

// ALIGNVLM: BRIDGING VISION AND LANGUAGE LATENT SPACES FOR MULTIMODAL UNDERSTANDING

Ahmed Masry^{1,2} Juan A. Rodriguez^{1,3,4} Tianyu Zhang^{1,3,5} Suyuchen Wang^{1,3,5}
 Chao Wang¹ Aarash Feizi^{1,3,6} Akshay Kalkunte Suresh¹ Abhay Puri¹
 Xiangru Jian^{1,7} Pierre-André Noël¹ Sathwik Tejaswi Madhusudhan¹
 Marco Pedersoli^{1,4} Bang Liu^{1,5,8} Nicolas Chapados¹ Yoshua Bengio^{3,5,8}
 Enamul Hoque² Christopher Pal^{1,3,8,9} Issam H. Laradji^{1,10} David Vazquez¹
 Perouz Taslakian¹ Spandana Gella¹ Sai Rajeswar^{1,3}

¹ServiceNow ²York University ³Mila ⁴École de Technologie Supérieure

⁵Université de Montréal ⁶McGill University ⁷University of Waterloo

⁸CIFAR AI Chair ⁹Polytechnique Montréal ¹⁰University of British Columbia

ABSTRACT

Aligning visual features with language embeddings is a key challenge in vision-language models (VLMs). The performance of such models hinges on having a good connector that maps visual features generated by a vision encoder to a shared embedding space with the LLM while preserving semantic similarity. Existing connectors, such as multilayer perceptrons (MLPs), often produce out-of-distribution or noisy inputs, leading to misalignment between the modalities. In this work, we propose a novel vision-text alignment method, ALIGNVLM, that maps visual features to a weighted average of LLM text embeddings. Our approach leverages the linguistic priors encoded by the LLM to ensure that visual features are mapped to regions of the space that the LLM can effectively interpret. ALIGNVLM is particularly effective for document understanding tasks, where scanned document images must be accurately mapped to their textual content. Our extensive experiments show that ALIGNVLM achieves state-of-the-art performance compared to prior alignment methods. We provide further analysis demonstrating improved vision-text feature alignment and robustness to noise.

1 INTRODUCTION

Vision-Language Models (VLMs) have gained significant traction in recent years as a powerful framework for multimodal document understanding tasks that involve interpreting both the visual and textual contents of scanned documents (Kim et al., 2022; Lee et al., 2023; Liu et al., 2023a; 2024; Hu et al., 2024; Wang et al., 2023a; Rodriguez et al., 2024b). Such tasks are common in real-world commercial applications, including invoice parsing (Park et al., 2019), form reading (Jaume et al., 2019), and document question answering (Mathew et al., 2021b). VLM architectures typically consist of three components: (i) a vision encoder to process raw images, (ii) a Large Language Model (LLM) pre-trained on text, and (iii) a connector module that maps the visual features from the vision encoder into the LLM’s semantic space.

A central challenge in this pipeline is to effectively map the continuous feature embeddings of the vision encoder into the latent space of the LLM while preserving the semantic properties of visual concepts. Existing approaches can be broadly categorized into *deep fusion* and *shallow fusion* methods. *Deep fusion* methods, such as NVLM (Dai et al., 2024), Flamingo (Alayrac et al., 2022), CogVLM (Wang et al., 2023b), and LLama 3.2-Vision (Grattafiori et al., 2024), integrate visual and textual features by introducing additional cross-attention and feed-forward layers at each layer of the LLM. While effective at enhancing cross-modal interaction, these methods substantially increase the parameter count of the VLM compared to the base LLM, resulting in high computational overhead and reduced efficiency.

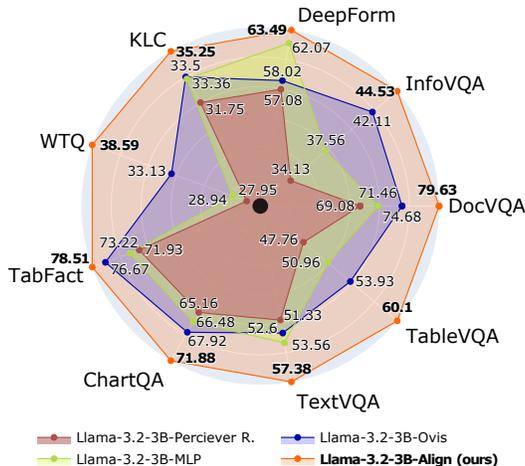


Figure 1: **Performance of Different VLM Connectors.** The proposed **ALIGN** connector outperforms other methods across benchmarks using the same training configuration. Radial distance is proportion of maximal score, truncated at 0.7 (black dot).

In contrast, *shallow fusion* methods project visual features from the vision encoder into the LLM input embedding space using either multilayer perceptrons (MLPs) (Liu et al., 2023b; 2024) or attention-based mechanisms such as the Perceiver Resampler (Li et al., 2023; Laurençon et al., 2024; Alayrac et al., 2022), before concatenating them with the textual prompt’s input embeddings. This approach is more parameter-efficient and computationally lighter than *deep fusion* methods, but it lacks a mechanism to ensure the projected embeddings remain within the region spanned by the LLM’s text embeddings – i.e. regions the LLM was pretrained to understand. As a result, unconstrained visual features can produce out-of-distribution (OOD) and noisy inputs, leading to misalignment between modalities and often degrading overall performance. Recent methods like Ovis (Lu et al., 2024) attempt to alleviate these issues by introducing separate visual embeddings indexed from the vision encoder outputs and combined together to construct the visual inputs to the LLM. However, this approach significantly increases parameter count due to the massive embedding matrix and requires extensive training to learn a new embedding space without guaranteeing alignment with the LLM’s input latent space.

To address these limitations, this paper introduces **ALIGNVLM**, a novel framework that sidesteps direct projection of visual features into the LLM embedding space. Instead, our proposed connector, **ALIGN**, maps visual features into probability distributions over the LLM’s *existing* pretrained vocabulary embeddings, which are then combined into a weighted representation of the text embeddings. By constraining each visual feature as a convex combination of the LLM text embeddings, our approach leverages the linguistic priors already encoded in the LLM’s text space. This ensures that the resulting visual features lie within the convex hull of the LLM’s embedding space, reducing the risk of noisy or out-of-distribution inputs and improving alignment between modalities. Our experimental results show that this approach improves performance on various document understanding tasks, outperforming prior connector methods by effectively fusing visual and linguistic content. We summarize our main contributions as follows:

- We propose a novel connector, **ALIGN**, to bridge the representation gap between vision and text modalities.
- We introduce a family of Vision-Language Models, **ALIGNVLM**, that achieves state-of-the-art performance on multimodal document understanding tasks by leveraging **ALIGN**.
- We conduct extensive experiments demonstrating the robustness and effectiveness of **ALIGN** across different model sizes ranging from 1B to 8B parameters.

2 RELATED WORK

2.1 VISION-LANGUAGE MODELS

Over the past few years, Vision-Language Models (VLMs) have achieved remarkable progress, largely due to advances in Large Language Models (LLMs). Initially demonstrating breakthroughs in text understanding and generation (Brown et al., 2020; Raffel et al., 2023; Achiam et al., 2023; Grattafiori et al., 2024; Qwen et al., 2025; Team, 2024), LLMs are now increasingly used to effectively interpret visual inputs (Liu et al., 2023b; Li et al., 2024; Wang et al., 2024; Chen et al., 2024b; Dai et al., 2024; Drouin et al., 2024; Rodriguez et al., 2022). This progress has enabled real-world applications across diverse domains, particularly in multimodal document understanding for tasks like form reading (Svetlichnaya, 2020), document question answering (Mathew et al., 2021b), and chart question answering (Masry et al., 2022). VLMs commonly adopt a three-component architecture: a pretrained vision encoder (Zhai et al., 2023; Radford et al., 2021), a LLM, and a connector module. A key challenge for VLMs is effectively aligning visual features with the LLM’s semantic space to enable accurate and meaningful multimodal interpretation.

2.2 VISION-LANGUAGE ALIGNMENT FOR MULTIMODAL MODELS

Existing vision-language alignment approaches can be classified into *deep fusion* and *shallow fusion*. Deep fusion methods integrate visual and textual features by modifying the LLM’s architecture, adding cross-attention and feed-forward layers. For example, Flamingo (Alayrac et al., 2022) employs the Perceiver Resampler, which uses fixed latent embeddings to attend to vision features and fuses them into the LLM via gated cross-attention layers. Similarly, NVLM (Dai et al., 2024) adopts cross-gated attention while replacing the Perceiver Resampler with a simpler MLP. CogVLM (Wang et al., 2023b) extends this approach by incorporating new feed-forward (FFN) and QKV layers for the vision modality within every layer of the LLM. While these methods improve cross-modal alignment, they significantly increase parameter counts and computational overhead, making them less efficient.

On the other hand, shallow fusion methods are more computationally efficient, mapping visual features into the LLM’s embedding space without altering its architecture. These methods can be categorized into three main types: (1) *MLP-based mapping*, such as LLaVA (Liu et al., 2023b) and PaliGemma (Beyer et al., 2024), which use multilayer perceptrons (MLP) to project visual features but often produce misaligned or noisy features due to a lack of constraints (Rodriguez et al., 2024b); (2) *cross-attention mechanisms*, BLIP-2 (Li et al., 2023) uses Q-Former, which utilizes a fixed set of latent embeddings to cross-attend to visual features, but that may still produce noisy or OOD visual features; and (3) *visual embeddings*, such as those introduced by Ovis (Lu et al., 2024), which use embeddings indexed by the vision encoder’s outputs to produce the visual inputs. While this regularizes feature mapping, it adds substantial parameter overhead and creates a new vision embedding space, risking misalignment with the LLM’s text embedding space. Encoder-free VLMs, like Fuyu-8B¹ and EVE (Diao et al., 2024), eliminate dedicated vision encoders but show degraded performance (Beyer et al., 2024).

In contrast, ALIGNVLM maps visual features from the vision encoder into probability distributions over the LLM’s text embeddings, using them to compute a convex combination. By leveraging the linguistic priors encoded in the LLM’s vocabulary, ALIGNVLM ensures that visual features remain within the convex hull of the text embeddings, mitigating noisy or out-of-distribution inputs and enhancing alignment, particularly for tasks that require joint modalities representation like multimodal document understanding.

3 METHODOLOGY

3.1 MODEL ARCHITECTURE

The overall model architecture, shown in Figure 2, consists of three main components:

(1) Vision Encoder. To handle high-resolution images of different aspect ratios, we divide each input image into multiple tiles according to one of the predefined aspect ratios (e.g., 1:1, 1:2, . . . , 9:1)

¹<https://www.adept.ai/blog/fuyu-8b>

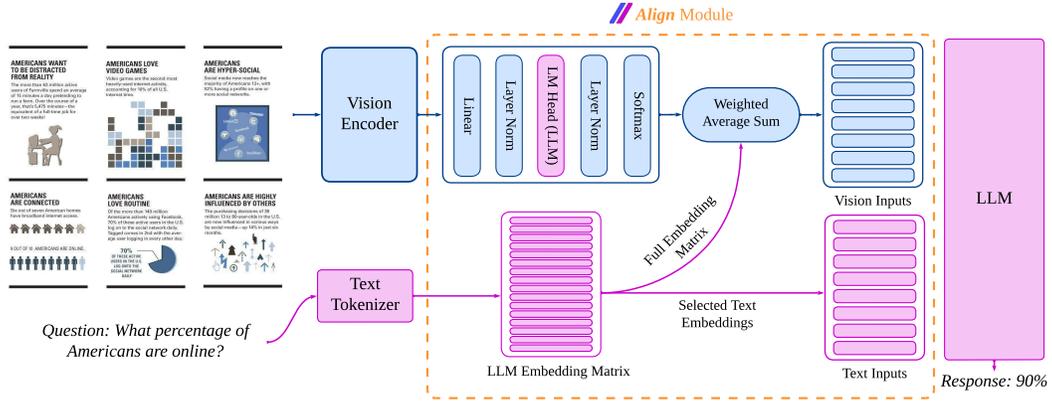


Figure 2: **ALIGNVLM Model Architecture.** The vision encoder extracts image features, which are processed to produce probabilities over the LLM embeddings. A weighted average combines these probabilities with embeddings to generate vision input vectors. Text inputs are tokenized, and the corresponding embeddings are selected from the embedding matrix, which is then used as input to the LLM. We display the vision layers in blue, and the text layers in purple.

chosen via a coverage ratio (Lu et al., 2024; Chen et al., 2024a). Due to limited computational resources, we set the maximum number of tiles to 9. Each tile is further partitioned into 14×14 patches, projected into vectors, and processed by a SigLip-400M vision encoder (Zhai et al., 2023) to extract contextual visual features.

Each tile $t \in \{1, \dots, T\}$ is divided into N_t patches

$$\mathbf{P}_t = \{\mathbf{p}_{t,1}, \dots, \mathbf{p}_{t,N_t}\},$$

where $\mathbf{p}_{t,i}$ is the i -th patch of tile t . The vision encoder maps these patches to a set of visual feature vectors

$$\begin{aligned} \mathbf{F}_t &= \text{VisionEncoder}(\mathbf{P}_t) \\ \mathbf{F}_t &= \{\mathbf{f}_{t,1}, \dots, \mathbf{f}_{t,N_t}\}, \quad \mathbf{f}_{t,i} \in \mathbb{R}^d. \end{aligned}$$

Finally, we concatenate the feature sets across all tiles into a single output

$$\mathbf{F} = \text{concat}(\mathbf{F}_1, \mathbf{F}_2, \dots, \mathbf{F}_T).$$

(2) ALIGN Module. This module aligns the visual features with the LLM. A linear layer $\mathbf{W}_1 \in \mathbb{R}^{D \times d}$ first projects the visual features $\mathbf{F} \in \mathbb{R}^{T \cdot N_t \times d}$ to the LLM’s token embedding space: one \mathbb{R}^D vector per token. A second linear layer $\mathbf{W}_2 \in \mathbb{R}^{V \times D}$ (initialized from the LLM’s language-model head) followed by a softmax, produces a probability simplex $\mathbf{P}_{\text{vocab}}$ over the LLM’s vocabulary (V tokens)

$$\mathbf{P}_{\text{vocab}} = \text{softmax}(\text{LayerNorm}(\mathbf{W}_2 \text{LayerNorm}(\mathbf{W}_1 \mathbf{F}))) \quad (1)$$

We then use the LLM text embeddings $\mathbf{E}_{\text{text}} \in \mathbb{R}^{V \times D}$ to compute a weighted sum

$$\mathbf{F}'_{\text{align}} = \mathbf{P}_{\text{vocab}}^T \mathbf{E}_{\text{text}}. \quad (2)$$

Finally, we concatenate $\mathbf{F}'_{\text{align}}$ with the tokenized text embeddings to form the LLM input

$$\mathbf{H}_{\text{input}} = \text{concat}(\mathbf{F}'_{\text{align}}, \mathbf{E}_{\text{text}}(\mathbf{x})),$$

where $\mathbf{E}_{\text{text}}(\mathbf{x})$ is obtained by tokenizing the input text $\mathbf{x} = (x_1, \dots, x_M)$ and selecting the corresponding embeddings from \mathbf{E}_{text} such that

$$\mathbf{E}_{\text{text}}(\mathbf{x}) = [\mathbf{E}_{\text{text}}(x_1), \dots, \mathbf{E}_{\text{text}}(x_M)]. \quad (3)$$

(3) Large Language Model. We feed the concatenated vision and text vectors, $\mathbf{H}_{\text{input}}$, into the LLM, which then generates output text auto-regressively. To demonstrate the effectiveness of our alignment technique, we experiment with the Llama 3.1 model family (Grattafiori et al., 2024). These models offer state-of-the-art performance and permissive licenses, making them suitable for commercial applications. In particular, we utilize Llama 3.2-1B, Llama 3.2-3B, and Llama 3.1-8B.

3.2 MOTIVATION AND RELATION WITH EXISTING METHODS

By construction, each \mathbb{R}^D representation in $\mathbf{F}'_{\text{align}}$ is constrained to the convex hull of the points \mathbb{E}_{text} , thus concentrating the visual features in the part of latent space that the LLM can effectively interpret. Moreover, we argue that our initialization of \mathbf{W}_2 to the language model head is an inductive bias toward *recycling* some of the semantics of these text tokens into visual tokens. This contrasts with past methods that have been proposed to adapt the vision encoder outputs $\mathbf{F} \in \mathbb{R}^{T \cdot N_t \times d}$ to an $\mathbf{F}' \in \mathbb{R}^{T \cdot N_t \times D}$ to be fed to the LLM. Here, we consider two examples in more detail, highlighting these contrasts.

(1) *MLP Connector* Liu et al. (2023b) applies a linear projection with parameters $\mathbf{W}_{\text{MLP}} \in \mathbb{R}^{D \times d}$ and $\mathbf{b}_{\text{MLP}} \in \mathbb{R}^D$, followed by an activation function σ (e.g., ReLU)

$$\mathbf{F}'_{\text{MLP}} = \sigma(\mathbf{W}_{\text{MLP}}\mathbf{F} + \mathbf{b}_{\text{MLP}}).$$

These parameters are all learned from scratch, with no particular bias aligning them to text embeddings.

(2) *Visual Embedding Table* Lu et al. (2024) introduces an entire new set of visual embeddings $\mathbf{E}_{\text{VET}} \in \mathbb{R}^{K \times D}$ which, together with the weights $\mathbf{W}_{\text{VET}} \in \mathbb{R}^{K \times d}$, specifies

$$\mathbf{F}'_{\text{VET}} = \text{softmax}(\mathbf{W}_{\text{VET}}\mathbf{F})^\top \mathbf{E}_{\text{VET}}.$$

When $D < d$, our $\mathbf{W}_2\mathbf{W}_1$ amounts to a low-rank version of \mathbf{W}_{VET} . There is thus much more to learn to obtain \mathbf{F}'_{VET} , and there is again no explicit pressure to align it with the text embeddings.

3.3 TRAINING DATASETS & STAGES

We train our model in three stages:

Stage 1. This stage focuses on training the ALIGN Module to map visual features to the LLM’s text embeddings effectively. We use the CC-12M dataset Changpinyo et al. (2021), a large-scale web dataset commonly used for VLM pretraining Liu et al. (2023b), which contains 12M image-text pairs. However, due to broken or unavailable links, we retrieved 8.1M pairs. This dataset facilitates the alignment of visual features with the text embedding space of the LLM. During this stage, we train the full model, as this approach improves performance and stabilizes the training of the ALIGN Module.

Stage 2. The goal is to enhance the model’s document understanding capabilities, such as OCR, document structure comprehension, in-depth reasoning, and instruction-following. We leverage the BigDocs-7.5M dataset Rodriguez et al. (2024a), a curated collection of license-permissive datasets designed for multimodal document understanding. This dataset aligns with the Accountability, Responsibility, and Transparency (ART) principles Bommasani et al. (2023); Vogus & Llansóe (2021), ensuring compliance for commercial applications. As in Stage 1, we train the full model during this stage.

Stage 3. To enhance the model’s instruction-tuning capabilities, particularly for downstream tasks like question answering, we further train it on the DocDownstream Rodriguez et al. (2024a); Hu et al. (2024) instruction tuning dataset. In this stage, the vision encoder is frozen, focusing training exclusively on the LLM and ALIGN module.

4 EXPERIMENTAL SETUP

Setup. We conduct all experiments using 8 nodes of H100 GPUs, totaling 64 GPUs. For model training, we leverage the MS-Swift framework (Zhao et al., 2024) for its flexibility. Additionally, we

Table 1: **Main Results on General Document Benchmarks.** We compare ALIGNVLM (ours) with state-of-the-art (SOTA) open and closed-source instructed models, and with base models that we trained using the process described in Section 3.3. ALIGNVLM models outperform all Base VLM models trained in the same data regime. Our models also perform competitively across document benchmarks even compared with SOTA models, in which the data regime is more targeted and optimized. Color coding for comparison: closed-source models, open-source models below 7B parameters, open-source models between 7-12B parameters.

Model	DocVQA VAL	InfoVQA VAL	DeepForm TEST	KLC TEST	WTQ TEST	TabFacet TEST	ChartQA TEST	TextVQA VAL	TableVQA TEST	Avg. Score
Closed-Source VLMs (Opaque Training Data)										
Claude-3.5 Sonnet	88.48	59.05	31.41	24.82	47.13	53.48	51.84	71.42	81.27	56.54
GeminiPro-1.5	91.23	73.94	32.16	24.07	50.29	71.22	34.68	68.16	80.43	58.46
GPT-4o 20240806	92.80	66.37	38.39	29.92	46.63	81.10	85.70	70.46	72.87	64.91
Open-Source Instruct VLMs (Semi-Opaque Training Data)										
Janus-1.3B (Wu et al., 2024a)	30.15	17.09	0.62	15.06	9.30	51.34	57.20	51.97	18.67	27.93
Qwen2-VL-2B (Wang et al., 2024)	89.16	64.11	32.38	25.18	38.20	57.21	73.40	79.90	43.07	55.84
InternVL-2.5-2B (Chen et al., 2024b)	87.70	61.85	13.14	16.58	36.33	57.26	74.96	76.85	42.20	51.87
DeepSeek-VL2-Tiny-3.4B (Wu et al., 2024b)	88.57	63.88	25.11	19.04	35.07	52.15	80.92	80.48	56.30	55.72
Phi3.5-Vision-4B (Abdin et al., 2024)	86.00	56.20	10.47	7.49	17.18	30.43	82.16	73.12	70.70	48.19
Qwen2-VL-7B (Wang et al., 2024)	93.83	76.12	34.55	23.37	52.52	74.68	83.16	84.48	53.97	64.08
LLaVA-NeXT-7B (Xu et al., 2024)	63.51	30.90	1.30	5.35	20.06	52.83	52.12	65.10	32.87	36.00
DocOwl1.5-8B (Hu et al., 2024)	80.73	49.94	68.84	37.99	38.87	79.67	68.56	68.91	52.60	60.68
InternVL-2.5-8B (Chen et al., 2024b)	91.98	75.36	34.55	22.31	50.33	74.75	82.84	79.00	52.10	62.58
Ovis-1.6-Gemma2-9B (Lu et al., 2024)	88.84	73.97	45.16	23.91	50.72	76.66	81.40	77.73	48.33	62.96
Llama3.2-11B (Grattafiori et al., 2024)	82.71	36.62	1.78	3.47	23.03	58.33	23.80	54.28	22.40	34.04
Pixtral-12B (Agrawal et al., 2024)	87.67	49.45	27.37	24.07	45.18	73.53	71.80	76.09	67.13	58.03
Document Understanding Instructed Models (Instruction Tuned on BigDocs-7.5M + DocDownStream (Rodriguez et al., 2024a; Hu et al., 2024))										
Qwen2-VL-2B (base+) (Qwen et al., 2025)	57.23	31.88	49.31	34.39	31.61	64.75	68.60	61.01	47.53	49.59
ALIGNVLM-Llama-3.2-1B (ours)	72.42	38.16	60.47	33.71	28.66	71.31	65.44	48.81	50.29	52.14
ALIGNVLM-Llama-3.2-3B (ours)	79.63	44.53	63.49	35.25	38.59	78.51	71.88	57.38	60.10	58.81
DocOwl1.5-8B (base+) (Hu et al., 2024)	78.70	47.62	64.39	36.93	35.69	72.65	65.80	67.30	49.03	57.56
Llama3.2-11B (base+) (Grattafiori et al., 2024)	78.99	44.27	67.05	37.22	40.18	78.04	71.40	68.46	56.73	60.26
ALIGNVLM-Llama-3.1-8B (ours)	81.18	53.75	63.25	35.50	45.31	83.04	75.00	64.60	64.33	62.88

utilize the DeepSpeed framework (Aminabadi et al., 2022), specifically the ZeRO-3 configuration, to optimize efficient parallel training across multiple nodes. Detailed hyperparameters are outlined in Appendix A.1.

Baselines. Our work focuses on architectural innovations, so we ensure that all baselines are trained on the same datasets. To enable fair comparisons, we evaluate our models against a set of **Base VLMs** fine-tuned on the same instruction-tuning tasks (Stages 2 and 3) as our models, using the BigDocs-7.5M and BigDocs-DocDownstream datasets. This approach ensures consistent training data, avoiding biases introduced by the **Instruct** versions of VLMs, which are often trained on undisclosed instruction-tuning datasets. Due to the scarcity of recently released publicly available Base VLMs, we primarily compare our model against the following Base VLMs of varying sizes: Qwen2-VL-2B (Wang et al., 2024), DocOwl1.5-8B (Hu et al., 2024), and Llama 3.2-11B (Grattafiori et al., 2024).

For additional context, we also include results from the Instruct versions of recent VLMs of different sizes: Phi3.5-Vision-4B (Abdin et al., 2024), Qwen2-VL-2B and 7B (Wang et al., 2024), LLaVA-NeXT-7B (Liu et al., 2024), InternVL2.5-2B and 8B (Chen et al., 2024b), Janus-1.3B (Wu et al., 2024a), DeepSeek-VL2-Tiny (Wu et al., 2024b), Ovis1.6-Gemma-9B (Lu et al., 2024), Llama3.2-11B (Grattafiori et al., 2024), DocOwl1.5-8B (Hu et al., 2024), and Pixtral-12B (Agrawal et al., 2024).

Evaluation Benchmarks. We evaluate our models on a diverse range of document understanding benchmarks that assess the model’s capabilities in OCR, chart reasoning, table processing, or form comprehension. In particular, we employ the VLMEvalKit (Duan et al., 2024) framework and report the results on the following popular benchmarks: DocVQA (Mathew et al., 2021b), InfoVQA (Mathew et al., 2021a), DeepForm (Svetlichnaya, 2020), KLC (Stanisławek et al.,

Table 2: **Impact of Connector Designs on VLM Performance:** We present the results of experiments evaluating different connector designs for conditioning LLMs on visual features. Our proposed **ALIGN** connector is compared against a basic Multi-Layer Perceptron (**MLP**), the **Perceiver Resampler**, and **Ovis**. The results demonstrate that ALIGN consistently outperforms these alternatives across all benchmarks.

Model	DocVQA VAL	InfoVQA VAL	DeepForm TEST	KLIC TEST	WTQ TEST	TabFact TEST	ChartQA TEST	TextVQA VAL	TableVQA TEST	Avg. Score
Llama-3.2-3B-MLP	71.46	37.56	62.07	33.36	28.94	73.22	66.48	53.56	50.96	53.06
Llama-3.2-3B-Perceiver R.	69.08	34.13	57.08	31.75	27.95	71.93	65.16	51.33	47.76	50.68
Llama-3.2-3B-Ovis	74.68	42.11	58.02	33.50	33.13	76.67	67.92	52.60	53.93	54.72
Llama-3.2-3B-ALIGN (ours)	79.63	44.53	63.49	35.25	38.59	78.51	71.88	57.38	60.10	58.81

2021), WTQ (Pasupat & Liang, 2015), TabFact (Chen et al., 2020), ChartQA (Masry et al., 2022), TextVQA (Singh et al., 2019), and TableVQA (Kim et al., 2024).

5 RESULTS

5.1 MAIN RESULTS

Table 1 presents the performance of ALIGNVLM compared to state-of-the-art (SOTA) open- and closed-source instructed models, as well as baseline Base VLMs fine-tuned in the same instruction-tuning setup. The results demonstrate that ALIGNVLM consistently outperforms all Base VLMs within the same size category and achieves competitive performance against SOTA Instruct VLMs despite being trained on a more limited data regime. Below, we provide a detailed analysis.

ALIGNVLM vs. Base VLMs. Our ALIGNVLM models, based on Llama 3.2-1B and Llama 3.2-3B, significantly outperform the corresponding Base VLM, Qwen2-VL-2B, by up to 9.22%. Notably, ALIGNVLM-Llama-3.2-3B surpasses DocOwl1.5-8B, which has 4B more parameters, demonstrating the effectiveness of ALIGN in enhancing multimodal capabilities compared to traditional *shallow fusion* methods (e.g., MLPs). Furthermore, our 8B model achieves a 2.62% improvement over Llama3.2-11B despite sharing the same Base LLM, Llama3.1-8B. Since all models in this comparison were trained on the same instruction-tuning setup, this experiment provides a controlled evaluation, isolating the impact of architectural differences rather than dataset biases. Consequently, these results suggest that ALIGNVLM outperforms VLMs with shallow fusion techniques and surpasses parameter-heavy *deep fusion* VLMs, such as Llama3.2-11B, while maintaining a more efficient architecture.

ALIGNVLM vs. Instruct VLMs. Even as open-source Instruct models are trained on significantly larger, often undisclosed instruction-tuning datasets, ALIGNVLM achieves superior performance. For instance, ALIGNVLM-Llama-3.2-3B (58.81%) outperforms all instructed VLMs in its size category, surpassing its closest competitor, Qwen2-VL-2B (55.84%), by 2.97%. Additionally, our 8B model outperforms significantly larger models such as Llama 3.2-11B and PixTral-12B by substantial margins. It also surpasses InternVL-2.5-8B and performs competitively with Qwen2-VL-7B, though a direct comparison may not be entirely fair since Qwen2-VL-7B was trained on an undisclosed instruction-tuning dataset. Finally, ALIGNVLM also exhibits comparable performance to closed-source models like GeminiPro-1.5 and GPT4o.

Overall, these results validate the effectiveness of ALIGN and establish ALIGNVLM as a state-of-the-art model for multimodal document understanding.

5.2 IMPACT OF CONNECTOR DESIGNS ON VLM PERFORMANCE

To assess the effectiveness of our ALIGN module, we compare it against three different and widely used *shallow fusion* VLM connectors: MLP, Perceiver Resampler, and Ovis. The results in Table 2 show that ALIGN consistently outperforms all alternatives, demonstrating its superiority both in aligning visual and textual modalities and in multimodal document understanding. MLP and Perceiver Resampler achieve the lowest performance, 53.06% and 50.68%, respectively, due to their

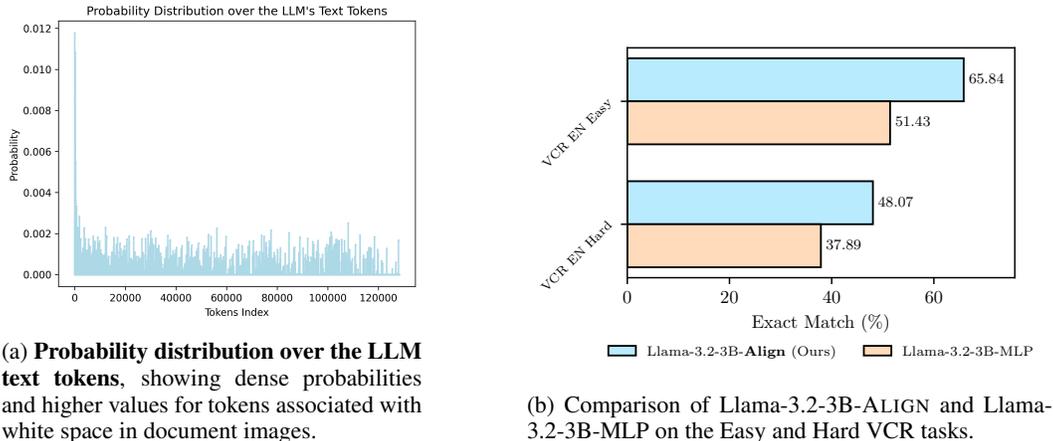


Figure 3: Probability distribution over the LLM text tokens and VCR Task Analysis.

direct feature projection, which lacks an explicit mechanism to align visual features with the LLM’s text space, leading to misalignment. Ovis introduces a separate visual embedding table, but this additional complexity does not significantly improve alignment, yielding only 54.72% accuracy. In contrast, ALIGN ensures that visual features remain within the convex hull of the LLM’s text latent space, leveraging the linguistic priors of the LLM to enhance alignment and mitigate noisy embeddings. This design leads to the highest performance (58.81%), establishing ALIGN as the most effective connector for integrating vision and language in multimodal document understanding. We provide some example outputs of the Llama-3.2-3B models with different connector designs in Appendix A.3.

5.3 PROBABILITY DISTRIBUTION OVER TEXT TOKENS ANALYSIS

To better understand the behavior of ALIGN, we examine the probability distribution, $\mathbf{P}_{\text{vocab}}$ in Eq equation 1, over the LLM’s text vocabulary generated from visual features. Specifically, we process 100 document images through the vision encoder and ALIGN, then average the resulting probability distributions across all image patches. The final distribution is shown in Figure 3a. As illustrated, the distribution is *dense* (rather than sparse), with the highest probability assigned to a single token being *0.0118*. This can be explained by the vision feature space being continuous and of much higher cardinality than the discrete text space. Indeed, while the LLM has 128K distinct vocabulary tokens, an image patch (e.g., 14×14 pixels) contains continuous, high-dimensional information that cannot be effectively mapped to a single or a few discrete tokens.

Furthermore, we observe that tokens on the left side of the distribution in Figure 3a have higher probabilities than the rest. Upon investigation, we found that these tokens correspond to patches that are predominantly white – a common feature in document images. Further analysis of the associated text tokens reveals that they predominantly consist of punctuation marks, as illustrated further in Appendix A.2. This suggests that the model repurposes punctuation marks to represent whitespaces. This may be attributed to the fact that both punctuation and whitespaces act as structural cues and separators. Other possibilities include whitespaces being rarely directly-required to perform a task, and LLMs may pay less specific attention to common tokens such as punctuation.

5.4 PIXEL-LEVEL TASKS ANALYSIS

To rigorously evaluate the ability of vision-language models to integrate fine-grained visual and textual pixel-level cues, we test our model on the VCR benchmark (Zhang et al., 2024), which requires the model to recover partially occluded texts with pixel-level hints from the revealed parts of the text. This task challenges VLM’s alignment of text and image in extreme situations. Current state-of-the-art models like GPT-4V OpenAI et al. (2023), Claude 3.5 Sonnet Anthropic (2024), and Llama-3.2 Dubey et al. (2024) significantly underperform humans on *hard* VCR task due to their inability to process subtle pixel-level cues in occluded text regions. These models frequently

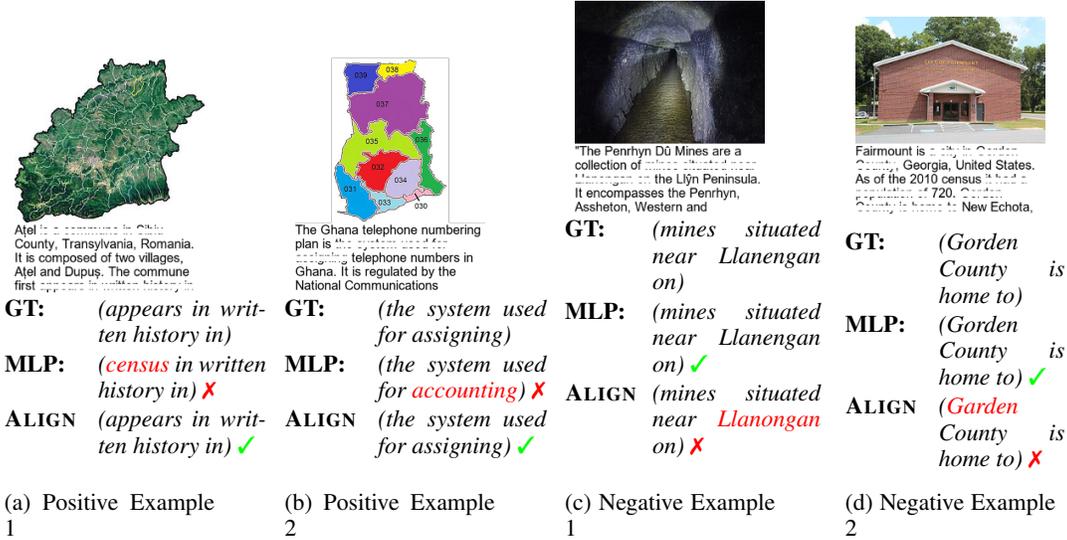


Figure 4: **Case Study for Pixel-Level Tasks.** We provide examples of our proposed **ALIGN** connector compared with a the Multi-Layer Perceptron (MLP) connector. The **ALIGN** connector tends to better map visual elements to common words. GT is the ground truth.

discard critical visual tokens during image tokenization on semantic priors, overlooking the interplay between partial character strokes and contextual visual scenes. To evaluate performance on VCR, we modify our Stage 3 SFT dataset composition by replacing the exclusive use of DocDownstream with a 5:1 blended ratio of DocDownstream and VCR training data. This adjustment enables direct evaluation of our architecture **ALIGN**'s ability to leverage pixel-level character cues.

From the experimental outcomes, it is evident that **ALIGNVLM** consistently outperforms the MLP Connector Model across both easy and hard settings of the pixel-level VCR task (see Figure 3b), with improvements ranging from 10.18% on the hard setting to 14.41% on the easy setting.

We provide a case study on VCR in Figure 4, featuring four representative examples. In Figure 4a, it is evident that the MLP connector model fails to capture semantic consistency as effectively as **ALIGNVLM**. The phrase “The commune first *census in written history in*” (where the words in italics are generated by the model while the rest are in the image) is not as semantically coherent as the phrase generated by **ALIGN** “The commune first *appears in written history in*”.

Beyond the issue of semantic fluency, in Figure 4b we also observe that **ALIGNVLM** successfully identifies the uncovered portion of the letter “g” in “accounting” and uses it as a pixel-level hint to infer the correct word. In contrast, the MLP model fails to effectively attend to this crucial detail.

Figures 4c and 4d show examples where **ALIGNVLM** fails on the VCR task. These carefully picked instances show that our method mistakes names of landmarks with common words when the two are very similar. As seen in the examples, **ALIGNVLM** mistakes “Llanengan” for “Llanongan” and “Gorden” for “Garden”. In both instances, the pairs differ by one character, indicating perhaps that **ALIGNVLM** tends to align vision representations to more common tokens in the vocabulary. One approach that would potentially mitigate such errors would be to train **ALIGNVLM** with more contextually-relevant data.

5.5 ROBUSTNESS TO NOISE ANALYSIS

To evaluate the robustness of our **ALIGN** connector to noisy visual features, we conduct an experiment where random Gaussian noise is added to the visual features produced by the vision encoder before passing them into the connector. Specifically, given the visual features $\mathbf{F} \in \mathbb{R}^{N \times d}$ output by the vision encoder (where N is the number of feature vectors and d is their dimensionality), we perturbed them as

$$\tilde{\mathbf{F}} = \mathbf{F} + \mathbf{N}, \quad \mathbf{N} \sim \mathcal{N}(0, \sigma = 3).$$

Table 3: **Robustness to Noise.** Comparison of Avg. Scores with and without Gaussian noise ($\sigma = 3$), including performance drop (Δ).

Model	Without Noise	With Noise	Drop (Δ)
Llama-3.2-3B-MLP	53.06	27.52	$\downarrow 25.54$
Llama-3.2-3B-ALIGN (ours)	58.81	57.14	$\downarrow 1.67$

As shown in Table 3, our **ALIGN** connector demonstrates high robustness to noise, with only a *1.67%* average drop in performance. In contrast, the widely adopted MLP connector suffers a significant performance degradation of *25.54%*, highlighting its vulnerability to noisy inputs. These empirical results support our hypothesis that leveraging the knowledge encoded in the LLM’s text embeddings and constraining the visual features within the convex hull of the text latent space act as a regularization mechanism, reducing the model’s sensitivity to noisy visual features.

6 CONCLUSION

We introduce **ALIGN**, a novel connector designed to align vision and language latent spaces in vision-language models (VLMs), specifically enhancing multimodal document understanding. By improving cross-modal alignment and minimizing noisy embeddings, our models, **ALIGNVLM**, which leverage **ALIGN**, achieve state-of-the-art performance across diverse document understanding tasks. This includes outperforming base VLMs trained on the same datasets and open-source instruct models trained on undisclosed data. Extensive experiments and ablations validate the robustness and effectiveness of **ALIGN** compared to existing connector designs, establishing it as a significant contribution to vision-language modeling. Future work will explore training on more diverse instruction-tuning datasets to generalize beyond document understanding to broader domains.

REFERENCES

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatiakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, Albert Q. Jiang, Kartik Khandelwal, Timothée Lacroix,

- Guillaume Lample, Diego Las Casas, Thibaut Lavril, Teven Le Scao, Andy Lo, William Marshall, Louis Martin, Arthur Mensch, Pavankumar Muddireddy, Valera Nemychnikova, Marie Peltat, Patrick Von Platen, Nikhil Raghuraman, Baptiste Rozière, Alexandre Sablayrolles, Lucile Saulnier, Romain Sauvestre, Wendy Shang, Roman Soletskyi, Lawrence Stewart, Pierre Stock, Joachim Studnia, Sandeep Subramanian, Sagar Vaze, Thomas Wang, and Sophia Yang. Pixtral 12b, 2024. URL <https://arxiv.org/abs/2410.07073>.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. URL <https://arxiv.org/abs/2204.14198>.
- Reza Yazdani Aminabadi, Samyam Rajbhandari, Minjia Zhang, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Jeff Rasley, Shaden Smith, Olatunji Ruwase, and Yuxiong He. Deepspeed inference: Enabling efficient inference of transformer models at unprecedented scale, 2022. URL <https://arxiv.org/abs/2207.00032>.
- Anthropic. The claude 3 model family: Opus, sonnet, haiku. 2024.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and Xiaohua Zhai. Paligemma: A versatile 3b vlm for transfer, 2024. URL <https://arxiv.org/abs/2407.07726>.
- Rishi Bommasani, Kevin Klyman, Shayne Longpre, Sayash Kapoor, Nestor Maslej, Betty Xiong, Daniel Zhang, and Percy Liang. The foundation model transparency index, 2023. URL <https://arxiv.org/abs/2310.12941>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts, 2021. URL <https://arxiv.org/abs/2102.08981>.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. In *International Conference Learning Representations*, 2020.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiaoyi Dong, Hang Yan, Hwei Guo, Conghui He, Botian Shi, Zhenjiang Jin, Chao Xu, Bin Wang, Xingjian Wei, Wei Li, Wenjian Zhang, Bo Zhang, Pinlong Cai, Licheng Wen, Xiangchao Yan, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites, 2024a. URL <https://arxiv.org/abs/2404.16821>.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24185–24198, 2024b.
- Wenliang Dai, Nayeon Lee, Boxin Wang, Zhuoling Yang, Zihan Liu, Jon Barker, Tuomas Rintamaki, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Nvlm: Open frontier-class multimodal llms. *arXiv preprint arXiv: 2409.11402*, 2024.

Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. Unveiling encoder-free vision-language models. *arXiv preprint arXiv:2406.11832*, 2024.

Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H. Laradji, Manuel Del Verme, Tom Marty, Léo Boisvert, Megh Thakkar, Quentin Cappart, David Vazquez, Nicolas Chapados, and Alexandre Lacoste. Workarena: How capable are web agents at solving common knowledge work tasks?, 2024. URL <https://arxiv.org/abs/2403.07718>.

Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 11198–11201, 2024.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, and et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie DelPierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papanikos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein,

Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding, 2024. URL <https://arxiv.org/abs/2403.12895>.

Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form understanding in noisy scanned documents, 2019. URL <https://arxiv.org/abs/1905.13538>.

- Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeongyeon Nam, Jinyoung Park, Jinyeong Yim, Won-seok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer, 2022. URL <https://arxiv.org/abs/2111.15664>.
- Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. Tablevqa-bench: A visual question answering benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205*, 2024.
- Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models?, 2024. URL <https://arxiv.org/abs/2405.02246>.
- Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding, 2023. URL <https://arxiv.org/abs/2210.03347>.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024. URL <https://arxiv.org/abs/2408.03326>.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023. URL <https://arxiv.org/abs/2301.12597>.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b.
- 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/>.
- Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. Ovis: Structural embedding alignment for multimodal large language model, 2024. URL <https://arxiv.org/abs/2405.20797>.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- Minesh Mathew, Viraj Bagal, Rubèn Pérez Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V Jawahar. Infographicvqa, 2021a. URL <https://arxiv.org/abs/2104.12756>.
- Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2021b. URL <https://arxiv.org/abs/2007.00398>.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. Cord: A consolidated receipt dataset for post-ocr parsing. *Document Intelligence Workshop at Neural Information Processing Systems*, 2019.
- Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *Annual Meeting of the Association for Computational Linguistics*, 2015.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,

- Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- 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, 2023. URL <https://arxiv.org/abs/1910.10683>.
- Juan Rodriguez, Xiangru Jian, Siba Smarak Panigrahi, Tianyu Zhang, Aarash Feizi, Abhay Puri, Akshay Kalkunte, François Savard, Ahmed Masry, Shravan Nayak, Rabiul Awal, Mahsa Masoud, Amirhossein Abaskohi, Zichao Li, Suyuchen Wang, Pierre-André Noël, Mats Leon Richter, Saverio Vadicchino, Shubham Agarwal, Sanket Biswas, Sara Shanian, Ying Zhang, Noah Bolger, Kurt MacDonald, Simon Fauvel, Sathwik Tejaswi, Srinivas Sunkara, Joao Monteiro, Krishnamurthy DJ Dvijotham, Torsten Scholak, Nicolas Chapados, Sepideh Kharagani, Sean Hughes, M. Özsu, Siva Reddy, Marco Pedersoli, Yoshua Bengio, Christopher Pal, Issam Laradji, Spandana Gella, Perouz Taslakian, David Vazquez, and Sai Rajeswar. Bigdocs: An open and permissively-licensed dataset for training multimodal models on document and code tasks, 2024a. URL <https://arxiv.org/abs/2412.04626>.
- Juan A. Rodriguez, David Vazquez, Issam Laradji, Marco Pedersoli, and Pau Rodriguez. Ocrvqgan: Taming text-within-image generation, 2022. URL <https://arxiv.org/abs/2210.11248>.
- Juan A. Rodriguez, Abhay Puri, Shubham Agarwal, Issam H. Laradji, Pau Rodriguez, Sai Rajeswar, David Vazquez, Christopher Pal, and Marco Pedersoli. Starvector: Generating scalable vector graphics code from images and text, 2024b. URL <https://arxiv.org/abs/2312.11556>.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *IEEE Conference Computer Vision Pattern Recognition*, 2019.
- Tomasz Stanisławek, Filip Graliński, Anna Wróblewska, Dawid Lipiński, Agnieszka Kaliska, Paulina Rosalska, Bartosz Topolski, and Przemysław Biecek. Kleister: key information extraction datasets involving long documents with complex layouts. In *International Conference on Document Analysis and Recognition*, 2021.
- S Svetlichnaya. Deepform: Understand structured documents at scale, 2020.
- Gemini Team. Gemini: A family of highly capable multimodal models, 2024. URL <https://arxiv.org/abs/2312.11805>.
- Caitlin Vogus and Emma Llansóe. Making transparency meaningful: A framework for policymakers. *Center for Democracy and Technology*, 2021.
- Dongsheng Wang, Natraj Raman, Mathieu Sibue, Zhiqiang Ma, Petr Babkin, Simerjot Kaur, Yulong Pei, Armineh Nourbakhsh, and Xiaomo Liu. Doclm: A layout-aware generative language model for multimodal document understanding, 2023a. URL <https://arxiv.org/abs/2401.00908>.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution, 2024. URL <https://arxiv.org/abs/2409.12191>.
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023b.

Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, and Ping Luo. Janus: Decoupling visual encoding for unified multimodal understanding and generation, 2024a. URL <https://arxiv.org/abs/2410.13848>.

Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding, 2024b. URL <https://arxiv.org/abs/2412.10302>.

Ruyi Xu, Yuan Yao, Zonghao Guo, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan Liu, Maosong Sun, and Gao Huang. Llava-uhd: an lmm perceiving any aspect ratio and high-resolution images. *European Conference on Computer Vision*, 2024. doi: 10.48550/arXiv.2403.11703.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. URL <https://arxiv.org/abs/2303.15343>.

Tianyu Zhang, Suyuchen Wang, Lu Li, Ge Zhang, Perouz Taslakian, Sai Rajeswar, Jie Fu, Bang Liu, and Yoshua Bengio. Vcr: Visual caption restoration. *arXiv preprint arXiv: 2406.06462*, 2024.

Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang, Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. Swift: a scalable lightweight infrastructure for fine-tuning, 2024. URL <https://arxiv.org/abs/2408.05517>.

A APPENDIX

A.1 EXPERIMENTAL SETUP

We provide detailed hyperparameters of our experiments in Table 4.

Table 4: Detailed hyperparameters for each training stage across different LLM backbones.

LLM Backbone	Llama 3.2-1B			Llama 3.2-3B			Llama 3.1-8B		
	Stage-1	Stage-2	Stage-3	Stage-1	Stage-2	Stage-3	Stage-1	Stage-2	Stage-3
Trainable Parameters	Full Model	Full Model	LLM & ALIGN	Full Model	Full Model	LLM & ALIGN	Full Model	Full Model	LLM & ALIGN
Batch Size	512	512	512	512	256	256	512	256	256
Text Max Length	1024	2048	2048	1024	2048	2048	1024	2048	2048
Epochs	1	1	5	1	1	5	1	1	5
Learning Rate	1×10^{-5}	5×10^{-5}	5×10^{-5}	1×10^{-5}	5×10^{-5}	5×10^{-5}	1×10^{-5}	1×10^{-5}	1×10^{-5}

A.2 VISION-TO-TEXT

In this experiment, we analyze how ALIGN maps visual features to the LLM’s text tokens. To do so, we manually curate a small dataset of image crops, each containing either a single word or a small set of visual text elements. Unlike the processing of high-resolution images described earlier (Section 3.1), these image crops are not divided into tiles. Instead, the backbone image encoder processes each crop as a single tile, producing 14×14 features from the input image. The resulting features pass through the Softmax operation (Equation 1), yielding a probability distribution over the LLM’s text tokens for each feature (region). We examine the decoded text tokens from specific image regions to better understand how visual features are mapped to textual representations.

As shown in Figure 5, white regions in the images tend to assign higher probabilities to punctuation tokens, such as commas or periods. Since punctuation structures written text, while white space separates document components like paragraphs, tables, and sections, ALIGN appears to leverage these implicit patterns to align visual structures with semantically meaningful representations in the LLM’s embedding space.

A.3 CASE STUDIES

In this section, we provide case studies for the experiments in Section 5.1. Specifically, we provide examples of our Llama-3.2-3B-ALIGN, and its counterpart model with alternative connectors Llama-3.2-3B-MLP and Llama-3.2-3B-Ovis on three different datasets: KLC (Stanisławek et al., 2021), DocVQA (Mathew et al., 2021b), and TextVQA (Singh et al., 2019). The examples are shown in Figure 6, 7, and 8.

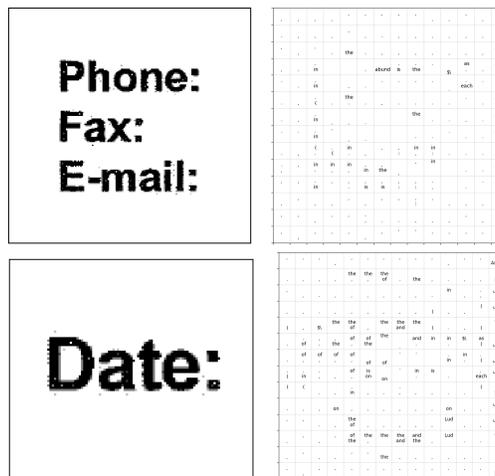
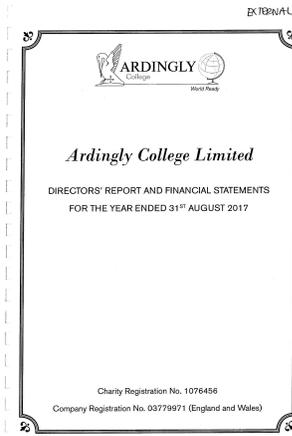
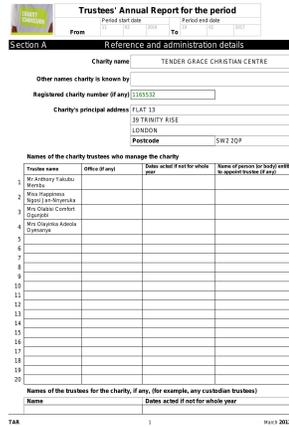


Figure 5: **Mapping Visual-to-Text tokens.** The left column shows the visual input to the model. In contrast, the right column visualizes the decoded tokens on a 14x14 grid, displaying the top k=2 tokens corresponding to the most likely LLM tokens predicted for the respective visual feature in each cell.



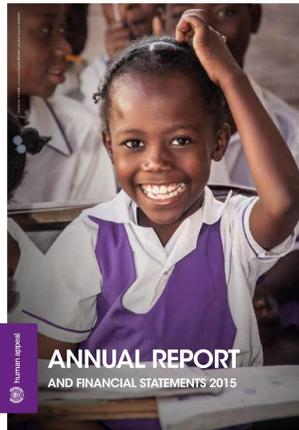
Question: What is the value for the charity name?
GT: (Ardingly College Ltd.)
MLP: (Ardington College Ltd.) ✗
Ovis: (Ardington College Ltd.) ✗
ALIGN: (Ardingly College Ltd.) ✓

(a) Positive Example #1



Question: What is the value for the address postcode?
GT: (SW2 2QP)
MLP: (SW22 0PQ) ✗
Ovis: (SW2 2OP) ✗
ALIGN: (SW2 2QP) ✓

(b) Positive Example #2



Question: What is the value for the charity name?
GT: (Human Appeal)
MLP: (Humanitarian Agenda) ✗
Ovis: (Human Appeal) ✓
ALIGN: (Human Rightsappeal) ✗

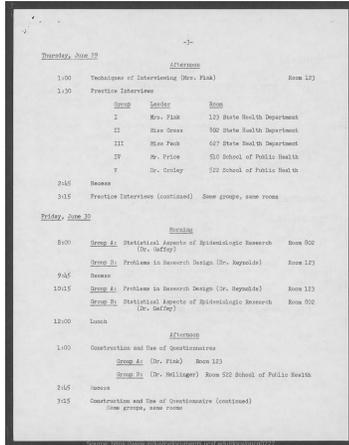
(c) Negative Example #1



Question: What is the value for the post town address?
GT: (Bishop's Stortford)
MLP: (Stortford) ✗
Ovis: (Bishop's Stortford) ✓
ALIGN: (Stortford) ✗

(d) Negative Example #2

Figure 6: Case Study for Connector Comparison on the KLC dataset (Stanisławek et al., 2021). We show four qualitative examples (including two correct and two incorrect examples) comparing Llama-3.2-3B-ALIGN to the same architecture with different connectors, Llama-3.2-3B-MLP and Llama-3.2-3B-Ovis. “GT” denotes the ground truth.



Question: What does the afternoon session begin on June 29?

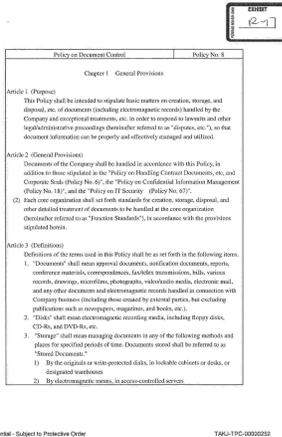
- GT:** (1:00)
MLP: (2:45) ✗
Ovis: (3:30) ✗
ALIGN: (1:00) ✓

(a) Positive Example #1

Question: What levels does the second table indicate?

- GT:** (hematocrit data - Massachusetts)
MLP: (SATISFACTORY) ✗
Ovis: (Females) ✗
ALIGN: (hematocrit data - Massachusetts) ✓

(b) Positive Example #2



Question: What type of policy is described in this document?

- GT:** (Policy on Document Control)
MLP: (Policy on Document Control) ✓
Ovis: (General Provisions) ✗
ALIGN: (Document Control) ✗

(c) Negative Example #1

Question: What was the diet fed to the #1 group?

- GT:** (basal diet)
MLP: (basal diet) ✓
Ovis: (Whole blood) ✗
ALIGN: (control diet) ✗

(d) Negative Example #2

Figure 7: Case Study for Connector Comparison on the DocVQA dataset (Mathew et al., 2021b). We show four qualitative examples (including two correct and two incorrect examples) comparing Llama-3.2-3B-ALIGN to the same architecture with different connectors, Llama-3.2-3B-MLP and Llama-3.2-3B-Ovis. “GT” denotes the ground truth.



Question: What greeting is written on the letter?

GT: (good bye)
MLP: (good) ✗
Ovis: (good buy) ✗
ALIGN: (good bye) ✓

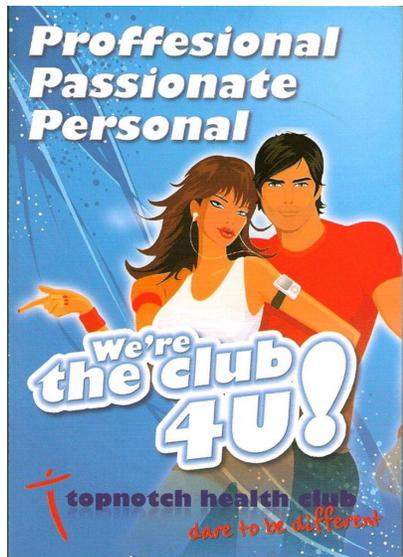
(a) Positive Example #1



Question: What indoor temperature is shown?

GT: (68.4)
MLP: (68 F) ✗
Ovis: (40.0) ✗
ALIGN: (68.4) ✓

(b) Positive Example #2



Question: What type of club is advertised?

GT: (health club)
MLP: (topnote health club) ✗
Ovis: (health club) ✓
ALIGN: (professional passionate personal) ✗

(c) Negative Example #1



Question: What credit card is this?

GT: (hadiah plus)
MLP: (hadiah plus) ✓
Ovis: (american big loyalty program) ✗
ALIGN: (hadia plus) ✗

(d) Negative Example #2

Figure 8: **Case Study for Connector Comparison on the TextVQA dataset (Singh et al., 2019).** We show four qualitative examples (including two correct and two incorrect examples) comparing Llama-3.2-3B-ALIGN to the same architecture with different connectors, Llama-3.2-3B-MLP and Llama-3.2-3B-Ovis. “GT” denotes the ground truth.