A Bounding Box is Worth One Token: Interleaving Layout and Text in a Large Language Model for Document Understanding

Anonymous Author(s) Affiliation Address email

Abstract

Recently, many studies have demonstrated that exclusively incorporating OCR-1 2 derived text and spatial layouts with large language models (LLMs) can be highly effective for document understanding tasks. However, existing methods that in-З tegrate spatial layouts with text have limitations, such as producing overly long 4 text sequences or failing to fully leverage the autoregressive traits of LLMs. In 5 this work, we introduce Interleaving Layout and Text in a Large Language Model 6 (LayTextLLM) for document understanding. In particular, LayTextLLM projects 7 each bounding box to a single embedding and interleaves it with text, efficiently 8 avoiding long sequence issues while leveraging autoregressive traits of LLMs. 9 LayTextLLM not only streamlines the interaction of layout and textual data but 10 also shows enhanced performance in Key Information Extraction (KIE) and Visual 11 Question Answering (VQA). Comprehensive benchmark evaluations reveal signifi-12 cant improvements, with a 27.0% increase on KIE tasks and 24.1% on VQA tasks 13 compared to previous state-of-the-art document understanding MLLMs, as well as 14 a 15.5% improvement over other SOTA OCR-based LLMs on KIE tasks. 15

16 **1 Introduction**

Recent research has increasingly focused on applying Large Language Models (LLMs) [1–17] to 17 document-oriented Visual Question Answering (VQA) and Key Information Extraction (KIE) scenar-18 ios. Efforts to build a text-sensitive MultiModal Large Language Models (MLLMs) based on existing 19 LLMs, particularly aimed at enhancing Visually Rich Document Understanding (VRDU), have made 20 significant progress [6, 12, 18]. Although existing MLLMs show promising results in document 21 understanding, they often encounter challenges related to image resolution. When the input image is 22 of low resolution, it is too blurry to extract visual features effectively. Conversely, high-resolution 23 images require additional computational resources to capture detailed textual information [12]. 24

Concurrently, another line of research employs off-the-shelf OCR tools to extract text and spatial layouts, which are then combined with LLMs to address VRDU tasks. These approaches assume that *most valuable information for document comprehension can be derived from the text and its spatial layouts, viewing spatial layouts as "lightweight visual information" [19]*. Following this premise, several studies [12, 20–23] have explored various approaches that integrate spatial layouts with text for LLMs, achieving results that are competitive with, or even surpass, those of MLLMs.

The most natural method to incorporate layout information is by treating spatial layouts as tokens, which allows for the seamless interleaving of text and layout into a unified text sequence [20, 22, 23]. For example, Perot et al. [20] employ format such as *"HARRISBURG 78*/09" to represent OCR text

and corresponding layout, where "*HARRISBURG*" is OCR text and "78/09" indicates the mean of

Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.

the horizontal and vertical coordinates, respectively. Similarly, He et al. [23] use " $[x_min, y_min, x_max, y_max]$ " to represent layout information. These approaches can effectively take advantage of autoregressive characteristics of LLMs and is known as the "*coordinate-as-tokens*" scheme [20]. In contrast, DocLLM [19] explores interacting spatial layouts with text through a disentangled spatial attention mechanism that captures cross-alignment between text and layout modalities.

However, we argue that both of the previous ap-40 proaches have limitations. As shown in Fig. 1, 41 coordinate-as-tokens significantly increases the 42 number of tokens. Additionally, to accurately 43 comprehend coordinates and enhance zero-shot 44 capabilities, this scheme often requires few-45 shot in-context demonstrations and large-scale 46 language models, such as ChatGPT Davinci-47 003 (175B) [23], which exacerbates issues re-48 lated to sequence length and GPU resource de-49 mands. Meanwhile, although DocLLM does not 50 increase sequence length and integrates spatial 51 layouts through attention, its generalizability is 52 limited. We believe that spatial cross attention 53 and masked span tasks in DocLLM cannot fully 54 utilize the autoregressive traits of LLMs. 55

56 To address these problems, this paper explores a 57 simple yet effective approach to enhance the in-

⁵⁸ teraction between spatial layouts and text — *In*-

59 terleaving Layout and Text in a Large Language

60 *Model (LayTextLLM)* for document understand-

61 ing. Adhering to the common practice of inter-

- leaving any modality with text [15, 24, 25], we
- specifically apply this principle to spatial lay outs. In particular, we maps each bounding box



Figure 1: The performance against input sequence length of different datasets across various OCR-based methods where data is from Tab. 2 and 5.

to a single embedding, which is then interleaved with its corresponding text. Then we propose a
tailored pre-training task—Layout-aware Next Token Prediction—a completely self-supervised task
that enhances the alignment between layout and textual modalities without using synthetic data.
Finally, through the proposed Shuffled-OCR Supervised Fine-tuning, LayTextLLM significantly
improves performance on downstream document-related VQA and KIE tasks. As shown in Fig. 1,
LayTextLLM significantly outperforms the 175B models, while only slightly increasing or even
reducing the sequence length compared to DocLLM. Our contributions can be listed as follows:

- We propose LayTextLLM for document understanding. To the best of the authors' knowledge, this is the first work to employ a unified embedding approach (Sec. 3.1.1) that interleaves spatial layouts directly with textual data within a LLM. By representing each bounding box with one token, LayTextLLM efficiently addresses sequence length issues brought by coordiante-as-tokens while fully leveraging autoregressive traits for enhanced document understanding.
- We propose two tailored training tasks: (1) Layout-aware Next Token Prediction (Sec. 3.2.1),
 a completely self-supervised training task to enhance the alignment between layout and
 textual modality; (2) Shuffled-OCR Supervised Fine-tuning task (Sec. 3.2.2) to better elicit
 the model generalizability in downstream tasks.
- Comprehensive experimental results demonstrate quantitatively that LayTextLLM significantly outperforms previous state-of-the-art (SOTA) OCR-free MLLMs by a large margin in zero-shot scenarios, particularly in KIE tasks with an improvement of 27.0%. Additionally, we illustrate that LayTextLLM competes effectively or even surpasses previous SOTA OCR-based methods in both zero-shot and SFT scenarios. Specifically, it surpasses DocLLM by 19.8% on VQA and 15.5% on KIE tasks (Sec. 4).
- Extensive ablations demonstrate the utility of the proposed component, with analysis showing that LayTextLLM not only improves performance but also reduces input sequence length compared to current OCR-based models.

91 2 Related Work

92 2.1 OCR-based LLMs for Document Understanding

Early document understanding methods [26-30] tend to solve the task in a two-stage manner, *i.e.*, first 93 reading texts from input document images using off-the-shelf OCR engines and then understanding 94 the extracted texts. Considering the advantages of LLMs (e.g., high generalizability), some recent 95 methods endeavor to combine LLMs with OCR-derived results to solve document understanding. 96 97 For example, inspired by the "coordinate-as-tokens" scheme [20], He et al. [23] propose to use " $[x_{min}, y_{min}, x_{max}, y_{max}]$ " to introduce the layout information, which can fuse the layout 98 information and texts into a unified text sequence and fully exploit the autoregressive merit of LLMs. 99 To reinforce the layout information while avoiding increasing the number of tokens, DocLLM [19] 100 designs a disentangled spatial attention mechanism to capture cross-alignment between text and 101 layout modalities. Recently, LayoutLLM [21] utilizes the pre-trained layout-aware model [31], to 102 insert the visual information, layout information and text information. However, the aforementioned 103 methods neither suffer from the computational overhead leading by the increasing tokens or hardly 104 take advantage of autoregressive characteristics of LLMs. Thus, it is an urgent problem to address 105 how to better incorporate layout information without significantly increasing the number of tokens. 106

107 2.2 OCR-free MLLMs for Document Understanding

Another approach to solve document understanding tasks is the OCR-free method. Benefiting from 108 the end-to-end training framework, it involves processing the text content of documents directly, 109 without relying on OCR engines. Donut [32] first presents an OCR-free method through mapping 110 a text-rich document image into the desired answers. Pix2Struct [33] is trained to parse masked 111 screenshots of web pages into simplified HTML, where variable resolution inputs are supported. 112 While these approaches eliminate the need for OCR tools, they still necessitate task-specific fine-113 tuning. With the increasing popularity of LLMs/MLLMs [10–17], various methods are proposed to 114 solve the document understanding task through explicitly training models on visual text understanding 115 datasets and fine-tuning them with instructions to perform a zero-shot prediction. LLaVAR [34] 116 and UniDoc [10] are notable examples that expand upon the document-oriented VQA capabilities 117 of LLaVA [35] by incorporating document-based tasks. These models pioneer the use of MLLMs 118 for predicting texts and coordinates from document images, enabling the development of OCR-119 free document understanding methods. Additionally, DocPedia [9] operates document images in 120 the frequency domain, allowing for higher input resolution without increasing the input sequence 121 length. Recent advancements in this field, including mPLUG-DocOwl [18], Qwen-VL [6], and 122 TextMonkey [12], leverage publicly available document-related VQA datasets to further enhance 123 the document understanding capability. Although these OCR-free methods have exhibited their 124 advantages, they still struggle with the high-resolution input to reserve more text-related details. 125

126 **3** Method

In this section, we present our LayTextLLM. First, we introduce a innovative Spatial Layout Projector (Sec. 3.1.1) converts four-dimensional layout coordinates into a single-token embedding. To reduce parameter overhead, we apply Partial Low-Rank Adaptation (Sec. 3.1.2). We also introduce two specific training tasks: Layout-aware Next Token Prediction (Sec. 3.2.1) to align layouts with text during pre-training, and Shuffled-OCR Supervised Fine-tuning (Sec. 3.2.2) to enhance the generalizability of the model. An illustration of our approach is shown in Fig. 2.

133 **3.1 Model Architecture**

LayTextLLM is built on the Llama2-7B-base model, which was originally designed to accept only
 text inputs [36, 37]. To enable the model to interleave spatial layouts with text, we introduce a
 novel Spatial Layout Projector. This projector converts OCR-derived coordinates into bounding box
 tokens. We also adopt the Partial Low-Rank Adaptation, a minimally invasive method to incorporate
 additional modalities while preserving the LLM's inherent knowledge intact.



Figure 2: An overview of LayTextLLM incorporates interleaving bounding box tokens (b^i) with text tokens (t^i) , where the superscripts represent the sequence positions of the tokens.

139 3.1.1 Spatial Layout Projector (SLP)

A key innovation in LayTextLLM is the Spatial Layout Projector (SLP), which transforms a spatial 140 layout into a singular bounding box token. This enhancement enables the model to process both 141 spatial layouts and textual inputs simultaneously. To be specifically, each OCR-derived spatial layout 142 is represented by a bounding box defined by four-dimensional coordinates $[x_1, y_1, x_2, y_2]$, these 143 coordinates represent the normalized minimum and maximum horizontal and vertical extents of the 144 box, respectively. The SLP maps these coordinates into a high-dimensional space that the language 145 model can process as a single token. The process can be computed as $z = W \cdot c + b$, where $c \in \mathbb{R}^4$ 146 is the vector of the bounding box coordinates. $W \in \mathbb{R}^{d \times 4}$ is a weight matrix with d represents 147 the dimension of the embedding, $b \in \mathbb{R}^{d \times 1}$ is a bias vector, z is the resulting bounding box token 148 represented as an *d*-dimensional embedding. As illustrated in Fig. 2, the resulting bounding box 149 token z will be interleaved with corresponding textual embeddings to put into LLMs. Note that the 150 SLP is shared by all bounding box tokens so very limited number of parameters are introduced. 151

Compared to the coordinate-as-tokens scheme, the SLP represents each bounding box with a single token. This approach significantly reduces the number of input tokens and adheres to the practice of interleaving any modality with text, effectively integrating layout and textual information into a unified sequence. This allows the model to process both modalities simultaneously and coherently, fully leveraging the autoregressive traits of LLMs.

157 3.1.2 Layout Partial Low-Rank Adaptation



¹⁷³ Figure 3: The illustration of P-LoRA, adapted from [15].

After using the SLP to generate bounding box tokens and a tokenizer to produce text tokens, these two modalities are then communicated using a Layout Partial Low-Rank Adaptation (P-LoRA) module in LLMs. P-LoRA, introduced in InternLM-XComposer2 [15], is originally used to adapt LLMs to visual modality. It applies plug-in low-rank modules specified to the visual tokens, which adds minimal parameters while preserving the LLMs inherent knowledge.

Formally, as shown in Fig. 3 for a linear layer in the LLM, the original weights $W_O \in \mathbb{R}^{C_{out} \times C_{in}}$ and bias $B_O \in \mathbb{R}^{C_{out}}$ are specified for input and output dimensions C_{in} and C_{out} . P-LoRA modifies this setup by incorporating two additional matrices, $W_A \in \mathbb{R}^{C_r \times C_{in}}$ and $W_B \in \mathbb{R}^{C_{out} \times C_r}$. These matrices are lower-rank, with C_r being considerably smaller than both C_{in} and C_{out} , and are specifically designed to interact with new modality tokens, which in our case are bounding box tokens. For example, given an input x =



Figure 4: Comparison of Layout-aware Next Token Prediction and normal Next Token Prediction.

175 $[x_b, x_t]$ comprising of bounding box tokens (x_b) and textual tokens (x_t) is fed into the system, the 176 forward process is as follows, where \hat{x}_t, \hat{x}_b and \hat{x} are outputs:

$$\hat{x}_{t} = W_{0}x_{t} + B_{0}
\hat{x}_{b} = W_{0}x_{b} + W_{B}W_{A}x_{b} + B_{0}
\hat{x} = [\hat{x}_{b}, \hat{x}_{t}]$$
(1)

177 3.2 Training Procedure

LayTextLLM is trained with innovative layout-aware training procedure, which consists of two stages:
 Layout-aware Next Token Prediction pre-training and Shuffled-OCR Supervised Fine-tuning.

180 3.2.1 Layout-aware Next Token Prediction

Inspired by the next token prediction commonly used in current LLM pre-training [1-7], we propose 181 the Layout-aware Next Token Prediction (LNTP). Fig. 4 presents the contrast of the proposed Layout-182 aware Next Token Prediction and the conventional next token prediction task. The traditional next 183 token prediction (Fig. 4(a)) relies solely on the textual content, predicting each subsequent token 184 based on the prior sequence of tokens without considering their spatial layouts. Layout-aware next 185 token prediction (Fig. 4(b)), however, interleaves the spatial information encoded by SLP (*i.e.*, b^i) 186 with the text tokens (*i.e.*, t^i). This integration considers both the content and its layout within the 187 document, leading to a richer, more precise understanding of both the structure and the content. 188

189 190

191

CASH	I (MYR)		20.00
Char	nge	1.30		
	IST%	Amt(RM)	GST(RM)	Total(RM)
SR	6	17.64	1.06	18.70
 GC)ODS S	old are non	-cash refu	NDABLE.
	EXCHA	NGE OF GOOD	S WITHIN 1	4 DAYS
	ACCOM	PANIED BY O	RIGINAL RE	CEIPT.

Figure 5: Receipt layout example.

197 198 199

192

193

194

195

196

200 the parameters of SLP and P-LoRA.

Similarly, primary objective of LNTP is to maximize the likelihood of its predictions for the next token. Thus the loss function is defined as

$$\mathcal{L} = -\frac{1}{T} \sum_{i=1}^{T} \log P\left(t^{i} \mid t^{1}, t^{2}, \dots, t^{i-1}\right)$$
(2)

where $P(t^i | t^1, t^2, \ldots, t^{i-1})$ represents the probability of i^{th} token t^i given the sequence of preceding tokens $t^1, t^2, \ldots, t^{i-1}$, as predicted by the model. Note that we compute the loss only for text tokens, excluding bounding box tokens. During pre-training, our goal is to enhance the alignment between spatial layouts and textual modality, while preserving the LLM's inherent knowledge as much as possible. Thus, we freeze the LLMs and only update

It is important to note that the proposed Layout-aware Next Token Prediction is a completely selfsupervised pre-training procedure, unlike previous works that require human annotations of document structure data or synthetic data generated by larger LLMs such as GPT-4 [21]. Thus, LNTP facilitates the creation of large-scale, high-fidelity pre-training datasets at minimal cost.

205 3.2.2 Shuffled-OCR Supervised Fine-tuning

OCR engines typically process text from top to bottom and left to right. This order is also adopted as the input sequence for current OCR-based LLMs [19, 21]. However, modern LLMs often exhibit a strong inductive bias toward the positions of input tokens, influenced by designs such as Rotary Position Embeddings (RoPE) [38]. Specifically, tokens that are close together in the input sequence are likely to receive higher attention scores, which is advantageous for processing standard text sequences. Such inductive bias brings cons and pros.

Consider the example illustrated in Fig. 5, where the OCR input text reads: " ... Change, 1.30, GST%, 212 Amt(RM), GST(RM), Total(RM), SR, 6, 17.64, 1.06, 18.70 ... ". If the question posed is "What is the 213 value of the field Change?" (highlighted in a blue box), the model easily identifies "1.30" as it is 214 closely positioned to the word "Change" in the sequence. However, for a more challenging query like 215 "What is the value of the field Total(RM)?" (highlighted in a red box), the model struggles to determine 216 the correct answer due to the presence of multiple subsequent numbers closed to "Total(RM)". 217 LayTextLLM integrates spatial layouts with textual data, reducing reliance on input sequence order. 218 Thus, we posit that shuffling the OCR input order could enhance the resilience of LayTextLLM in 219 discerning relevant information irrespective of token proximity in the sequence. 220

Specifically, we propose Shuffled-OCR Supervised Fine-tuning (SSFT) that randomly shuffles the order of OCR-derived text in a certain proportion of examples. The range of exploration for the shuffling ratio can be found in Tab. 7 and 20% shuffled ratio is applied. The training objective is equivalent to predicting the next tokens, but in this scenario, only the tokens of the response are used to compute loss. During SSFT, we unfreeze all parameters including those of LLMs. Experimental results in Section 4.6 demonstrate that utilizing SSFT can further enhance model performance, making it more robust to disruptions in input token order.

228 4 Experiments

229 4.1 Datasets

Pre-training data In our training process, we exclusively use open-source data to facilitate replication. We collect data from two datasets for pre-training: (1) IIT-CDIP Test Collection 1.0 [39] and (2) DocBank [40]. The IIT-CDIP Test Collection 1.0 comprises an extensive repository of more than 16 million document pages. DocBank consists of 500K documents, each presenting distinct layouts with a single page per document. For training efficiency, we choose to utilize the entire DocBank dataset and only subsample 5 million pages from the IIT-CDIP collection 1.0.

SFT data For document-oriented VQA, we select Document Dense Description (DDD) and
Layout-aware SFT data used in Luo et al. [21], which are two synthetic datasets generated by
GPT-4. Besides, DocVQA [41], InfoVQA [42], ChartQA [43], VisualMRC [44] is included
following [12]. For KIE task, we select SROIE [45], CORD [46], FUNSD [47], POIE [48] datasets
following [12, 19, 21].

241 4.2 Implementation Detail

The LLM component of LayTextLLM is initialized from the Llama2-7B-base [36], which is a widelyused backbone. Other parameters including SLP and P-LoRA are randomly initialized. During pre-training, the LLM is frozen, and the parameters of SLP and P-LoRA modules are updated. During SFT, all parameters are fine-tuned. Other detailed setup can be found in Appendix B.

We have configured the model with three versions of LayTextLLM for a side-by-side comparison 246 under different settings. Aligned with Luo et al. [21], the first version, LayTextLLM_{zero}, is 247 trained exclusively with DDD and Layout-aware SFT data. Building upon this, and in alignment 248 with the setting of Liu et al. [12], we introduce the DocVQA, InfoVQA, and ChartQA training 249 sets to the dataset pool for our second version, termed LayTextLLM_{vga}. Finally, we incorporate 250 a comprehensive suite of KIE datasets-FUNSD, CORD, POIE, SROIE, and VisualMRC-as 251 described by Wang et al. [19], creating our most extensive version, LayTextLLM_{all}. Note that all 252 versions are based on the same pre-trained LayTextLLM weight. 253

	Document-Oriented VQA				KIE		
	DocVQA	InfoVQA	Avg	FUNSD	SROIE	POIE	Avg
Metric			Acc	curacy %			
OCR-free							
UniDoc [10]	7.7	14.7	11.2	1.0	2.9	5.1	3.0
DocPedia [9]	47.1*	15.2*	31.2	29.9	21.4	39.9	30.4
Monkey [49]	50.1*	25.8^{*}	38.0	24.1	41.9	19.9	28.6
InternVL [50]	28.7*	23.6^{*}	26.2	6.5	26.4	25.9	19.6
InternLM-XComposer2 [15]	39.7	28.6	34.2	15.3	34.2	49.3	32.9
TextMonkey [12]	64.3*	28.2^{*}	46.3	32.3	47.0	27.9	35.7
TextMonkey ₊ [12]	66.7*	28.6^{*}	47.7	42.9	46.2	32.0	40.4
text + polys							
LayTextLLM _{zero} (Ours)	71.8	33.8	52.8	49.4	86.7	66.1	67.4
LayTextLLM _{vqa} (Ours)	77.4*	66.1 *	71.8	48.9	74.6	70.6	64.7

Table 1: Comparison with SOTA OCR-free MLLMs. * indicates the training set used.

	Document-Oriented VQA				KIF	2	
	DocVQA	VisualMRC	Avg	FUNSD	CORD	SROIE	Avg
Metric	AN	LS % / CIDEr			F-score	e %	
Text							
Llama2-7B-base	34.0	182.7	108.3	25.6	51.9	43.4	40.3
Llama2-7B-chat	20.5	6.3	13.4	23.4	51.8	58.6	44.6
Text + Polys							
Llama2-7B-base _{coor} [23]	8.4	3.8	6.1	6.0	46.4	34.7	29.0
Llama2-7B-chat _{coor} [23]	12.3	28.0	20.1	14.4	38.1	50.6	34.3
Davinci-003-175B _{coor} [23]	-	-	-	-	92.6	95.8	-
DocLLM [19]	69.5*	264.1*	166.8	51.8*	67.4*	91.9*	70.3
LayTextLLM _{zero} (Ours)	65.4	260.7	163.0	48.6	74.5	86.4	69.8
LayTextLLM _{vqa} (Ours)	75.7*	260.2^{*}	168.0	52.7	70.9	78.6	67.4
LayTextLLM _{all} (Ours)	77.3*	295.9*	186.6	64.2*	96.5 *	95.8 *	85.8

Table 2: Comparison with other OCR-based methods. * indicates the training set used.

254 4.3 Baselines

OCR-free baselines In the category of OCR-free MLLMs, we have chosen the following SOTA
 models as our strong baselines due to their superior performance in both document-oriented VQA
 and KIE tasks. These include UniDoc [10], DocPedia [9], Monkey [49], InternVL [50], InternLM-

258 XComposer2 [15], TextMonkey, and TextMonkey₊ [12].

OCR-based baselines For OCR-based baseline models, we implemented a basic approach using 259 only OCR-derived text as input. This was done using two versions: Llama2-7B-base and Llama2-260 7B-chat. We also adapted the coordinate-as-tokens scheme from He et al. [23] for these models, 261 resulting in two new variants: Llama2-7B-base coor and Llama2-7B-chat coor. It's important to note 262 that we did not employ the ICL strategy with these models, as it would significantly exceed their 263 maximum sequence length constraints. Additionally, we included results from a stronger baseline 264 using the ChatGPT Davinci-003 (175B) model [23], termed Davinci-003-175B_{coor}. One other recent 265 SOTA OCR-based approach, DocLLM [19] is also considered in our analysis. Finally, LayoutLLM 266 and LayoutLLM_{CoT} [21], which integrates visual cues, text and layout is also included. 267

268 **4.4 Evaluation Metrics**

To ensure a fair comparison with OCR-free methods, we adopted the accuracy metric, where a response from the model is considered correct if it fully captures the ground truth. This approach aligns with the evaluation criteria described by [9, 10, 12]. To further enhance the comparability with other OCR-based methods, we conducted additional evaluations using original metrics specific to certain datasets, such as F1 score [19, 23], ANLS [19, 21, 51] and CIDEr [19, 52].

	Docume	nt-Oriented V	QA		KIE		
	DocVQA	VisualMRC	Avg	FUNSD ⁻	$CORD^{-}$	SROIE ⁻	Avg
Metric				ANLS %			
Visual + Text + Polys							
LayoutLLM [21]	72.3	-	-	74.0	-	-	-
LayoutLLM $_{CoT}$ [21]	74.2	55.7	64.9	79.9	63.1	72.1	71.7
Text							
Llama2-7B-base	34.0	25.4	29.7	42.1	46.7	60.6	49.8
Llama2-7B-chat	20.5	9.9	15.2	15.1	20.0	35.6	23.5
Text + Polys							
Llama2-7B-base _{coor} [23]	8.4	6.7	7.5	4.3	33.0	47.2	28.1
Llama2-7B-chat _{coor} [23]	12.3	12.2	12.2	11.9	6.4	39.4	19.2
LayTextLLM _{zero} (Ours)	65.4	36.2	50.8	71.0	66.9	89.2	75.7
LayTextLLM _{all} (Ours)	77.3*	41.7^{*}	59.5	81.3 *	82.6 *	96.2 *	86.7

Table 3: Comparison with LayoutLLM. - indicates that the cleaned test set used in Luo et al. [21].

274 4.5 Quantitative Results

275 **Comparison with SOTA OCR-free Methods** The experimental results shown in Tab. 1 demon-276 strate the outstanding performance of the LayTextLLM series across various tasks. Note that the results for ChartQA are reported in Appendix E due to concerns about fairness in comparison, as 277 the dataset does not include OCR-derived results and we used in-house OCR tools instead. Firstly, 278 LayTextLLM_{zero} significantly outperforms previous SOTA OCR-free methods, such as TextMon-279 key [12], in zero-shot capabilities, even when these methods use the training set of the dataset. For 280 example, in the DocVQA and InfoVQA datasets, LayTextLLM_{zero} achieves accuracies of 71.8% and 281 282 33.8%, respectively, which are markedly higher than existing OCR-free methods such as TextMonkey and InternLM-XComposer2. When fine-tuned with corresponding datasets, LayTextLLM shows even 283 greater performance improvements, particularly in document-oriented VOA datasets. Specifically, 284 its accuracies on DocVQA and InfoVQA increase to 77.4% and 66.1%, respectively, demonstrating 285 the model's strong ability to leverage task-specific data. Additionally, LayTextLLM_{zero} excels in 286 KIE datasets, particularly on the SROIE and POIE datasets, achieving accuracies of 86.7% and 287 66.1%, respectively. These results significantly surpass those of previous SOTA OCR-free model (*i.e.*, 288 TextMonkey $_{+}$) by margins of 40.5% and 34.1%, respectively. This significant performance gain is 289 likely due to these datasets containing low-resolution images that are too blurred for current MLLMs 290 to extract visual features, whereas LayTextLLM shows robustness in such challenging scenarios. 291

Comparison with SOTA OCR-based Methods For comprehensive comparison, we have also 292 conducted correspinding experiments to align with OCR-based methods [19, 21]. The experimental 293 results presented in Tab. 2 showcase significant performance improvements achieved by LayTextLLM 294 models compared to pure OCR-based SOTA methods such as DocLLM [19]. Specifically, when 295 comparing with DocLLM, LayTextLLM_{zero} demonstrates notably superior performance, with even its 296 zero-shot capabilities being competitive with supervised SFT approaches. We believe that the subpar 297 performance of DocLLM is likely due to its use of cross-attention and the masked span pre-training 298 tasks [53], which fail to leverage the autoregressive features of LLMs effectively. Similarly, when 299 contrasting with coordinate-as-tokens employed in Llama2-7B, LayTextLLM_{zero} again outperforms 300 significantly. This disparity in performance can be attributed to the following three reasons: (1) The 301 coordinate-as-tokens approach tends to introduce an excessive number of tokens, often exceeding the 302 pre-defined maximum length of Llama2-7B (i.e., 4096). Consequently, this leads to a lack of crucial 303 304 OCR information, resulting in hallucination and subpar performance. (2) When re-implementing the coordinate-as-tokens method with Llama2-7B, we did not introduce the ICL strategy, as it would 305 contribute additional length to the input sequence. (3) The coordinate-as-tokens approach necessitates 306 a considerably larger-sized LLM to comprehend the numerical tokens effectively. 307

In comparison to LayoutLLM [21], our approach exhibits discrepant performance in different tasks, as shown in Tab. 3. In zero-shot scenarios, we outperform LayoutLLM in most KIE datasets, validating our capability to leverage OCR-based results effectively. However, we fall short on documentoriented VQA tasks since answering some questions that are strongly related to vision information may challenge our approach. Two main reasons may well explain this performance discrepancy: (1) The visual encoder in LayoutLLM provides additional visual information. (2) LayoutLLM
incorporates the Chain-of-Thought (CoT) mechanism to model contextual information while it is
not used in our approach. However, when fine-tuned with tailored data, LayTextLLM significantly
outperforms LayoutLLM, showcasing its strong ability to utilize task-specific data. More qualitative
example demonstrates can be found in Appendix A.

318 4.6 Analysis

Document-Oriented VQA								KIE		
LNTP	SSFT	DocVQA	InfoVQA	VisualMRC	Avg	FUNSD	CORD	SROIE	POIE	Avg
\checkmark	\checkmark	72.3 74.5	35.7 38.0	24.4 23.9	44.2 45.5	50.6 56.4	91.6 95.8	92.8 93.2	58.6 59.6	73.4 76.3
\checkmark	\checkmark	78.8	42.7	35.1	52.2	62.9	95.9	95.2	61.7	78.9

Table 4: Ablations on pre-training and SFT component of LayTextLLM (Accuracy).

Ablations To better assess the utility of Layout-aware Next Token Prediction and Shuffled-OCR Supervised Fine-tuning in LayTextLLM, an ablation study was performed (see Tab. 4). Details on the training setup for all variants are provided in Appendix B. It is evident that both LNTP and SSFT significantly enhance the utility of LayTextLLM. Specifically, disabling LNTP results in an 8% decrease in performance on VQA tasks and a 5.5% decrease on KIE tasks. Disabling SSFT leads to a decrease in average accuracy by 6.7% and 2.6% for VQA and KIE tasks, respectively.

Sequence Length Tab. 5 presents statistics on the average input sequence length across different 325 datasets. Intriguingly, despite interleaving bounding box tokens, LayTextLLM consistently exhibits 326 the shortest sequence length in three out of four datasets, even surpassing DocLLM, which is coun-327 terintuitive. We attribute this to the tokenizer mechanism. For example, using tokenizer.encode(), a 328 single word from the OCR engine, like "International" is encoded into a single ID [4623]. Conversely, 329 when the entire OCR output is processed as one sequence, such as "... CPC, International, Inc...", the 330 word "International" is split into two IDs [17579, 1288], corresponding to "Intern" and "ational" 331 respectively. This type of case occurs frequently, more discussion in Appendix C. 332

Dataset	LayTextLLM	DocLLM [19]	Coor-as-tokens [23]
DocVQA	664.3	827.5	4085.7
CORD	137.9	153.2	607.3
FUNSD	701.9	847.5	4183.4
SROIE	529.2	505.1	1357.7

Table 5: Average sequence length of each data for different methods using Llama2 tokenizer.

333 5 Limitation

Although LayTextLLM has shown significant capabilities in text-rich VQA and KIE tasks, this alone does not suffice for all real-world applications. There are some instances, particularly in chart analysis, where reasoning must be based solely on visual cues (*e.g.* size, color)—a challenge that remains unmet. Questions such as "*What is the difference between the highest and the lowest green bar?*" illustrate this gap. The ChartQA results, detailed in Appendix E, also underscore these limitations. Addressing these challenges highlights the urgent need for future enhancements that integrate visual cue within the capabilities of LayTextLLM.

341 6 Conclusion

We propose LayTextLLM for various VRDU tasks, in which spatial layouts and textual data are
seamlessly interleaved to make more accurate prediction by introducing a innovative Spatial Layout
Projector. Two tailored training tasks — Layout-aware Next Token Prediction and Shuffled-OCR
Supervised Fine-tuning — are designed to improve the comprehension of document layouts. Extensive
experiments confirm the effectiveness of LayTextLLM.

347 **References**

- [1] OpenAI: Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-348 349 ciaLeoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red 350 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, 351 Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim 352 Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany 353 Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek 354 Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, 355 HyungWon Chung, Dave Cummings, and Jeremiah Currier. Gpt-4 technical report. arXiv 356 preprint arXiv:2303.08774, Dec 2023. 357
- [2] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and
 Lijuan Wang. The dawn of Imms: Preliminary explorations with gpt-4v(ision). *arXiv preprint arXiv:2309.17421*, Sep 2023.
- [3] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [4] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.
- [5] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al.
 Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [6] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [7] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng
 Ren, Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language
 understanding. *arXiv preprint arXiv:2403.05525*, 2024.
- [8] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li,
 Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.
- [9] Hao Feng, Qi Liu, Hao Liu, Wengang Zhou, Houqiang Li, and Can Huang. Docpedia: Un leashing the power of large multimodal model in the frequency domain for versatile document
 understanding. *arXiv preprint arXiv:2311.11810*, 2023.
- [10] Hao Feng, Zijian Wang, Jingqun Tang, Jinghui Lu, Wengang Zhou, Houqiang Li, and Can Huang.
 Unidoc: A universal large multimodal model for simultaneous text detection, recognition,
 spotting and understanding. *arXiv preprint arXiv:2308.11592*, 2023.
- [11] Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin
 Jin, Fei Huang, et al. mPLUG-DocOwl 1.5: Unified structure learning for ocr-free document
 understanding. *arXiv preprint arXiv:2403.12895*, 2024.
- [12] Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai.
 Textmonkey: An ocr-free large multimodal model for understanding document. *arXiv preprint arXiv:2403.04473*, 2024.
- [13] Jingqun Tang, Chunhui Lin, Zhen Zhao, Shu Wei, Binghong Wu, Qi Liu, Hao Feng, Yang Li,
 Siqi Wang, Lei Liao, et al. Textsquare: Scaling up text-centric visual instruction tuning. *arXiv preprint arXiv:2404.12803*, 2024.
- [14] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong,
 Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to
 commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*,
 2024.

- [15] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei,
 Songyang Zhang, Haodong Duan, Maosong Cao, et al. Internlm-xcomposer2: Mastering
 free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024.
- [16] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu,
 Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision
 language models. *arXiv preprint arXiv:2403.18814*, 2024.
- [17] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae
 Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL https:
 //llava-vl.github.io/blog/2024-01-30-llava-next/.
- [18] Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai
 Xu, Chenliang Li, Junfeng Tian, et al. mPLUG-DocOwl: Modularized multimodal large
 language model for document understanding. *arXiv:2307.02499*, 2023.
- [19] Dongsheng Wang, Natraj Raman, Mathieu Sibue, Zhiqiang Ma, Petr Babkin, Simerjot Kaur,
 Yulong Pei, Armineh Nourbakhsh, and Xiaomo Liu. Docllm: A layout-aware generative
 language model for multimodal document understanding. *arXiv preprint arXiv:2401.00908*,
 2023.
- [20] Vincent Perot, Kai Kang, Florian Luisier, Guolong Su, Xiaoyu Sun, Ramya Sree Boppana,
 Zilong Wang, Jiaqi Mu, Hao Zhang, and Nan Hua. Lmdx: Language model-based document
 information extraction and localization. *arXiv preprint arXiv:2309.10952*, 2023.
- ⁴¹⁶ [21] Chuwei Luo, Yufan Shen, Zhaoqing Zhu, Qi Zheng, Zhi Yu, and Cong Yao. Layoutllm: Layout ⁴¹⁷ instruction tuning with large language models for document understanding. *CVPR* 2024, 2024.
- [22] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra:
 Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.
- Izabang He, Lei Wang, Yi Hu, Ning Liu, Hui Liu, Xing Xu, and Heng Tao Shen. Icl-d3ie:
 In-context learning with diverse demonstrations updating for document information extraction.
 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19485–19494, 2023.
- [24] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao
 Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning
 perception with language models. *arXiv:2302.14045*, 2023.
- [25] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu
 Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv:2306.14824*, 2023.
- [26] Wonseok Hwang, Jinyeong Yim, Seung-Hyun Park, Sohee Yang, and Minjoon Seo. Spatial
 dependency parsing for semi-structured document information extraction. *Cornell University arXiv, Cornell University arXiv*, May 2020.
- [27] Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. Layoutlm:
 Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Aug 2020.
 doi: 10.1145/3394486.3403172. URL http://dx.doi.org/10.1145/3394486.3403172.
- [28] Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei
 Florencio, Cha Zhang, Wanxiang Che, Min Zhang, and Lidong Zhou. Layoutlmv2: Multimodal pre-training for visually-rich document understanding. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Jan 2021. doi: 10.18653/
 v1/2021.acl-long.201. URL http://dx.doi.org/10.18653/v1/2021.acl-long.201.

- Teakgyu Hong, DongHyun Kim, Mingi Ji, Wonseok Hwang, Daehyun Nam, and Sungrae Park.
 Bros: A pre-trained language model focusing on text and layout for better key information
 extraction from documents. *Proceedings of the AAAI Conference on Artificial Intelligence*,
 page 10767–10775, Jul 2022. doi: 10.1609/aaai.v36i10.21322. URL http://dx.doi.org/
 10.1609/aaai.v36i10.21322.
- [30] Zineng Tang, Zhenfeng Yang, Guoxin Wang, Yuwei Fang, Yang Liu, Zhu C, Michael Zeng,
 Zhang Cha, and Mohit Bansal. Unifying vision, text, and layout for universal document
 processing. *Cornell University arXiv,Cornell University arXiv*, Dec 2022.
- [31] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. Layoutlmv3: Pre-training for
 document ai with unified text and image masking. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4083–4091, 2022.
- [32] Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeong Yeon Nam, Jinyoung Park, Jinyeong Yim,
 Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document
 understanding transformer. In *European Conference on Computer Vision*, pages 498–517, 2022.
- [33] Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisensch los, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct:
 Screenshot parsing as pretraining for visual language understanding. In *International Confer- ence on Machine Learning*, pages 18893–18912. PMLR, 2023.
- [34] Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun.
 Llavar: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023.
- [35] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances
 in neural information processing systems, 36, 2024.
- [36] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [37] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan
 Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction
 model. *arXiv:2304.15010*, 2023.
- [38] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer:
 Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [39] D. Lewis, G. Agam, S. Argamon, O. Frieder, D. Grossman, and J. Heard. Building a test
 collection for complex document information processing. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*,
 Aug 2006. doi: 10.1145/1148170.1148307. URL http://dx.doi.org/10.1145/1148170.
 1148307.
- [40] Minghao Li, Yiheng Xu, Lei Cui, Shaohan Huang, Furu Wei, Zhoujun Li, and Ming Zhou.
 Docbank: A benchmark dataset for document layout analysis. In *Proceedings of the 28th International Conference on Computational Linguistics*, Jan 2020. doi: 10.18653/v1/2020.
 coling-main.82. URL http://dx.doi.org/10.18653/v1/2020.coling-main.82.
- [41] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on
 document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021.
- [42] Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawa har. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1697–1706, 2022.
- [43] Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA:
 A benchmark for question answering about charts with visual and logical reasoning. In
 Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Findings of the Association*

- for Computational Linguistics: ACL 2022, pages 2263–2279, Dublin, Ireland, May 2022.
 Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.177. URL
 https://aclanthology.org/2022.findings-acl.177.
- [44] Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. Visualmrc: Machine reading comprehen sion on document images. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 volume 35, pages 13878–13888, 2021.
- [45] Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawahar. Icdar2019 competition on scanned receipt ocr and information extraction. In 2019 *International Conference on Document Analysis and Recognition (ICDAR)*, pages 1516–1520. IEEE, 2019.
- [46] Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and
 Hwalsuk Lee. Cord: a consolidated receipt dataset for post-ocr parsing. In *Workshop on Document Intelligence at NeurIPS 2019*, 2019.
- [47] Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form
 understanding in noisy scanned documents. In 2019 International Conference on Document
 Analysis and Recognition Workshops (ICDARW), volume 2, pages 1–6. IEEE, 2019.
- [48] Jianfeng Kuang, Wei Hua, Dingkang Liang, Mingkun Yang, Deqiang Jiang, Bo Ren, and Xiang
 Bai. Visual information extraction in the wild: practical dataset and end-to-end solution. In
 International Conference on Document Analysis and Recognition, pages 36–53. Springer, 2023.
- [49] Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang
 Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large
 multi-modal models. *arXiv preprint arXiv:2311.06607*, 2023.
- [50] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Zhong Muyan, Qinglong
 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
 for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023.
- [51] Liangcai Gao, Yilun Huang, Herve Dejean, Jean-Luc Meunier, Qinqin Yan, Yu Fang, Florian
 Kleber, and Eva Lang. Icdar 2019 competition on table detection and recognition (ctdar). In
 International Conference on Document Analysis and Recognition, 2019.
- [52] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image
 description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015.
- [53] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified
 text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [54] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words
 with subword units. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi:
 10.18653/v1/P16-1162. URL https://aclanthology.org/P16-1162.

532 Appendix

533 A Qualitative Examples

Qualitative examples of document-oriented VQA (upper row) and KIE (bottom row) are shown in Fig. 6. The results indicate that LayTextLLM is highly effective in utilizing spatial layout information to make more accurate predictions for these challenging examples. For example, in the upper right figure, many numeric texts in the receipt act as noise for the baseline method. In contrast, LayTextLLM integrates layout information to accurately predict the total price, as demonstrated by the other examples, underscoring the utility of LayTextLLM.



Figure 6: Qualitative comparison with the baseline method.

540 **B** Implementation Detail

All training and inference procedures are conducted on eight NVIDIA A100 GPUs.

Training LayTextLLM is initialized with Llama2-7B-Base model, the pre-training, SFT, and other model hyper-parameters can be seen in Tab. 6. Please note that all variants of LayTextLLM, including those utilized in ablation studies, are trained in accordance with the SFT settings. All baseline results are sourced from their respective original papers, with the exception of the Llama2-7B series and the Llama2-7B_{coor} series. These were re-implemented and can be referenced in [21, 23].

			mux lengen	Trecision	11 am params	Fix parallis
Llama2-7B-base Llama2-7B-base	256 256	128 256	2048 4096	bf16 bf16	648 M 7.4 B	6.7 B 0B
Learning rate	Weight decay	Scheduler	Adam betas	Adam epsilon	Warm up	Epoch
1.0e-04 2.0e-05	0.01 0.01	cosine cosine	[0.9, 0.999] [0.9, 0.999]	1.0e-08 1.0e-08	0.005 0.005	2 2
L L 1. 2	earning rate .0e-04 .0e-05	Iama2-7B-base 256 earning rate Weight decay .0e-04 0.01 .0e-05 0.01	Initial-7B-base 250 128 lama2-7B-base 256 256 earning rate Weight decay Scheduler .0e-04 0.01 cosine .0e-05 0.01 cosine	Iama2-7B-base 256 128 2048 lama2-7B-base 256 256 4096 earning rate Weight decay Scheduler Adam betas .0e-04 0.01 cosine [0.9, 0.999] .0e-05 0.01 cosine [0.9, 0.999]	Image: 7B-base 256 256 2048 b110 lama2-7B-base 256 256 4096 bf16 earning rate Weight decay Scheduler Adam betas Adam epsilon .0e-04 0.01 cosine [0.9, 0.999] 1.0e-08 .0e-05 0.01 cosine [0.9, 0.999] 1.0e-08	Iama2-7B-base 256 256 2048 b110 048 M lama2-7B-base 256 256 4096 bf16 7.4 B earning rate Weight decay Scheduler Adam betas Adam epsilon Warm up .0e-04 0.01 cosine [0.9, 0.999] 1.0e-08 0.005 .0e-05 0.01 cosine [0.9, 0.999] 1.0e-08 0.005

able	6:	Lav	TextL	LM	trainng	Hy	per-	parame	ters

Inference For the document-oriented VQA test set, we use the original question-answer pairs as 547 the prompt and ground truth, respectively. For Key Information Extraction (KIE) tasks, we reformat 548 the key-value pairs into a question-answer format, as described in [12, 19, 21]. Additionally, for the 549 FUNSD dataset, we focus our testing on the entity linking annotations as described in [21]. 550

To eliminate the impact of randomness on evaluation, no sampling methods are employed during 551 testing for any of the models. Instead, beam search with a beam size of 1 is used for generation across 552 all models. Additionally, the maximum number of new tokens is set to 512, while the maximum 553 number of input tokens is set to 4096. 554

Discussion of Input Sequence Length С 555

As mentioned in Section 4.6, it is intriguing that LayTextLLM has fewer input sequences than Do-556 cLLM, which is counterintuitive given that LayTextLLM interleaves bounding box tokens, typically 557 resulting in longer sequence lengths. We attribute this to the Byte Pair Encoding (BPE) tokenizers [54] 558 prevalently used in modern LLMs such as Llama2. 559

BPE operates by building a vocabulary of commonly occurring subwords (or token pieces) derived 560 from the training data. Initially, it tokenizes the text at the character level and then progressively 561 merges the most frequent adjacent pairs of characters or sequences. The objective is to strike a 562 balance between minimizing vocabulary size and maximizing encoding efficiency. 563

Thus, when tokenizing a single word like "International" on its own, the tokenizer might identify it 564 as a common sequence in the training data and encode it as a single token. This is especially likely if 565 "International" frequently appears as a standalone word in the training contexts. However, when the 566 word "International" is part of a larger sequence of words such as including in a long sequence of 567 OCR-derived texts like "...335 CPC, International, Inc ... ", the context changes. The tokenizer might 568 split "International" into sub-tokens like "Intern" and "ational" because, in various contexts within 569 the training data, these subwords might appear more frequently in different combinations or are more 570 useful for the model to understand variations in meaning or syntax. 571

When using LayTextLLM, we input word-level OCR results into the tokenizer, typically resulting in 572 the former situation, where words are encoded as single tokens. Conversely, with DocLLM, the entire 573 OCR output is processed as one large sequence, leading to the latter situation and a longer sequence 574 length than in LayTextLLM. This difference underscores the utility of LayTextLLM in achieving 575 both accuracy and inference efficiency due to its shorter sequence length. 576

D **Shuffle Ratio Exploration** 577

Tab. 7 presents the results of exploring training and testing shuffling ratios on the FUNSD dataset 578 using two different models: Llama2-7B-base and LayTextLLM. The table shows the performance of 579 these models at various shuffling ratios (100%, 50%, 20%, and 0%). 580

LayTextLLM consistently outperforms Llama2-7B-base across all levels of shuffling, which further 581 underscores the significance of interleaving spatial layouts with text. Particularly at the 100% shuffle 582 level, Llama2-7B-base demonstrates limited accuracy at only 20.3, while LayTextLLM maintains a 583 relatively higher performance. It is also interesting to note that Llama2-7B-base generally improves as 584 the shuffling percentage decreases, whereas LayTextLLM performs best when 20% of the examples 585 with OCR-derived text are shuffled. This observation suggests that LayTextLLM effectively utilizes 586 spatial layouts and is less dependent on the sequence of input tokens. Therefore, a certain proportion 587 of shuffled examples can serve as adversarial examples to enhance the model's robustness, addressing 588

situations such as errors in the text order from the OCR engine, which are caused by subtle differences

⁵⁹⁰ in horizontal or vertical coordinates.

	FUN	SD
Ratio	Llama2-7B-base	LayTextLLM
100%	20.3	44.7
50%	49.1	62.1
20%	50.2	65.4
0%	52.3	65.1

	_							
Table 7	7:	Shuffling	ratio o	explo	ration	in	FUNSD	dataset.

591 E Results of ChartQA

As shown in Fig. 7, the question-answer pairs in ChartQA [43] tend to involve the visual cues for 592 reasoning. However, with only text and layout information as input, the proposed LayTextLLM 593 inevitably have difficulties in reasoning visual-related information. Thus, on the ChartQA dataset, 594 LayTextLLM can hardly achieve better performance than previous methods that include visual inputs. 595 Although the visual information is not used in LayTextLLM, it can still exhibit better zero-shot ability 596 than UniDoc [10]. After incorporating the training set of ChartQA, the performance of LayTextLLM 597 can be boosted to 35.7%. Considering the importance of visual cues in ChartQA-like tasks, we will 598 try to involve the visual information into LayTextLLM in future work. 599



Figure 7: A	A failure case	of LayTextL	LM on	CharQA.
		-		

	ChartQA
OCR-free	
UniDoc [10]	10.9
DocPedia [9]	46.9*
Monkey [49]	54.0*
InternVL [50]	45.6*
InternLM-XComposer2 [15]	51.6*
TextMonkey [12]	58.2*
TextMonkey ₊ [12]	59.9 *
text + polys	
LayTextLLM _{zero} (Ours)	21.4
LayTextLLM _{vqa} (Ours)	29.8*
LayTextLLM _{all} (Ours)	35.7*

Table 8: Comparison with SOTA OCR-free MLLMs on ChartQA. * indicates the training set used.

600 NeurIPS Paper Checklist

 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? Answer: [Yes] Justification: We have detailed the contributions accurately in the abstract and introduction. Guidelines: The answer NA means that the abstract and introduction do not include the claims made in the paper. The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. 2 Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer; [Yes] Justification: A discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Cuidelines: The authors are encouraged to create a separate "Limitations, noiseless settings, model vell-specification, asymptotic approximations only holding locally). The authors should reflect on the whese assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results are to address problems of privacy and fairness. The authors should deflect on the scope of the claims made, e.g., if the approach was only tested on a few data	601	1.	Claims
 paper's contributions and scope? Answer: [Yes] Justification: We have detailed the contributions accurately in the abstract and introduction. Guidelines: The answer NA means that the abstract and introduction do not include the claims made in the paper. The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include appirational goals as motivation as long as it is clear that these goals are not attained by the paper. 2. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The authors are encouraged to create a separate "Limitations" section in their paper. The authors are encouraged to reate a separate "Limitations" section in their paper. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, who load approximations only holding locally). The authors and how they scale with dataset size. The authors should reflect on the factors that influence the performance of the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, who hould approximations on only holding locally). The authors and how they scale with dataset size.	602		Ouestion: Do the main claims made in the abstract and introduction accurately reflect the
 Answer: [Yes] Justification: We have detailed the contributions accurately in the abstract and introduction. Guidelines: The answer NA means that the abstract and introduction do not include the claims made in the paper. The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. 21 Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymption; e.g., independence assumptions, noiseless settings, model well-specification, asymption; which should be articulated. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should discuss the complete honesty about limitations of th	603		paper's contributions and scope?
 Justification: We have detailed the contributions accurately in the abstract and introduction. Guidelines: The answer NA means that the abstract and introduction do not include the claims made, in the paper. The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The anthors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions only holding locally). The authors should reflect on the tactors that influence the performance of the approach was only tested on implicit assumptions, which should be articulated.	604		Answer: [Yes]
 Guidelines: The answer NA means that the abstract and introduction do not include the claims made in the paper. The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The authors are encouraged to create a separate "Limitations" section in their paper. The authors are encouraged to create a separate "Limitations" section in their paper. The authors are encouraged to create a separate "Limitations" section in their paper. The authors should point out any strong assumptions and how robust the results are to violations of these assumptions might be violated in practice and what the implications would be. The authors should reflect on the factors that influence the performance of the approach was an only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, on algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle tec	605		Justification: We have detailed the contributions accurately in the abstract and introduction.
 The answer NA means that the abstract and introduction do not include the claims made in the paper. The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or which should be. The authors should reflect on the factors that influence the performance of the approach. The authors should discuss the computational efficiency of the rapproach to address problem so firvice assumptions of their approach was only tested on a few datasets or which aboud be articulated. The authors should fedect on the factors that influence the performance of the approach. The authors should discuss the computational efficiency of the rapproach to address problems of privacy	606		Guidelines:
 The answer NA means that the abstract and infroduction do not include the claims made in the paper. The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should diffect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the scope of nonelay because it fails to handle technical jargon. The authors should reflect on the scope of nonelay that the approach was only tested on a few datasets o	007		• The answer NA means that the abstract and introduction do not include the claims
 The abstract ad/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be precived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. 2. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer; [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions for online lectures because it fails to handle technical jargon. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a specch-to-text system might hot be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If appl	608		made in the paper.
 contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. 2. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions might be violated in practice and what the implications would reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how thy scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors should discuss tho complete honesty about limitations	609		• The abstract and/or introduction should clearly state the claims made, including the
 The claims made should match the optication and experimental results, and reflect how much the results can be expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. 2. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The apper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performace of the approach. For example, a facil recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairnees. While the authors might fear that compl	610 611		contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers
 The animate induction and the expected to generalize to other settings. It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. 2. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The anthors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiscless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performace of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations that are side approach to address problems of privacy an	612		 The claims made should match theoretical and experimental results, and reflect how
 It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations, socieless settings, model well-specification, asymptotic approximations only holding locally). The authors should point out any strong assumptions and how robust the results are to violations of these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The a	613		much the results can be expected to generalize to other settings.
are not attained by the paper. 2. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: ••••••••••••••••••••••••••••••••••••	614		• It is fine to include aspirational goals as motivation as long as it is clear that these goals
 2. Limitations Question: Does the paper discuss the limitations of the work performed by the authors? 	615		are not attained by the paper.
Guestion: Does the paper discuss the limitations of the work performed by the authors? Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: • The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. • The authors are encouraged to create a separate "Limitations" section in their paper. • The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. • The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. • The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. • If applicable, the authors should discuss possible limitations of their approach to address pro	616	2.	Limitations
 Answer: [Yes] Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparancy play an imporinating the used by reviewe	617		Question: Does the paper discuss the limitations of the work performed by the authors?
 Justification: As discussed in Section. 5, we have listed some limitations of our work and shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty	618		Answer: [Yes]
 shown corresponding failure cases in the supplementary material. Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and acomplete (and correct) encof?<	619		Justification: As discussed in Section 5, we have listed some limitations of our work and
 Guidelines: The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The authors are encouraged to create a separate "Limitations" section in their paper. The authors are encouraged to create a separate "Limitations" section in their paper. The authors of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should serie best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limit	620		shown corresponding failure cases in the supplementary material.
 The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. While the specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a computate (and accreact) wropfical result. 	621		Guidelines:
 the paper has limitations, but those are not discussed in the paper. The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a computer (a general). 	622		• The answer NA means that the paper has no limitation while the answer No means that
 The authors are encouraged to create a separate "Limitations" section in their paper. The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a commetite (and correct) wave? 	623		the paper has limitations, but those are not discussed in the paper.
 The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a computate (and correct) wave? 	624		• The authors are encouraged to create a separate "Limitations" section in their paper.
 Violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a completa (and correct) proof. 	625		• The paper should point out any strong assumptions and how robust the results are to
 and the wein-spectre atton, asymptote approximations only holding locarly). The authors should reflect on how these assumptions might be violated in practice and what the implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a cormet herefore. 	626		violations of these assumptions (e.g., independence assumptions, noiseless settings, model well specification, asymptotic approximations only holding locally). The authors
 implications would be. The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a completa (and correct) proof? 	628		should reflect on how these assumptions might be violated in practice and what the
 The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a completa (ord correct) proof2 	629		implications would be.
 only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	630		• The authors should reflect on the scope of the claims made, e.g., if the approach was
 depend on implicit assumptions, which should be articulated. The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an impor- tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	631		only tested on a few datasets or with a few runs. In general, empirical results often
 The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an impor- tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	632		depend on implicit assumptions, which should be articulated.
 For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an impor- tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	633		• The authors should reflect on the factors that influence the performance of the approach.
 ⁶³⁵ Is low of images are taken in low lighting. Of a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and corract) proof? 	634		For example, a facial recognition algorithm may perform poorly when image resolution is law or images are taken in law lighting. Or a speech to tart system might not be
 ⁶³⁶ technical jargon. The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an impor- tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	635		is low of images are taken in low lighting. Of a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle
 The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and each a complete (and correct) proof? 	635 637		technical jargon
 and how they scale with dataset size. If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and each a complete (and correct) proof? 	638		• The authors should discuss the computational efficiency of the proposed algorithms
 If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	639		and how they scale with dataset size.
 address problems of privacy and fairness. While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and proof? 	640		• If applicable, the authors should discuss possible limitations of their approach to
 While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	641		address problems of privacy and fairness.
 reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an impor- tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and 	642		• While the authors might fear that complete honesty about limitations might be used by
 limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an impor- tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and 	643		reviewers as grounds for rejection, a worse outcome might be that reviewers discover
 judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and error a complete (and correct) proof? 	644		limitations that aren't acknowledged in the paper. The authors should use their best
 tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations. 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	645		judgment and recognize that individual actions in favor of transparency play an impor-
 3. Theory Assumptions and Proofs Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof? 	646 647		tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize hopesty concerning limitations
Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?	047	2	Theory Assumptions and Proofs
649 Question: For each theoretical result, does the paper provide the full set of assumptions and	648	э.	Theory Assumptions and Froois
a complete (and confect) proof?	649 650		Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

651 Answer: [NA]

652	Justification: This work does not include theoretical results.
653	Guidelines:
654	• The answer NA means that the paper does not include theoretical results.
655	• All the theorems, formulas, and proofs in the paper should be numbered and cross-
656	referenced.
657	• All assumptions should be clearly stated or referenced in the statement of any theorems.
658	• The proofs can either appear in the main paper or the supplemental material, but if
659	they appear in the supplemental material, the authors are encouraged to provide a short
660	proof sketch to provide intuition.
661	• Inversely, any informal proof provided in the core of the paper should be complemented
662	by formal proofs provided in appendix or supplemental material.
663	• Theorems and Lemmas that the proof relies upon should be properly referenced.
664	4. Experimental Result Reproducibility
665	Question: Does the paper fully disclose all the information needed to reproduce the main ex-
666	perimental results of the paper to the extent that it affects the main claims and/or conclusions
667	of the paper (regardless of whether the code and data are provided or not)?
668	Answer: [Yes]
669	Justification: In Section 4.1, 4.2 and Appendix B, we have described the details of imple-
670	menting and training the proposed model to ensure the reproducibility of our work.
671	Guidelines:
672	• The answer NA means that the paper does not include experiments.
673	• If the paper includes experiments, a No answer to this question will not be perceived
674	well by the reviewers: Making the paper reproducible is important, regardless of
675	whether the code and data are provided or not.
676	• If the contribution is a dataset and/or model, the authors should describe the steps taken
677	to make their results reproducible or verifiable.
678	• Depending on the contribution, reproducibility can be accomplished in various ways.
679	For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and ampirical evaluation, it may
681	be necessary to either make it possible for others to replicate the model with the same
682	dataset, or provide access to the model. In general, releasing code and data is often
683	one good way to accomplish this, but reproducibility can also be provided via detailed
684	instructions for how to replicate the results, access to a hosted model (e.g., in the case
685	of a large language model), releasing of a model checkpoint, or other means that are
686	appropriate to the research performed.
687	• While NeurIPS does not require releasing code, the conference does require all submis-
688	sions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
689	(a). If the contribution is primerily a new algorithm, the paper should make it clear how
690 691	(a) If the control of its primarity a new argorithm, the paper should make it creat now to reproduce that algorithm
692	(b) If the contribution is primarily a new model architecture, the paper should describe
693	the architecture clearly and fully.
694	(c) If the contribution is a new model (e.g., a large language model), then there should
695	either be a way to access this model for reproducing the results or a way to reproduce
696	the model (e.g., with an open-source dataset or instructions for how to construct
697	the dataset).
698	(u) we recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility
099 700	In the case of closed-source models, it may be that access to the model is limited in
701	some way (e.g., to registered users), but it should be possible for other researchers
702	to have some path to reproducing or verifying the results.
703	5. Open access to data and code
704	Ouestion: Does the paper provide open access to the data and code with sufficient instruc-
705	tions to faithfully reproduce the main experimental results. as described in supplemental
706	material?

707		Answer: [No]
708		Justification: We would release our code after this paper is accepted.
709		Guidelines:
710		• The answer NA means that paper does not include experiments requiring code.
711		• Please see the NeurIPS code and data submission guidelines (https://nips.cc/
712		public/guides/CodeSubmissionPolicy) for more details.
713		• While we encourage the release of code and data, we understand that this might not be
714		possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
715		including code, unless this is central to the contribution (e.g., for a new open-source
716		benchmark).
717		• The instructions should contain the exact command and environment needed to run to
718		//nips.cc/public/guides/CodeSubmissionPolicy) for more details
720		• The authors should provide instructions on data access and preparation including how
721		to access the raw data, preprocessed data, intermediate data, and generated data, etc.
722		• The authors should provide scripts to reproduce all experimental results for the new
723		proposed method and baselines. If only a subset of experiments are reproducible, they
724		should state which ones are omitted from the script and why.
725 726		• At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
727		• Providing as much information as possible in supplemental material (appended to the
728		paper) is recommended, but including URLs to data and code is permitted.
729	6.	Experimental Setting/Details
730		Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
731		parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
732		results?
733		Answer: [Yes]
734 735		Justification: We have detailed the experimental setting and implementation details in Section 4.1, 4.2 and Appendix B of the supplementary material.
736		Guidelines:
737		• The answer NA means that the paper does not include experiments.
738		• The experimental setting should be presented in the core of the paper to a level of detail
739		that is necessary to appreciate the results and make sense of them.
740 741		• The full details can be provided either with the code, in appendix, or as supplemental material
741	7	Experiment Statistical Significance
749		Ouestion: Does the paper report error bars suitably and correctly defined or other appropriate
743		information about the statistical significance of the experiments?
745		Answer: [No]
746		Justification: Our deep learning model is designed for a complex task (requiring huge
747		computing resources) where traditional error bars are less informative due to the high
748		variability in model training and initialization. We ensured the robustness of our model
749 750		models demonstrated improvements in key performance areas underscoring the practical
751		effectiveness of our approach. We acknowledge the limitation of not using traditional
752		statistical tests and suggest that future work could explore statistical significance in more
753		controlled settings.
754		Guidelines:
755		• The answer NA means that the paper does not include experiments.
756		• The authors should answer "Yes" if the results are accompanied by error bars, confi-
757		dence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper
/ 30		une main cianns of the paper.

759 760	• The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions)
761	• The method for calculating the error bars should be explained (closed form formula,
763	call to a library function, bootstrap, etc.)
764	• The assumptions made should be given (e.g., Normally distributed errors).
765	• It should be clear whether the error bar is the standard deviation or the standard error
766	of the mean.
767	• It is OK to report 1-sigma error bars, but one should state it. The authors should
768	preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
709	• For asymmetric distributions, the authors should be careful not to show in tables or
770	figures symmetric error bars that would vield results that are out of range (e.g. negative
772	error rates).
773 774	• If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
775	8. Experiments Compute Resources
776	Question: For each experiment, does the paper provide sufficient information on the com-
777	puter resources (type of compute workers, memory, time of execution) needed to reproduce
778	the experiments?
779	Answer: [Yes]
780	Justification: We have reported the needed computer resources in Section B of the supple-
781	mentary material.
782	Guidelines:
783	 The answer NA means that the paper does not include experiments.
784 785	• The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
786	• The paper should provide the amount of compute required for each of the individual
787	experimental runs as well as estimate the total compute.
788	• The paper should disclose whether the full research project required more compute
789 790	than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
791	9. Code Of Ethics
792 793	Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
794	Answer: [Yes]
795 796	Justification: The research in this paper conforms with the NeurIPS Code of Ethics in every respect.
797	Guidelines:
798	• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
799	• If the authors answer No, they should explain the special circumstances that require a
800	deviation from the Code of Ethics.
801	• The authors should make sure to preserve anonymity (e.g., if there is a special consid-
802	eration due to laws or regulations in their jurisdiction).
803	10. Broader Impacts
804	Question: Does the paper discuss both potential positive societal impacts and negative
805	societal impacts of the work performed?
806	Answer: [NA]
807	Justification: There is no societal impact of the work performed.
808	Guidelines:
809	• The answer NA means that there is no societal impact of the work performed.

810 811		• If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
812		• Examples of negative societal impacts include potential malicious or unintended uses
813		(e.g., disinformation, generating fake profiles, surveillance), fairness considerations
815		groups), privacy considerations, and security considerations.
816		• The conference expects that many papers will be foundational research and not tied
817		to particular applications, let alone deployments. However, if there is a direct path to
818		any negative applications, the authors should point it out. For example, it is legitimate
819		to point out that an improvement in the quality of generative models could be used to
820		generate deepfakes for disinformation. On the other hand, it is not needed to point out
821		that a generic algorithm for optimizing neural networks could enable people to train models that generate Deenfakes faster
922		• The authors should consider possible harms that could arise when the technology is
824		being used as intended and functioning correctly, harms that could arise when the
825		technology is being used as intended but gives incorrect results, and harms following
826		from (intentional or unintentional) misuse of the technology.
827		• If there are negative societal impacts, the authors could also discuss possible mitigation
828		strategies (e.g., gated release of models, providing defenses in addition to attacks,
829		mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML)
830	1.1	ceduack over time, improving the efficiency and accessionity of ML).
831	11.	Sateguards Ouestion: Does the paper describe safeguards that have been put in place for responsible
833		release of data or models that have a high risk for misuse (e.g., pretrained language models,
834		image generators, or scraped datasets)?
835		Answer: [NA]
836 837		Justification: All the datasets used in this paper are publicly available and they contain no unsafe images.
838		Guidelines:
839		• The answer NA means that the paper poses no such risks.
840		• Released models that have a high risk for misuse or dual-use should be released with
841		necessary safeguards to allow for controlled use of the model, for example by requiring
842		that users adhere to usage guidelines or restrictions to access the model or implementing
843		• Datasets that have been screpted from the Internet could nose sofety risks. The authors
844 845		should describe how they avoided releasing unsafe images.
846		• We recognize that providing effective safeguards is challenging, and many papers do
847		not require this, but we encourage authors to take this into account and make a best
848		faith effort.
849	12.	Licenses for existing assets
850		Question: Are the creators or original owners of assets (e.g., code, data, models), used in
851		the paper, properly credited and are the license and terms of use explicitly mentioned and
852		property respected?
853		Answer: [Yes]
854 855		Justification: For the used datasets and pre-trained models, we have cited their corresponding works.
856		Guidelines:
857		• The answer NA means that the paper does not use existing assets.
858		• The authors should cite the original paper that produced the code package or dataset.
859		• The authors should state which version of the asset is used and, if possible, include a
860		URL.
861		• The name of the license (e.g., CC-BY 4.0) should be included for each asset.

862 863		• For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
864 865 866 867		• If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
868 869		• For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
870 871		• If this information is not available online, the authors are encouraged to reach out to the asset's creators.
872	13.	New Assets
873 874		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
875		Answer: [NA]
876		Justification: This paper does not release new assets
877		Guidelines:
878		• The answer NA means that the paper does not release new assets.
879 880		• Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
882		• The paper should discuss whether and how consent was obtained from people whose
883		asset is used.
884 885		• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
886	14.	Crowdsourcing and Research with Human Subjects
887 888 889		Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?
890		Answer: [NA]
891		Justification: This paper does not involve crowdsourcing nor research with human subjects.
892		Guidelines:
893 894		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
895 896 897		• Including this information in the supplemental material is fine, but if the main contribu- tion of the paper involves human subjects, then as much detail as possible should be included in the main paper.
898 899 900		• According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
901 902	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects
903 904 905 906		Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?
907		Answer: [NA]
908		Justification: This paper does not involve crowdsourcing nor research with human subjects.
909		Guidelines:
910 911		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

912	• Depending on the country in which research is conducted, IRB approval (or equivalent)
913	may be required for any human subjects research. If you obtained IRB approval, you
914	should clearly state this in the paper.
915	• We recognize that the procedures for this may vary significantly between institutions
916	and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
917	guidelines for their institution.
918	• For initial submissions, do not include any information that would break anonymity (if
919	applicable), such as the institution conducting the review.