

DOLMA: Visual Instruction Tuning for Document AI

Anonymous ACL submission

Abstract

The rapid expansion of Vision-Language Models (VLMs) has spurred research into their applicability across various domains. While VLMs excel in understanding environmental contexts, their effectiveness declines with visually-rich scanned documents. Although some VLMs use Optical Character Recognition (OCR) to mitigate this, OCR alone is insufficient for the complex textual and visual insights required. Developing tailored models for Document AI applications also demands substantial labeled data and high training costs. To address these challenges, we conducted experiments with various models, data types, architectures, and training methodologies. Based on our findings, we introduce DOLMA, an OCR-free vision-language model designed for diverse Document AI applications in a zero-shot setting. Despite having a moderate parameter count of 7 billion, DOLMA performs on par with models ten times larger on numerous Document AI benchmarks. The complete model, including weights, training data, and code, is publicly available.

1 Introduction

In recent years, there has been a notable surge in interest surrounding the understanding of visually-rich scanned documents (VRD). The latter encompasses PDFs and document images such as business forms, receipts, driving licenses and invoices. The understanding and digitization of those document images entails intricate tasks such as document visual question answering (DVQA), document classification (CLS), and key information extraction (KIE).

Traditional approaches address these challenges by employing Optical Character Recognition (OCR) alongside handcrafted rules or

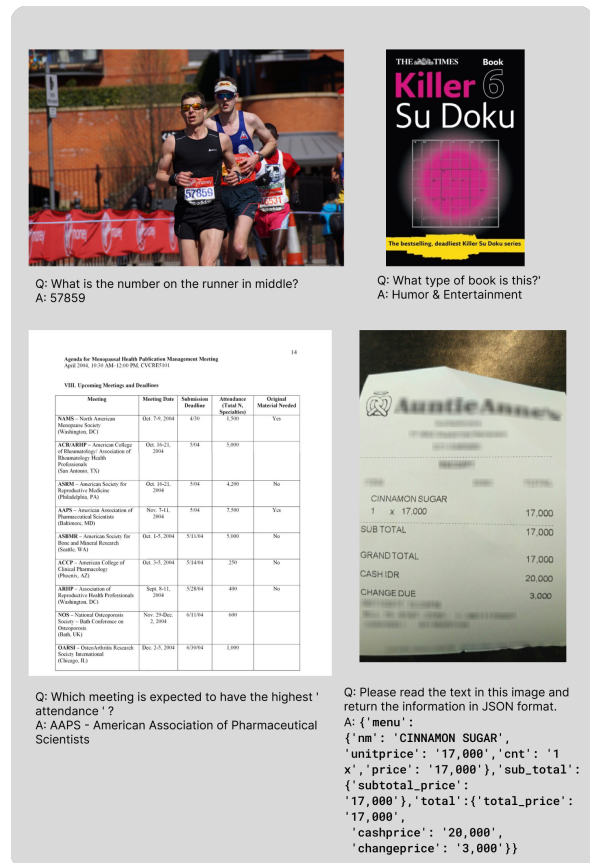


Figure 1: The training pipeline of DOLMA.

layout analysis. However, these methods often necessitate post-processing steps, potentially limiting the efficacy and use of those models. In recent years, the Document AI community has proposed various transformer-based architectures providing remarkable progress on VRD understanding (VRDU). Notably, Transformer-based models like LayoutLM and its variants have showcased advancements by integrating OCR, image, and layout information. Nevertheless, recent efforts in OCR-free, end-to-end document understanding from images indicate

053 a shift towards more versatile models, mini-
054 mizing task-specific engineering and reducing
055 reliance on external components during infer-
056 ence.

057 In this study, we aimed to explore various de-
058 sign choices to identify the optimal combination
059 of models, data, and architecture based on our
060 experiments. We also imposed constraints on
061 model size and resource usage to demonstrate
062 the most efficient and cost-effective approach
063 to developing a model that can perform on par
064 with other state-of-the-art models. To assess
065 the quality and utility of the model, we evaluate
066 it based on the following properties:

- 067 • **Property 1: Multitasking.** The model
068 is expected to perform the main Document
069 AI tasks such as document classification,
070 document question answering, and key in-
071 formation extraction.
- 072 • **Property 2: OCR-independence.** Key
073 information in documents is many times in-
074 corporated in non-optical characters such
075 as logos, images, charts and other visu-
076 als. OCR-dependent models do not have
077 the capability to extract this information.
078 Nonetheless, we consider the models that
079 do not necessarily rely on OCR yet can
080 improve the results using OCR informa-
081 tion. We call them OCR-enhanced models
082 as they can still perform without relying
083 on OCR.
- 084 • **Property 3: Instruction following.**
085 The typical usage of information extraction
086 from documents is related to structuring
087 image data into programmatically read-
088 able formats such as JSON, XML or CSV.
089 As the use cases of information extraction
090 can be different, the Document AI founda-
091 tion model should have the ability to
092 follow the user’s instruction and generate
093 extracted output in the required format
094 (including notation format such as JSON
095 and its internal structure such as key/value
096 hierarchy).
- 097 • **Property 4: Template independence.**
098 The Document AI foundation model
099 should be able to provide competitive per-
100 formance on the same documents even if
101 the templates are different.

102 We outline the following roadmap of experi-
103 ments, which will be discussed in subsequent
104 sections. The modalities we consider include
105 a Vision encoder, a Language decoder, and a
106 bridge connector between them. We establish
107 two stages for training: (1) pretraining and (2)
108 fine-tuning. Stage (1) is designed to enable
109 the model to acquire OCR capabilities, while
110 stage (2) focuses on task-specific supervised
111 instruction tuning.

112 During stage (1), we experiment with (a)
113 the design of the bridge connector and (b) the
114 choice of language model. For (a), we report
115 findings using design choices from LLAVA (Liu
116 et al., 2023) for the linear projection strat-
117 egy, QwenVL (Bai et al., 2023) for the cross-
118 attention strategy, and Idefics2 (Laurençon
119 et al., 2024) for the projection + perceiver-
120 resampler strategy. For (b), we evaluate Vi-
121 cuna (Zheng et al., 2023), LLAMA 3 (Team,
122 2024), and Phi 3 (Abdin et al., 2024). We select
123 Vicuna as a well-established instruction model,
124 LLAMA 3 as a state-of-the-art large language
125 model, and Phi to assess the impact of using
126 smaller models.

127 During stage (2), our primary focus is on
128 training strategies. We discovered that train-
129 ing all modalities yields the best results. Con-
130 sequently, the main variable is the strategic
131 approach to the largest modality, which in our
132 case is the LLM. We report on three strate-
133 gies: fine-tuning only the attention layers of
134 the LLM, full LLM fine-tuning, and applying
135 LoRA on top of the LLM. In all three scenar-
136 ios, we fully fine-tune the vision encoder and
137 the bridge connector. All the aforementioned
138 experiments are conducted using 8 H100 GPU
139 spot instances to ensure the fastest possible
140 training time. Building on our observations, we
141 propose DOLMA, "Document Optimized Lan-
142 guage Model for Automation," which adheres
143 to the four principles outlined above. DOLMA
144 is a 7-billion-parameter Vision-Language Model
145 (VLM) that achieves results on various Docu-
146 ment AI benchmarks on par with state-of-the-
147 art models, even matching the performance of
148 models that are ten times larger.

149 2 Related Work

150 The advent of ChatGPT represents a significant
151 advancement in the domain of Large Language

Models (LLMs). LLMs constitute a substantial area of study in natural language processing, specializing in processing and generating textual content for tasks like language translation, summarizing, question answering, and text completion.

2.1 LLMs

Through extensive pre-training on textual datasets, LLMs acquire proficiency in contextual relationships and linguistic patterns. The transformative impact of transformers, as introduced in "Attention is All You Need," has played a pivotal role in the success of LLMs, leading to the development of pre-trained models such as BERT, BART, and others. This success has spurred further exploration into LLMs like OPT (Zhang et al., 2022), BLOOM (Workshop et al., 2023), PaLM (Chowdhery et al., 2022), and LLaMA (Touvron et al., 2023a). Particularly noteworthy is LLaMA3 (Team, 2024), an open-source LLM demonstrating comparable or superior performance to both open and closed-source models. The open-source nature of LLaMA has encouraged numerous researchers to build models on top of it, employing diverse training strategies and architectural modifications, including models like Vicuna (Zheng et al., 2023) and Alpaca (Taori et al., 2023).

2.2 Multimodal LLM

In the realm of multimodal AI, Multimodal Language Models (MLMs) have emerged as a significant focus. Unlike text-to-text models, MLMs are designed to comprehend and generate content across multiple modalities, often integrating text and images. These models exhibit proficiency in tasks requiring a fusion of textual and visual understanding, such as generating image captions, image-text matching, visual question answering and contextualizing information in mixed-media environments. Training MLMs involves leveraging datasets encompassing both textual and visual information, facilitating the capture of intricate relationships between words and images. Notable MLMs include GPT-4V, LLaVA (Liu et al., 2023), Gemini Pro Vision, and others.

2.3 Document AI

Transformer-based architectures have found success in Visual Document Understanding (VRDU) and Document Visual Question Answering (DVQA) tasks (Wang et al., 2023b; Ye et al., 2023a,b; Kim et al., 2022; Hong et al., 2023; Bai et al., 2023). Recent works like LayoutLM (Huang et al., 2022) focus on pre-training a language model, such as BERT (Devlin et al., 2019), alongside an OCR-based engine to comprehend both textual content and layout information in document images. This approach extends traditional language models by incorporating positional embeddings that encode the spatial arrangement of words on a page, enabling the model to capture both structural relationships and contextual meanings. Recent works, such as DocLLM (Wang et al., 2023a), integrate lightweight visual information by utilizing spatial positions and dimensions of text tokens obtained through OCR. It employs separate vectors to represent vision and image modalities, extending the self-attention mechanism of the transformer architecture to compute their interdependencies in a disentangled manner. Alternative methods, exemplified by DONUT (Kim et al., 2022), leverage transformer architectures for document understanding tasks, focusing on extracting information directly from the document's content without relying on OCR. DONUT employs the Swin transformer (Liu et al., 2021) as the vision encoder and BART (Lewis et al., 2019) as the decoder model. A more general model, Qwen-VL (Bai et al., 2023), incorporates an adapter with cross-attention layers to attenuate vision encoder embeddings with language embeddings. Qwen-VL, trained on a large corpus of both regular and document images, demonstrates proficiency in tasks such as image captioning, question answering, visual grounding, and text reading.

As shown in the Table 1, among the models mentioned above, GPT-4V, Gemini Pro Vision, LLaVA and Qwen-VL are the only models that satisfy the 4 properties we seek in a Document AI foundation model.

In summary, our review highlights the significant strides in LLMs, the emergence of multimodal AI with MLMs, and the successful applications of transformer architectures in VrDu

Table 1: Comparison of different models across the 4 properties we seek

Model	Property 1 Multitasking	Property 2 OCR-free	Property 3 Instruction following	Property 4 Template independent
GPT4-V	✓	✓	✓	✓
Gemini-Pro-Vision	✓	✓	✓	✓
Donut	×	✓	×	×
LayoutLMV3	✓	×	×	✓
DocLLM	✓	×	×	✓
Qwen-VL	✓	✓	✓	✓
LLaVA	✓	✓	✓	✓
CogAgent	✓	✓	✓	✓
UReader	✓	✓	✓	✓
DocOwl	✓	✓	✓	✓
DOLMA (ours)	✓	✓	✓	✓

and DVQA tasks. These advancements lay the groundwork for versatile models like Qwen-VL, showcasing the evolving landscape of AI and machine learning.

3 Analysing the design possibilities for vision-language models

In this section, we will examine the various design choices for vision-language models in Document AI as documented in the open-source literature and present our findings. For our experiments, we will utilize the IIT-CDIP dataset and our own PDF-generated dataset for pretraining. Additionally, we will employ various benchmarking datasets, including DocVQA (Mathew et al., 2021b), *CORD-V2* (Park et al., 2019), *Infographics-VQA* (Mathew et al., 2021a), *ICDAR-SROIE* (ICDAR, 2019), *Chart-QA* (Masry et al., 2022), *OCR-VQA* (Mishra et al., 2019), *RVL-CDIP* (Harley et al., 2015), and *TextVQA* (Singh et al., 2019), for fine-tuning experiments.

3.1 The design of the bridge connector

Vision-language models comprise two modalities: vision and language. While there are numerous models available for these modalities, it is essential for them to effectively "communicate" with each other. We explore three types of connectors: linear projection, cross-attention, and projection + perceiver-resampler.

For our experiments, we fixed Swin Base (Liu et al., 2021) as the vision model and Vicuna (Zheng et al., 2023) as the language model across all three bridge connector designs. Dur-

ing the pretraining stage, we trained both the vision model and each connector while keeping the LLM model frozen. Utilizing a total of 3 million image-text pairs from the IIT-CDIP (Soboroff, 2022) dataset and our in-house PDF-generated data (details of which will be discussed in subsequent sections), we pretrained the model for the text extraction task, thereby imparting OCR capabilities. We trained each configuration for a total of one epoch. Among the three connectors, only the linear projector connector successfully converged. We hypothesize that extended training might enhance the performance of the other connectors; however, given our experimental setup and resource constraints, the linear projector layer demonstrated the best results.

Insight 1

The linear projector layer is the fastest and most straightforward method to connect the vision and language models, achieving training convergence with just one epoch on 3 million image-text pairs.

3.2 The design of the Vision model

For the vision model, we selected the SWIN Transformer (Liu et al., 2021) (base, large) and the Vision Transformer (Dosovitskiy et al., 2020) (CLIP ViT-L/14). Similar to the previous section, we conducted the pretraining stage by keeping the LLM frozen and training the vision model along with the connector. We fixed the bridge connector to the linear projection design and experimented with different vision

311 encoders.

312 Our findings revealed that 3 million image-
313 text pairs and one epoch of training were
314 insufficient for the vision models to acquire
315 text extraction capabilities. We utilized pre-
316 trained weights for each model, but none of
317 the trainings converged except for the Swin
318 Base model. For Swin Base, we used weights
319 from the DONUT (Kim et al., 2022) model,
320 which had been pretrained for the text extrac-
321 tion task using over 11 million image-text pairs
322 and trained for 200K steps with a batch size of
323 196. The DONUT model employed Swin Base,
324 and similarly, when we integrated Swin Base
325 into our architecture, the training eventually
326 converged.

Insight 2

Vision models require tens of millions of text extraction pretraining data and extended training sessions to develop OCR capabilities.

3.3 The design of the LLM model

327 We selected Vicuna 1.5 (Zheng et al., 2023),
328 LLama3 (Team, 2024), and Phi (Abdin et al.,
329 2024) for our experiments. Following the suc-
330 cess of LLaVA, we chose Vicuna as our initial
331 model. We included LLama3 to evaluate the
332 impact of a relatively newly released state-of-
333 the-art model. Additionally, we decided to use
334 the Phi-3 model to assess the performance of
335 a model with fewer than seven billion param-
336 eters. The vision encoder employed was Swin
337 Base, and the projection design was used as
338 the bridge connector. During the pretraining
339 stage, we kept the LLM frozen and only pre-
340 trained the vision model and bridge connector.
341 At the conclusion of the experiments, both Phi
342 and LLama3 failed to converge during train-
343 ing, whereas only Vicuna was able to achieve
344 near-zero loss for the text extraction task.

345 The objective of this experiment is to com-
346 pare different fine-tuning strategies within the
347 given setup, based on performance across vari-
348 ous benchmarks as well as computational com-
349 plexity. Our experimental configuration in-
350 cludes Swin Base (pretrained with Donut) as
351 the vision encoder, a projection as the bridge
352 connector, and Vicuna as the LLM decoder.
353 Swin Base and the projection were pretrained

354 as described in the previous sections. During
355 the fine-tuning stage, we continue training the
356 entire model (all three modalities) on bench-
357 mark datasets. While we fully train the Vision
358 Encoder and the bridge connector, we apply
359 three different strategies for training the LLM:
360 full LLM fine-tuning, fine-tuning only the at-
361 tention layers of the LLM, and applying LoRA
362 to the LLM. We focus these strategies on the
363 LLM because it constitutes 98% of the model’s
364 weights, making the full training of the Vision
365 Encoder and bridge connector less computa-
366 tionally intensive. For LoRA, we set $r=128$.
367 Each strategy involves training the model for a
368 total of four epochs.

369 We report the results of these three strate-
370 gies on selected benchmark datasets in Table 2.
371 For each benchmark dataset, we use the official
372 train, validation, and test splits. Evaluation
373 results are presented on the test split, except
374 for the TextVQA dataset, as the test set labels
375 are not available. For each dataset, we use
376 the corresponding evaluation metric commonly
377 employed in the literature. The experiments
378 indicate that LoRA training produces the low-
379 est results across all benchmarks compared to
380 full and attention-only fine-tuning. We also
381 experimented with changing the compression
382 dimension of LoRA to $r=256$ but obtained sim-
383 ilar, near-identical scores.

384 Attention-only and full fine-tuning yield sig-
385 nificantly better results, with each method out-
386 performing the other on different benchmarks.
387 For instance, the attention-only method outper-
388 forms full fine-tuning by 2 points on DocVQA,
389 whereas full fine-tuning scores 3 points higher
390 on SROIE compared to attention-only fine-
391 tuning. Overall, the average results are very
392 close, with attention-only fine-tuning being
393 marginally better than full fine-tuning.

Insight 3

Given our experimental setup, fine-tuning only the attention layers of the LLM is equivalent to full LLM fine-tuning in terms of performance. And both are better than LORA in terms of evaluation scores.

394 In the next section, we will take the best
395 model from our experiments and compare it
396 with other Document Language Models (Doc-
397 cLMs).

Model	DocVQA [ANLS]	CORD V2 [F1]	Info VQA [ANLS]	SROIE [F1]	Chart QA [Rel. EM]	OCR VQA [EM]	RVL-CDIP [Accuracy]	TextVQA [VQA Score]
Attention	0.75	0.76	0.3633	0.76	0.5952	0.722	0.94	0.4644
LORA	0.47	0.463	0.28	0.53	0.442	0.504	0.91	0.2498
Full	0.73	0.76	0.30	0.79	0.5948	0.741	0.94	0.4698

Table 2: Performance of fine-tuning strategy on various benchmarks.

4 DOLMA

4.1 Architecture

We constructed DOLMA by incorporating the insights derived from our previous experiments. The architecture consists of Swin Base as the visual encoder and Vicuna as the LLM decoder, connected via a projection layer serving as the bridge connector. Swin Base is a Swin Transformer with a patch size of 4 and a window size of 10, comprising fewer than 100 million parameters. This model is pretrained with Donut and was trained on 11 million image-text pairs.

Following the approach of LLaVA, we employ a 2-layer MLP as the projection layer, utilizing the GELU activation function between layers, resulting in a total of fewer than 30 million parameters. Vicuna 1.5 serves as the LLM, featuring 7 billion parameters. It is trained by fine-tuning Llama 2 (Touvron et al., 2023b) on user-shared conversations collected from ShareGPT.

4.2 Datasets

We utilized two datasets for pretraining and eight datasets for fine-tuning our model. Specifically, IIT-CDIP (Soboroff, 2022) and our own PDF-generated datasets were employed for pretraining the vision encoder and the projection layer. The other datasets—DocVQA, CORD V2, Infographics VQA, SROIE, Chart QA, OCR VQA, RVL-CDIP, and TextVQA—were used to train the full model, with fine-tuning applied only to the attention layers of the LLM. For each dataset, we constructed a unique instruction prompt to ensure that the model retains its instruction-following capabilities. The prompts can be found in the appendix.

IIT-CDIP (Soboroff, 2022). "CDIP" stands for "Complex Document Information Processing" and "IIT" stands for "Illinois Institute of Technology" who originally built the dataset. The dataset consists of documents from the states' lawsuit against the tobacco industry

in the 1990s. Labels are the text extracted from the dataset using Tesseract. Overall, the datasets consist of around 7 million documents. As the quality of the dataset is crucial for our task we applied some pre processing techniques and removed all the images that had almost no text and ha low quality OCR. The final short-listed number is around 2 million image-text pairs.

PDF-archive. We downloaded an additional 1 million pages of open-source archive PDF documents and extracted the text from them using PyPDF. To obtain these documents, we utilized the arXiv API to download various scientific papers. To get png images from PDFs we set the image zoom equal to 1.8. Given that the archive data is too clean and perfect, it would not adequately represent the everyday scanned document types that Document AI models typically encounter. Therefore, we employed Augraphy (Project) to augment the PDF data by adding random marks, paper folding effects, various colors, and blur effects. Example of an augmented images can be found in the appendix.

CORD V2 (Park et al., 2019). Public benchmark of 1000 receipts images. We follow the official split of 800 - train, 100 - validation and 100 - test samples. The text is fully in Latin characters. Each image may contain different fields with the total number of unique fields amounting to 30. Our data generation process imposes instruction to extract either all or a subset of those fields in a predefined structured format (e.g. JSON) or in unstructured, question-answering manner.

ICDAR SROIE (ICDAR, 2019). A dataset of 1000 whole scanned receipt images. The text is in English characters and each image contains around 4 main fields. The dataset comes with JSON structured annotation intended for KIE task. We separate 347 images for the testing set and utilize the rest in training.

485 *DocVQA* (Mathew et al., 2021b). Document
 486 question answering dataset consisting of 50k
 487 records sourced from the Industry Documents
 488 Library, maintained by the UCSF. The dataset
 489 includes mixture of printed, typewritten and
 490 handwritten documents that are letters, memos,
 491 notes, reports and other types of documents.
 492 We follow the official split with 40k - train, 5k
 493 - validation and 5k - test sets.

494 *RVL-CDIP* (Harley et al., 2015). Relatively
 495 larger dataset of 400k images used for docu-
 496 ment classification task. The dataset includes
 497 documents such as letter, memo, email and oth-
 498 ers. Overall, there are 16 unique classes with
 499 25k images per class. We follow the official split
 500 of 320k - training, 40k - validation and 40k -
 501 testing splits.

502 *Infographic VQA* (Mathew et al., 2021a).
 503 Similar to typical VQA task, task is to answer
 504 questions asked on a given infographic image.
 505 Similar to extractive QA framework popular in
 506 NLP, and the DocVQA dataset, here question-
 507 answers are primarily extractive type. But
 508 there are a small percentage of questions where
 509 answers are not extractive. There are 30 K
 510 questions and 5K Images in the dataset. Im-
 511 ages are collected from the Internet. Questions
 512 and answers are manually annotated.

513 *ChartQA* (Masry et al., 2022). A Benchmark
 514 for Question Answering about Charts with Vi-
 515 sual and Logical Reasoning. The datasets is
 516 split into 30K train, 2K validation and 2.5K
 517 test image-question-answer pairs.

518 *OCR VQA* (Mishra et al., 2019). OCR-VQA
 519 dataset contains 207572 images and associated
 520 question-answer pairs. They provide questions
 521 inquiring about title, author, edition, year and
 522 genre of the book and corresponding ground-
 523 truth answer. This dataset contains approxi-
 524 mately 1 million QA pairs.

525 *Text VQA* (Singh et al., 2019). TextVQA re-
 526 quires models to read and reason about text in
 527 images to answer questions about them. Specif-
 528 ically, models need to incorporate a new modal-
 529 ity of text present in the images and reason
 530 over it to answer TextVQA questions.

531 4.3 Training Details

532 The training process consists of multiple steps,
 533 as outlined in Figure 2.

534 First, we pre-trained the Swin Base and the
 535 MLP projector using the IIT-CDIP and PDF-

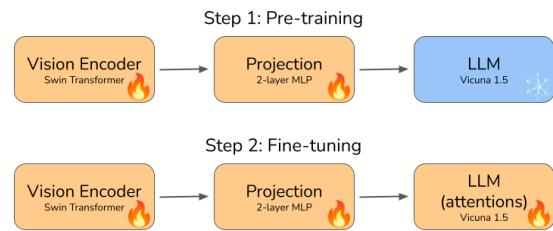


Figure 2: The training pipeline of DOLMA.

536 archive data. The objective of this pretraining
 537 stage is to enable the model to acquire OCR ca-
 538 pabilities and learn to project the visual embed-
 539 dings into the LLM embedding space. During
 540 this stage, the language model is kept frozen.
 541 We pre-trained the model for 1 epoch, with a
 542 batch size of 16 per device, a learning rate of
 543 $2e-4$, and a cosine learning rate scheduler with
 544 3% warmup steps. We used the AdamW opti-
 545 mizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 1e-8$.
 546 Given the importance of image resolution in
 547 Document AI, we increased the resolution of
 548 images to 1280x960 pixels and applied padding
 549 when necessary.

550 Second, we unfroze the entire model and
 551 continued with the fine-tuning process. As
 552 suggested by Insight 3, we fine-tuned only the
 553 attention layers of the LLM. The model was
 554 trained for 10 epochs, with a batch size of 10,
 555 a learning rate of $2e-5$, and a cosine learning rate
 556 scheduler with 3% warmup steps. Similarly, we
 557 used the AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 =$
 558 0.999 , and $\epsilon = 1e-8$. The image resolution was
 559 maintained at 1280x960 pixels.

560 The training was conducted using 8x H100
 561 80GB GPUs¹.

562 4.4 Qualitative analysis and 563 benchmark results

564 We compare DOLMA with models that sat-
 565 isfy the four properties outlined in Table 1.
 566 The evaluation scores are reported in Table 3.
 567 All scores are sourced from their respective
 568 papers. For scores that were not directly avail-
 569 able, we referenced other papers: specifically,
 570 the OCR VQA score of Qwen VL was taken
 571 from the CogAgent paper, and the Infographic-
 572 sQA, ChartQA, and TextVQA scores of Donut

¹Cloud resources were generously provided by AWS

Model	DocVQA [ANLS]	CORD V2 [F1]	Info VQA [ANLS]	Chart QA [Rel. EM]	OCR VQA [EM]	RVL-CDIP [Accuracy]	TextVQA [VQA Score]
DOLMA (ours)	0.75	0.76	0.363	0.595	0.722	0.94	0.464
Donut	0.675	0.841	0.116	0.418	-	0.95	0.435
Qwen-VL	0.651	-	0.354	0.657	0.757	-	0.638
UReader	0.654	-	0.422	0.593	0.411	-	0.576
DocOwl	0.622	-	0.382	0.574	-	-	0.526
CogAgent	0.816	-	0.445	0.684	0.75	-	0.761

Table 3: Comparison of document AI models on various Document AI tasks.

were sourced from the UReader paper. For each benchmark dataset, we used the official train, validation, and test splits. Evaluation results are reported on the test split, except for the TextVQA dataset, where test set labels are unavailable. We employed the evaluation metrics commonly used in the literature for each dataset.

DOLMA outperforms Donut in all tasks except for CORD V2 and RVL-CDIP. The reason for this discrepancy is the relative simplicity of these tasks and the fact that the evaluation used task-specific fine-tuned models, meaning that the models were fine-tuned on a single dataset for many epochs, as described in the Donut paper. Nevertheless, DOLMA managed to outperform Donut in the DocVQA tasks under the same training conditions.

Overall, DOLMA demonstrated performance on par with models such as Qwen VL and DocOwl, even though the vision encoders in these models are 20x and 5x larger in parameter size, respectively. For the DocVQA task, DOLMA outperforms all models except CogAgent. It is important to note that while the other models listed have fewer than 10 billion parameters, CogAgent has 17 billion parameters. As the scores illustrate, model size has a significant impact on performance in our case.

5 Conclusion and Future Work

In this paper, we conducted experiments to understand the requirements for building a Vision-Language model for Document AI tasks. Our findings highlight the effectiveness of different model architectures, model sizes, pretraining, and fine-tuning strategies. Based on these insights, we introduced DOLMA, an OCR-free, instruction-following vision-language model that can be utilized for various Document AI tasks. We demonstrated that DOLMA can per-

form on par with larger VLM models, despite being trained on fewer data samples and with fewer resources.

In future research, we plan to investigate the possibility of scaling DOLMA to handle multilingual and multi-page documents.

6 Limitations

While DOLMA demonstrates promising results in various Document AI tasks, several limitations must be acknowledged:

1. *Data Diversity*: Although we utilized a substantial amount of data for pretraining and fine-tuning, the datasets may not fully capture the diversity of real-world documents. This could limit the model’s generalizability to unseen document types and formats.

2. *Model Size*: Despite DOLMA’s competitive performance with a moderate parameter count of 7 billion, it remains computationally intensive. This may pose challenges for deployment in resource-constrained environments.

3. *OCR Capabilities*: While DOLMA is designed to be OCR-free, its performance in extracting text from highly complex or degraded documents may still lag behind specialized OCR systems. Further improvements are needed to enhance its robustness in such scenarios.

4. *Multilingual and Multi-page Documents*: Our current experiments focus primarily on single-page, monolingual documents. The model’s effectiveness in handling multilingual and multi-page documents remains unexplored and warrants further investigation.

5. *Training Costs*: Although we aimed to minimize training costs, the process still requires significant computational resources, particularly for fine-tuning. This could be a barrier for smaller research groups or organizations

with limited access to high-performance computing resources.

6. *Evaluation Metrics*: The evaluation metrics used in our experiments are standard in the literature, but they may not fully capture the nuanced performance of the model in practical applications. Future work should consider more comprehensive evaluation frameworks.

7. *Ethical Considerations*: As with any AI model, there are ethical considerations related to data privacy and potential biases in the training data. These issues need to be addressed to ensure the responsible deployment of DOLMA.

By acknowledging these limitations, we aim to provide a balanced view of our work and highlight areas for future research and improvement.

References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Qin Cai, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Yen-Chun Chen, Yi-Ling Chen, Parul Chopra, Xiyang Dai, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Victor Fragoso, Dan Iter, Mei Gao, Min Gao, Jianfeng Gao, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Ce Liu, Mengchen Liu, Weishung Liu, Eric Lin, Zeqi Lin, Chong Luo, Piyush Madan, Matt Mazzola, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Xin Wang, Lijuan Wang, Chunyu Wang, Yu Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Haiping Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Sonali Yadav, Fan Yang, Jianwei Yang, Ziyi Yang, Yifan Yang, Donghan Yu, Lu Yuan, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. *Phi-3 technical report: A highly capable language model locally on your phone*.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. *Palm: Scaling language modeling with pathways*. *Preprint*, arXiv:2204.02311.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. *Bert: Pre-training of deep bidirectional transformers for language understanding*. *Preprint*, arXiv:1810.04805.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2020. *An image is worth 16x16 words: Transformers for image recognition at scale*.

Adam W Harley, Alex Ufkes, and Konstantinos G Derpanis. 2015. Evaluation of deep convolutional nets for document image classification and retrieval. In *International Conference on Document Analysis and Recognition (ICDAR)*.

Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. 2023. *Cogagent: A visual language model for gui agents*.

Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. *Layoutlmv3: Pre-training for document ai with unified text and image masking*. *Preprint*, arXiv:2204.08387.

ICDAR. 2019. Sroie. Available at: <https://rrc.cvc.uab.es/?ch=13&com=introduction>.

768	Geewook Kim, Teakgyu Hong, Moonbin Yim,	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang,	820
769	JeongYeon Nam, Jinyoung Park, Jinyeong Yim,	Yann Dubois, Xuechen Li, Carlos Guestrin,	821
770	Wonseok Hwang, Sangdoo Yun, Dongyoon Han,	Percy Liang, and Tatsunori B. Hashimoto.	822
771	and Seunghyun Park. 2022. Ocr-free document	2023. Stanford alpaca: An instruction-	823
772	understanding transformer. In <i>European Confer-</i>	following llama model. https://github.com/	824
773	<i>ence on Computer Vision (ECCV)</i> .	tatsu-lab/stanford_alpaca .	825
774	Hugo Laurençon, Léo Tronchon, Matthieu Cord,	Meta LLaMA Team. 2024. Introducing meta llama	826
775	and Victor Sanh. 2024. What matters when	3: The most capable openly available llm to date.	827
776	building vision-language models?		
777	Mike Lewis, Yinhan Liu, Naman Goyal, Mar-	Hugo Touvron, Thibaut Lavril, Gautier Izacard,	828
778	jan Ghazvininejad, Abdelrahman Mohamed,	Xavier Martinet, Marie-Anne Lachaux, Timo-	829
779	Omer Levy, Ves Stoyanov, and Luke Zettle-	thée Lacroix, Baptiste Rozière, Naman Goyal,	830
780	moyer. 2019. Bart: Denoising sequence-to-	Eric Hambro, Faisal Azhar, Aurelien Rodriguez,	831
781	sequence pre-training for natural language gener-	Armand Joulin, Edouard Grave, and Guil-	832
782	ation, translation, and comprehension. <i>Preprint,</i>	laume Lample. 2023a. Llama: Open and ef-	833
783	arXiv:1910.13461 .	ficient foundation language models. <i>Preprint,</i>	834
		arXiv:2302.13971 .	835
784	Haotian Liu, Chunyuan Li, Qingyang Wu, and	Hugo Touvron, Louis Martin, Kevin Stone, Pe-	836
785	Yong Jae Lee. 2023. Visual instruction tuning.	ter Albert, Amjad Almahairi, Yasmine Babaei,	837
786	In <i>NeurIPS</i> .	Nikolay Bashlykov, Soumya Batra, Prajjwal	838
787	Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan	Bhargava, Shruti Bhosale, Dan Bikel, Lukas	839
788	Wei, Zheng Zhang, Stephen Lin, and Baining	Blecher, Cristian Canton Ferrer, Moya Chen,	840
789	Guo. 2021. Swin transformer: Hierarchical vision	Guillem Cucurull, David Esiobu, Jude Fernan-	841
790	transformer using shifted windows. <i>Preprint,</i>	des, Jeremy Fu, Wenyin Fu, Brian Fuller, Cyn-	842
791	arXiv:2103.14030 .	thia Gao, Vedanuj Goswami, Naman Goyal, An-	843
792	Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq	thony Hartshorn, Saghar Hosseini, Rui Hou,	844
793	Joty, and Enamul Hoque. 2022. Chartqa: A	Hakan Inan, Marcin Kardas, Viktor Kerkez,	845
794	benchmark for question answering about charts	Madian Khabsa, Isabel Kloumann, Artem Ko-	846
795	with visual and logical reasoning.	renev, Punit Singh Koura, Marie-Anne Lachaux,	847
796	Minesh Mathew, Viraj Bagal, Rubèn Pérez Tito,	Thibaut Lavril, Jenya Lee, Diana Liskovich, Ying-	848
797	Dimosthenis Karatzas, Ernest Valveny, and C. V	hai Lu, Yuning Mao, Xavier Martinet, Todor	849
798	Jawahar. 2021a. Infographicvqa .	Mihaylov, Pushkar Mishra, Igor Molybog, Yixin	850
799	Minesh Mathew, Dimosthenis Karatzas, and C. V.	Nie, Andrew Poulton, Jeremy Reizenstein, Rashi	851
800	Jawahar. 2021b. Docvqa: A dataset for vqa on	Rungta, Kalyan Saladi, Alan Schelten, Ruan	852
801	document images. <i>Preprint,</i> arXiv:2007.00398 .	Silva, Eric Michael Smith, Ranjan Subramanian,	853
802	Anand Mishra, Shashank Shekhar, Ajeet Kumar	Xiaoqing Ellen Tan, Binh Tang, Ross Taylor,	854
803	Singh, and Anirban Chakraborty. 2019. Ocr-	Adina Williams, Jian Xiang Kuan, Puxin Xu,	855
804	vqa: Visual question answering by reading text	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela	856
805	in images. In <i>ICDAR</i> .	Fan, Melanie Kambadur, Sharan Narang, Aure-	857
806	Seunghyun Park, Seung Shin, Bado Lee, Junyeop	lien Rodriguez, Robert Stojnic, Sergey Edunov,	858
807	Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk	and Thomas Scialom. 2023b. Llama 2: Open	859
808	Lee. 2019. Cord: A consolidated receipt dataset	foundation and fine-tuned chat models. <i>Preprint,</i>	860
809	for post-ocr parsing.	arXiv:2307.09288 .	861
810	The Augraphy Project. Augraphy: an augmenta-	Dongsheng Wang, Natraj Raman, Mathieu Sibue,	862
811	tion pipeline for rendering synthetic paper print-	Zhiqiang Ma, Petr Babkin, Simerjot Kaur, Yu-	863
812	ing, faxing, scanning and copy machine processes.	long Pei, Armineh Nourbakhsh, and Xiaomo Liu.	864
813	Amanpreet Singh, Vivek Natarajan, Meet Shah,	2023a. Docllm: A layout-aware generative lan-	865
814	Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,	guage model for multimodal document under-	866
815	and Marcus Rohrbach. 2019. Towards vqa mod-	standing. <i>Preprint,</i> arXiv:2401.00908 .	867
816	els that can read.		
817	Ian Soboroff. 2022. Complex document information	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi	868
818	processing (cdip) dataset, national institute of	Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,	869
819	standards and technology.	Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu,	870
		Juanzi Li, Yuxiao Dong, Ming Ding, and Jie	871
		Tang. 2023b. Cogvlm: Visual expert for pre-	872
		trained language models.	873
		BigScience Workshop, :, Teven Le Scao, Angela	874
		Fan, Christopher Akiki, Ellie Pavlick, Suzana	875
		Ilić, Daniel Hesslow, Roman Castagné, Alexan-	876
		dra Sasha Luccioni, François Yvon, Matthias	877
		Gallé, Jonathan Tow, Alexander M. Rush, Stella	878

879	Biderman, Albert Webson, Pawan Sasanka Amanamanchi, Thomas Wang, Benoît Sagot,	943
880	Niklas Muennighoff, Albert Villanova del Moral,	944
881	Olatunji Ruwase, Rachel Bawden, Stas Bek-	945
882	man, Angelina McMillan-Major, Iz Beltagy, Huu	946
883	Nguyen, Lucile Saulnier, Samson Tan, Pedro Or-	947
884	tiz Suarez, Victor Sanh, Hugo Laurençon, Yacine	948
885	Jernite, Julien Launay, Margaret Mitchell, Colin	949
886	Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa,	950
887	Alham Fikri Aji, Amit Alfassy, Anna Rogers,	951
888	Ariel Kreisberg Nitzav, Canwen Xu, Chenghao	952
889	Mou, Chris Emezue, Christopher Klamm, Colin	953
890	Leong, Daniel van Strien, David Ifeoluwa Ade-	954
891	lani, Dragomir Radev, Eduardo González Pon-	955
892	ferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar	956
893	Natan, Francesco De Toni, Gérard Dupont, Ger-	957
894	mán Kruszewski, Giada Pistilli, Hady Elshahar,	958
895	Hamza Benyamina, Hieu Tran, Ian Yu, Idris Ab-	959
896	dulmumin, Isaac Johnson, Itziar Gonzalez-Dios,	960
897	Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian	961
898	Zhu, Jonathan Chang, Jörg Frohberg, Joseph To-	962
899	bong, Joydeep Bhattacharjee, Khalid Almubarak,	963
900	Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon	964
901	Weber, Long Phan, Loubna Ben allal, Ludovic	965
902	Tanguy, Manan Dey, Manuel Romero Muñoz,	966
903	Maraim Masoud, María Grandury, Mario Šaško,	967
904	Max Huang, Maximin Coavoux, Mayank Singh,	968
905	Mike Tian-Jian Jiang, Minh Chien Vu, Mo-	969
906	hammad A. Jauhar, Mustafa Ghaleb, Nishant	970
907	Subramani, Nora Kassner, Nurulaqilla Khamis,	971
908	Olivier Nguyen, Omar Espejel, Ona de Gibert,	972
909	Paulo Villegas, Peter Henderson, Pierre Colombo,	973
910	Priscilla Amuok, Quentin Lhoest, Rheza Harli-	974
911	man, Rishi Bommasani, Roberto Luis López, Rui	975
912	Ribeiro, Salomey Osei, Sampo Pyysalo, Sebas-	976
913	tian Nagel, Shamik Bose, Shamsuddeen Hassan	977
914	Muhammad, Shanya Sharma, Shayne Longpre,	978
915	Somaieh Nikpoor, Stanislav Silberberg, Suhas	979
916	Pai, Sydney Zink, Tiago Timponi Torrent, Timo	980
917	Schick, Tristan Thrush, Valentin Danchev, Vas-	981
918	silina Nikoulina, Veronika Laippala, Violette Lep-	982
919	ercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat,	983
920	Arun Raja, Benjamin Heinzerling, Chenglei Si,	984
921	Davut Emre Taşar, Elizabeth Salesky, Sabrina J.	985
922	Mielke, Wilson Y. Lee, Abheesht Sharma, An-	986
923	drea Santilli, Antoine Chaffin, Arnaud Stiegler,	987
924	Debajyoti Datta, Eliza Szczechla, Gunjan Chh-	988
925	ablani, Han Wang, Harshit Pandey, Hendrik	989
926	Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao,	990
927	Lintang Sutawika, M Saiful Bari, Maged S. Al-	991
928	shaibani, Matteo Manica, Nihal Nayak, Ryan	992
929	Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-	993
930	David, Stephen H. Bach, Taewoon Kim, Tali	994
931	Bers, Thibault Fevry, Trishala Neeraj, Urmish	995
932	Thakker, Vikas Raunak, Xiangru Tang, Zheng-	996
933	Xin Yong, Zhiqing Sun, Shaked Brody, Yallow	997
934	Uri, Hadar Tojarieh, Adam Roberts, Hyung Won	998
935	Chung, Jaesung Tae, Jason Phang, Ofir Press,	999
936	Conglong Li, Deepak Narayanan, Hatim Bour-	1000
937	foune, Jared Casper, Jeff Rasley, Max Ryabinin,	1001
938	Mayank Mishra, Minjia Zhang, Mohammad	1002
939	Shoeybi, Myriam Peyrounette, Nicolas Patry,	1003
940	Nouamane Tazi, Omar Sanseviero, Patrick von	1004
941	Platen, Pierre Cornette, Pierre François Laval-	1005
942	lée, Rémi Lacroix, Samyam Rajbhandari, San-	1006
	chit Gandhi, Shaden Smith, Stéphane Requena,	
	Suraj Patil, Tim Dettmers, Ahmed Baruwa,	
	Amanpreet Singh, Anastasia Cheveleva, Anne-	
	Laure Ligozat, Arjun Subramonian, Aurélie	
	Névéol, Charles Lovering, Dan Garrette, Deepak	
	Tunuguntla, Ehud Reiter, Ekaterina Taktasheva,	
	Ekaterina Voloshina, Eli Bogdanov, Genta Indra	
	Winata, Hailey Schoelkopf, Jan-Christoph Kalo,	
	Jekaterina Novikova, Jessica Zosa Forde, Jor-	
	dan Clive, Jungo Kasai, Ken Kawamura, Liam	
	Hazan, Marine Carpuat, Miruna Clinciu, Na-	
	joung Kim, Newton Cheng, Oleg Serikov, Omer	
	Antverg, Oskar van der Wal, Rui Zhang, Ruochen	
	Zhang, Sebastian Gehrmann, Shachar Mirkin,	
	Shani Pais, Tatiana Shavrina, Thomas Scialom,	
	Tian Yun, Tomasz Limisiewicz, Verena Rieser,	
	Vitaly Protasov, Vladislav Mikhailov, Yada Pruk-	
	sachatkun, Yonatan Belinkov, Zachary Bam-	
	berger, Zdeněk Kasner, Alice Rueda, Amanda	
	Pestana, Amir Feizpour, Ammar Khan, Amy	
	Faranak, Ana Santos, Anthony Hevia, Antig-	
	ona Undreaaj, Arash Aghagol, Arezoo Abdol-	
	lahi, Aycha Tammour, Azadeh HajiHosseini,	
	Bahareh Behroozi, Benjamin Ajibade, Bharat	
	Saxena, Carlos Muñoz Ferrandis, Daniel Mc-	
	Duff, Danish Contractor, David Lansky, Davis	
	David, Douwe Kiela, Duong A. Nguyen, Edward	
	Tan, Emi Baylor, Ezinwanne Ozoani, Fatima	
	Mirza, Frankline Ononiwu, Habib Rezanejad,	
	Hessie Jones, Indrani Bhattacharya, Irene So-	
	laiman, Irina Sedenko, Isar Nejadgholi, Jesse	
	Passmore, Josh Seltzer, Julio Bonis Sanz, Livia	
	Dutra, Mairon Samagaio, Maraim Elbadri, Mar-	
	got Mieskes, Marissa Gerchick, Martha Akin-	
	lolu, Michael McKenna, Mike Qiu, Muhammed	
	Ghauri, Mykola Burynok, Nafis Abrar, Nazneen	
	Rajani, Nour Elkott, Nour Fahmy, Olanre-	
	waju Samuel, Ran An, Rasmus Kromann, Ryan	
	Hao, Samira Alizadeh, Sarmad Shubber, Silas	
	Wang, Sourav Roy, Sylvain Vignier, Thanh	
	Le, Tobi Oyebade, Trieu Le, Yoyo Yang, Zach	
	Nguyen, Abhinav Ramesh Kashyap, Alfredo	
	Palasciano, Alison Callahan, Anima Shukla, An-	
	tonio Miranda-Escalada, Ayush Singh, Benjamin	
	Beilharz, Bo Wang, Caio Brito, Chenxi Zhou,	
	Chirag Jain, Chuxin Xu, Clémentine Fourier,	
	Daniel León Perrián, Daniel Molano, Dian Yu,	
	Enrique Manjavacas, Fabio Barth, Florian Fuhri-	
	mann, Gabriel Altay, Giyaseddin Bayrak, Gully	
	Burns, Helena U. Vrabec, Imane Bello, Ishani	
	Dash, Jihyun Kang, John Giorgi, Jonas Golde,	
	Jose David Posada, Karthik Rangasai Sivaraman,	
	Lokesh Bulchandani, Lu Liu, Luisa Shinzato,	
	Madeleine Hahn de Bykhovetz, Maiko Takeuchi,	
	Marc Pàmies, Maria A Castillo, Marianna	
	Nezhurina, Mario Sängler, Matthias Samwald,	
	Michael Cullan, Michael Weinberg, Michiel De	
	Wolf, Mina Mihaljcic, Minna Liu, Moritz Frei-	
	dank, Myungsun Kang, Natasha Seelam, Nathan	
	Dahlberg, Nicholas Michio Broad, Nikolaus	
	Muellner, Pascale Fung, Patrick Haller, Ramya	
	Chandrasekhar, Renata Eisenberg, Robert Mar-	
	tin, Rodrigo Canalli, Rosaline Su, Ruisi Su,	

1007 Samuel Cahyawijaya, Samuele Garda, Shlok S
1008 Deshmukh, Shubhanshu Mishra, Sid Kiblawi,
1009 Simon Ott, Sinee Sang-aaronsiri, Srishti Ku-
1010 mar, Stefan Schweter, Sushil Bharati, Tanmay
1011 Laud, Théo Gigant, Tomoya Kainuma, Wojciech
1012 Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash
1013 Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu,
1014 Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras,
1015 Younes Belkada, and Thomas Wolf. 2023. [Bloom:
1016 A 176b-parameter open-access multilingual lan-
1017 guage model](#). *Preprint*, arXiv:2211.05100.

1018 Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye,
1019 Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu,
1020 Chenliang Li, Junfeng Tian, Qian Qi, Ji Zhang,
1021 and Fei Huang. 2023a. [mplug-docowl: Mod-
1022 ularized multimodal large language model for
1023 document understanding](#).

1024 Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye,
1025 Ming Yan, Guohai Xu, Chenliang Li, Junfeng
1026 Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He,
1027 Xin Alex Lin, and Fei Huang. 2023b. [Ure-
1028 ader: Universal ocr-free visually-situated lan-
1029 guage understanding with multimodal large lan-
1030 guage model](#).

1031 Susan Zhang, Stephen Roller, Naman Goyal, Mikel
1032 Artetxe, Moya Chen, Shuohui Chen, Christo-
1033 pher Dewan, Mona Diab, Xian Li, Xi Victoria
1034 Lin, Todor Mihaylov, Myle Ott, Sam Shleifer,
1035 Kurt Shuster, Daniel Simig, Punit Singh Koura,
1036 Anjali Sridhar, Tianlu Wang, and Luke Zettle-
1037 moyer. 2022. [Opt: Open pre-trained transformer
1038 language models](#). *Preprint*, arXiv:2205.01068.

1039 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng,
1040 Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
1041 Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing,
1042 Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
1043 2023. [Judging llm-as-a-judge with mt-bench and
1044 chatbot arena](#). *Preprint*, arXiv:2306.05685.

A Prompts for the pertaining stage 1045

System 1046

You are a helpful language and vision assistant. You are able to understand the visual content that the user provides. "

User 1047

Extract all the text from the document.

B Prompts for the fine-tuning stage 1048

Prompt template for CORD [task: KIE] 1049

Please read the text in this image and return the information in JSON format.

1050

The nested JSON should have the following keys: menu, void_menu, subtotal, total. Each key has subkeys as listed below (with descriptions in brackets):

menu:

- nm (name of menu)
- num (identification # of menu)
- unitprice (unit price of menu)
- menu.cnt (quantity of menu)
- discountprice (discounted price of menu)
- price (total price of menu)
- itemsubtotal (price of each menu after discount applied)
- vatyn (whether the price includes tax or not)
- etc (others)
- sub_nm (name of submenu)
- sub_unitprice (unit price of submenu)
- sub_cnt (quantity of submenu)
- sub_price (total price of submenu)
- sub_etc (others)

void_menu:

- nm (name of menu)
- price (total price of menu)

subtotal:

- subtotal_price (subtotal price)
- discount_price (discounted price in total)
- service_price (service charge)
- othersvc_price (added charge other than service charge)
- tax_price (tax amount)
- etc (others)

total:

- total_price (total price)
- total_etc (others)
- cashprice (amount of price paid in cash)
- changeprice (amount of change in cash)
- creditcardprice (amount of price paid in credit/debit card)
- emoneyprice (amount of price paid in emoney, point)
- menutype_cnt (total count of type of menu)
- menuqty_cnt (total count of quantity)

1051

1052

Prompt template for SROIE [task: KIE]

Please read the text in this image and return the information in JSON format.
The JSON should have the following keys: company, date, address, total.

1053

Prompt template for DocVQA [task: VQA]

1054

"Please read the text in this image and answer to the question: {question}\n

<image>"	1055
Prompt template for InfographicVQA [task: VQA]	1056
"Given this infographic image, {question}\n<image>"	
Prompt template for TextVQA [task: VQA]	1057
"Given the image, {question}\n<image>"	
Prompt template for OCRVQA [task: VQA]	1058
"Here is an image of a book cover, {question}\n<image>"	
Prompt template for RVL-CDIP [task: classification]	1059
'Please classify the given image to one of the following classes: ["letter", "memo", "email", "filefolder", "form", "handwritten", "invoice", "advertisement", "budget", "news article", "presentation", "scientific publication", "questionnaire", "resume", "scientific report", "specification"].'	
Prompt template for CartQA [task: VQA]	1060
Given this image of a chart, {question}\n<image>"	

C Samples from our PDF-arxiv dataset

4 JOURNAL BROWSE

in particular that $R_1(A) = R_1(B)$. Examining the pre-orderings of $\text{abd}(A)$ (3.3) in our example in the last section, we see that each P_i has the same $R_i(\tau)$ for each factor τ of $\text{abd}(A)$. It will be helpful to determine the exact structure of the $R_i(\tau)$ in the shelling obtained by listing the factors of $\text{abd}(A)$ in the reverse lexicographical order.

Let τ be a face of A , and O an ordering of V . Then let $\text{full}(\tau) = \{i \mid \tau(i) = |P_i| - 1\}$, and for $i \in \text{full}(\tau)$ let $\text{omit}(\tau, i)$ be the element of P_i not in τ (the notation is meant to suggest that $\text{full}(\tau)$ reflects the indices of the sets P_i such that $\tau \cap P_i$ is full in the sense that no further elements of P_i could be added without leaving A , and $\text{omit}(\tau, i)$ is the element of P_i missing from τ). Let $\text{sort}(\tau)$ be the first element of $V \setminus \text{full}(\tau)$ (omit τ, i) not appearing in τ (with respect to order O) if such an element exists, otherwise set $\text{sort}(\tau) = \infty$. Let $\tau_{\text{succ}} = \{\tau \in \mathcal{F} \mid \tau \supset \tau_0\}$, and $R_i(\tau) = \{p \mid p \in P_i \text{ and } p \supset \text{omit}(\tau, i) \text{ for some } \tau \in \text{full}(\tau)\}$. Finally, let $R_i(\tau) = \tau_{\text{succ}} \cap R_i(\tau)$.

Example 4.1. Let $A = (A_1, A_2, A_3)$, with vertex ordering O as above.

FIGURE 2. Vertex set of $A(1,5,4,3)$ with ordering O

Consider the face $\tau = (v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8)$. Then $\text{full}(\tau) = \{2\}$, $\text{omit}(\tau, 2) = \{v_2\}$, $O(\tau) = (v_1, v_3, v_4, v_5, v_6, v_7, v_8)$, $\text{sort}(\tau) = v_2$, and $\tau_{\text{succ}} = \{(v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8), (v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8, v_9)\}$. So $R_2(\tau) = \{(v_1, v_2, v_3, v_4)\}$.

FIGURE 3. τ

FIGURE 4. $O(\tau)$

FIGURE 5. $\tau_{\text{succ}}(\tau)$

Now, if τ is a face of $\text{abd}(A)$ and γ is a face of \mathcal{T} , $\tau \cap \gamma = \emptyset$ for some $g \in \tau$. Then γ appears as a face of a facet occurring before τ in the reverse lexicographical order determined by the ordering O on the vertices, if and only if g is a face of $\text{abd}(A)$. For some property P , we say τ is *regularly shelled* if the reverse lexicographically first face of $\text{abd}(A)$ containing $R_2(\tau)$, so if $R_2(\tau) \subseteq \gamma$, γ occurs as an earlier facet. On the other hand, if there is $g \in R_2(\tau)$ such that $g \notin \gamma$, either $g \in O(\tau)$, in which case γ contains τ , or $g \in P_i$ is a reverse lexicographical order face of $\text{abd}(A)$ containing τ and g is not in γ . In this case, γ occurs as an earlier facet of $\text{abd}(A)$ containing τ . Thus, if $\tau \in \mathcal{T}$, $\tau_{\text{succ}}(\tau)$ is the reverse lexicographical order of the facets of $\text{abd}(A)$, $\mathcal{T} \setminus \{\tau_0, \tau\} = \{\gamma \in \mathcal{T} \mid R_2(\tau) \subseteq \gamma\}$. Our inductive shellings will share this structure.

arXiv:0812.4554v2 [physics.chem-ph] 24 Dec 2008

Laser-induced atomic fragment fluorescence spectroscopy: A facile technique for molecular spectroscopy of spin-forbidden states

Qun Zhang,^{1,2*} Yang Chen,^{3,*} and Mark Kapp^{1,†}

¹Hebei National Laboratory for Physical Sciences at the Microscale, University of Science and Technology of China, Hefei, Anhui 230026, People's Republic of China

²Hebei, Anhui 230026, People's Republic of China

³Department of Physics, Ben-Gurion University, Be'er Sheva 84105, Israel

(Dated: November 5, 2008)

Spectra of spin-forbidden and spin-allowed transitions in the mixed $^1\Pi_u - ^3\Sigma_u^-$ state of N_2 are measured separately by two-photon excitation using a single tunable dye laser. The two-photon excitation produces N_2^+ by photoionization, which is easily and rapidly detected by atomic fluorescence. At low laser power, only the $^3\Sigma_u^-$ state is excited, completely free of triplet excitation. At high laser power, photoionization on the $^1\Pi_u$ triplet state intermediate becomes non-negligible, effectively "switching" the observations from singlet spectroscopy to triplet spectroscopy with only minor apparatus changes. This technique of perturbation-assisted laser-induced atomic fragment fluorescence (LIAFF) may therefore be especially useful as a general vehicle for investigating perturbation-induced physics pertinent to the spin-forbidden states, as well as for studying allowed and forbidden states of other molecules.

PACS numbers: 33.20.+v, 33.30.+z, 33.50.+z, 33.50.Wi, 33.50.Wz

Keywords: fluorescence spectroscopy, spin-forbidden molecular transition, perturbation, photoionization, atomic molecule, laser-induced fluorescence, atomic fragments

I. INTRODUCTION

The well-known $^1\Pi_u - ^3\Sigma_u^-$ intersystem crossing-based system of the N_2 molecule has long served as a prototype for developing new spectroscopic techniques [1, 2, 3, 4]. In particular, the recently developed laser-induced atomic fragment fluorescence (LIAFF) spectroscopy [5, 6] has been used as a prototype for developing new spectroscopic techniques [7, 8, 9, 10]. LIAFF is a technique that uses optical triplet resonance (PMOTR) spectroscopy [20] and triplet energy transfer to generate atomic fragment fluorescence [5]. To date, the three Π components of the $^1\Pi_u$ state have been fully resolved only by the continuous wave (cw) PMOTR technique which demands three excitation beams, all of which must be "locked" on the transition [20].

We report here a newly developed and simple technique of perturbation-assisted laser-induced atomic fragment fluorescence (LIAFF) spectroscopy that is suitable for detailed studies of spin-forbidden electronic states, and in particular, we have used the perturbation-assisted LIAFF technique to fully resolve all three Π components of the $^1\Pi_u$ state in N_2 . Based on a single conventional pulsed dye laser, the LIAFF technique does not attain the high resolution of the sophisticated techniques mentioned above instead we focus here on demonstrating the robustness and experimental utility for easily obtaining information on spin-forbidden states. We expect that the technique will be widely applicable, as shown by our example studies herein on the $^1\Pi_u - ^3\Sigma_u^-$ intersystem crossing band system of the N_2 molecule.

II. EXPERIMENTAL

*Corresponding authors. Electronic address: qunzhang@ustc.edu.cn

†Corresponding authors. Electronic address: markkapp@ustc.edu.cn

The perturbation-assisted LIAFF technique is illustrated schematically in Fig. 1 by referring to the four

available routes of fixation (e.g. [10, 2, 5, 23]). It is a delicate, a delicate problem. However, even in other related gels, plasticity in the areas where our existing transport systems (e.g. photonic fiber, polymer-based) are not able to provide the required level of mobility, it is still a challenge. In addition, showing the protein solubility in D-gels is still a challenge, and needs to be improved.

3. Improvement of the standard 2D electrophoresis technique

The constraints imposed by the IEF separation are rather strict: low sample amount, no modification of the gel, and low amounts of toxic detergent. In this case, we recently discovered two parameters which the experimental setup comply with: a) increase the solubility of the proteins in the gel, and b) concentrate the proteins in the narrow lanes, showing the protein solubility in D-gels is still a challenge, and needs to be improved.

The situation was improved by the introduction of the so-called "microgel" droplets in a different area [1]. While this initial effort was largely qualitative and related to membrane proteins, the use of water-in-oil emulsions (W/O emulsions) in combination with a microfluidic chip were considered as promising alternatives. As a result, the use of emulsions dramatically widened the solubility range of even hydrophobic proteins, as shown by the work carried out by Letourneau et al. (Fig. 2) [1]. While the combination of emulsions and microfluidics is not particularly a combination of microfluidics and TIRF, X100s field of view and focal length, the combination of both is a promising alternative.

This electrophoresis experiment was carried out using a 2D gel electrophoresis apparatus, and the data were analyzed using the software [21]. For the case of detection, a 488 nm laser line was used as a light source [21]. The visualization of the bands in the membrane plot was done using a 488 nm (PRL) with a seven-transverse mode (TM7) fiber optic probe beam, placed at various transverse positions, including below the protein solution. The same mode (TM7) which was used for the visualization so far [21] or as a probe for protein case [21] [23] or for individual transmitters [25].

FIGURE 4. Comparison of the exact solution and the solution for $N=4$ for randomly generated β_0 and λ_0 .

The LIAFF technique is also been observed for the spin-forbidden model with a damped Jaynes-Cummings Hamiltonian [26]. At the point where the Kramers gap becomes non-invertible, the LIAFF solution deviates from the exact solution (see Fig. 9). We verified that both β_0 and λ_0 were not the same.

III. CONCLUSION

In this subsection, we compare the exact solution to TCLA, NZ (and the solution of the optically bistable system) with the solution of the exact solution. The exact solution is $\beta_0 = 0$ and $\lambda_0 = 0$ when $\beta_0 = 1$. Here we observe that while the short-time behavior of the exact solution is approximated well by the optically bistable system, the long-time behavior is approximated well only by PM.