

VISA: Retrieval Augmented Generation with Visual Source Attribution

Anonymous ACL submission

Abstract

001 Generation with source attribution is impor- 042
002 tant for enhancing the verifiability of retrieval- 043
003 augmented generation (RAG) systems. How- 044
004 ever, existing approaches in RAG primarily 045
005 link generated content to document-level ref- 046
006 erences, making it challenging for users to 047
007 locate evidence among multiple content-rich 048
008 retrieved documents. To address this chal- 049
009 lenge, we propose *Retrieval-Augmented Gener- 050*
010 *ation with Visual Source Attribution* (VISA), a 051
011 novel approach that combines answer genera- 052
012 tion with visual source attribution. Leveraging 053
013 large vision-language models (VLMs), VISA 054
014 identifies the evidence and highlights the ex- 055
015 act regions that support the generated answers 056
016 with bounding boxes in the retrieved document 057
017 screenshots. To evaluate its effectiveness, we 058
018 curated two datasets: Wiki-VISA, based on 059
019 crawled Wikipedia webpage screenshots, and 060
020 Paper-VISA, derived from PubLayNet and tai- 061
021 lored to the medical domain. Experimental re- 062
022 sults demonstrate the effectiveness of VISA for 063
023 visual source attribution on documents’ origi- 064
024 nal look, as well as highlighting the challenges 065
025 for improvement. Code, data, and model check- 066
026 points will be released. 067

027 1 Introduction 028

028 Retrieval-augmented generation (RAG) has be- 029
029 come a key technique for enhancing the reliabil- 030
030 ity in information-seeking processes (Gao et al., 031
031 2024). Traditional RAG pipeline directly gen- 032
032 erates an answer to a user query from retrieved 033
033 candidate documents (Chen et al., 2017; Lewis 034
034 et al., 2020). Yet, it is hard for users to verify 035
035 the sources and appropriately trust generated an- 036
036 swers, given that models could produce halluci- 037
037 nated content (Min et al., 2023; Malaviya et al., 038
038 2024). Recent works have introduced the genera- 039
039 tion with citation paradigm (Gao et al., 2023; Ye 040
040 et al., 2024), prompting the model to not only gen- 041
041 erate answers but also directly cite the identifiers

of the source documents. Such source attribution 042
approaches make it possible for users to check the 043
reliability of the outputs (Asai et al., 2024). 044

045 However, text-based generation with source attri- 046
046 bution faces several issues: First, citing the source 047
047 at the document level could impose a heavy cogni- 048
048 tive burden on users (Foster, 1979; Sweller, 2011), 049
049 where users often struggle to locate the core evi- 050
050 dence at the section or passage level within the 051
051 dense and multi-page document. Despite such 052
052 granularity mismatch could be addressed through 053
053 passage-citation-based generation methods — link- 054
054 ing answers to specific text chunks, it requires non- 055
055 trivial extra engineering efforts to match the chunk 056
056 in the document source. Moreover, visually high- 057
057 lighting text chunks in the source document is more 058
058 intuitive for users, but it remains challenging as it 059
059 requires control over document rendering, which is 060
060 not always accessible, such as in PDF scenarios. 061

061 Inspired by the recent document screenshot em- 062
062 bedding retrieval paradigm — dropping the docu- 063
063 ment processing module and directly using VLM 064
064 to preserve the content integrity and encoding docu- 065
065 ment screenshots for retrieval (Ma et al., 2024), 066
066 we ask whether source attribution can also be in- 067
067 tegrated into such a unified visual paradigm to es- 068
068 tablish a fully visual, end-to-end verifiable RAG 069
069 pipeline that is both user-friendly and effective? 070

070 To this end, we propose *Retrieval Augmented 071*
071 *Generation with Visual Source Attribution* (VISA). 072
072 In our approach, a large vision-language model 073
073 (VLM) processes single or multiple retrieved docu- 074
074 ment images and not only generates an answer to 075
075 the user query but also returns the bounding box of 076
076 the relevant region within the evidence document. 077
077 As illustrated in Figure 1, this method enables di- 078
078 rect attribution by visually pinpointing the exact 079
079 position within the document, allowing users to 080
080 quickly check the supporting evidence within the 081
081 original context for the generated answer. VLMs 082
082 are not restricted by document format or rendering,

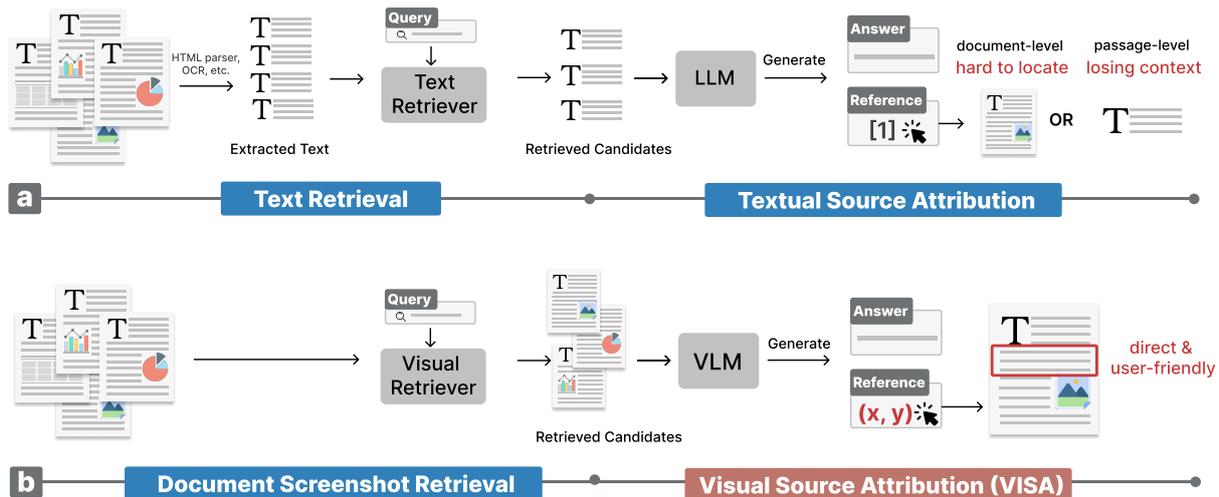


Figure 1: Comparison between (a) Text-based generation with source attribution in a RAG pipeline. and (b) Visual-based generation with source attribution in a V-RAG pipeline. VISA directly pinpoint the source evidence of the answer for user query in the original document with a bounding box.

083 making them more versatile for diverse use cases.
 084 Moreover, this task serves as a meaningful evaluation
 085 of VLMs, assessing their ability to provide
 086 self-explanations and accurately localize support-
 087 ing information within their original input in an
 088 RAG paradigm.

089 To train and evaluate VISA, we curated two
 090 datasets: Wiki-VISA and Paper-VISA. Wiki-
 091 VISA is derived from the Natural Questions
 092 dataset (Kwiatkowski et al., 2019). It reconstructs
 093 the original Wikipedia webpages, using short an-
 094 swers as generation targets and corresponding long
 095 answer’s HTML bounding box as source attribution
 096 targets. This dataset supports the test of model’s
 097 ability to attribute sources across multi-document,
 098 multi-page, and multi-modal content. On the other
 099 hand, Paper-VISA, built from PubLayNet (Zhong
 100 et al., 2019) with synthetic query generation, fo-
 101 cuses on the biomedical domain by evaluating per-
 102 formance on multi-modal scientific paper PDFs.
 103 Together, they provide diverse and challenging
 104 benchmarks for assessing the granularity and ac-
 105 curacy of source attribution in RAG systems. Our
 106 experiments, spanning both in-domain training and
 107 zero-shot evaluation, revealed existing state-of-the-
 108 art models like QWen2-VL-72B (Wang et al., 2024)
 109 struggle with precise visual source attribution in
 110 zero-shot prompting. Fine-tuning VISA on our cu-
 111 rated datasets significantly improved model perfor-
 112 mance in visual attribution accuracy. Further anal-
 113 ysis highlights key areas for improvement, such
 114 as enhancing bounding box precision for long im-
 115 age documents, multi-documents, and zero-shot

generalization capabilities.

2 Related Work 117

2.1 RAG attribution 118

119 Open-domain question answering with LLMs often
 120 suffer from two key issues: hallucinations and out-
 121 dated internal knowledge. Retrieval-Augmented
 122 Generation (RAG) has been recognized as an ef-
 123 fective solution to these problems (Lewis et al.,
 124 2020; Gao et al., 2024; Ovidia et al., 2024). In
 125 RAG, relevant documents are first retrieved from
 126 an external database and then fed into LLMs along-
 127 side the question. This allows LLMs to reference
 128 the retrieved documents during answer generation.
 129 Furthermore, RAG can generate a list of citations
 130 attached to the generated answers, linking them
 131 to the retrieved documents so users can verify the
 132 accuracy of the output. This process is known as
 133 source attribution (Rashkin et al., 2023; Bohnet
 134 et al., 2023; Khalifa et al., 2024).

135 Typically, RAG with source attribution follows
 136 a text-only pipeline where all inputs and outputs,
 137 such as questions, retrieved documents, generated
 138 answers, and citations, are in textual form. Re-
 139 cently, vision-based RAG pipelines have emerged,
 140 where the retrieved documents are represented
 141 as screenshot images (Ma et al., 2024; Faysse
 142 et al., 2024), and VLMs process both textual ques-
 143 tions and these document images to generate an-
 144 swers (Riedler and Langer, 2024; Xia et al., 2024;
 145 Yu et al., 2024). Compared to traditional text-only
 146 RAG, vision-based RAG can leverage structured

147 and visual information from documents, such as
148 tables, graphs, and images, which are often chal-
149 lenging to extract through text-only pipelines.

150 Our VISA attribution method proposed in this
151 paper is a novel approach for vision-based RAG
152 pipelines: directly drawing bounding boxes around
153 the content in retrieved document screenshots that
154 potentially supports the generated answers. This
155 approach differs from existing attribution methods
156 in two ways: (1) Granularity: Existing attribution
157 methods often operate at the document level, re-
158 quiring users to read entire documents to locate
159 supportive content. In contrast, our method directly
160 attributes the answer to specific content within the
161 document, such as a passage, table, or image in
162 the screenshot. (2) Presentation: Traditional attri-
163 bution methods provide a list of textual citations,
164 whereas our method uses bounding boxes, offering
165 a visually-oriented form of attribution. This can
166 help users quickly locate the relevant information.

167 2.2 Bounding Box Drawing with VLM

168 Bounding box-based object detection is a well-
169 established task in computer vision (CV) (Zhao
170 et al., 2019; Zou et al., 2023). Traditional ap-
171 proaches rely on convolutional neural networks
172 (CNNs) (LeCun et al., 2015) or Vision Trans-
173 formers (ViTs) (Dosovitskiy et al., 2021) to extract fea-
174 tures and predict bounding boxes alongside object
175 classification (Ren et al., 2015; Dai et al., 2016;
176 Redmon et al., 2016; Carion et al., 2020).

177 Recent vision-language models (VLMs) like
178 GPT4o (OpenAI et al., 2024), QWen2-VL (Wang
179 et al., 2024), and PaliGemma (Steiner et al., 2024)
180 have shown the ability to generate bounding box co-
181 ordinates in an image-to-text manner, taking input
182 images and generate the top-left and bottom-right
183 coordinates of target objects. Unlike traditional
184 object detection that focuses on natural images,
185 our method applies bounding box drawing to text-
186 intensive document screenshots.

187 Additionally, grounding elements on screenshots
188 has been explored in GUI agent systems (Cheng
189 et al., 2024; Lin et al., 2024), where bounding
190 boxes are used to localize UI elements like but-
191 tons. While these approaches focus on GUI con-
192 texts, our work targets visual source attribution in
193 vision-based RAG processes, grounding bounding
194 boxes to locate evidence within document images.

3 Method 195

3.1 Task Definition 196

197 Our VISA is a novel source attribution method pri-
198 marily designed for vision-based RAG systems. To
199 formally define the task of RAG with VISA: given
200 a textual user query q as the RAG system input, the
201 retrieval component of the system needs to retrieve
202 a set of candidate documents $D = \{d_1, \dots, d_n\}$
203 from corpus \mathcal{C} . Then the generation component of
204 the system needs to return three outputs: an answer
205 a that answers the query q , the identifier i of the
206 most relevant document d^* in D , and a bounding
207 box coordinates $B_{d^*} = [(x_1, y_1), (x_2, y_2)]$ within
208 d^* that highlight the content supporting the gener-
209 ated answer a .

210 In a vision-based RAG setup, user queries are
211 textual, while all documents in the corpus \mathcal{C} are
212 screenshots of documents (e.g., webpages or PDF
213 pages) provided as image inputs.

3.2 Generation with Visual Source Attribution 214

215 This paper focuses on VISA within the generation
216 component of vision-based RAG systems. As dis-
217 cussed in the previous section, VISA must handle
218 multimodal input. To achieve this, we leverage
219 VLMs for implementing VISA. Specifically, for a
220 given query and a set of retrieved candidate docu-
221 ments (i.e., screenshots of documents), the system
222 processes the inputs as follows: query tokens are
223 directly input into the language model, while docu-
224 ment screenshots are first processed by the image
225 encoder to extract image representations, which are
226 then fed into the language model.

227 The language model subsequently generates the
228 answer, the identifier of the relevant document, and
229 the xy-coordinates of the bounding box’s top-left
230 and bottom-right corner on the content that sup-
231 ports the generated answer. Notably, this entire
232 process can be framed as a next-token prediction
233 task. Finally, the generated identifier and bounding
234 box coordinates are used to draw the bounding box
235 on the target document screenshot, which is pre-
236 sented to the user along with the generated answer.

237 Technically, existing instruction-tuned VLMs,
238 such as QWen2-VL-72B (Wang et al., 2024), can
239 potentially be prompted to perform VISA in a zero-
240 shot manner. However, we find that VISA remains
241 a challenging task. Consequently, further super-
242 vised fine-tuning on a dedicated VISA task dataset
243 is necessary. In the next section, we introduce the
244

245 datasets we crafted specifically for training and
246 evaluating VISA.

247 3.3 Dataset Acquisition

248 The training and evaluation data suitable for the
249 VISA task needs to be formatted as follows: the
250 input consists of a textual query and document
251 screenshot images as multimodal inputs, while the
252 target outputs include the textual short answer, the
253 relevant document identifier, and the coordinates
254 of the bounding box. To create datasets that meet
255 these requirements, we craft existing publicly avail-
256 able datasets to support the training and evaluation
257 of our proposed VISA method.

258 **Wiki-VISA** is derived from the Natural Questions (NQ) dataset (Kwiatkowski et al., 2019). The
259 original NQ dataset provides natural questions,
260 along with short and long answers sourced from
261 Wikipedia webpages. We use the short answers
262 as answer targets. However, the original dataset
263 does not contain the original webpage screenshots.
264 We use the Selenium Python toolkit¹ to access and
265 render the webpage with the original URL with a
266 history version stamp. And take a screenshot with
267 980 pixels width and up to 3920 pixels (4 pages)
268 height. Using the long answer, we identify the cor-
269 responding element in the HTML from which the
270 long answer is derived. We then draw a bounding
271 box around this element to obtain the coordinates.
272 Notably, the answers in this dataset can come from
273 various elements, such as passages, tables, lists,
274 or images within the webpage. Since the ques-
275 tions and answers in Wiki-VISA are human-judged,
276 we consider this dataset a high-quality, supervised
277 dataset and evaluation for VISA on general knowl-
278 edge, with Wikipedia webpage.

280 **Paper-VISA** is derived from PubLayNet (Zhong
281 et al., 2019), a dataset originally designed for doc-
282 ument layout analysis of single page PubMed PDF
283 documents. PubLayNet provides bounding box
284 coordinates and class labels (e.g., title, text, table,
285 figure, etc.) for each element in a paper’s PDF
286 screenshot. However, the dataset does not include
287 queries or answers associated with each document.
288 To address this limitation, we leverage instruction-
289 tuned VLMs (e.g. Qwen2-VL-72B) to syntheti-
290 cally generate queries and answers. Specifically,
291 for each paper screenshot sample in the PubLayNet
292 training data, we select a bounding box within the
293 sample and overlay it on the screenshot. The mod-

294 ified screenshot is then input to the VLM with a
295 prompt designed to instruct the model to generate
296 a question and a short answer based on the content
297 within the bounding box. See Appendix A.1 for the
298 prompt details and generation example. By aug-
299 menting the original PubLayNet in this way, we
300 create synthetic queries and answers, enabling it
301 to support VISA training. We consider the result-
302 ing Paper-VISA dataset as synthetic training and
303 evaluation for scientific paper PDFs in the medical
304 domain.

305 **FineWeb-VISA** is based on the FineWeb-edu
306 corpus (Penedo et al., 2024), a high-quality text
307 corpus of crawled webpages. We sampled 60k web-
308 page URLs and used Selenium to capture screen-
309 shots of diverse, content-rich webpages. A passage
310 containing more than 50 words was randomly se-
311 lected as the target source. A bounding box was
312 drawn around the selected content, and a VLM
313 was prompted to generate a query and short answer
314 supported by the target content, similar as Paper-
315 VISA. Although Fineweb-VISA provides diverse
316 layout, it do not guaranteed to high quality data
317 has human annotated in Wiki-VISA or Paper-VISA
318 that assessing a specific domain, we only leverage
319 Fineweb-VISA as training data to analysis zeroshot
320 and data augmentation effectiveness.

321 **Multi-Candidates** By now, each query is paired
322 with the triplet of a positive document, target short
323 answer, and target evidence bounding box. To set
324 up a RAG experimental environment for evaluat-
325 ing VISA, we in addition need to let the generator
326 take multiple candidates as input, simulating the
327 scenario that the generator is taking multiple re-
328 trieval candidates and attributing the evidence in
329 most relevant documents. Given the query q , we
330 use a retriever R to retrieve top- k candidates. And
331 randomly sampled $m - 1$ candidates that are not
332 ground truth as hard negative candidates. The hard
333 negative candidates are mixed with the one ground
334 truth document together as the input for the multi-
335 document VISA. The reason we did not directly
336 take top- m documents as the retrieval candidate
337 is that we do not want VISA biased on a specific
338 retriever and position of the candidate docs. Gener-
339 ally, our VISA does not rely on the type of retriever.
340 It can be either a traditional text-based retriever that
341 indexes the document with extracted text or a recent
342 document screenshot retriever that directly indexes
343 the original document screenshot. However, inte-
344 grating with those visual-based retrievers enables

¹<https://pypi.org/project/selenium/>

Dataset	# Train	# Test
Wiki-VISA	87k	3,000
Paper-VISA	100k	2,160
Fineweb-VISA	60k	-

Table 1: Datasets statistics for train and test splits.

us to build an end-to-end RAG solution without the necessity of explicit document content processes such as HTML parsing or OCR. Thus, we leverage an off-the-shelf Document Screenshot Embedding (DSE) model (Ma et al., 2024) to serve as the retrieval component of the RAG system. When encoding queries and documents, the model directly encodes textual queries and document screenshot images into single vector embeddings and performs cosine similarity search during inference. In this work, we set $k = 20$ and $m = 3$.

Additionally, an RAG pipeline may have the chance of having no ground truth document returned from the retriever. We use a probability of 20% to randomly replace the ground truth document in the candidates, to access the model’s capability to detect no-answer situations. After these operations, the data statistics are shown in Table 1.

4 Experiment Setup

4.1 Evaluation

Evaluation metrics assessed both the generated answers and bounding box predictions. For answer generation, relaxed exact match (EM) was used to measure accuracy. If the golden answer and predicted answer have a sub-sequence relationship and the difference in string length is within 20 characters. The predicted answer is considered as correct. For bounding boxes, Intersection over Union (IoU) was calculated to determine localization precision, with an IoU threshold of 0.5 indicating a correct prediction.

To analyze performance across varying content types, test samples were categorized by the modality and location of the evidence. For Wiki-VISA, categories included first-page passages, passages beyond the first page, and non-passage content such as tables and figures. For Paper-VISA, since it is a single-page document, categories were divided into passage and non-passage content. The overall accuracy for each dataset was computed as a macro average across these categories.

We evaluate the effectiveness of VISA in two dif-

ferent settings: *Single oracle candidate* and *Multi-candidate*. *Single oracle candidate* setting solely evaluates the generation and visual attribution component. We conduct controlled experiments by training and testing the VLMs using only a single ground truth relevant document screenshot as input. In this setup, it is guaranteed that the answer can be found within the input document. The VLMs do not need to predict the relevant document identifier and can focus exclusively on answer generation and bounding box prediction.

In a *Multi-candidate* setting, the model is evaluated on its ability to distinguish relevant documents from irrelevant ones, in addition to generating accurate answers and bounding boxes. This setup better reflects the RAG scenarios in which multiple candidate documents are retrieved, and the model must not only generate a correct response but also attribute it to the correct source document. For the *Multi-candidate* evaluation, we assess two configurations: *Multi-candidate*, *Oracle in Candidates* which has ground truth in candidates, this setting has the same query set as the single setting, hence directly comparable. *Multi-candidate*, *Full* contains the additional 20% cases where ground truth has no answer.

4.2 Training details

To train vision-language models (VLMs) for answer generation with VISA, we initialized the models using the open-source Qwen2-VL-2B and Qwen2-VL-7B (Wang et al., 2024), finetuning on the training datasets described in Section 3.3.

We first trained the models in a single-candidate setup, where the input was limited to a single oracle document image. In this setup, the model was trained to generate both the answer and its corresponding bounding box. We used the prompt template provided in Appendix A.2 to format the model’s input and output.

Next, we trained the models in a multi-candidate setup. Here, the model received three document candidates and the task was to generate the identifier of the relevant document (if present), the answer, and the bounding box for the evidence. For cases where no relevant document was present (20% of the training samples), the model was trained to generate “No answer.” We used the prompt template provided in Appendix A.3 to format the model’s input and output.

The training objective for both setups was next-token prediction with cross-entropy loss. We fine-

Method	Wiki-VISA								Paper-VISA					
	Average		[<1] Passage		[>1] Passage		Non-Passage		Average		Passage		Non-Passage	
	bbx	ans	bbx	ans	bbx	ans	bbx	ans	bbx	ans	bbx	ans	bbx	ans
<i>Zeroshot Prompt</i>														
QWen2-VL-72B	1.5	60.4	3.4	58.5	0.1	54.9	0.9	67.9	1.5	43.1	0.5	40.2	2.5	45.9
<i>Fine-tune, Single Oracle Candidates</i>														
VISA-2B-single	37.5	57.1	70.0	61.1	18.7	44.9	23.8	65.3	63.0	38.3	50.6	34.4	75.3	42.1
VISA-7B-single	54.2	65.2	75.6	66.5	50.1	56.0	36.8	73.1	68.2	43.8	58.1	41.6	78.2	45.9
<i>Fine-tune, Multi Candidates, Oracle in Candidates</i>														
VISA-2B-multi	22.5	37.9	46.5	46.1	6.4	27.2	14.6	40.5	51.3	33.8	41.1	30.1	61.4	37.4
VISA-7B-multi	37.7	41.8	58.1	49.2	32.8	32.0	22.2	44.1	59.9	39.2	47.7	35.9	72.0	42.4
<i>Fine-tune, Multi Candidates, Full</i>														
VISA-2B-full	32.1	46.9	51.0	53.6	18.9	38.0	26.5	49.1	59.8	44.7	51.6	42.6	67.9	46.7
VISA-7B-full	41.6	51.1	56.6	57.1	34.4	43.2	33.9	53.1	66.8	50.3	57.1	47.5	76.5	53.0

Table 2: Effectiveness of VISA on Wiki-VISA and Paper-VISA datasets for bounding box accuracy (bbx) and answer accuracy (ans). Fine-tuned models are trained individually on in-domain data. The *Multi-Candidate, Oracle in Candidates* setting uses the same query set as the Single Oracle Candidates setting, allowing direct comparison. The full setting has an additional 20% queries with no ground truth documents in candidates.

tuned the models for two epochs in the single-candidate setting, using LoRA with a learning rate of $1e-4$, a batch size of 64, and $4 \times H100$ GPUs. For the multi-candidate setting, we initialized the models with weights from the single-candidate setup and trained for one epoch with the same learning rate. We froze the image encoder to reduce GPU memory usage in the multi-candidate setting.

During the training, random cropping was applied outside of the bounding box. This augmentation exposed the model to varying input sizes, which enhanced its zero-shot effectiveness on unseen document layouts. Bounding box targets were represented using absolute coordinate values. We also explored normalizing the scale of bounding box coordinates to values in the range[0-1]. Details can be found in Section 6.3.

5 Experimental Results

Table 2 presents the performance of VISA on the Wiki-VISA and Paper-VISA datasets across different experimental settings. Zero-shot prompting results reveal the difficulty of directly applying state-of-the-art VLMs to the visual source attribution task. QWen2-VL-72B achieves a reasonable answer generation accuracy of 60.4% on average on Wiki-VISA but fails to deliver effective bounding box predictions, with only 1.5% accuracy. This observation is consistent on Paper-VISA. These highlight the limitations of existing VLMs in pinpointing the source evidence in original documents with proper location and granularity.

Fine-tuning on our crafted training data enables

the model to effectively execute the task. In the single-candidate setup, where the model processes only the relevant document, fine-tuned models demonstrate substantial gains compared to zero-shot prompting a much larger model. On Wiki-VISA, the 7B variant achieves 54.2% bounding box accuracy and 65.2% answer accuracy, while on Paper-VISA, the corresponding scores reach 68.2% and 43.8%. Performance in the multi-candidate setting, which more closely mirrors real-world retrieval-augmented generation (RAG) systems, shows similar trends. The 7B model achieves 41.6% bounding box accuracy and 51.1% answer accuracy when handling three candidate documents, including cases where no relevant document is present. This demonstrates the model’s capability to identify relevant sources among multiple documents while enabling fine-grained attribution. However, when comparing the multi-candidates, oracle in candidates setting, We can see the model facing challenges when handling multiple candidates compared to just handling a single relevant document. E.g. on Wiki-VISA, bounding box accuracy for 7B model is 37.7% on average which is 17 points lower than the corresponding single candidate setting. Showing that visual source attribution among multi-candidates is much harder than just locating the source element in a single one.

It further demonstrates that the effectiveness of VISA is influenced by document characteristics, such as content location and modality. For Wiki-VISA, bounding box accuracy is significantly higher for passages on the first page ([<1] passage) compared to passages beyond the first page ([>1]

Train Data	Wiki-VISA								Paper-VISA					
	Average		[<1] Passage		[>1] Passage		Non-Passage		Average		Passage		Non-Passage	
	bbx	ans	bbx	ans	bbx	ans	bbx	ans	bbx	ans	bbx	ans	bbx	ans
Wiki	54.2	65.2	75.6	66.5	50.1	56.0	36.8	73.1	27.8	36.2	20.5	32.6	35.1	39.7
Paper	0.2	42.6	0	46.3	0.4	33.5	0.1	48.1	68.2	43.8	58.1	41.6	78.2	45.9
FineWeb	37.6	50.2	48.9	45.1	57.3	52.3	6.6	53.1	22.0	43.3	26.5	41.7	17.4	44.9
Wiki+Fineweb	58.2	65.3	68.7	66.6	61.7	57.1	44.1	72.1	21.0	43.1	18.5	42.2	23.4	43.9
Paper+Fineweb	36.1	48.7	51.8	49.6	49.6	44.2	6.8	52.4	66.5	44.6	56.1	42.2	76.9	47.0
Wiki+Paper+Fineweb	58.1	64.8	69.9	65.0	58.7	56.7	45.8	72.7	67.6	44.3	55.9	41.5	79.3	47.1

Table 3: Effectiveness of VISA trained on different combinations training data for bounding box accuracy (bbx) and answer accuracy (ans) in the single oracle candidate setting.

passage). For example, the 2B variant achieves 70.0% accuracy for [<1] passages but only 18.7% for [>1] passages, indicating the challenges posed by long, multi-page documents. The larger model, the 7B variant, narrows this gap, reflecting the better handling of long-context inputs. Non-passage content, such as tables and figures, also have obviously a different level of grounding effectiveness, indicating the difference of effectiveness in different visual elements.

6 Analysis

6.1 Out-of-Domain Zeroshot

Table 3 shows the effectiveness of VISA while trained with different data combinations in the single candidate setting. It enables us to study the effectiveness of out-of-domain transfer and augmentation. First, we highlight the challenges of zero-shot generalization in VISA. Training and evaluating on in-domain achieves an effective bounding box accuracy, e.g. 54.2% on average for Wiki-VISA. However, significant performance drops are observed when models are tested on out-of-domain datasets. For instance, a model trained on Wiki-VISA achieves only 27.8% bounding box accuracy on Paper-VISA, while a model trained on Paper-VISA achieves near-zero performance (0.2%) on Wiki-VISA. This gap underscores the difficulty of transferring visual source attribution capabilities across datasets with differing document structures, layouts, and content modalities. Interestingly, Wiki-VISA appears to transfer better to Paper-VISA compared to the reverse. This may be because of the multi-page nature of Wiki-VISA, which provides richer training signals that generalize better to simpler single-page setting in Paper-VISA.

FineWeb-VISA shows as a promising resource for training models with improved zero-shot capabilities. When trained on FineWeb-VISA alone, the model achieves 37.6% bounding box accu-

racy on Wiki-VISA and 22.0% on Paper-VISA. Notably, FineWeb-VISA outperforms Wiki-VISA training on [>1] passage bbx accuracy for Wiki-VISA (57.3% vs. 50.1%), suggesting its effectiveness in handling long and complex document structures. However, FineWeb-VISA does not perform as well on non-passage content, likely due to its training focus on passage-level targets.

6.2 Data Augmentation

The results also demonstrate the benefits of augmenting training data with FineWeb-VISA. On Wiki-VISA, combining Wiki and FineWeb training data improves bounding box accuracy from 54.2% to 58.2% and improves performance on [>1] passages from 50.1% to 61.7%, indicating that FineWeb complements Wiki by enhancing the model’s ability to attribute evidence in multi-page contexts. For Paper-VISA, however, augmenting with FineWeb does not significantly improve in-domain performance. Training on Paper+FineWeb achieves a comparable bounding box accuracy to Paper alone, but it enhances zero-shot performance on Wiki-VISA (from 0.2% to 36.1%).

Training on the full combination of datasets (Wiki+Paper+FineWeb) yields strong results across both domains, with 58.1% bbx accuracy on Wiki-VISA and 67.6% on Paper-VISA. This shows the importance of diverse training data for building generalizable models capable of handling different document types, layouts, and evidence modalities. Future work should focus on expanding the dataset diversity to further improve generalization and enable robust visual source attribution for a wide range of document structures.

6.3 Bounding Box Target

Table 4 shows the impact of different bounding box target representations and cropping strategies during training. Training with random cropping

Error Type	Type-I: Wrong source attribution	Type-II: Position misalignment	Type-III: Granularity mismatch
Question	Where is the energy released from when food is metabolized?	Who is the movie phantom thread based on?	Who played skeleton in the movie masters of the universe?

Document

Ground Truth VISA Output

Figure 2: Type of errors in the evaluation of Wiki-VISA.

Train Data	Wiki-VISA		Paper-VISA	
	bbx	ans	bbx	ans
Crop, Absolute	54.2	65.2	27.8	36.2
No Random Crop	58.8	65.6	1.7	36.9
Normalized Value	56.4	64.4	0.1	37.2
No Bounding Box	0	67.6	0	35.2

Table 4: Impact of bounding box target representation and cropping strategies during training on Wiki-VISA in the single oracle candidate setting.

and absolute coordinate values achieves a balance between in-domain performance on Wiki-VISA (54.2%) and zero-shot generalization to Paper-VISA (27.8%) in bounding box accuracy. Removing random cropping slightly improves Wiki performance but drastically reduces zero-shot generalization, indicating that random cropping enhances the model’s robustness to varied input sizes. Normalizing coordinate values achieves moderate performance on Wiki-VISA but fails on Paper-VISA, suggesting that absolute bounding box values are better suited to our experiments.

The “No Bounding Box” row represents a vanilla visual retrieval-augmented generation setup without visual source attribution, where models generate answers without bounding box predictions. VISA enables visual source attribution capability while the effectiveness of answer generation is preserved at about the same level of effectiveness.

6.4 Error Analysis

We conducted an error analysis on 50 randomly sampled cases from Wiki-VISA to better understand the limitations of VISA. Errors were categorized into three main types as demonstrated in Figure 2. The first type, wrong source attribution,

occurred in 43 cases where the model attributed the source to an incorrect section of the document, failing to identify the precise region containing the evidence. The second type, position misalignment, was observed in 4 cases where the model appeared to have the correct intent but drew the bounding box inaccurately, either slightly off position or incorrectly sized. The third type, granularity mismatch, appeared in 3 cases where the model’s attributed source, such as a specific cell in a table or an item in a list, did not match the ground truth granularity. While these cases could potentially be considered false negatives, we leave it in error analysis to emphasize the challenge in real-world use cases where user preferences for granularity may differ from the model’s output.

7 Conclusion

In this paper, we introduced VISA, a visual source attribution approach for retrieval-augmented generation pipeline. By leveraging vision-language models, VISA not only generates answers to user queries but also provides bounding boxes that visually attribute the supporting evidence within document screenshots. This capability enhances transparency and supports users in verifying the generated information effectively. Through the development of curated datasets, we demonstrated the effectiveness of VISA across diverse document types and layouts, including complex multi-page documents and multimodal content. Our experimental results highlight the potential of VISA to bridge the gap between information retrieval and answer generation by offering finer-grained, visually grounded evidence attribution. Moving forward, we hope VISA represents a pioneering step for more verifiable and user-friendly RAG systems.

8 Limitations

While VISA demonstrates promising results for answer generation and content grounding in vision-based RAG systems, it has several limitations. First, it focuses on generating short answers, which may not suffice for scenarios requiring detailed or explanatory responses, highlighting the need for enhancements in generating richer context. Second, it assumes answers are derived from a single, localized region within a document, which limits its effectiveness for cases where evidence spans multiple sections or modalities (e.g., combining text and tables). Third, while our evaluation spans web and medical scientific papers with various content modalities (e.g., passages, tables, figures), it does not fully capture the diversity of real-world documents such as scanned or handwritten content. Additionally, as VISA aims to make it intuitive for users to verify answers, conducting user studies could further confirm its practical utility.

References

- Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh, Luke Zettlemoyer, Hannaneh Hajishirzi, and Wen tau Yih. 2024. [Reliable, adaptable, and attributable language models with retrieval](#). *Preprint*, arXiv:2403.03187.
- Bernd Bohnet, Vinh Q. Tran, Pat Verga, Roei Aharoni, Daniel Andor, Livio Baldini Soares, Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, Tom Kwiatkowski, Ji Ma, Jianmo Ni, Lierni Sestorain Saralegui, Tal Schuster, William W. Cohen, Michael Collins, Dipanjan Das, Donald Metzler, Slav Petrov, and Kellie Webster. 2023. [Attributed question answering: Evaluation and modeling for attributed large language models](#). *Preprint*, arXiv:2212.08037.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *Computer Vision – ECCV 2020*, pages 213–229, Cham. Springer International Publishing.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading Wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Li YanTao, Jianbing Zhang, and Zhiyong Wu. 2024. [SeeClick: Harnessing GUI grounding for advanced visual GUI agents](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational*

Linguistics (Volume 1: Long Papers), pages 9313–9332, Bangkok, Thailand. Association for Computational Linguistics.

- Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. 2016. [R-fcn: Object detection via region-based fully convolutional networks](#). In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.

- 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. 2021. [An image is worth 16x16 words: Transformers for image recognition at scale](#). In *International Conference on Learning Representations*.

- Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. 2024. [Colpali: Efficient document retrieval with vision language models](#). *Preprint*, arXiv:2407.01449.

- Jeremy J. Foster. 1979. *The Use of Visual Cues in Text*, pages 189–203. Springer US, Boston, MA.

- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. [Enabling large language models to generate text with citations](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6465–6488, Singapore. Association for Computational Linguistics.

- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. [Retrieval-augmented generation for large language models: A survey](#). arXiv:2312.10997.

- Muhammad Khalifa, David Wadden, Emma Strubell, Honglak Lee, Lu Wang, Iz Beltagy, and Hao Peng. 2024. [Source-aware training enables knowledge attribution in language models](#). In *First Conference on Language Modeling*.

- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural Questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.

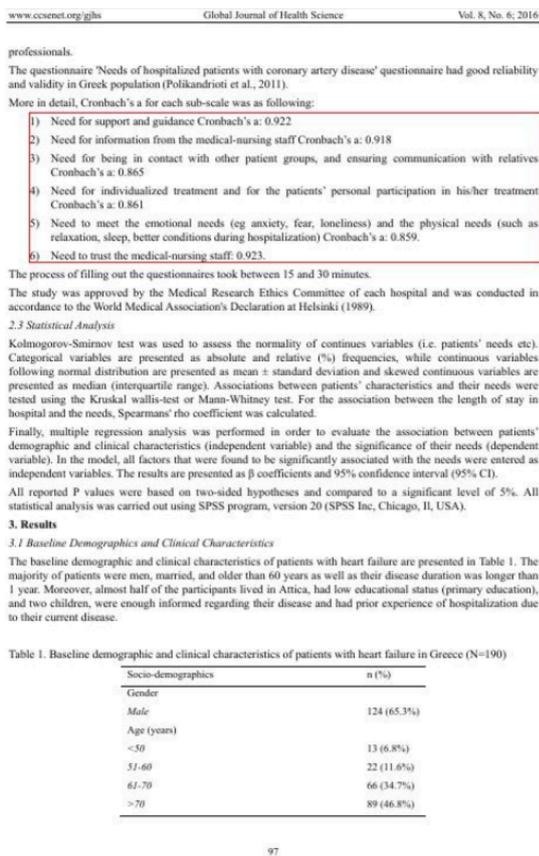
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature*, 521(7553):436–444.

- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020.

747	Retrieval-augmented generation for knowledge-intensive nlp tasks. In <i>Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20</i> , Red Hook, NY, USA. Curran Associates Inc.	807
748		808
749		809
750		810
751		811
752	Kevin Qinghong Lin, Linjie Li, Difei Gao, Zhengyuan Yang, Zechen Bai, Weixian Lei, Lijuan Wang, and Mike Zheng Shou. 2024. Showui: One vision-language-action model for generalist gui. In <i>NeurIPS 2024 Workshop on Open-World Agents</i> .	812
753		813
754		814
755		815
756		816
757	Xueguang Ma, Sheng-Chieh Lin, Minghan Li, Wenhu Chen, and Jimmy Lin. 2024. Unifying multimodal retrieval via document screenshot embedding. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 6492–6505, Miami, Florida, USA. Association for Computational Linguistics.	817
758		818
759		819
760		820
761		821
762		822
763		823
764	Chaitanya Malaviya, Subin Lee, Sihao Chen, Elizabeth Sieber, Mark Yatskar, and Dan Roth. 2024. ExpertQA: Expert-curated questions and attributed answers. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 3025–3045, Mexico City, Mexico. Association for Computational Linguistics.	824
765		825
766		826
767		827
768		828
769		829
770		830
771		831
772		832
773	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100, Singapore. Association for Computational Linguistics.	833
774		834
775		835
776		836
777		837
778		838
779		839
780		840
781	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichen, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varava, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jiang, Joaquin Quinero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lillian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeih, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatabaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix,	841
782		842
783		843
784		844
785		845
786		846
787		847
788		848
789		849
790		850
791		851
792		852
793		853
794		854
795		855
796		856
797		857
798		858
799		859
800		860
801		861
802		862
803		863
804		864
805		865
806		866

871	Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Pappay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyl Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. 2024. Gpt-4o system card . <i>arXiv:2410.21276</i> .	931
872		932
873		933
874		934
875		935
876		
877		936
878		937
879		938
880		
881		939
882		940
883		941
884		942
885		943
886		944
887		945
888		946
889		
890		947
891		948
892		949
893		
894		950
895		951
896		952
897		953
898		954
899		955
900		956
901		
902		957
903		958
904		959
905		960
906		961
907		
908	Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. 2024. Fine-tuning or retrieval? comparing knowledge injection in LLMs . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 237–250, Miami, Florida, USA. Association for Computational Linguistics.	962
909		963
910		964
911		965
912		966
913		967
914		968
915		969
916		
917		970
918		971
919		972
920		973
921		974
922		975
923		
924		976
925		977
926		978
927		
928		979
929		980
930		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

Input document screenshot with bounding box



Generated question and answer

Question: What is the Cronbach's alpha for the need for support and guidance sub-scale?

Short Answer: 0.922

Figure 3: An example of synthetic data from Paper-VISA.

A Appendix

A.1 Prompt for synthetic data generation

The following prompt was used for prompting QWen2-VL-72B to generate synthetic questions and answers for Paper-VISA and Fineweb-VISA datasets.

System:

Ask a question that can be specifically answered by the content in the red bounding box area and give a short answer. The question can be a wh- question, a yes/no question, or a how question, that can be answered in a few words.

Output format:

Question: <question>

Short Answer: <short answer>

Or simply return 'Empty' if the bounding box area is not visible or informative.

User: {image}

Figure 3 shows an example of synthetic data from Paper-VISA.

A.2 Prompt for Single Oracle candidate VISA

994

The following prompt template was used to format the model’s inputs and outputs for training the *Single Oracle Candidate VISA*.

995

996

```
Model Input:
System:
Given a document image, your task is to answer the question and locate the source of the answer
via a bounding box.

User:
{image} Image Size: {image.size}
Question: {question}

Model Output:
Assistant:
Answer: {answer}
Bounding Box: {bounding_box}
```

997

A.3 Prompt for Multi-candidate VISA

998

The following prompt template was used to format the model’s inputs and outputs for training the *Multi-candidate VISA*.

999

1000

```
Model Input:
System:
Given document images, your task is to answer the question and locate the source of the answer
via a bounding box.

User:
{image1} Image Size: {image1.size}
{image2} Image Size: {image2.size}
{image3} Image Size: {image3.size}
Question: {question}

Model Output:
Assistant:
Answer: {answer}
Evidence Document: {index}
Bounding Box: {bounding_box}
```

1001

A.4 Dataset Licenses

1002

- **NQ:** Apache License 2.0 1003
- **Wikipedia:** Creative Commons Attribution Share Alike, GNU Free Documentation License family. 1004
- **Fineweb-edu:** Open Data Commons License Attribution family. 1005
- **PubLayNet:** Community Data License Agreement – Permissive, Version 1.0. 1006
- **VISA Datasets:** Our crafted datasets follow the same license as the source of the documents. 1007

1008 **A.5 Model Backbone Licenses**

- 1009 • **QWen2-VL-72B**: Qwen LICENSE AGREEMENT.
- 1010 • **QWen2-VL-2B**: Apache License.
- 1011 • **QWen2-VL-7B**: Apache License.
- 1012 • **VISA Models**: Our fine-tuned models follow the same licenses as the original model backbone.

1013 **A.6 AI Assistant Usage**

1014 GPT4o is used during the writing to correct grammar errors and format tables.