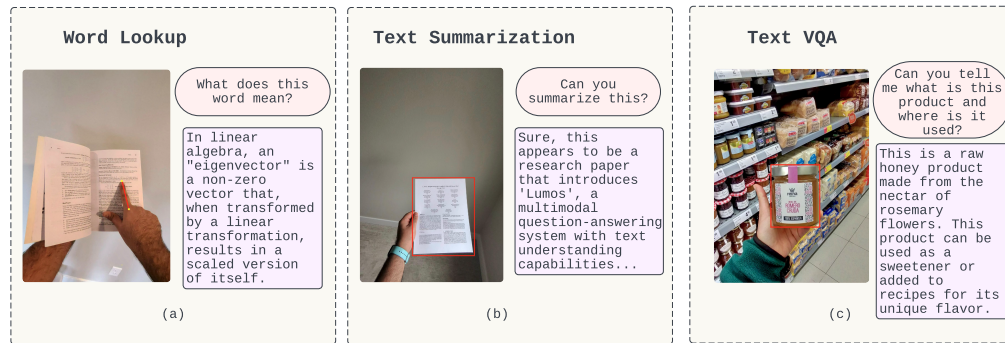
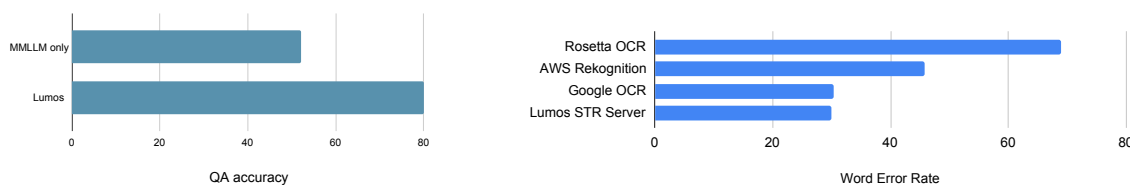


Figure 1: Text based use cases that LUMOS supports.



(a) LUMOS achieved 80% QA accuracy, while adding the (b) LUMOS STR has lowest word error rate compared with other STR solutions

Figure 2: LUMOS Quality metrics

one might choose to run a separate Scene Text Recognition (STR), another mature technique, on the image and send the recognized texts to the MM-LLM as prompt input to facilitate text understanding. We will now discuss in detail why such implementations are inadequate and the challenges we solve within LUMOS.

The first and key challenge we encounter is *latency*: just transferring a high-resolution image from device to cloud cost significant time resulting in a poor user experience. For instance, transmitting an image of size  $3k \times 4k$  (standard resolution for today's devices) from a device to the cloud may take several seconds before even running any AI models. And the end-to-end time to get a response would be even longer making for a poor experience.

Alternatively, if we transfer only a low-resolution thumbnail, the transfer time can be significantly reduced (e.g., transferring a thumbnail of size  $450 \times 600$  pixels takes only a few hundred ms). However, this results in significantly degraded quality on text recognition. As shown in Table 2, the accuracy of question answering relying solely on running MM-LLM over thumbnails is only 52%. A separate cloud-based STR can barely recognize texts on the thumbnails either, since the size is too small, illegible even for humans.

Now assuming we choose an on-device STR solution, the second challenge is the constrained compute and memory resources on devices. Although running STR models on-device may seem like a viable solution to address latency concerns, current state-of-the-art STR models are not readily suitable for on-device usage; for example, Google's recent work [19] features a text detection model that alone has a size of 240MB, impractical for on-device use where several other processes might be running and sharing memory.

The final set of challenges arise with doing STR on in-the-wild text images, which are different from common web images, scanned documents or zoomed-in images. Images taken on-the-go and outdoors can amplify the challenges of STR. 1) The cameras are typically wide angle, and thus the text of interest occupies only a small portion of the image; furthermore, there is often additional background text which can be irrelevant to the user query (see Figure 1(c)). 2) The text in the scene may not have any uniformity: rotated, diverse orientations and font sizes. 3) The image quality might be poor owing to sub-optimal lighting condition, user movement, and the camera angle. For all of these reasons, traditional OCR (Optical Character Recognition) systems, despite their strong performance on scanned documents and screenshots, can fall short on a STR task in an in-the-wild text setting. As an example, the cloud-based OCR solution Rosetta [7] exhibits a surprising 53% Word Error Rate (WER) on our in-the-wild text STR benchmark (see Section 6 for details).

In this paper, we discuss our results overcoming these three challenges. (1) In our tests, our proposed system has an average end-to-end latency of  $\leq 5$  seconds, including photo capture, image transfer, on-device STR execution, and on-cloud MM-LLM execution. (2) Our on-device STR models have a total size of  $\leq 8$ Mb, a peak memory footprint of  $\leq 200$ Mb, an average latency of  $\leq 1$ sec, and 0.4 mWh power usage. (3) Despite the low cost, our STR solution achieves competitive quality on public STR benchmarks when compared to state-of-the-art STR solutions from other cloud service providers (Figure 2b). On our own in-the-wild text benchmarks, it achieves a 14.6% WER and enables an average accuracy of 80% on complex

text-based QA tasks, improving over vanilla MM-LLM solution by 28% (see Figure 2a).

There are three key innovations in LUMOS: First, a hybrid approach to multimodal text-understanding with *an architecture leveraging components across on-device and on-cloud*. In particular, we conducted on-device STR, such that we can achieve high-quality text recognition results on the full-resolution image; we then send the recognized texts, together with the low-resolution image to the MM-LLM on cloud for question answering; as we run STR in parallel to image transfer, which is the main latency bottleneck, the on-device STR does not add additional latency most of the time (see Section 3). Running STR on the full-resolution image can still be computationally expensive on device, hence our second innovation is an *ROI (Region Of Interest) detection solution* that allows the STR to focus on the area of interest and thus reduce the computational overhead. Our ROI detection solution first effectively detects salient areas in the visual, and then crops the salient area as STR input (see Section 4.1). Third, we developed *a state-of-the-art on-device and resource-preserving STR model*. We optimized our models to run with hardware acceleration resulting in a smaller memory and compute footprint, and efficient battery usage, with minimum sacrifice on quality (see Section 4.2-5).

To the best of our knowledge, we are the first to propose a multimodal assistant with text understanding capabilities that heavily leverages on-device computation. We summarize our key contributions as follows:

- We propose LUMOS, an end-to-end (E2E) multimodal assistant system with text understanding capabilities; through careful placement of components on-device or on-cloud, we are able to achieve high quality, low latency, and minimize on-device resource usage.
- We present an on-device STR pipeline with a set of models for ROI detection, text detection, text recognition, and reading order reconstruction that together achieved high quality (WER=14.6%) and low cost (latency=0.9s, peak runtime memory=200 Mb, power=0.4 mwh on testing device).
- Through a comprehensive evaluation of our system on QA benchmarks, we validated the high effectiveness and efficiency of our system.

## 2 PREVIOUS WORK

**OCR and STR.** The field of OCR has been a focal point of research for many years. However, the spectrum of difficulty in recognizing text in natural environments is notably broad. At one end, OCR’s application to scanned documents containing well-structured printed text is widely recognized as one of the most successful implementations of computer vision [13, 21]. Conversely, STR focuses on recognizing text in the wild, which still represent a significant challenge due to the larger variance of wild text objects [7, 14, 22, 25, 26, 33]. The STR problem we are solving in this paper considers in-the-wild text images (so the area of interest is considerably smaller), and needs to be tackled on device, thus is much harder and requires better model designs and tuning.

**On-device STR.** When it comes to STR on-device, in [9] an extremely lightweight OCR system with a size of only 3.5Mb is proposed; the model achieves impressive latency on GPUs but still

falls short when it comes to CPUs. Munjal et al. [22] describes an efficient lightweight STR system, which has only 0.88M parameters and performs real-time text recognition at a speed of 2.44 ms per word crop of size  $16 \times 64$ . In comparison, the STR solution described in this paper takes 0.29 ms per word crop of size  $48 \times 320$ .

**Multimodal LLMs and Text Recognition Ability** More recently, MM-LLMs have demonstrated potential in addressing a variety of tasks, including text recognition [3, 4, 10, 16, 29, 35, 37]. While the current trend leans towards the use of all-modality LLMs, they have limitations particularly in handling text-in-the-wild scenarios. Furthermore, the challenges associated with high transfer latency as described in Section 1 makes these models impractical for immediate use [18, 27]. A different approach, the *Flamingo* models [4, 6], have shown impressive performance on tasks such as generic VQA and captioning, but fall short when compared to [12] on text rich VQA. Both sets of models are sub-optimal compared to OCR-assisted VQA as we discussed in this paper and are not optimized for memory and compute at inference time.

## 3 OVERALL ARCHITECTURE

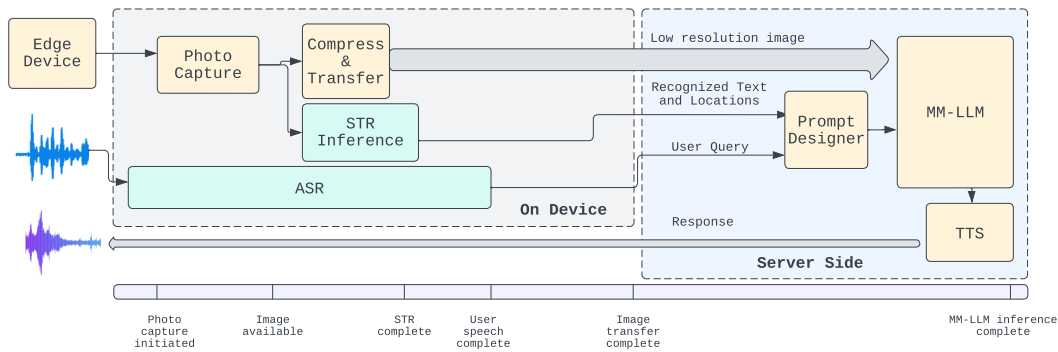
We now describe the overall architecture of LUMOS (see Figure 3). To simplify, we focus on multimodal use cases, assuming a picture will be taken once the user triggers the flow, and the device will provide the image at two resolutions  $3K \times 4K$  (full resolution), and  $450 \times 600$  (thumbnail).

**Device-side:** At the device side, when a user gives a voice query, three components will start in parallel. First, Automatic Speech Recognition (ASR) starts processing the query after a wake word. Second, the photo capture, compression (e.g., from a  $3k \times 4k$  full-resolution image to a  $450 \times 600$  thumbnail) and transfer to cloud will begin in parallel to the voice query completion (to reduce overall system latency). Third, the STR component will start as soon as the full-resolution image is ready. As indicated by in Figure 3, we carefully design the system to parallelize the time consuming components, STR inference and image transfer, to reduce latency.

**Cloud-side:** The cloud side hosts a *MM-LLM model*, which takes as input the low-resolution thumbnail, a prompt composed of the recognized texts and their coordinates from STR, and the user query from ASR, and generates the answer response. An illustrative prompt to MM-LLM can be found in Appendix Table 7. Subsequently, the **TTS (Text-to-Speech)** component translates the response to voice signal and sends back to the user.

This architecture incorporates three design choices we have carefully made.

- *Where to do STR?* As discussed in detail in Section 1, to reduce latency, we transfer *only* a low-resolution image to the cloud. However, neither an MM-LLM nor an STR model can achieve desired quality on such a low-resolution image, especially given that the text area is typically quite small in the in-the-wild text image. We thus apply STR on device with the full-resolution image, and only on the region of interest (see section 4.1 for details).
- *How to cut the STR latency?* Running STR on device can be time-consuming. To reduce this latency, we took two actions: 1) use hardware acceleration (see section 5), 2) execute STR



**Figure 3: Overall architecture of LUMOS. The width of the blocks on device are roughly represents runtime latency. The arrow width roughly represents to the size of the payload being transferred. Blue blocks indicate models using hardware acceleration.**

and image transfer in parallel (see Figure 3). With this design, for the majority of the cases STR does not add extra latency.

- **How to extend to MM-LLM use cases where STR is not necessary to answer the user question?** Ideally, we wish to build a *single* multimodal assistant that can answer text-heavy questions as well as generic questions where text comprehension is not necessary. Determining whether a user question is based on the text in the scene requires an NLU (Natural Language Understanding) component, which can take extra time and may have limited quality with the limited computation power on device. We instead conduct STR in all cases and defer the decision to the MM-LLM on the cloud. This approach is feasible only because of our significant reduction of latency (particularly through parallelization) and optimization of hardware efficiency for STR.

It is worth mentioning that placing STR on-device poses significant constraints on the model’s architecture, latency, memory, and battery consumption, in addition to the quality challenges for in-the-wild text STR discussed in Section 1. Despite these limitations, our on-device STR model achieves strong performance compared to three state-of-the-art cloud STR solutions that do not have such constraints (see Table 3 for details). In the next section, we describe how we achieve this.

## 4 SCENE-TEXT RECOGNITION

We now describe our core technique—the on-device STR. This pipeline contains four sub-components as depicted in Figure 4.

- **Region of Interest (ROI) detection** takes an image as input (at both  $3k \times 4k$  resolution and a thumbnail resolution), outputs a cropped image (about  $1k \times 1.3k$  size) that contains all the text likely needed to answer the user query. This component plays a key role to ensure that we run the rest of the STR pipeline only on the relevant portion of the input image, reducing both computational cost and background noise.
- **Text detection** takes the cropped image from ROI detection as input, detects words, and outputs the identified bounding box coordinates for each word.

- **Text recognition** takes the cropped image from ROI detection and the word bounding box coordinates from Text detection as input, returns the recognized words.
- **Reading-order reconstruction** organizes recognized words into paragraphs and in reading order within each paragraph based on the layout. It outputs text paragraphs as well as their location coordinates.

We note that in most previous works STR refers to only the Text detection and Text recognition parts. We included two additional components—ROI detection and Reading order reconstruction—in our STR system to address LUMOS specific challenges. The primary challenges we face include the limited hardware for inference and the large variation of texts in the wild. We address these challenges through *careful model architecture selection and training data curation and augmentation*, as we discuss in detail next.

### 4.1 ROI Detection

**Motivation** ROI detection plays a key role for on-device STR and there are three motivations behind it. First, as shown in Figure 1(b), because of the nature of in-the-wild text images, the text area of interest often occupies only a small fraction of the image, even if the object is only an arm length away from the device. Running STR directly on the full-resolution image can be prohibitively expensive with the limited computational power of the device, whereas downsizing the image can make the texts too small to be legible even to humans. Second, as shown in Figure 1(c), the image may contain a lot of background text that are irrelevant to the user query, such as text from products on the shelves. Recognizing these texts consumes the limited hardware resources, increases the latency, and confuses the MM-LLM at the downstream. Third, users often hold the paper or the object of interest like in Figure 1(c), or point to the particular words or phrases like in Figure 1(a), where those gestures provide critical clues for ROI detection. These motivations underscore the importance of identifying the ROI before proceeding with other steps in STR.

**Problem definition and challenges** The ROI detection module uses a low resolution thumbnail  $450 \times 600$  to detect the ROI, and returns the cropped area from the raw image  $3k \times 4k$  containing the

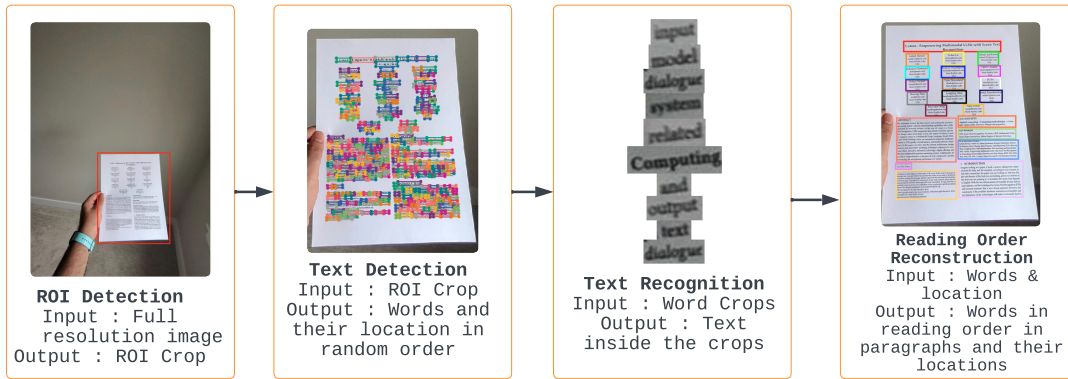


Figure 4: On-device STR component flow of LUMOS.

ROI. A major challenge for ROI is the non-holding or non-pointing hands in the picture, which can lead to wrong detection results (see example in Figure 8 in the Appendix).

**Solution and modeling** We treat ROI detection as an object (salient area) detection problem, facilitated with keypoint detection in presence of a *pointing finger*. For finger pointing, we detect two key points—the last joint and the tip of index finger; the two points formulate a pointing vector, as shown in Figure 1(a). We train a model that jointly detects both the ROI and the two keypoints (when present). If the keypoints are detected, we include an additional prompt to the downstream MM-LLM, describing the pointing event as well as the words and the paragraphs closest to the tip of the index finger in the direction of the pointing vector. We use the Mask-rcnn [11] model since it can provide a unified framework for both object and keypoint detection. We apply inference on the  $450 \times 600$  thumbnail.

**Training data** We trained the model using 80K in-the-wild text images annotated with salient regions, and 20K images with hand holding or finger pointing. To reduce false positives caused by accidental hands, we included 10K images with a hand that is neither holding nor pointing as hard negatives in our training data.

## 4.2 Text Detection

**Problem definition and challenges** Text detection takes the cropped image in full-resolution as input, predicts location of each word as bounding boxes. There are three challenges for detecting text in the wild: C1. the text size can be very small (e.g., ingredients on a coke can at arm length) or very big (e.g., storefront); C2. text can often be tilted with the nature of the image; C3. we are not able to use state-of-the-art text detection model architectures like [15, 19] with the on-device constraint.

**Solution and modeling** To account for the tilted text (C2), our detector predicts rotated bounding box as mentioned in [20]. To be computationally efficient (C3), we use an anchor-free single-stage detector as described in [30] (instead of a two-stage detector). We use FBNetv2 (with 1.1 million parameters) [31] with PAN neck [17] for the backbone of the detector. FBNetv2 is a CNN model designed for transforming input images into feature maps; this backbone not only is computationally efficient (C3) but also provides strong

image features at different scales (C1). For the loss, we use a variant of the well-known focal loss [36] as classification loss, and the KLD loss [34] as our box regression loss for its state-of-the-art performance on rotated box (C2).

**Training data** Our training data consist of 140K images with 6 million annotated bounding boxes, combining public STR datasets like text OCR [28] and in-house annotations on in-the-wild text images. To address the challenge of text scale variation (C1), we applied aggressive *scale jittering*, data augmentation that increases or reduces input image sizes, to create variational sizes of bounding boxes in training data.

## 4.3 Text Recognition

**Problem definition and challenges** Text recognition takes the image crop from ROI detection and the word bounding box coordinates, and outputs the recognized words for each box. There are three key challenges we need to address: C1. huge diversity in the widths of bounding boxes (e.g., URLs tend to be longer, price tags tend to be extremely small); C2. diversity of text appearances in terms of font, size, orientation, and background; C3. existence of (quite some) text detection errors; C4. hardware constraints.

**Solution and modeling** We transform the problem of recognizing a word into the problem of recognizing a sequence of characters. Because of hardware acceleration constraints (C4) as we will describe in Section 5, we are limited to using fixed width and height for each bounding box. Therefore, we scale each bounding box to a fixed height of 48 pixels and a fixed width of 320 pixels to ensure that the input to the model is consistent and can be processed efficiently. Based on statistics we assume that each individual character has a width of 8 pixels. Thus, we recognize a maximum of 40 characters ( $320/8$ ) per bounding box; a word rarely exceeds this limit. The final recognizer output is a posterior of shape  $40 \times |\text{alphabets}|$  and the size of the alphabets in our model is top-150 most frequently used Latin characters obtained from the training data.

We again use the FBNetv2 backbone and train the model using CTC (Connectionist Temporal Classification) loss, as it can handle variable-length input sequences (C1) and has lower latency and computational complexity (C4), critical in dense text scenarios.



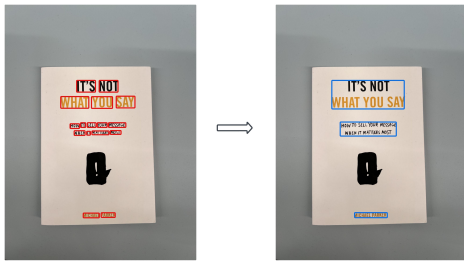


Figure 5: Left: Word bounding boxes. Right: Paragraphs from out Reading Order Reconstruction component

**Training data** During training, to handle the extreme variations in bounding box lengths (C1), we employ *curriculum learning*; that is, we gradually increase the complexity of the input images. We begin with words containing a maximum of 16 characters and progressively increase the character limit up to a maximum of 40 characters. This helps the model learn the necessary features and patterns more effectively.

Overall, the recognizer model is trained on 3M word bounding boxes, with 15% being synthetically generated to increase the robustness of the model. To be more robust against detector errors (C3), we introduce random cropping around the boundaries of the bounding boxes based on error patterns we have observed in detector evaluation, combined with jittering. We incorporated RandAug [8], which applies random combinations of image transformations such as rotation, shearing, brightness adjustment, and contrast adjustment to input images. By exposing the model to a wide range of transformed images, it learns to be more robust to these transformations and generalizes better to new, unseen data (C2).

#### 4.4 Reading Order Reconstruction

**Problem definition** The Reading Order Reconstruction module connects the words to paragraphs, returns the words in the paragraph in reading order, together with the coordinates of each paragraph. Figure 5 shows sample paragraphs.

**Solutions** We identify paragraphs in three steps. First, we connect the words to paragraphs. We expand the word bounding boxes both vertically and horizontally by predefined ratios, as shown in Figure 9. The expansion ratios are selected to fill the gaps between words within a line and lines within a paragraph and are the same for all bounding boxes. We then group bounding boxes that have significant overlap after expansion as a paragraph. For each paragraph, we then apply *raster scan* (i.e., sorting by Y coordinate then X) to the words to generate the paragraph in reading order. Finally, we compute the location of the paragraph by finding the minimum area rectangle enclosing all words in the paragraph. See Algorithm 1 in the Appendix for detailed description of the Reading order reconstruction module.

We found this simple heuristic approach achieves a good quality most of the time with low computation cost. The accuracy for this module is 92% using metrics defined in [?].

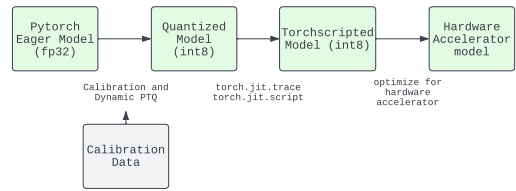


Figure 6: Model Export Pipeline

Table 1: Evaluation dataset details

Name	Size
In-house wild text benchmark	968 images, 47K words
Public wild text benchmark	1.7K images, 146K words
Summarization question set	245 images, 735 questions
Word Lookup question set	200 images, 600 questions
Other question set	200 images, 600 questions

## 5 ON-DEVICE EXPORT

As mentioned in the introduction, LUMOS need to be compatible with devices to make our smart assistant more accessible. We evaluated our on-device system’s performance with on our testing devices, which is equipped with hardware accelerators for deep learning models. We describe the process of exporting our models to the testing device as well as the memory/latency in this setting.

- (1) **Quantization to int8** We first quantize the float32 models to int8 models to save inference latency and runtime memory. We use Post Training Quantization (PTQ) [23] to do this, because the sizes of our models were relatively small and PTQ requires a calibration step only after models are full trained.
- (2) **On-device CPU models** We next transfer the models to TorchScript models using packages provided by PyTorch. This provides a model that is executable on CPU of the device.
- (3) **On-device hardware accelerated models** Modern devices often comes with a hardware accelerator for deep learning models. To utilize this, we take a further step making our model hardware accelerator compatible, and evaluate the latency of our system on hardware accelerator.

We emphasize that the model execution efficiency is achieved with cost. First, we are constrained to use quantization and hardware accelerator friendly models, limited our modeling choices as stated in Section 4. Second, quantization and model export would cause accuracy drops in the ML models. Regardless, our system still achieves competitive performance compared to other STR services as we show soon in Section 6.

## 6 EXPERIMENTAL RESULTS

We answer two questions with our experiments: 1) How good is LUMOS as an end-to-end text visual question answering system? 2) What is the quality, efficiency, and hardware usage for our on-device STR solution?

**Table 2: QA accuracy of LUMOS variants on task-specific benchmarks. On-device STR boosts QA accuracy by 28%.**

System	Summarization	Word Lookup	Others	Avg
MMLLM only	53.0%	43.0%	60.1%	52.0%
+ STR	87.7%	65.0%	81.3%	78.0%
+ STR + Pos	<b>88.3%</b>	<b>67.2%</b>	<b>83.3%</b>	<b>79.6%</b>

## 6.1 Experiment Setup

**Datasets** Table 1 summarizes the datasets we used for evaluation. We have two benchmarks: *In-house wild text benchmark* and *Public wild text benchmark*. *In-house wild text benchmark* contains 968 in-the-wild text images taken from an edge device and contains 47K word boxes. The benchmark contains annotations for the word boxes and transcriptions, and in addition annotations for salient areas for ROI evaluation. *Public wild text benchmark* is a broadly-used STR benchmark, containing 1.7K images and 146K word boxes. We then created task-specific datasets to evaluate end-to-end quality of summarization, word lookup and a few other tasks on the *In-house wild text benchmark*. We first sampled text-heavy images from the benchmark, and then our annotators created  $\sim 3$  task-related questions for each image.

**Metrics definition** We have two major metrics. To understand the end-to-end question answering quality, we measure *QA accuracy* as the percentage of successful responses among all answers. A group of raters manually decided the correctness of each response judging from the image, user query and generated response, based on the relevancy, fluency and factual accuracy of the response.

To understand the quality of STR solutions, we measured the *Word Error Rate* (WER), a standard metric extensively used in the domain of speech and text recognition. WER considers 3 types of errors: 1) Deletion: a ground truth word that is not detected; 2) Insertion: a prediction that is not matched to any ground truth word box; 3) Substitution: a prediction that matches a ground truth box, but the word recognized is different from the ground truth. WER is the sum of Deletion, Insertion, Substitution errors divided by the total number of words in the ground truth. With the existence of insertion errors, WER can be higher than 1. A lower WER is indicative of higher quality of the models.

## 6.2 End-to-End Quality

We evaluated the overall quality of three variants of LUMOS: 1) MMLLM only: we provide only the  $450 \times 600$  thumbnail and user query to the MM-LLM; 2) MM-LLM+STR: we in addition provide the text output from the on-device STR to MM-LLM; 3) MM-LLM+STR+Positions: we in addition provide the paragraph location (from reading order reconstruction module). See Table 7 for detailed input formats of these variants.

Table 2 compares the QA accuracy of the three variants on the task-specific E2E datasets. We have four observations. First, LUMOS obtains a high average QA accuracy, 80%, in question answering. Second, the on-device STR significantly improves QA accuracy on all three tasks over MM-LLM only (80% vs. 52%). The improvement is particularly large for the summarization task (+35%), where LUMOS needs to comprehend dense texts. Third, sending positions to

**Table 3: WER comparison on public wild text benchmarks. LUMOS STR obtains the lowest WER with a small size, and the on-device model sacrifices quality only slightly.**

Model	WER	Del	Ins	Sub	#Params
<i>Public wild text benchmark</i>					
Rosetta OCR	68.9%	58.1%	2.3%	8.5%	15Mb
AWS Rekognition [1]	45.8%	38.1%	1.6%	6.1%	-
Google OCR [2]	30.4%	9.4%	9.5%	11.5%	240Mb <sup>+</sup>
<b>LUMOS STR Server</b>	<b>29.9%</b>	<b>17.7%</b>	<b>2.5%</b>	<b>9.7%</b>	30Mb
LUMOS STR Device	32.4%	18.5%	2.7%	11.2%	<b>8Mb</b>
<i>In-house wild text benchmark</i>					
Rosetta OCR	53%	46.0%	1.1%	5.9%	15Mb
<b>LUMOS STR Server</b>	<b>13%</b>	<b>4.7%</b>	<b>1.4%</b>	<b>6.9%</b>	30Mb
LUMOS STR Device	14.6%	5.1%	1.8%	7.7%	<b>8Mb</b>

MM-LLM further improves the performance on all tasks (+1.6%), as it allows the model to better handle the spatial relationships between words in the scene. Finally, among different tasks, we observe the best quality on summarization (88%), which has higher tolerance on small recognition errors; the quality on word lookup is lowest (67%), as we observe a large variety of hand-word positions, making the problem much more difficult.

## 6.3 STR quality

**LUMOS STR quality** We next compare quality of 5 STR Systems: 1) Rosetta [7], a well known STR system from the research community; 2) Google Cloud OCR [2]; 3) AWS Rekognition [1]; 4) LUMOS STR Cloud: LUMOS STR running on cloud; 5) LUMOS STR Device: LUMOS STR running on our device hardware. For a fair comparison, we removed punctuations from the benchmarks since different baseline STR systems treat them differently, as a separate word or part of a word. We also removed words smaller than 8 pixels high since it is hard for humans to read.

Table 3 shows the WER of each solution, together with error breakdowns in terms of deletion, insertion, substitution errors. We have four observations. 1) LUMOS STR has a reasonably low WER, 30% on the public benchmark and 13% on the in-house benchmark. 2) LUMOS STR outperforms Rosetta, AWS, and Google, despite never trained on the public wild text benchmark (we do not know if Google and AWS were trained on the public wild text benchmark). Rosetta made a lot of deletion errors as it missed small texts and has a low word limit per image. Similarly, AWS has a low word limit per image, leading to high deletion errors. 3) LUMOS STR Device is smallest in model size with only  $\sim 8$ Mb parameters; nevertheless, it sacrifices WER by only 1-2% comparing with the on-server model and still has a competitive performance. 4) Finally, among different types of errors, Substitution errors is only a small portion ( $<10\%$ ), showing that word detection is a much bigger challenge than word recognition for STR tasks.

**Ablation study** We now listed the main drivers for WER improvements. We compared with Rosetta, a two-step STR system (faster-rnn [24] word detector and CNN + CTC recognizer) on the *In-house*

<sup>+</sup>Estimated based on [19], using the size of MaX-DeepLab-S [32]

**Table 4: WER gains from each component**

Component	Reason	WER	Comp. to baseline
Baseline (Rosetta OCR)	-	53%	
+ROI detection	avoid aggressive input image downsizing	42%	-11%
+Text Detection	stronger model, data augmentation, more in domain training data, increased word limit	26%	-16%
+Text Recognition	synthetic data on rare/hard symbols det error simulation, RandAug	13%	-13%
+on-device export	model quantization error	14.6%	+1.6%

**Table 5: Recall for ROI detection. On average our ROI method is able to reduce image size by 25% while including 99% words of interest.**

Method	Recall	Improvement
Center Crop	65.9%	
ROI detection	97.7%	+31.8%
ROI detection with Hand cues	99.0%	+1.3%

*wild text benchmark*. There are three contributors for quality improvements as shown in Table 4.

- ROI detection allows us to run our detection and recognition on a text-dense cropped region in original size, instead of on an aggressively downsized (3x-4x) full image, thus reducing WER by 11%, and especially reducing WER on small-font texts.
- Our detection model uses additional in-domain data and data augmentation for training to improve robustness, and increases word limit per image, thus reducing WER by 16%. In particular, we increased recall of detecting word boxes, thus reducing deletion errors, in detection of small text (<15 pixels tall) by 14% and of large text (>120 pixels tall) by 20%.
- Our recognition model used data augmentation to accommodate more variations for text in the wild, thus reducing WER by 13%.

Finally, these improvements are well preserved in model quantization and export, which increased WER by only 1.6% but achieved huge efficiency gains as we discuss soon in Section 6.4.

**ROI detection recall** To illustrate the effectiveness of the ROI detection component, we compared the performance of 3 image cropping methods: 1) Center Crop: heuristic-rule baseline that crops the 1500\*2000 center region (similar as the ROI output size); 2) ROI detection: use an object detection model to detect the region; 3) ROI detection with hand cues: use object detection together with the holding and pointing gestures to detect the region.

We measured ROI quality by word-level *recall*—how many words of interest are included in the ROI output region. Table 5 shows the results on the *in house wild text benchmark*. We are able to reach 99% recall with our ROI detection component while reducing image size by 25% on average. Our model achieves much higher recall (+32%) than the Center Crop baseline, and including hand cues further improves the recall (+1.3%).

**Table 6: Model execution metrics. Running the models on hardware accelerator (HA) saved latency by 9X and energy by 3X comparing with running on CPU.**

Metrics	CPU	HA	Saving
Overall on device latency (100 words)	8390ms	940ms	8.9X
Text Detection latency	750ms	66ms	11.4X
Text Recognition latency	238ms	29ms	8.2X
ROI detection latency	300ms	30ms	10X
Model size	-	8Mb	-
Peak memory footprint	-	200Mb	-
Overall on device energy cost	1.1mwh	0.4mwh	2.8X

## 6.4 STR Efficiency

Finally, we show the efficiency of our STR models in Table 6 when running on testing devices. The model export steps generated on-device compatible models with the total size around 8Mb. Running the models on hardware accelerator provided huge gain in terms of both latency (9x) and battery usage (3x).

## 7 CONCLUSION

This paper presented LUMOS, one of the first smart multimodal assistant with strong text understanding capabilities which is also device compatible. Our comprehensive evaluation demonstrates the effectiveness of our proposed method, outperforming existing approaches in terms of accuracy. Additionally, we have shown that our system meets the stringent latency, size, memory, power, and compute requirements for on-device deployment. Overall, our work represents a significant step towards enabling MM-LLMs to read in real-world scenarios, paving the way for more advanced applications in the fields of computer vision and natural language processing. Future work includes further optimizations to our on-device models, and research on end-to-end text recognition and visual translation with multimodal large language models.

## ACKNOWLEDGMENTS

The authors would like to thank Mei-Yuh Hwang, Praveen Krishnan, Guan Pang, Becka Silvert, Renato Sanchez, Crystal Nakatsu, Lucas Kabela, Frank Seide, Samyak Datta, Peyman Heidari, Shashank Jain, Nish Gupta, Kate Ovchinnikova, Rongzhou Shen, Saumya Mukul, Shane Moon, David Strauss, Lintao Cui, Sofiane Djeflal, Megha Tiwari, Vitaly Berov, Shanying Luo for their valuable inputs and contributions.



## REFERENCES

- [1] [n. d.]. AWS Rekognition. <https://aws.amazon.com/rekognition/>
- [2] [n. d.]. Google Cloud OCR. <https://cloud.google.com/vision/docs/ocr>
- [3] OpenAI (2023). 2023. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL]
- [4] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. Flamingo: a Visual Language Model for Few-Shot Learning. arXiv:2204.14198 [cs.CV]
- [5] Stanislaw Antol, Aishwarya Agrawal, Jiaseen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual Question Answering. *CoRR* abs/1505.00468 (2015). arXiv:1505.00468 <http://arxiv.org/abs/1505.00468>
- [6] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. OpenFlamingo: An Open-Source Framework for Training Large Autoregressive Vision-Language Models. arXiv:2308.01390 [cs.CV]
- [7] Fedor Borisjuk, Albert Gordo, and Viswanath Sivakumar. 2018. Rosetta: Large Scale System for Text Detection and Recognition in Images. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (London, United Kingdom) (KDD '18). Association for Computing Machinery, New York, NY, USA, 71–79. <https://doi.org/10.1145/3219819.3219861>
- [8] Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V. Le. 2019. RandAugment: Practical data augmentation with no separate search. *CoRR* abs/1909.13719 (2019). arXiv:1909.13719 <http://arxiv.org/abs/1909.13719>
- [9] Yuning Du, Chenxia Li, Ruoyu Guo, Xiaoting Yin, Weiwei Liu, Jun Zhou, Yifan Bai, Zilin Yu, Yehua Yang, Qingqing Dang, and Haoshuang Wang. 2020. PP-OCR: A Practical Ultra Lightweight OCR System. arXiv:2009.09941 [cs.CV]
- [10] Hao Feng, Zijian Wang, Jingqun Tang, Jinghui Lu, Wengang Zhou, Houqiang Li, and Can Huang. 2023. UniDoc: A Universal Large Multimodal Model for Simultaneous Text Detection, Recognition, Spotting and Understanding. arXiv:2308.11592 [cs.AI]
- [11] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2018. Mask R-CNN. arXiv:1703.06870 [cs.CV]
- [12] Wenbo Hu, Yifan Xu, Yi Li, Weiye Li, Zeyuan Chen, and Zhuowen Tu. 2023. BLIVA: A Simple Multimodal LLM for Better Handling of Text-Rich Visual Questions. arXiv:2308.09936 [cs.CV]
- [13] Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and C. V. Jawahar. 2019. ICDAR2019 Competition on Scanned Receipt OCR and Information Extraction. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. 1516–1520. <https://doi.org/10.1109/ICDAR.2019.00244>
- [14] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2016. Reading Text in the Wild with Convolutional Neural Networks. *Int. J. Comput. Vision* 116, 1 (jan 2016), 1–20. <https://doi.org/10.1007/s11263-015-0823-z>
- [15] Minghui Liao, Pengyuan Lyu, Minghang He, Cong Yao, Wenhao Wu, and Xiang Bai. 2019. Mask TextSpotter: An End-to-End Trainable Neural Network for Spotting Text with Arbitrary Shapes. arXiv:1908.08207 [cs.CV]
- [16] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning. arXiv:2304.08485 [cs.CV]
- [17] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. 2018. Path Aggregation Network for Instance Segmentation. arXiv:1803.01534 [cs.CV]
- [18] Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Yang Liu, Biao Yang, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Xucheng Yin, Cheng lin Liu, Lianwen Jin, and Xiang Bai. 2023. On the Hidden Mystery of OCR in Large Multimodal Models. arXiv:2305.07895 [cs.CV]
- [19] Shangbang Long, Siyang Qin, Dmitry Pantelev, Alessandro Bissacco, Yasuhisa Fujii, and Michalis Raptis. 2022. Towards End-to-End Unified Scene Text Detection and Layout Analysis. arXiv:2203.15143 [cs.CV]
- [20] Jianqi Ma, Weiyuan Shao, Hao Ye, Li Wang, Hong Wang, Yingbin Zheng, and Xiangyang Xue. 2018. Arbitrary-Oriented Scene Text Detection via Rotation Proposals. *IEEE Transactions on Multimedia* 20, 11 (Nov. 2018), 3111–3122. <https://doi.org/10.1109/tmm.2018.2818020>
- [21] Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. 2021. DocVQA: A Dataset for VQA on Document Images. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*. 2199–2208. <https://doi.org/10.1109/WACV48630.2021.00225>
- [22] Rachit S Munjal, Arun D Prabhu, Nikhil Arora, Sukumar Moharana, and Gopi Ramena. 2021. STRIDE: Scene Text Recognition In-Device. In *2021 International Joint Conference on Neural Networks (IJCNN)*. 1–8. <https://doi.org/10.1109/IJCNN52387.2021.9534319>
- [23] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tijmen Blankevoort. 2021. A White Paper on Neural Network Quantization. arXiv:2106.08295 [cs.LG]
- [24] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2016. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. arXiv:1506.01497 [cs.CV]
- [25] Baoguang Shi, Xiang Bai, and Cong Yao. 2017. An End-to-End Trainable Neural Network for Image-Based Sequence Recognition and Its Application to Scene Text Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 39, 11 (nov 2017), 2298–2304. <https://doi.org/10.1109/TPAMI.2016.2646371>
- [26] Baoguang Shi, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. 2016. Robust Scene Text Recognition with Automatic Rectification. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 4168–4176. <https://doi.org/10.1109/CVPR.2016.452>
- [27] Yongxin Shi, Dezhi Peng, Wenhui Liao, Zening Lin, Xinhong Chen, Chongyu Liu, Yuyi Zhang, and Lianwen Jin. 2023. Exploring OCR Capabilities of GPT-4V(ision): A Quantitative and In-depth Evaluation. arXiv:2310.16809 [cs.CV]
- [28] Amanpreet Singh, Guan Pang, Mandy Toh, Jing Huang, Wojciech Galuba, and Tal Hassner. 2021. TextOCR: Towards large-scale end-to-end reasoning for arbitrary-shaped scene text. arXiv:2105.05486 [cs.CV]
- [29] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, and et al. 2023. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805 [cs.CL]
- [30] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. 2019. FCOS: Fully Convolutional One-Stage Object Detection. arXiv:1904.01355 [cs.CV]
- [31] Alvin Wan, Xiaoliang Dai, Peizhao Zhang, Zijian He, Yuandong Tian, Saining Xie, Bichen Wu, Matthew Yu, Tao Xu, Kan Chen, Péter Vajda, and Joseph Gonzalez. 2020. FBNetV2: Differentiable Neural Architecture Search for Spatial and Channel Dimensions. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020)*, 12962–12971. <https://api.semanticscholar.org/CorpusID:215744832>
- [32] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. 2021. MaX-DeepLab: End-to-End Panoptic Segmentation with Mask Transformers. arXiv:2012.00759 [cs.CV]
- [33] Kai Wang and Serge Belongie. 2010. Word Spotting in the Wild. In *Proceedings of the 11th European Conference on Computer Vision: Part I* (Heraklion, Crete, Greece) (ECCV'10). Springer-Verlag, Berlin, Heidelberg, 591–604.
- [34] Xue Yang, Xiaojiang Yang, Jirui Yu, Qinghao Ye, Ming Yan, Guohai Xu, Qian Tian, and Junchi Yan. 2022. Learning High-Precision Bounding Box for Rotated Object Detection via Kullback-Leibler Divergence. arXiv:2106.01883 [cs.CV]
- [35] Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He, Xin Alex Lin, and Fei Huang. 2023. UReader: Universal OCR-free Visually-situated Language Understanding with Multimodal Large Language Model. arXiv:2310.05126 [cs.CV]
- [36] Haoyang Zhang, Ying Wang, Feras Dayoub, and Niko Sünderhauf. 2021. VarifocalNet: An IoU-aware Dense Object Detector. arXiv:2008.13367 [cs.CV]
- [37] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models. arXiv:2304.10592 [cs.CV]

## A APPENDIX



Figure 7: Sample input image for LUMOS



Figure 8: Example of a hand that is not indicative of ROI

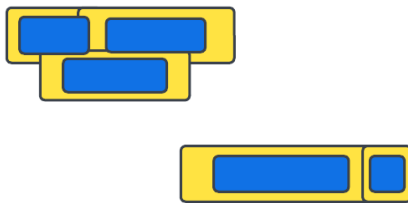


Figure 9: How we construct paragraphs: Blue boxes are from text detection which is then expanded (yellow). There are two paragraphs based on how the yellow boxes overlap

**Algorithm 1** Reading Order Reconstruction Algorithm

- 1: **Input** Recognized word  $w_1, w_2 \dots w_n$ , their corresponding bounding box  $b_1, b_2 \dots b_n$ , vertical expansion ratio  $r_v$ , horizontal expansion ratio  $r_h$ , iou threshold  $t$
- 2: **Step 1** Expand all boxes  $b_1, b_2 \dots b_n$  by  $r_v$  vertically and  $r_h$  horizontally to get expanded boxes  $\hat{b}_1, \hat{b}_2 \dots \hat{b}_n$  as shown in figure 9
- 3: **Step 2** Compute pairwise IOU of  $\hat{b}_1, \hat{b}_2 \dots \hat{b}_n$ , construct adjacency matrix  $A \in [0, 1]^{n \times n}$  where each element of  $A$  is defined as

$$a_{i,j} = \begin{cases} 1, & \text{if } iou(\hat{b}_i, \hat{b}_j) \geq t \\ 0, & \text{otherwise} \end{cases}$$

- 4: **Step 3** Finding connected components

$$C_1 = \{1_1, 1_2, \dots\}, \dots, C_m = \{m_1, m_2, \dots\}$$

from the graph defined by the adjacency matrix  $A$  (each word box is a vertex) where  $k_l$  is the index of the  $l^{th}$  box in the  $k^{th}$  connected component. Each connect component stands for a paragraph.

- 5: **Step 4** Applying raster scan to words within each paragraph (connected components) to generate words in reading order for each paragraph.
- 6: **Step 5** Compute the location of the paragraph by finding the minimum area rectangle enclosing all words in each connects components
- 7: **Output** Words in each paragraphs in reading order:  $P_1 = \{w_{1_1}, w_{1_2}, \dots\}, \dots, P_m = \{w_{m_1}, w_{m_2}\},$
- 8: Locations of the rectangular area for each paragraph in the format of (*top, left, height, width*).

**Table 7: Sample prompt to different variants of LUMOS for Figure 7**

System	Prompt	Response	Relevant
MM-LLM only	What does this sign say	The sign says, "Haz Matado" which translates to "You've killed it."	irrelevant
MM-LLM+STR	Here is a transcription of the recognized text in the image, which may contain misspelled words. [HAZ VOLUNTARIADO en esta tienda R OKFAM intermon]. My question is what does this sign sayin	The sign says, "HAZ VOLUNTARIADO en esta tienda" which translates to "DO VOLUNTEER WORK this store"	relevant
MM-LLM+STR+Position	I have included the coordinates of the text within the image. The coordinates (0,0) is indicative of the top-left corner of the image and (3024, 4032) is indicative of the bottom right corner of the image. For example, the tuple (zzzz, 10, 20) means the paragraph zzzz is centered at the point (10, 20) within the image. Here is a transcription of the recognized text in the image according to this coordinate system. It may contain misspelled words. [(HAZ VOLUNTARIADO en esta tienda , 1260, 1532) (R OKFAM intermon, 1282, 1939)]. My question is what does this sign say.	The sign says, "HAZ VOLUNTARIADO en esta tienda" which translates to "DO VOLUNTEER WORK in this store"	relevant