

Look As You Think: Unifying Reasoning and Visual Evidence Attribution for Verifiable Document RAG via Reinforcement Learning

Shuochen Liu, Pengfei Luo, Chao Zhang, Yuhao Chen, Haotian Zhang, Qi Liu, Xin Kou, Tong Xu*, Enhong Chen

University of Science and Technology of China
shuochenliu@mail.ustc.edu.cn, tongxu@ustc.edu.cn

Abstract

Aiming to identify precise evidence sources from visual documents, visual evidence attribution for visual document retrieval-augmented generation (VD-RAG) ensures reliable and verifiable predictions from vision-language models (VLMs) in multimodal question answering. Most existing methods adopt end-to-end training to facilitate intuitive answer verification. However, they lack fine-grained supervision and progressive traceability throughout the reasoning process. In this paper, we introduce the **Chain-of-Evidence (CoE)** paradigm for VD-RAG. CoE unifies Chain-of-Thought (CoT) reasoning and visual evidence attribution by grounding reference elements in reasoning steps to specific regions with bounding boxes and page indexes. To enable VLMs to generate such evidence-grounded reasoning, we propose **Look As You Think (LAT)**, a reinforcement learning framework that trains models to produce verifiable reasoning paths with consistent attribution. During training, LAT evaluates the attribution consistency of each evidence region and provides rewards only when the CoE trajectory yields correct answers, encouraging process-level self-verification. Experiments on vanilla Qwen2.5-VL-7B-Instruct with Paper- and Wiki-VISA benchmarks show that LAT consistently improves the vanilla model in both single- and multi-image settings, yielding average gains of 8.23% in soft exact match (EM) and 47.0% in IoU@0.5. Meanwhile, LAT not only outperforms the supervised fine-tuning baseline, which is trained to directly produce answers with attribution, but also exhibits stronger generalization across domains.

1 Introduction

With the development of multimodal understanding capabilities in vision-language models (VLMs) (Chen et al. 2025c; Bai et al. 2025a), visual document retrieval-augmented generation (VD-RAG) has emerged as a critical research frontier. Nevertheless, current VLMs remain susceptible to hallucinations, whereby their outputs may deviate from the source document content (Bai et al. 2025b). Without reliable visual evidence attribution mechanisms to identify sources in documents, users cannot intuitively trace back the specific information employed by the model, thereby reducing the reliability of VD-RAG systems in applications.

*Corresponding author.

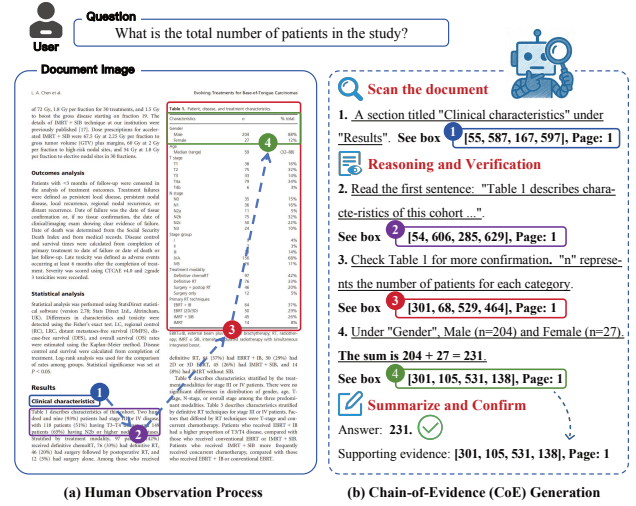


Figure 1: (a) Humans infer information by observing and locating supporting evidence in the document. (b) Each element in the reasoning step is linked to a visual attribution via a bounding box during Chain-of-Evidence generation.

Along this line, recent works (Shao et al. 2024a; Wu et al. 2025; Qi et al. 2025) mitigate hallucination in single-image Chain-of-Thought (CoT) (Kojima et al. 2022) reasoning by attending to critical visual regions as evidence, leveraging large-scale annotations of stepwise evidence regions. VISA (Ma et al. 2024b) further incorporates visual evidence attribution into the VD-RAG framework by linking answers with supporting sources via bounding boxes. Despite these advances, effective evidence attribution in visual documents still faces challenges: **(1) Lack of progressive reasoning mechanisms for verifiable attribution.** Existing methods, such as VISA, directly associate answers with final evidence, but do not reveal the intermediate process that would help users clearly understand the content and trace how the evidence is located. In contrast, when addressing complex problems, humans do not directly locate evidence but progressively search for relevant information. As shown in Figure 1a, such an observation pathway is structured as “chapter-paragraph summary-specific element”, localizing evidence from coarse to fine. However, current VLMs strug-

gle to replicate this progressive observation process, especially in multi-image scenarios. **(2) Limited supervision for learning multimodal reasoning.** Training models to learn CoT reasoning requires extensive annotated data (Xu et al. 2025; Wang et al. 2025), especially for stepwise evidence attribution. However, manually constructing such datasets is a costly process. Therefore, how to effectively learn stepwise visual evidence attribution and generalize reasoning abilities under limited annotated data remains a critical challenge.

Given these challenges, we introduce **Chain-of-Evidence (CoE)**, a reasoning paradigm that integrates Chain-of-Thought (CoT) reasoning with visual evidence attribution, designed for VD-RAG. Unlike existing methods that perform reasoning within a single image or directly link the final answer to the source, CoE models the intermediate reasoning trajectory by grounding each step to a supporting source from documents. This coarse-to-fine attribution process mirrors human problem-solving and enhances the reliability of VD-RAG. To this end, we introduce **Look-As-You-Think (LAT)**, a two-stage training approach designed to implement the proposed CoE paradigm. Specifically, in the first stage, we fine-tune VLM on a set of few human-verified annotated CoE data to learn reasoning patterns. Then we adopt reinforcement learning (RL) under the Group Relative Policy Optimization (GRPO) (Shao et al. 2024b) through a tailored reward design. The model is guided by a stepwise reward based on the semantic alignment between predicted visual evidence and corresponding context, encouraging faithful reasoning without requiring stepwise annotations. Furthermore, the reward is issued only when the CoE trajectory yields the correct answer. We further incorporate the outcome reward to ensure answer accuracy and guide the model in defining the scope of answer localization. This combined reward scheme enhances the ability to generate CoE reasoning, thereby linking sub-step verification to end-task performance. Our contributions are summarized as follows:

- 1) In the VD-RAG scenario, we formalize the **Chain-of-Evidence (CoE)** paradigm by modeling multimodal reasoning as a sequence of grounded steps, where reference elements (e.g., figures, tables, or factual information) are linked to their source through a bounding box and page index.
- 2) Building upon CoE, we propose **LAT**, an RL-based approach that jointly optimizes reasoning and visual grounding through a stepwise reward. By aligning each reference element with its visual evidence, LAT enables attribution-aware reasoning with few CoE-annotated samples.
- 3) LAT balances traceability and performance. Compared to the vanilla model, it shows average improvements of 8.23% EM and 47.0% IoU@0.5 in both single- and multi-image scenarios. Meanwhile, LAT not only outperforms the supervised fine-tuning baseline, which is trained to directly produce answers with attribution without CoE reasoning, but also exhibits stronger cross-domain generalization.

2 Related Work

Visual Evidence Attribution

Early end-to-end grounding methods integrate object detection into generated text by using markdown hyperlinks to

generate bounding-box tokens (Chen et al. 2023; Peng et al. 2023). They are trained on large corpora of grounded images and texts. Building on this foundation, multi-step visually grounded CoT frameworks (Shao et al. 2024a; Li et al. 2025; Wu et al. 2025; Xia et al. 2025) interleave reasoning and localization by predicting bounding boxes as evidence within the reasoning process, thereby yielding interpretable reasoning traces. However, these approaches have been validated exclusively on general visual perception tasks and rely on large-scale annotated evidence regions. They also do not explore the visual evidence attribution task in VD-RAG involving heterogeneous layouts and multi-page retrieval.

Moreover, existing text-based evidence attribution methods (Gao et al. 2023; Ye et al. 2024) in document RAG often operate at the document level, requiring users to read entire documents to locate supportive content. VISA (Ma et al. 2024b) first adapts visual evidence attribution to document screenshots by aligning the final answer with its evidence. Despite enabling intuitive verification of correctness, it fails to explicate the intermediate reasoning steps through which the model arrives at the answer. These limitations motivate the CoE reasoning paradigm, which generates faithful reasoning steps and validates each reference element against its corresponding source, as shown in Figure 1b.

RL for VLM Reasoning

Recent research has shown that RL-based policy optimization (Zhang et al. 2024) can improve the reasoning capabilities of large language models (LLMs) (Chen et al. 2025a,b; Zhang et al. 2025). DeepSeek-R1 (Guo et al. 2025) demonstrated that RL training elicits emergent CoT behaviors in LLMs, revealing the hallmark “aha moment”. Inspired by this phenomenon, several methods (Peng et al. 2025) have extended R1-style RL strategies to VLMs, leveraging rule-based reward functions to boost performance on mathematical reasoning and visual perception tasks. Unlike supervised fine-tuning (SFT), RL-based approaches achieve deeper reasoning and stronger generalization without relying on extensive human-annotated data (Chu et al. 2025).

However, existing RL frameworks for multimodal tasks are primarily optimized for answer accuracy as the reward signal (Yang et al. 2025; Shen et al. 2025), with no explicit supervision for verifying intermediate reasoning (Ni et al. 2025; Cao et al. 2025), and without design considerations for the visual evidence attribution task in VD-RAG. Drawing on human reading strategies, we introduce a stepwise, process-level reward that aligns each reasoning step with verifiable evidence. Leveraging the extracted CoE, we explicitly reward trajectories that are evidentially consistent and culminate in the correct answer, thereby ensuring valid reasoning and fostering faithful, attribution-based explanations.

3 Proposed Approach

To remedy the lack of verifiable progressive reasoning for visual evidence attribution, we first formalize the **Chain-of-Evidence (CoE)** paradigm. Building upon CoE, we present **LAT**, a two-stage RL-based framework shown in Figure 2. Stage I performs supervised fine-tuning to align annotated

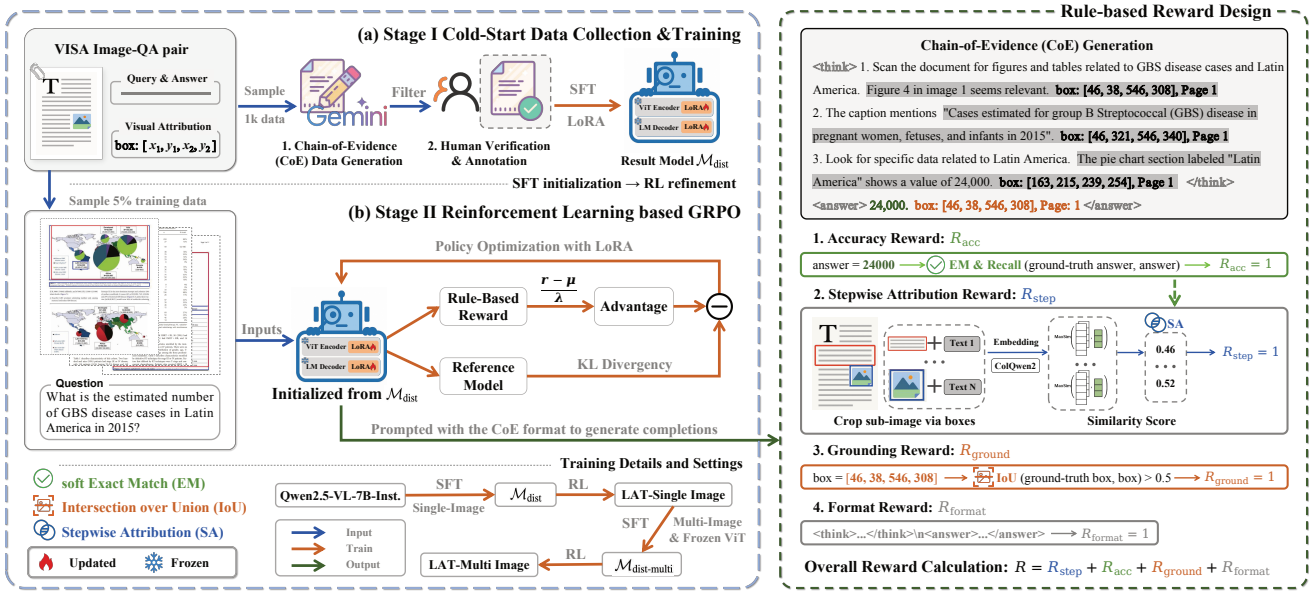


Figure 2: Overview of the proposed LAT framework. Left: A two-stage training pipeline. Stage I generates and filters the CoE data for fine-tuning. Stage II: The model undergoes refinement via RL under the GRPO algorithm. Right: Rule-based reward design. In GRPO training, the model generates CoE reasoning to guide policy updates through the reward signals.

CoE traces, and Stage II conducts fine-grained reinforcement learning with an answer-conditioned attribution reward that refines step-level grounding.

CoE Formalization and Notations

We formalize the generation with Chain-of-Evidence (CoE) reasoning as follows. Specifically, we define a textual query q and a set of document pages $\mathcal{P} = \{p_n\}_{n=1}^N$, which are pre-retrieved from the corpus. Given (q, \mathcal{P}) , CoE requires a VLM ϕ to perform CoT reasoning with stepwise visual attribution and then produce both an answer and its supporting evidence for evaluation, formulated as:

$$\mathcal{R}, \mathcal{B}, \mathcal{A} = \phi(q, \mathcal{P}). \quad (1)$$

Here, $\mathcal{R} = \{r_t\}_{t=1}^T$ denotes the textual reasoning steps, while $\mathcal{B} = \{(i_t, B_t)\}_{t=1}^T$ is the corresponding evidence chain, where $i_t \in [1, N]$ indicates the page index and $B_t = [(x_1^t, y_1^t), (x_2^t, y_2^t)]$ specifies the bounding box of the visual evidence for r_t on i_t -th page. After CoE reasoning, the final output $\mathcal{A} = \{a, (i^*, B_{\text{ans}})\}$ consists of the answer a , the most relevant page $p_{i^*} \in \mathcal{P}$, and the bounding box $B_{\text{ans}} \in \mathcal{B}$ confirming the evidence that supports a . It should be noted that not all reasoning steps in the response require a bounding box over images, such as calculation or conclusion steps. We omit these cases in our formulation for mathematical conciseness.

Stage I: Cold-Start Data Collection and Training

To prime the model for CoE reasoning, we first sample 1,000 instances from each training dataset and prompt a stronger proprietary model Gemini 2.5 Pro (Comanici et al. 2025) using two CoE exemplars as in-context demonstrations. The model produces stepwise rationales, each anno-

tated with bounding-box evidence, yielding a cold-start corpus that reflects the desired reasoning format.

For each query q and its corresponding response y , we assess answer quality with a recall metric,

$$\text{Recall}(a, a_{gt}) = \frac{|a \cap a_{gt}|}{|a_{gt}|}, \quad (2)$$

where a_{gt} is the ground-truth answer from the dataset, $|a|$ denotes the number of words in the extracted answer portion a from the response y , and $|a \cap a_{gt}|$ is the number of overlapping words between a and a_{gt} . Only samples with recall above a threshold γ are retained to ensure sufficient answer accuracy in the initial CoE traces.

We manually verify and correct any bounding-box drift, retaining verified samples with correct answers ($\sim 30\%$) to ensure data quality. Details of the resulting dataset $\mathcal{D}_{\text{final}}$ used in the cold-start training, including its split distribution, are provided in Table 6 in the appendix. Next, we fine-tune the VLM on $\mathcal{D}_{\text{final}}$ using LoRA (Hu et al. 2022), aiming to minimize the cross-entropy loss between the generated output and the annotated reasoning sequences through SFT. The resulting model is defined as $\mathcal{M}_{\text{dist}}$.

Stage II: Unified Reasoning and Visual Attribution via Reinforcement Learning

To emulate the human observation process shown in Figure 1, the model needs to deliver both an accurate answer and attribution-aware CoE reasoning to identify evidence supporting the final answer. We decompose the overall objective into four sub-goals: answer accuracy, stepwise visual attribution quality, evidence grounding precision, and adherence to the structured output format. Accordingly, we design four reward functions for rule-based RL training.

Accuracy Reward (R_{acc}). We reward the model based on soft exact match (EM), considering a response correct and $\text{EM}(a, a_{gt})$ equal to 1 if the normalized predicted answer a is a substring of the ground truth a_{gt} , or vice versa. To prevent the reward from becoming too sparse, we enhance this signal by including the recall metric. Formally,

$$R_{\text{acc}} = \frac{\mathbb{I}(\text{EM}(a, a_{gt}) = 1) + \text{Recall}(a, a_{gt})}{2}, \quad (3)$$

where $\mathbb{I}(\cdot)$ denotes the indicator function. By incorporating the recall metric, we prevent the wholesale rejection of semantically relevant yet not exact-match samples, while assigning higher rewards to perfectly matched outputs, thereby enhancing the model’s generalization capability.

Stepwise Attribution Reward (R_{step}). To ensure each reasoning step is grounded in the associated evidence attribution, we design a stepwise reward that measures the semantic alignment between each step and the evidence indicated by the bounding box. For each textual reasoning r_t involving reference element, we crop the raw document image of the i_t -th page according to the predicted bounding box B_t and encode both the cropped sub-image and r_t using a multimodal retriever ColQwen2 (Faysse et al. 2024), which enables efficient document indexing through visual features and handles dynamic resolutions. Formally,

$$e_{\text{img}}^{(t)} = \text{Norm}(\text{enc}_{\text{img}}(\text{crop}(B_t))), \quad (4)$$

$$e_{\text{txt}}^{(t)} = \text{Norm}(\text{enc}_{\text{txt}}(r_t)), \quad (5)$$

where $\text{Norm}(\cdot)$ denotes L_2 normalization for consistent cosine similarity evaluation. $\text{enc}_{\text{img}}(\cdot)$ and $\text{enc}_{\text{txt}}(\cdot)$ represent the image and text encoders of the retriever, respectively. $\text{crop}(\cdot)$ denotes the extraction of a sub-image from the original document page, delimited by the bounding box B_t . The stepwise attribution reward is defined as:

$$S = \min_{1 \leq k \leq K} \cos(e_{\text{img}}^{(k)}, e_{\text{txt}}^{(k)}), \quad (6)$$

$$I = \max_{1 \leq i, j \leq K, i \neq j} \text{IoU}(B_i, B_j), \quad (7)$$

$$R_{\text{step}} = \frac{\mathbb{I}(S \geq \tau) + \mathbb{I}(I \leq \delta)}{2} \cdot \mathbb{I}(R_{\text{acc}} \geq \epsilon), \quad (8)$$

where K is the number of context-evidence pairs and $\cos(\cdot)$ denotes the cosine similarity measure. We require all pairs in the CoE reasoning process to exceed the similarity threshold τ to ensure proper alignment. Nevertheless, the model may exploit this mechanism by repeating bounding boxes to satisfy the reward, which undermines coarse-to-fine evidence attribution and rich visual cues.

To mitigate such a phenomenon, we introduce a constraint on bounding box overlap by computing the maximum pairwise intersection over union (IoU) I among bounding boxes from the reasoning steps (Equation 7) and enforcing $I \leq \delta$ in Equation 8. This constraint promotes attribution diversity and encourages progressive grounding across reasoning steps. Furthermore, the reward is constrained by R_{acc} , ensuring that only faithful CoE processes contributing to correct answers are reinforced. This formulation guides the model to treat visual grounding as an internal retrieval problem, supporting stepwise and fine-grained visual attribution.

Grounding Reward (R_{ground}). The model needs to select the relevant source from the CoE reasoning process as evidence to support the final answer. We measure its precision by computing the IoU between the predicted bounding box (i^*, B_{ans}) and the ground-truth evidence (i_{gt}, B_{gt}) :

$$R_{\text{ground}} = \mathbb{I}(\text{IoU}(B_{\text{ans}}, B_{gt}) > 0.5), \text{ s.t. } i^* = i_{gt}. \quad (9)$$

By penalizing misaligned evidence attribution, R_{ground} drives the model to identify correct content within the source page i_{gt} instead of blindly collecting irrelevant information.

Format Reward (R_{format}). We prompt the model to conduct CoE reasoning on visual documents and generate answers. Well-formatted outputs contain CoE reasoning within `<think>...</think>` tags and answers with supporting evidence within `<answer>...</answer>` tags.

$$R_{\text{format}} = \begin{cases} 1, & \text{if the format is correct,} \\ -1, & \text{otherwise.} \end{cases} \quad (10)$$

By assigning negative rewards to incorrectly formatted outputs, we impose stricter output constraints and accelerate convergence toward the desired format, enabling more targeted optimization of the remaining reward components.

Training Algorithm. As shown in Figure 2, we sample 5% training data from the raw dataset, initialize the policy model π_θ from the cold-start checkpoint $\mathcal{M}_{\text{dist}}$, and adopt the GRPO (Shao et al. 2024b) algorithm, combined with the fine-grained reward design (Shen et al. 2025). GRPO, which supports rule-based rewards and optimizes the policy π_θ by sampling a group of candidate outputs for each query and computing a group-relative advantage, eliminates the need for training a critic model. We iteratively update π_θ , encouraging it to produce responses with higher relative advantages within the sampled group under the GRPO objective.

4 Experiment Setup

Datasets

We conducted the experiments on the VISA benchmark (Ma et al. 2024b), which is the first dataset designed for visual evidence attribution in real-world VD-RAG scenarios and comprises three subsets: (1) **Wiki-VISA** is derived from the Natural Questions (NQ) (Kwiatkowski et al. 2019). VISA renders the Wikipedia pages and identifies the HTML element containing the answer with a bounding box. (2) **Paper-VISA** builds upon PubLayNet (Zhong, Tang, and Yepes 2019). VISA synthesizes QA pairs grounded in annotated layouts. (3) **FineWeb-VISA** extends FineWeb-edu (Penedo et al. 2024) by selecting passages longer than 50 tokens and synthesizing grounded QA pairs. In the single-image setup, each query is paired with a source document page, a short answer extracted from it, and an evidence bounding box.

When extended to multiple images, for each query q , two additional document images are randomly sampled from the top K screenshots retrieved by Document Screenshot Embedding (DSE) (Ma et al. 2024a) within the VISA dataset, then merged with the source document page as input. For no-answer scenarios, the source document page is replaced with an irrelevant page from the dataset. Example cases are shown in Figure 5 in Appendix B.

Models	Param.	Wiki-VISA (Single)		Paper-VISA (Single)		Wiki-VISA (Multi)		Paper-VISA (Multi)	
	Size	EM	IoU@0.5	EM	IoU@0.5	EM	IoU@0.5	EM	IoU@0.5
Proprietary and Open-source Models, Direct Answer									
Gemini2.0 flash †	-	73.3	-	46.8	-	28.2	-	26.7	-
Qwen-VL max †	-	63.1	-	46.6	-	48.7	-	32.1	-
mPLUG-DocOwl2	8B	23.6	-	20.3	-	19.3	-	21.5	-
Qwen2.5-VL	7B	67.7	-	38.2	-	54.1	-	34.6	-
Qwen2.5-VL	32B	66.3	-	36.0	-	60.8	-	38.8	-
InternVL2.5	8B	45.9	-	42.4	-	37.7	-	32.1	-
LLaVA-OneVision	7B	48.4	-	28.7	-	38.5	-	30.8	-
LLaVA-CoT	11B	45.3	-	28.8	-	-	-	-	-
R1-OneVision	7B	58.0	-	24.4	-	57.2	-	31.4	-
Proprietary Model, Full Data Training, Direct Answer & Attribution									
VISA †	7B	74.1	28.3	49.6	62.2	58.7	31.6	54.1	65.6
Qwen2.5-VL (DA)	7B	69.4	0.80	43.8	2.87	52.5	0.97	37.2	5.56
LAT-Ind.	7B	73.6	53.7	45.4	49.9	64.5	38.0	51.2	46.9
LAT-Full	7B	73.1	57.8	46.2	48.4	64.8	41.4	50.6	49.3
△ Vanilla model	-	+4.20	+57.0	+2.40	+47.0	+12.3	+40.4	+14.0	+43.7

Table 1: Performance comparison on Paper- and Wiki-VISA in both single- and multi-image settings. Bold indicates the best score in each column. † denotes the proprietary model, which is fine-tuned on the full in-domain dataset, serving as an upper-bound baseline. “DA” refers to the zero-shot prompting setting for direct answer & attribution without training.

Baseline

To evaluate the effectiveness of LAT, we compared it against three categories of baselines: **(1) Proprietary and open-source models**, including Gemini (Anil et al. 2025), Qwen-VL Max (Bai et al. 2023), mPLUG-DocOwl2 (Hu et al. 2024), InternVL2.5 (Chen et al. 2025c), LLaVA-OneVision (Li et al. 2024), and Qwen2.5-VL (Bai et al. 2025a). **(2) Reasoning models**, including LLaVA-CoT (Xu et al. 2025), which is trained via SFT on CoT data, and R1-OneVision (Yang et al. 2025), which is optimized with RL in general scenarios. **(3) Attribution-supervised model:** VISA-7B (Ma et al. 2024b), trained to generate answers with direct attribution. These models serve as baselines to evaluate whether LAT balances answer accuracy with attribution and achieves CoE reasoning under limited supervision.

Training Details

We used Qwen2.5-VL-7B-Instruct (Bai et al. 2025a) as the backbone and applied LoRA (Hu et al. 2022) for parameter-efficient fine-tuning, with rank $r=64$ and scaling factor $\alpha=64$. Following the pipeline in Figure 2, we performed SFT on $\mathcal{D}_{\text{final}}$ using a learning rate of $1e-4$, followed by RL with $5e-5$. During RL training, LAT was trained on 5% of QA pairs sampled from the raw dataset. After each stage, we merged the LoRA parameters for subsequent training.

In the multi-image setting, each query is paired with three retrieved documents provided in the *candidates* field, as included in the VISA dataset. When no relevant information is available, the model is trained to output “No answer”. We initialized the multi-image model from the single-image trained version and further performed SFT using multi-image CoE data in $\mathcal{D}_{\text{final}}$, fine-tuning the LoRA adapter of the LM while keeping the vision transformer (ViT) frozen to reduce GPU memory usage. Additional experimental de-

tails are provided in Appendix A. Our code is available at <https://github.com/PolarisLiu1/LAT>.

Evaluation

We measured the performance across three dimensions: answer accuracy, evidence grounding, and stepwise attribution quality. Specifically, we reported answer accuracy using **soft Exact Match (EM)**, and evaluated grounding precision by computing **IoU@0.5**, which measures the proportion of the predicted box B_{ans} whose IoU with the ground truth evidence exceeds 0.5. To evaluate the quality of stepwise visual attribution, we employed the **Stepwise Attribution (SA)** reward function. We adopted the default threshold $\tau=0.3$ in Equation 8 to determine whether a step is correctly verified.

Following the evaluation of VISA (Ma et al. 2024b), we assessed the performance on both the Paper- and Wiki-VISA datasets under two settings. (1) **Single-image**: The model is provided solely with the source document and evaluated across three dimensions. (2) **Multi-image**: The model is additionally required to identify the source document from a set of retrieved candidates. For reproducibility, we adopt greedy decoding as the decoding strategy during evaluation.

5 Results and Analysis

Main Results

Attribution-aware performance under in-domain and cross-domain settings. To assess both answer accuracy and attribution precision, we evaluated models under direct answer and attribution-aware settings. We define an in-domain setup as evaluation on datasets from the same distribution as the training data. As shown in Table 1, LAT trained on in-domain data (**LAT-Ind.**) outperforms open-source models. Compared to the vanilla model, LAT enhances both answer correctness (EM, **+7.95%**) and evidence

Models	Param.	Wiki-VISA (Single)			Paper-VISA (Single)			Wiki-VISA (Multi)			Paper-VISA (Multi)		
	Size	EM	IoU@0.5	SA	EM	IoU@0.5	SA	EM	IoU@0.5	SA	EM	IoU@0.5	SA
Proprietary Models, CoE Reasoning													
Gemini2.0 flash †	-	52.0	4.47	24.0	46.4	4.72	14.1	32.1	1.20	5.20	27.0	3.24	20.0
Qwen-VL max †	-	55.2	0.17	10.3	47.8	2.87	8.6	35.8	0.10	5.90	32.6	1.85	12.5
Open-source Models, CoE Reasoning													
Qwen2.5-VL	7B	60.4	2.20	13.9	36.9	3.80	12.4	54.9	11.0	12.5	39.6	16.2	29.5
+one shot ICL	7B	61.8	1.37	12.2	37.3	8.30	29.3	54.2	9.20	11.9	39.5	18.7	30.9
Qwen2.5-VL	32B	62.8	9.27	0.16	35.2	2.50	0.11	62.5	18.7	0.60	43.3	19.0	0.40
+one shot ICL	32B	61.7	9.13	5.23	35.0	4.12	1.25	64.1	22.1	1.87	42.6	19.6	0.51
LAT-Ind.	7B	73.6	53.7	64.6	45.4	49.9	35.5	64.5	38.0	71.8	51.2	46.9	46.3
LAT-Full	7B	73.1	57.8	59.6	46.2	48.4	33.8	64.8	41.4	75.2	50.6	49.3	53.2
Δ Vanilla model	-	+13.2	+55.6	+50.7	+9.3	+46.1	+23.1	+9.9	+30.4	+62.7	+11.6	+33.1	+23.7

Table 2: Performance comparison for Chain-of-Evidence (CoE) reasoning processes on Paper- and Wiki-VISA datasets. Bold indicates the best score in each column. Results include both result accuracy (EM, IoU@0.5) and process quality (SA) metrics.

grounding (IoU@0.5, **+44.6%**) by optimizing attribution-aware reasoning through CoE-guided RL. The improvements are evident in multi-page scenarios, where the complexity of evidence selection emphasizes the advantages of our approach. In particular, for unanswerable cases, LAT outputs “No answer” and achieves an average precision of 65%, highlighting its robustness in handling such scenarios.

Table 2 highlights the effectiveness of LAT in improving CoE reasoning quality. To assess the impact of in-context learning (ICL) (Brown et al. 2020), we included an annotated CoE example as a prompt. While this yields a moderate improvement (SA, +4.0%), suggesting that demonstrations can partially guide the generation of structured reasoning, it has a limited effect on the result. LAT achieves a substantial SA gain of **37.5%**, demonstrating its ability to ground each reasoning step accurately. Moreover, LAT maintains high answer accuracy, indicating an alignment between faithful reasoning and correct outcomes. To assess generalization across domains, we trained **LAT-Full** on the sampled subset from all datasets. Compared to the in-domain variant, **LAT-Full** shows further improvements (e.g., Wiki-VISA Multi, IoU@0.5: 38.0%→**41.4%**; SA: 71.8%→**75.2%**), exhibiting generalization across diverse document distributions.

CoE reasoning performance with limited supervision.

Unlike VISA, which relies on large-scale (100k) supervised data and directly links the final answer to supporting evidence without reasoning, LAT is trained on only 5% of raw QA pairs during the RL stage. In low-resource settings, LAT achieves comparable performance to VISA-7B in answer accuracy and attribution precision, while maintaining traceable CoE reasoning. Notably, on high-resolution Wiki-VISA images, LAT outperforms VISA-7B, demonstrating robustness in visually complex scenarios under limited supervision.

To ensure fair comparison, we established an SFT baseline by training the model with VISA’s supervision protocol on the same data subset and experimental setup used for our approach. Figure 3a demonstrates LAT’s superior performance in both answer accuracy and evidence precision, with greater improvements observed in the multi-image set-

Train→Eval	Method	EM	IoU@0.5
Paper→Wiki (Single)	SFT	66.0	29.4
	LAT-Ind.	67.7 _{↑1.7}	35.6 _{↑6.2}
Paper→Wiki (Multi)	SFT	48.7	10.3
	LAT-Ind.	57.3 _{↑8.6}	21.4 _{↑11.1}

Table 3: LAT demonstrates robust generalization with cross-domain transfer between Paper-VISA and Wiki-VISA.

ting (Figure 3b). Meanwhile, Table 3 highlights LAT’s generalization across datasets, outperforming SFT by 1.7% in EM and 6.2% in IoU@0.5 in the “Paper→Wiki” transfer setting, while preserving the vanilla model’s performance on Wiki- (EM: 67.7%) and Paper-VISA (EM: 38.2%). This demonstrates that LAT improves transfer effectiveness without sacrificing adaptability to diverse data types.

Ablation Study

Effectiveness of reward components. We conducted ablation studies on both datasets to examine the contributions of individual components in our reward formulation (Table 4). We first analyzed the impact of distillation based on annotated CoE reasoning trajectories. The model $\mathcal{M}_{\text{dist}}$, obtained through fine-tuning during the cold start stage, shows an average improvement of **6.48%** in EM and **20.9%** in IoU@0.5 compared to the vanilla model. This indicates that distillation improves adherence to CoE reasoning formats.

Next, we assess the impact of stepwise attribution by ablating the process reward R_{step} . The model fails to align intermediate reasoning steps with visual evidence, resulting in a 15.5% reduction in SA. Since the evidence grounding of the answer is inherently linked to the quality of intermediate visual attribution ($B_{\text{ans}} \in \mathcal{B}$), we also observed a decline in IoU@0.5, confirming the necessity of step-level attribution. Meanwhile, without the overlap constraint $I \leq \delta$ in R_{step} , the model tends to reuse large and redundant regions across reasoning steps. This behavior undermines the coarse-to-fine grounding strategy and reduces attribution fidelity.

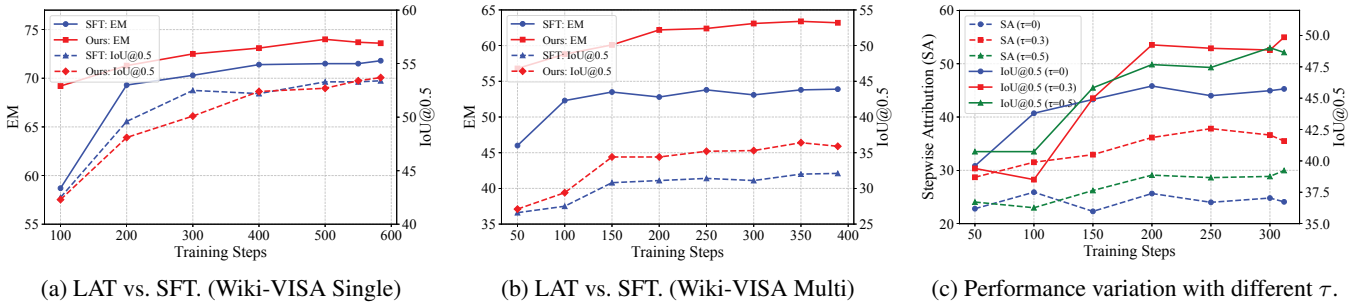


Figure 3: Comparison of LAT and SFT performance across different settings and ablation study on threshold τ .

Models	Wiki-VISA (Single)			Paper-VISA (Single)			Wiki-VISA (Multi)			Paper-VISA (Multi)		
	EM	IoU@0.5	SA	EM	IoU@0.5	SA	EM	IoU@0.5	SA	EM	IoU@0.5	SA
Vanilla model (Qwen2.5-VL-7B)	60.4	2.20	13.9	36.9	3.80	12.4	54.9	11.0	12.5	39.6	16.2	29.5
$\mathcal{M}_{\text{dist}}$ (Stage I)	67.1	25.6	32.8	44.7	27.0	13.0	59.4	27.8	44.0	46.5	36.3	28.2
LAT-Ind. (Stage II)	73.6	53.7	64.6	45.4	49.9	35.5	64.5	38.0	71.8	51.2	46.9	46.3
w/o. R_{step}	73.1	49.8	43.8	45.0	45.7	24.1	61.8	35.0	48.0	49.0	44.8	40.5
w/o. $R_{\text{acc}}, R_{\text{ground}}$	72.7	30.4	59.8	44.9	28.7	44.2	59.0	31.8	67.1	50.7	41.3	54.7
w/o. $R_{\text{acc}}, R_{\text{ground}}, (R_{\text{acc}} \geq \epsilon)$	67.0	28.9	56.0	45.0	22.6	60.7	31.7	22.5	81.7	37.6	30.3	63.2

Table 4: Ablation study results on LAT. “w/o.” denotes the results of models trained without the corresponding reward function.

Aligning reasoning and result objectives. We further analyzed the role of jointly optimizing process- and result-level rewards, where we removed the supervision from the final answer (w/o. $R_{\text{acc}}, R_{\text{ground}}$) and used it solely as a filter for valid reasoning paths. The model maintains moderate EM performance, indicating that consistency between the reasoning process and the answer inherently provides useful training signals. However, when supervision $R_{\text{acc}} \geq \epsilon$ is removed from R_{step} , the training objective is reduced to aligning only the reasoning format, resulting in significant degradation in both answer accuracy and evidence grounding precision. This suggests that R_{step} and result-level rewards work synergistically, where process rewards guide reasoning coherence while answer supervision ensures factual accuracy and proper grounding to source content.

Further Analysis

Sensitivity analysis of attribution threshold τ . To evaluate the alignment quality between visual evidence and textual references in CoE reasoning, we introduced a similarity threshold τ in the stepwise attribution reward function to distinguish positive from negative samples. Based on the synthetic CoE dataset $\mathcal{D}_{\text{final}}$, we computed semantic similarity scores across different answer types. As shown in Figure 8 in Appendix E, we reported the range of similarity values for each category, excluding pairs with zero scores. We set the default threshold to 0.3 based on the distribution analysis.

Given the parameter sensitivity, we further conducted experiments on both datasets to evaluate the robustness of τ across different document types. Specifically, we compared a high-threshold setting ($\tau=0.5$), a no-step variant (equivalent to $\tau=0$ or 1, where all steps are uniformly rewarded), and the default setting. As shown in Figure 3c and Figure 4

(Appendix E), $\tau=0.3$ achieves a better balance, yielding consistent improvements in both IoU@0.5 and SA throughout training. In contrast, the high threshold fails to sustain performance gains, as such strict criteria make it difficult to sample sufficient positive instances, while the no-step variant underperforms due to the lack of fine-grained attribution guidance. These results suggest that attribution supervision at $\tau=0.3$ offers relatively effective guidance for stepwise grounding. Additional analyses are reported in Appendix E.

Traceable Reasoning with the CoE Paradigm. The CoE paradigm achieves stepwise visual attribution, generating a traceable reasoning process toward the final answer. Through the LAT framework, we leverage stepwise rewards to achieve an average SA of 57.1%. Meanwhile, by penalizing repetitive reasoning processes in Equation 7, we encourage the model to generate diverse and fine-grained reasoning. As illustrated in Figure 12–17 in Appendix F, LAT accurately identifies the answer regions and generates faithful reasoning paths that closely align with the visual evidence.

6 Conclusion

In this paper, we introduce Chain-of-Evidence (CoE), a reasoning paradigm that unifies CoT with stepwise visual evidence attribution. To achieve CoE, we propose Look As You Think (LAT), a reinforcement learning framework that aligns the intermediate process to mitigate ungrounded reasoning for the visual evidence attribution task in VD-RAG. By incorporating stepwise rewards under the GRPO algorithm, LAT facilitates verification at each reasoning step. Experiments on Paper- and Wiki-VISA show that LAT outperforms baselines. We hope this work inspires further research on enhancing the verifiability of VD-RAG systems.

7 Acknowledgements

This work was supported in part by the grants from National Science and Technology Major Project (No. 2023ZD0121104), and National Natural Science Foundation of China (No.62222213, 62072423).

References

- Anil, R.; Borgeaud, S.; Alayrac, J.-B.; Yu, J.; Soricut, R.; Schalkwyk, J.; Dai, A. M.; Hauth, A.; Millican, K.; et al. 2025. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805.
- Bai, J.; Bai, S.; Yang, S.; Wang, S.; Tan, S.; Wang, P.; Lin, J.; Zhou, C.; and Zhou, J. 2023. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond. arXiv:2308.12966.
- Bai, S.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; Song, S.; Dang, K.; Wang, P.; Wang, S.; Tang, J.; Zhong, H.; Zhu, Y.; Yang, M.; Li, Z.; Wan, J.; Wang, P.; Ding, W.; Fu, Z.; Xu, Y.; Ye, J.; Zhang, X.; Xie, T.; Cheng, Z.; Zhang, H.; Yang, Z.; Xu, H.; and Lin, J. 2025a. Qwen2.5-VL Technical Report. arXiv:2502.13923.
- Bai, Z.; Wang, P.; Xiao, T.; He, T.; Han, Z.; Zhang, Z.; and Shou, M. Z. 2025b. Hallucination of Multimodal Large Language Models: A Survey. arXiv:2404.18930.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Cao, M.; Zhao, H.; Zhang, C.; Chang, X.; Reid, I.; and Liang, X. 2025. Ground-R1: Incentivizing Grounded Visual Reasoning via Reinforcement Learning. arXiv:2505.20272.
- Chen, K.; Zhang, Z.; Zeng, W.; Zhang, R.; Zhu, F.; and Zhao, R. 2023. Shikra: Unleashing Multimodal LLM’s Referential Dialogue Magic. arXiv:2306.15195.
- Chen, Y.; Liu, S.; Lyu, Y.; Zhang, C.; Shi, J.; and Xu, T. 2025a. Xiangqi-R1: Enhancing Spatial Strategic Reasoning in LLMs for Chinese Chess via Reinforcement Learning. arXiv:2507.12215.
- Chen, Y.; Lyu, Y.; Liu, S.; Zhang, C.; Lv, J.; and Xu, T. 2025b. Think Wider, Detect Sharper: Reinforced Reference Coverage for Document-Level Self-Contradiction Detection. In Christodoulopoulos, C.; Chakraborty, T.; Rose, C.; and Peng, V., eds., *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 1273–1288. Suzhou, China: Association for Computational Linguistics. ISBN 979-8-89176-332-6.
- Chen, Z.; Wang, W.; Cao, Y.; Liu, Y.; Gao, Z.; Cui, E.; Zhu, J.; Ye, S.; Tian, H.; Liu, Z.; Gu, L.; Wang, X.; Li, Q.; Ren, Y.; Chen, Z.; Luo, J.; Wang, J.; Jiang, T.; Wang, B.; He, C.; Shi, B.; Zhang, X.; Lv, H.; Wang, Y.; Shao, W.; Chu, P.; Tu, Z.; He, T.; Wu, Z.; Deng, H.; Ge, J.; Chen, K.; Zhang, K.; Wang, L.; Dou, M.; Lu, L.; Zhu, X.; Lu, T.; Lin, D.; Qiao, Y.; Dai, J.; and Wang, W. 2025c. Expanding Performance Boundaries of Open-Source Multimodal Models with Model, Data, and Test-Time Scaling. arXiv:2412.05271.
- Chu, T.; Zhai, Y.; Yang, J.; Tong, S.; Xie, S.; Schuurmans, D.; Le, Q. V.; Levine, S.; and Ma, Y. 2025. SFT Memorizes, RL Generalizes: A Comparative Study of Foundation Model Post-training. arXiv:2501.17161.
- Comanici, G.; Bieber, E.; Schaekermann, M.; Pasupat, I.; et al. 2025. Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities. arXiv:2507.06261.
- Dao, T. 2023. FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning. arXiv:2307.08691.
- Faysse, M.; Sibille, H.; Wu, T.; Omrani, B.; Viaud, G.; Hudelot, C.; and Colombo, P. 2024. Colpali: Efficient document retrieval with vision language models. In *The Thirteenth International Conference on Learning Representations*.
- Gao, T.; Yen, H.; Yu, J.; and Chen, D. 2023. Enabling Large Language Models to Generate Text with Citations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; et al. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. arXiv:2501.12948.
- Hu, A.; Xu, H.; Zhang, L.; Ye, J.; Yan, M.; Zhang, J.; Jin, Q.; Huang, F.; and Zhou, J. 2024. mPLUG-DocOwl2: High-resolution Compressing for OCR-free Multi-page Document Understanding. arXiv:2409.03420.
- Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; Chen, W.; et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2): 3.
- Kojima, T.; Gu, S. S.; Reid, M.; Matsuo, Y.; and Iwasawa, Y. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213.
- Kwiatkowski, T.; Palomaki, J.; Redfield, O.; Collins, M.; Parikh, A.; Alberti, C.; Epstein, D.; Polosukhin, I.; Devlin, J.; Lee, K.; et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7: 453–466.
- Li, B.; Zhang, Y.; Guo, D.; Zhang, R.; Li, F.; Zhang, H.; Zhang, K.; Zhang, P.; Li, Y.; Liu, Z.; and Li, C. 2024. LLaVA-OneVision: Easy Visual Task Transfer. arXiv:2408.03326.
- Li, Z.; Luo, R.; Zhang, J.; Qiu, M.; Huang, X.-J.; and Wei, Z. 2025. VoCoT: Unleashing Visually Grounded Multi-Step Reasoning in Large Multi-Modal Models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 3769–3798.
- Ma, X.; Lin, S.-C.; Li, M.; Chen, W.; and Lin, J. 2024a. Unifying Multimodal Retrieval via Document Screenshot Embedding. In Al-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Proceedings of the 2024 Conference on Empirical*

- Methods in Natural Language Processing*, 6492–6505. Miami, Florida, USA: Association for Computational Linguistics.
- Ma, X.; Zhuang, S.; Koopman, B.; Zuccon, G.; Chen, W.; and Lin, J. 2024b. VISA: Retrieval Augmented Generation with Visual Source Attribution. arXiv:2412.14457.
- Ni, M.; Yang, Z.; Li, L.; Lin, C.-C.; Lin, K.; Zuo, W.; and Wang, L. 2025. Point-RFT: Improving Multimodal Reasoning with Visually Grounded Reinforcement Finetuning. arXiv:2505.19702.
- Penedo, G.; Kydlíček, H.; Lozhkov, A.; Mitchell, M.; Raffel, C. A.; Von Werra, L.; Wolf, T.; et al. 2024. The fineweb datasets: Decanting the web for the finest text data at scale. *Advances in Neural Information Processing Systems*, 37: 30811–30849.
- Peng, Y.; Zhang, G.; Zhang, M.; You, Z.; Liu, J.; Zhu, Q.; Yang, K.; Xu, X.; Geng, X.; and Yang, X. 2025. LMM-R1: Empowering 3B LMMs with Strong Reasoning Abilities Through Two-Stage Rule-Based RL. arXiv:2503.07536.
- Peng, Z.; Wang, W.; Dong, L.; Hao, Y.; Huang, S.; Ma, S.; and Wei, F. 2023. Kosmos-2: Grounding Multimodal Large Language Models to the World. arXiv:2306.14824.
- Qi, J.; Ding, M.; Wang, W.; Bai, Y.; Lv, Q.; Hong, W.; Xu, B.; Hou, L.; Li, J.; Dong, Y.; and Tang, J. 2025. CogCoM: A Visual Language Model with Chain-of-Manipulations Reasoning. arXiv:2402.04236.
- Rajbhandari, S.; Ruwase, O.; Rasley, J.; Smith, S.; and He, Y. 2021. Zero-infinity: Breaking the gpu memory wall for extreme scale deep learning. In *Proceedings of the international conference for high performance computing, networking, storage and analysis*, 1–14.
- Rasley, J.; Rajbhandari, S.; Ruwase, O.; and He, Y. 2020. DeepSpeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 3505–3506.
- Shao, H.; Qian, S.; Xiao, H.; Song, G.; Zong, Z.; Wang, L.; Liu, Y.; and Li, H. 2024a. Visual cot: Advancing multimodal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning. *Advances in Neural Information Processing Systems*, 37: 8612–8642.
- Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Bi, X.; Zhang, H.; Zhang, M.; Li, Y. K.; Wu, Y.; and Guo, D. 2024b. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. arXiv:2402.03300.
- Shen, H.; Liu, P.; Li, J.; Fang, C.; Ma, Y.; Liao, J.; Shen, Q.; Zhang, Z.; Zhao, K.; Zhang, Q.; Xu, R.; and Zhao, T. 2025. VLM-R1: A Stable and Generalizable R1-style Large Vision-Language Model. arXiv:2504.07615.
- Wang, J.; Kang, Z.; Wang, H.; Jiang, H.; Li, J.; Wu, B.; Wang, Y.; Ran, J.; Liang, X.; Feng, C.; and Xiao, J. 2025. VGR: Visual Grounded Reasoning. arXiv:2506.11991.
- Wu, Q.; Yang, X.; Zhou, Y.; Fang, C.; Song, B.; Sun, X.; and Ji, R. 2025. Grounded Chain-of-Thought for Multimodal Large Language Models. arXiv:2503.12799.
- Xia, J.; Tong, B.; Zang, Y.; Shao, R.; and Zhou, K. 2025. Bootstrapping Grounded Chain-of-Thought in Multimodal LLMs for Data-Efficient Model Adaptation. arXiv:2507.02859.
- Xu, G.; Jin, P.; Wu, Z.; Li, H.; Song, Y.; Sun, L.; and Yuan, L. 2025. LLaVA-CoT: Let Vision Language Models Reason Step-by-Step. arXiv:2411.10440.
- Yang, Y.; He, X.; Pan, H.; Jiang, X.; Deng, Y.; Yang, X.; Lu, H.; Yin, D.; Rao, F.; Zhu, M.; Zhang, B.; and Chen, W. 2025. R1-Onevision: Advancing Generalized Multimodal Reasoning through Cross-Modal Formalization. arXiv:2503.10615.
- Ye, X.; Sun, R.; Arik, S.; and Pfister, T. 2024. Effective Large Language Model Adaptation for Improved Grounding and Citation Generation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 6237–6251.
- Zhang, C.; Wang, Y.; Xu, D.; Zhang, H.; Lyu, Y.; Chen, Y.; Liu, S.; Xu, T.; Zhao, X.; Gao, Y.; Hu, Y.; and Chen, E. 2025. TeaRAG: A Token-Efficient Agentic Retrieval-Augmented Generation Framework. arXiv:2511.05385.
- Zhang, H.; Shen, S.; Xu, B.; Huang, Z.; Wu, J.; Sha, J.; and Wang, S. 2024. Item-difficulty-aware learning path recommendation: From a real walking perspective. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 4167–4178.
- Zhong, X.; Tang, J.; and Yepes, A. J. 2019. Publaynet: largest dataset ever for document layout analysis. In *2019 International conference on document analysis and recognition (ICDAR)*, 1015–1022. IEEE.

Look As You Think: Unifying Reasoning and Visual Evidence Attribution for Verifiable Document RAG via Reinforcement Learning

Supplementary Material

Visual evidence attribution in visual document retrieval-augmented generation (VD-RAG) requires models to both generate accurate answers and identify supporting evidence within external visual documents.

We propose the LAT framework, which enables Chain-of-Evidence (CoE) reasoning through a two-stage training paradigm. In Stage I, the model is fine-tuned on a set of human-verified CoE annotations to acquire effective reasoning patterns. Next, it undergoes reinforcement learning (RL) with tailored reward functions that guide stepwise attribution. This design enhances both reasoning capability and attribution quality, thereby achieving robust stepwise attribution and consistent performance across single- and multi-image scenarios without relying on process-level ground-truth annotations.

A Experiment Details

Since CoE reasoning requires flexible grounding capabilities and adaptability in low-resource conditions, we conducted preliminary experiments to identify a suitable backbone model. Specifically, we evaluated Qwen2.5-VL and InternVL2.5 by prompting them to answer questions with direct evidence attribution. As shown in Table 8, Qwen2.5-VL achieved a 2.1% improvement in EM under zero-shot prompting, whereas InternVL2.5 exhibited performance degradation. The weaker performance of InternVL2.5 may stem from its reliance on special tokens (e.g., `<box>`, `<ref>`) for grounding, which constrains the flexible grounding necessary for CoE reasoning in zero-shot settings and impedes effective convergence when training under low-resource conditions.

In contrast, Qwen2.5-VL represents bounding boxes in JSON format, allowing more flexible and direct visual evidence attribution without the constraints of special tokens. This design aligns well with the requirements of CoE reasoning. Building on this foundation, we introduce the two-stage training framework to enhance evidence grounding and reasoning consistency.

Stage I Cold Start: We sampled instances from each dataset (Paper-VISA, Wiki-VISA, FineWeb-VISA) and filtered them using the recall metric defined in Section 3 to retain correct samples. Manual correction was then performed to address bounding box drift by adjusting their positions and sizes to ensure alignment with the correct content regions. The final dataset splits are summarized in Table 6 (cf. Table 5 for data usage comparison). We employed LoRA for parameter-efficient fine-tuning and maintained these configurations throughout training.

Stage II RL Training: We employed GRPO as the policy optimization framework, with the reward function defined in Section 3. To ensure training stability, the learning rate was

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

Table 5: Statistics for the entire VISA dataset, where the construction quantity for multi-image settings is consistent with that for single-image settings.

Dataset	# Train-LAT		# Test
	Cold-Start Stage	RL Stage	
Single-Image / Multi-Image			
Wiki-VISA	264 / 463	4,676 / 4,676	3,000
Paper-VISA	271 / 355	5,000 / 3,562	2,160
Fineweb-VISA	108 / 124	2,000 / 2,000	-

Table 6: Training subsets sampled from the Paper- and Wiki-VISA datasets under single- and multi-image settings.

set to $5e-5$ for both Paper- and Wiki-VISA, and training was conducted for one epoch. We further configured 8 and 6 rollouts for Paper- and Wiki-VISA, respectively, to enhance sampling diversity. The maximum prompt and completion lengths were limited to 16,384 and 600 tokens.

For the stepwise attribution reward R_{step} , we required $R_{\text{acc}} \geq \epsilon$, which zeroes the step rewards for incorrect answers and thereby reduces the advantage of erroneous samples. As specified in Equation 3, if the predicted answer a exactly matches the ground truth a_{gt} ($\text{EM} = 1$), the recall is 1 and the reward is set to 1. For non-exact matches, recall was computed as the proportion of overlapping tokens between a and a_{gt} , normalized by the ground truth length (Equation 2). Manual evaluation of a representative subset indicates that answers require at least 80% token overlap to preserve semantic completeness ($\gamma = 0.8$). This corresponds to the accuracy reward R_{acc} of approximately 0.4 ($\frac{\mathbb{I}(\text{EM}(a, a_{gt})=1) + \text{Recall}(a, a_{gt})}{2} = \frac{0+0.8}{2}$). Accordingly, we set $\epsilon = 0.4$ as the threshold, ensuring that faithful reasoning towards the answer is rewarded.

To promote attribution diversity across reasoning steps, we applied IoU-based constraints with a threshold $\gamma = 0.5$. A predicted box was regarded as valid evidence if its IoU with the corresponding ground truth exceeded 0.5, indicating sufficient relevance. Conversely, pairwise IoU among predicted boxes was required to remain below 0.5, ensuring that distinct reasoning steps relied on independent, non-overlapping evidence regions.

Implementation Details: All experiments are conducted using Qwen2.5-VL-7B-Instruct on 4xA800 (80 GB) GPUs,

Stage	Parameter	Value
Cold Start	γ	0.8
	learning rate	1e-4
	batch Size	8
	epoch(s)	2
	lora rank (r)	64
	lora scaling (α)	64
	lora dropout	0.05
RL Training	learning rate	5e-5
	batch size	16 / 8
	epoch(s)	1
	rollout	8 / 6
	τ	0.3
	δ	0.5
	ϵ	0.4
	lora rank (r)	64
	lora scaling (α)	64
	lora dropout	0.05

Table 7: Hyperparameter Settings for Training. Values marked with A / B correspond to different configurations used for Paper- and Wiki-VISA, respectively.

requiring approximately 20 GPU hours for single-domain training and 48 GPU hours for multi-domain settings. Training leverages mixed precision with DeepSpeed (Rasley et al. 2020; Rajbhandari et al. 2021) and CPU offloading for memory efficiency. To mitigate computational overhead, we employed Flash-attention2 (Dao 2023) and gradient checkpointing techniques. Detailed hyperparameters for different training stages are presented in Table 7. The pseudocode for LAT is provided in Algorithm 1. For multi-image scenarios, initialize the model from the single-image RL-trained parameters, then further apply SFT on multi-image CoE data in $\mathcal{D}_{\text{final}}$ following the above procedures.

For reproducibility, we fixed the random seed to 3407 during training and disabled sampling at evaluation by setting `do_sample=False`. In contrast to VISA (Ma et al. 2024b), we maintained original image resolutions since our CoE reasoning framework requires comprehensive visual understanding across the entire image.

B Dataset Construction

VISA Datasets

1) Wiki-VISA: Selenium renders Wikipedia pages for Natural Questions (NQ) (Kwiatkowski et al. 2019) query-answer pairs, with HTML elements (containing answers) and their bounding boxes annotated.

2) Paper-VISA: Based on PubLayNet (Zhong, Tang, and Yepes 2019). Vision-language models (VLMs) generate QA pairs from layout-annotated scientific documents, with answer-region bounding boxes extracted.

3) FineWeb-VISA: Sampled from FineWeb-edu (Penedo et al. 2024). VLMs generate queries/short answers for educational webpage passages longer than 50 tokens, with screenshots and answer-region bounding boxes.

Algorithm 1: LAT Framework

Input: Sampled query set $\mathcal{Q}_{\text{cold_start}}$ and \mathcal{Q}_{RL}
Parameter: Learning rates η_1, η_2 ; Reward thresholds $\tau, \delta, \epsilon, \gamma$
Output: LAT model π_θ

```

1: Stage I: Cold-start Supervised Fine-tuning
2: for each query  $q \in \mathcal{Q}_{\text{cold\_start}}$  do
3:   Generate CoE reasoning traces using Gemini2.5 pro
4:   Filter outputs by the recall metric
5:   if  $\text{Recall}(a, a_{gt}) \geq \gamma$  then
6:      $\mathcal{D}_{\text{final}} \leftarrow \mathcal{D}_{\text{final}} \cup \{q\}$ 
7:   end if
8: end for
9: Manually verify bounding boxes and format
10: Train model  $\pi_\theta$  via SFT on verified data  $\mathcal{D}_{\text{final}}$  with learning rate  $\eta_1$ 
11: Stage II: Reinforcement Learning with GRPO
12: for each training step do
13:   Sample a batch of queries  $q \sim \mathcal{Q}_{\text{RL}}$ 
14:   Model  $\pi_\theta$  generates CoE reasoning steps  $\{r_1, \dots, r_T\}$  and bounding boxes  $\{B_1, \dots, B_T\}$  for each query  $q$ 
15:   Compute reward  $R = R_{\text{acc}} + R_{\text{step}} + R_{\text{ground}} + R_{\text{format}}$ 
16:   Update policy  $\pi_\theta$  using GRPO with  $R$  and learning rate  $\eta_2$ 
17: end for
18: return final LAT model  $\pi_\theta$ 

```

Multi-Image Data Construction

To simulate real-world VD-RAG scenarios, VISA constructs a multi-image document experimental environment after obtaining the query-document-answer-bounding box triplets. Specifically, given a query q , VISA employs a retriever to retrieve top- k candidate documents, then randomly samples $m-1$ hard negative candidates that do not contain the ground truth. These negative samples are combined with one source document containing the correct answer to serve as input for the multi-image scenarios. VISA deliberately avoids directly utilizing the top- m retrieval results to prevent bias toward specific retrievers or candidate document positions, thereby ensuring methodological generalizability.

To evaluate the model’s capability in handling no-answer scenarios, VISA randomly replaces the source document in the candidate set with a 20% probability, simulating realistic cases where the retriever fails to return documents containing the correct answer. In the specific experimental setup, VISA leverages the Document Screenshot Embedding (DSE) model as the retriever. The parameters are set to $k = 20$ and $m = 3$. Example cases are shown in Figure 5.

Training Data Preprocess

We employed Gemini2.5 Pro (Comanici et al. 2025) to generate CoE data. We used the prompt shown in Figure 9 to guide stepwise evidence attribution during generation.

We subsequently applied a recall metric to measure response against reference answers for preliminary sample filtering. To ensure data quality, we manually reviewed and corrected instances of irrelevant attributions, spatial coordinate offsets, and format inconsistencies. The resulting dataset $\mathcal{D}_{\text{final}}$ is then employed for cold-start training. Examples are illustrated in Figure 6 and 7.

Models	Param.	Wiki-VISA (Single)		Paper-VISA (Single)		Wiki-VISA (Multi)		Paper-VISA (Multi)	
	Size	EM	IoU@0.5	EM	IoU@0.5	EM	IoU@0.5	EM	IoU@0.5
Open-source Models, Direct Answer									
InternVL2.5	8B	45.9	-	42.4	-	37.7	-	32.1	-
Qwen2.5-VL	7B	67.7	-	38.2	-	54.1	-	34.6	-
Open-source Models, Direct Answer & Attribution (DA)									
InternVL2.5	8B	45.1	0	42.0	0.60	35.0	1.67	31.4	1.94
Qwen2.5-VL	7B	69.4	0.80	43.8	2.87	52.5	0.97	37.2	5.56

Table 8: Performance comparison of various models for direct answer generation and attribution tasks. “DA” refers to the zero-shot prompting setting for direct answer & attribution without training.

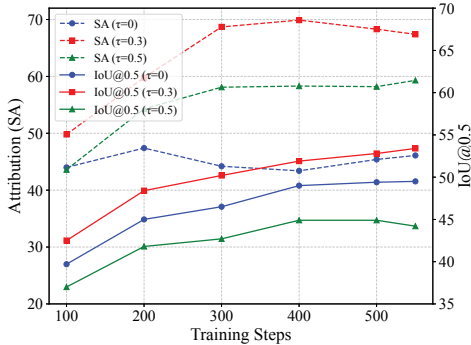


Figure 4: Performance variation with different τ (Wiki).

C Prompt

We designed task-specific prompts adapted to reasoning strategies and model architectures. For models such as LLaVA-CoT (Xu et al. 2025) and R1-OneVision (Yang et al. 2025), which are trained to take questions directly as input, with no additional prompting. Their performance is evaluated by extracting content from the generated outputs, specifically from segments enclosed within the `<CONCLUSION>` and `<answer>` tags. Since both Gemini and InternVL generate bounding boxes in normalized coordinates. During evaluation, these normalized values are converted to absolute coordinates based on the size of the source images. For other models, we employed prompts that initiate responses with “The answer is:”, which encourages concise answers without irrelevant explanations.

Each bounding box must be paired with its source image index using the JSON format `{"bbox_2d": [x_1, y_1, x_2, y_2], "image_index": i}`. As illustrated in Figure 10, we employed a consistent prompt format across both the cold-start and RL training stages.

D Baseline

Our results demonstrate that LAT consistently outperforms open-source models of comparable or larger scales, including InternVL2.5-8B (Chen et al. 2025c), Qwen2.5-VL-7/32B-Instruct (Bai et al. 2025a), LLaVA-OneVision (Li et al. 2024), mPLUG-DocOwl2 (Hu et al. 2024), LLaVA-CoT (Xu et al. 2025), and R1-OneVision (Yang et al. 2025). We note that LLaVA-CoT is trained exclusively on single-

image datasets without multi-image reasoning supervision and is therefore excluded from the multi-image evaluation.

LAT achieves superior performance compared to certain closed-source models, including Qwen-VL Max (Bai et al. 2023) and Gemini 2.0 Flash (Anil et al. 2025). In our experimental setting, VISA-7B is performed with the open-source implementation available on HuggingFace¹. GCoT (Wu et al. 2025) and VisCoT (Shao et al. 2024a) rely on single-image datasets with costly step annotations (e.g., 438k) and are not trained in multi-image settings. The EM of VisCoT on Paper/Wiki-VISA is 10.5/8.8%, illustrating the domain gap.

E Further Analysis

Sensitivity analysis of attribution threshold τ .

To ensure the alignment quality of visual evidence-text pairs in CoE reasoning, we set a threshold τ in the stepwise attribution recall reward function to distinguish positive from negative samples. Using the synthetic CoE data $\mathcal{D}_{\text{final}}$, we computed the values of semantic similarity on manually corrected annotations. As shown in Figure 8, we recorded the maximum and minimum similarities for samples grouped by answer type in each dataset, excluding pairs with zero similarity. Based on this analysis, we selected 0.3 as the threshold, slightly above the minimum similarity observed in human-corrected annotations.

To examine the effect of different thresholds, we conducted ablation experiments on single-image scenarios from Paper-VISA (Figure 3c) and Wiki-VISA (Figure 4). These experiments reveal that the choice of τ has an impact on reasoning outcomes. For example, a relatively high threshold, such as 0.5, leads to notable declines in both accuracy and IoU@0.5, and the overall process traceability quality (SA) also fails to improve, because such strict thresholds make it difficult to sample sufficient positive instances.

Generalization and Ablation

The results in Table 3 demonstrate LAT’s cross-domain transferability, achieving superior performance over SFT with improvements of 1.7% in EM and 6.2% in IoU@0.5 when transferring from Paper- to Wiki-VISA datasets. In the reverse transfer (Wiki→Paper), a slight drop in IoU@0.5 is

¹<https://huggingface.co/collections/MrLight/visa-rag-with-visual-source-attribution>

Train→Eval	Method	EM	IoU@0.5
Wiki→Paper (Single)	SFT	36.9	11.5
	LAT-Ind.	43.6 \uparrow 6.7	10.4 \downarrow 1.1
Wiki→Paper (Multi)	SFT	34.3	5.5
	LAT-Ind.	44.3 \uparrow 10.0	21.9 \uparrow 16.4

Table 9: LAT demonstrates robust generalization with cross-domain transfer between Paper-VISA and Wiki-VISA.

Dateset	Method	Total	Correct	Acc.
Paper	LAT-Ind.	425	356	0.84
	LAT-Full	425	343	0.81
Wiki	LAT-Ind.	578	267	0.46
	LAT-Full	578	283	0.49

Table 10: The accuracy of correctly detecting no-answer cases in multi-image settings.

observed. This can be attributed to the training bias on Wiki-VISA, which contains high-resolution images with relatively dispersed layouts, whereas Paper-VISA features compact medical documents. Such dense layouts often lead to incomplete evidence localization. Nevertheless, the model maintains strong reasoning consistency and answer accuracy, preserving the vanilla model’s performance on both Wiki-VISA (EM: 67.7%) and Paper-VISA (EM: 38.2%). The generalization is even more apparent in multi-image scenarios (Table 3 and 9). These results indicate that LAT improves transfer effectiveness while retaining adaptability across diverse document types.

We also conducted further ablation studies. Without the cold-start stage, the performance achievable by RL is fundamentally capped (EM, 43.9%; IoU, 24.1%; SA, 33.4% on Paper-VISA). Moreover, the learning curves indicate that RL training for a single epoch has not yet reached its performance ceiling. We therefore extended the training to 2 epochs on Paper-VISA, achieving improved results (EM, 46.5 \uparrow 1.1%; IoU, 50.1 \uparrow 0.2%; SA, 38.1 \uparrow 2.6%).

F Case Study

Figure 12–17 present representative examples of model responses. Overall, LAT improves the traceability of the reasoning process. Based on the evaluation of several model responses, we find that LAT generates coherent Chain-of-Evidence (CoE) reasoning traces while maintaining general QA performance and attribution precision. The model is guided by the prompt and CoE training data to directly locate content relevant to the query, following a coarse-to-fine observation process as shown in Figures 12 and 13. In Wiki-VISA, where reasoning requires searching across different layouts, CoE enables the model to verify answer correctness by progressively narrowing down evidence. Figures 16 and 17 show the results under multi-image scenarios. Table 10 reports the accuracy of correctly detecting no-answer cases in multi-image settings.

G Limitation and Future Work

While LAT substantially improves stepwise attribution and overall performance, several practical aspects remain open for future enhancement. First, the thresholds used in the stepwise attribution rewards are manually set based on annotated examples, a choice that offers training stability. Exploring adaptive thresholding mechanisms, such as confidence-aware scaling, may improve flexibility and enhance alignment robustness. Second, LAT is currently trained only on the VISA dataset, whose evidence annotations predominantly follow a single-hop structure. As a result, LAT’s ability to generalize to multi-hop or cross-source evidence reasoning is not yet fully evaluated or supervised. Incorporating datasets with explicit multi-hop evidence chains would provide the necessary supervision to assess and strengthen LAT’s capacity for longer reasoning trajectories and cross-page evidence composition.

Paper-VISA

Diagnosis 2017, 5, 25

4 of 8

conformity (SDOC), where the manufacturer is responsible for ensuring that the product complies with the relevant requirement and then produces a written self-declaration statement [17].

3. Results

3.1. Prostate Volume and Pain Score

Figure 2 shows the changes in PV in CP/CPPS patients. In the control group, the mean prostate volume increased from 30.77 ± 6.436 to 31.58 ± 7.138 mL at the end of the study period, whereas in the treatment group the mean PV decreased from 31.75 ± 7.009 to 27.07 ± 4.522 mL. For the treatment group, the z value was -5.592 at the significance level p value < 0.001 . These data indicated that the therapeutic device reduced the prostate volume significantly, whereas in the control group the prostate volume increased.

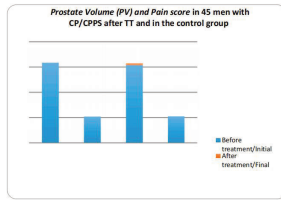


Figure 2. Dynamics of Pain scores and Prostate Volume (PV) mL in 45 men with CP/CPPS after TT and in the control group measured by the National Institute of Health-Chronic Prostatitis Symptom Index (NIH-CPSI) and ultrasound.

Figure 2 also shows the changes in pain score in CP/CPPS patients at the beginning and at the end of the study. In the control group, the mean of pain score decreased from 10.49 to 9.71 at the end of the study period, whereas in the treatment group the mean of pain score decreased from 10.38 to 8.58. In the treatment group, the z value was -5.725 at the significance level p value < 0.001 . These data suggest that the pain score decreased in groups. However, pain in the 'treatment' group decreased considerably while in the 'no treatment' group it only decreased slightly.

3.2. Quality of Life and Maximum Urinary Flow Rate

We assessed the QoL according to NIH-CPSI (see Figure 3). In the control group, the mean QoL decreased slightly from 8.47 to 8.33, whereas in the treatment group the mean QoL decreased from 8.11 to 2.98. For the treatment group, the z value was -5.661 at the significance level p value < 0.001 . These results indicated that the treatment with therapeutic device decreased the QoL significantly while in the control group it decreased slightly.

Wiki-VISA

SEC Championship Game

The SEC Championship Game (officially The SEC Championship on CBS presented by Dr Pepper for broadcast and sponsorship purposes) is an annual American football game that has determined the Southeastern Conference's season champion since 1992. The championship game pits the SEC Western Division regular season champion against the Eastern Division regular season champion. It is typically played on the first Saturday of December.

The SEC Championship Game (officially The SEC Championship on CBS presented by Dr Pepper for broadcast and sponsorship purposes) is an annual American football game that has determined the Southeastern Conference's season champion since 1992. The championship game pits the SEC Western Division regular season champion against the Eastern Division regular season champion. It is typically played on the first Saturday of December.

While the SEC members have played in the game, only six have won: Florida, Georgia, and Tennessee of the Eastern Division, and Auburn, Arkansas, and LSU of the Western Division. Each of these teams has won the championship multiple times. Auburn is the current SEC champion.

The SEC Championship Game has been aired on CBS since 1992, the current broadcast home is ESPN+ as the play is being streamed, only Auburn as the color analyst, and Mike LeFlore as the sideline reporter. Since 2007, the game is held the first Saturday of December at 4:00 PM Eastern.

History

The SEC was the first NCAA conference in any division to hold a football championship game that was awarded from NCAA under its regular season contests. Such a game was made possible by two separate developments. The first came in 1985, when the NCAA membership approved a proposal sponsored by the Division I Football Fairs, Athletic Conference and Central Intercollegiate Athletic Association that allowed any conference with 12 football members to split into divisions and stage a championship game between the division winners. The SEC took advantage of this rule by adding the University of Arkansas and the University of South Carolina in 1985, bringing the conference membership to 12, and adding two new football divisions. The first game was played at the Georgia Dome in 1985, where the conference decided that football champion (the first being the Big 12 in 1983).

The first two SEC Championship Games were held at Legion Field in Birmingham, Alabama. From 1984 until 2018, the game has been played at the Georgia Dome in Atlanta. With the Georgia Dome scheduled to be demolished after the 2018 season, the SEC chose to move the game to Auburn at the Jordan-Hare Stadium, described here. Auburn has been the site of the game since 2019.

In 2018, Auburn and Florida met in the SEC Championship Game for the first time in nearly half a century. The game was the first to be played at the Jordan-Hare Stadium. Auburn won 35–13 and earned a berth in the 2019 SEC Football Championship Game, which it won on its own field.

Auburn and Florida met in the 2019 SEC Championship Game on November 24th year in the SEC. Auburn won the game 35–13, breaking the previous record of 44 points for most points scored by a single team in the SEC Championship Game (previously set by Auburn in 2015).

In recent years, the game has been nicknamed the "national semifinals" for college football. The 2009, 2009, and 2013 games were exclusively national semifinal games, as both participating teams were usually guaranteed a berth in the BCS national championship game as well. The 2013 game was the first of an era of no such at the time of the game, but results of other games later this year meant that it effectively was such a semifinal.

Between 2008 and 2013 the venue of the SEC Championship Game went on to play in the BCS National Championship Game eight straight years, making it a BCS game. Since 2014, the SEC Championship Game venue has gone on to appear in the College Football Playoff every season, and as such, the SEC Championship Game has been labeled a "national semifinal", since the winner (Alabama each time) has advanced to the CFP semifinals.

Results

Results from all SEC Championship games that have been played. Ratings are from the AP 100 released prior to each game.

Year	Eastern Division	Western Division	Site	Attendance	TV Rating	MVP
1992	#12 Florida	#2 Alabama	Legion Field + Birmingham, AL	83,001	5.8	QB Dennis Longhorn, Auburn
1993	#9 Florida	#16 Auburn		78,345		QB Terry Davis, Florida
1994	#6 Florida	#3 Auburn		74,751	10.5	OT Josh Jennings, Florida
1995	#2 Florida	#3 Auburn		71,325	7.2	QB Cherry Wuerffels, Florida
1996	#6 Florida	#1 Auburn		74,122	7.0	QB Cherry Wuerffels, Florida
1997	#3 Tennessee	#1 Auburn		74,586		QB Peyton Manning, Tennessee
1998	#1 Tennessee	#4 Mississippi State		74,790		QB Thomas Price, Tennessee
1999	#5 Florida	#7 Arkansas		71,900		WR Freddie Milers, Auburn
2000	#7 Florida	#8 Auburn		75,427		QB Tim Lincecum, Florida
2001	#2 Tennessee	#21 LSU		74,803	7.0	QB Matt Music, LSU
2002	#4 Georgia	#22 Auburn		75,895	3.2	QB David Connor, Georgia
2003	#5 Georgia	#3 LSU		74,863	4.1	RB Justin Vincent, LSU
2004	#25 Tennessee	#3 Auburn		74,892	4.8	QB Jason Campbell, Auburn
2005	#12 Georgia	#3 LSU		75,717	3.9	QB J. J. Shockey, Georgia
2006	#4 Florida	#6 Auburn	Georgia Dome + Atlanta, GA	75,374	4.7	WR Percy Harvey, Florida
2007	#14 Tennessee	#1 LSU		75,832	6.0	QB Matt Fretwell, Florida
2008	#2 Florida	#1 Auburn		75,850	10.4	QB Tim Lincecum, Florida
2009	#1 Florida	#2 Auburn		75,824	12.8	QB Long McRedd, Auburn
2010	#3 South Carolina	#1 Auburn		75,862	6.3	QB Tim Lincecum, Auburn
2011	#12 Georgia	#1 LSU		74,635	7.4	QB Tyrone Matthews, LSU
2012	#3 Auburn	#2 Alabama		75,824	10.0	RB Blake Lery, Auburn
2013	#5 Missouri	#2 Auburn		75,832	8.7	RB Tim Mason, Auburn
2014	#14 Missouri	#3 Alabama		72,126	7.7	QB Blake Sims, Auburn
2015	#28 Florida	#2 Auburn		76,320	6.95	RB Derrick Henry, Auburn
2016	#25 Florida	#1 Auburn		74,432	7.0	LB Justin Feltz, Auburn

Results by team

Appearances	School	Wins	Losses	Win %	Years Won
12	Florida	7	5	.583	1992, 1994, 1996, 1998, 2000, 2008, 2009
12	Alabama	7	5	.583	1992, 1999, 2004, 2005, 2006, 2010, 2011, 2013
6	LSU	4	1	.800	2005, 2006, 2007, 2013
5	Auburn	3	2	.600	2004, 2010, 2013
5	Georgia	2	3	.400	2002, 2003
5	Tennessee	2	3	.400	1997, 1998
3	Arkansas	0	3	.000	
2	Missouri	0	2	.000	
1	Mississippi State	0	1	.000	
1	South Carolina	0	1	.000	

ID: PMC5750536_00004.jpg

Question: What is the significance level of the increase in Qmax mL/s in the treatment group?

Short Answer: $p < 0.001$

Candidates: ["PMC5750536_00004.jpg",
"PMC3279134_00007.jpg",
"PMC4937636_00005.jpg"]

Pos_idx: 0

Bounding_box: [76, 326, 518, 381]

ID: -7360365691130648166

Question: Who won the first sec championship in football

Short Answer: Alabama

Candidates: ["9005458971896989335",
"-8782797140655775746",
"5099534616456568876"]

Pos_idx: -1 (No answer for multi-candidate setting)

Bounding_box: [24, 2157, 940, 2215]

Figure 5: Data examples from Paper-VISA (left) and Wiki-VISA (right). Each image is assigned a unique identifier, with every dataset entry containing a reference image paired with a specific question, ground-truth answer, and answer source localized by a bounding box. In multi-image scenarios, a retriever selects two images and appends their IDs to the reference image, forming a candidate list. For example, the red bounding box (left) indicates the answer source, where $\text{pos_idx}=0$ signifies that the reference image occupies the first position in the candidate list. For entries lacking ground-truth answers (right), the reference image is substituted with an irrelevant image in the candidate list ($\text{pos_idx}=-1$).

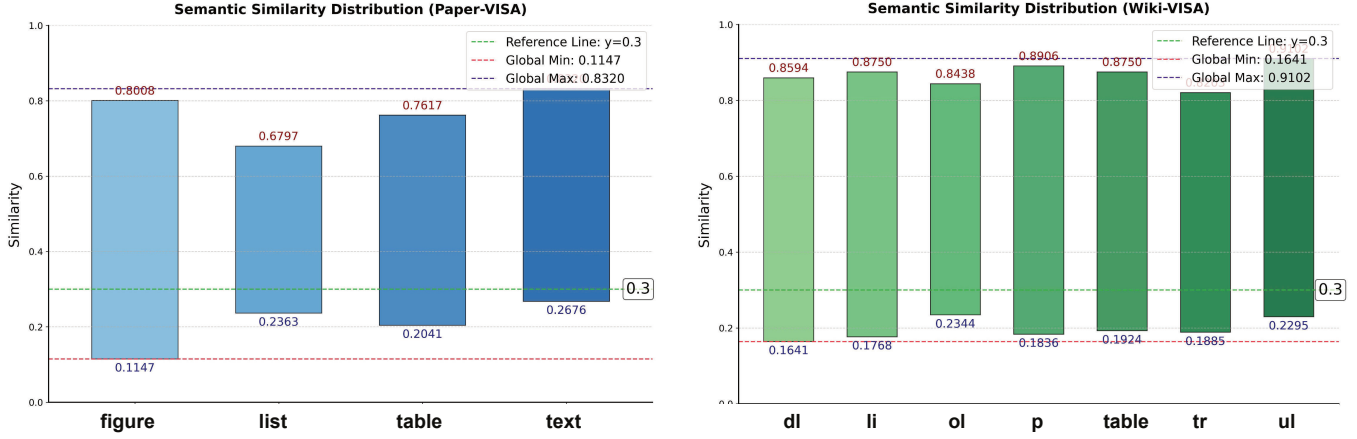


Figure 8: Semantic similarity distribution of synthetic CoE data. We computed the mean upper and lower bounds of semantic similarity between manually annotated evidence attributions and context-evidence pairs across different categories of synthetic CoE data in the cold-start stage. To ensure the effectiveness of the reward function, we selected a threshold $\tau = 0.3$ in Equation 8 to enforce well-aligned CoE reasoning steps during RL training.

<image> Image Size: ()

Task Description:

Given images of document pages, your task is to solve a document reasoning problem. Please first think about the long reasoning process, and then provide the user with the answer.

Restriction:

1. The reasoning process moves from the whole to the details, aligning with human observation.
2. While reasoning, the assistant needs to detect precise evidence for ****each item**** by using a 2D bounding box: For each item (such as figure, table, or factual information), the assistant should attach a 2D bounding box in “box_2d” and the image index number in “Image_index”.
3. Finally, provide the final answer and locate the source of the answer via a 2D bounding box with the corresponding image index. If no answer is found in the images, please indicate “No answer.” Here is an example:

Example: <example>

Question: <question>

Figure 9: Prompt template used for CoE Generation with Gemini2.5 pro.

```

<|im_start|>user
<image> Image Size: ()
Task Description:
Given a document image and a relevant question, analyze the image to extract information relevant to the question, then
provide the final answer. Finally, please locate the source of the final answer via a bounding box with an image index.
Restriction:
1. For each identified element (e.g., figures, tables, or factual text) during analysis, provide a bounding box and include
its image index to highlight the visual evidence.
2. The analysis and final answer are enclosed within <think> </think> and <answer> </answer> tags, re-
spectively, i.e., <think> analysis with visual evidence. </think><answer> the final answer and the corresponding
bounding box as its source. </answer>
3. Each bounding box must be formatted as:
Bounding box: {bbox_2d: [x1, y1, x2, y2], image_index: i}
Question: <question> <|im_end|>
<|im_start|>assistant

```

Figure 10: Prompt template used for training and inference (Single-image) with Qwen2.5-VL. The structured prompt encourages attribution-aware reasoning.

```

<|im_start|>user
<image> Image Size: ()
<image> Image Size: ()
<image> Image Size: ()
Task Description:
Given document images and a relevant question, analyze the images to extract information relevant to the question, then
provide the final answer. Finally, please locate the source of the final answer via a bounding box with an image index.
Restriction:
1. For each identified element (e.g., figures, tables, or factual text) during analysis, provide a bounding box and include
its image index to highlight the visual evidence.
2. The analysis and final answer are enclosed within <think> </think> and <answer> </answer> tags, re-
spectively, i.e., <think> analysis with visual evidence. </think><answer> the final answer and the corresponding
bounding box as its source. </answer>
3. Each bounding box must be formatted as:
Bounding box: {bbox_2d: [x1, y1, x2, y2], image_index: i}
Question: <question> <|im_end|>
<|im_start|>assistant

```

Figure 11: Prompt template used for training and inference (Multi-image) with Qwen2.5-VL. The structured prompt encourages attribution-aware reasoning.

Table 3. Univariate factors associated with antibiotic prescribing for ARIs (excluding pneumonia), Minya District, Egypt, 2011.

Characteristic	Child Patients (under 18 years) Prescribed Antibiotics No. (%)	χ^2 p	Adult Patients Prescribed Antibiotics No. (%)	χ^2 p
Healthcare facility type				
Tertiary care hospitals	109 (90.8)	<0.01	53 (88.3)	>0.05
Primary healthcare units	70 (72.2)		45 (84.9)	
Age groups				
<5 y	114 (79.2)	>0.05	—	—
6–20 y	65 (89)		15 (78.9)	>0.05
21–40 y	—		40 (81.6)	
41–60 y	—		30 (93.8)	
≥60 y	—		12 (100)	
Gender				
Male	87 (80.6)	>0.05	36 (92.3)	>0.05
Female	92 (84.4)		62 (83.8)	
Educational level				
No formal education	75 (83.3)	>0.05	40 (92.3)	>0.05
Elementary/primary school	36 (76.6)		19 (81.8)	
Secondary school	62 (84.9)		24 (82.8)	
University	6 (83.7)		8 (80)	
Presenting symptoms				
Fever				
Yes	118 (88.7)	<0.01	22 (88)	>0.05
No	61 (72.6)		76 (86.4)	
Cough				
Yes	119 (77.3)	<0.05	58 (89.2)	>0.05
No	60 (95.2)		40 (83.3)	
Loss of appetite				
Yes	114 (88.4)	<0.01	—	—
No	65 (73.9)		—	—
Sore throat				
Yes	106 (93)	<0.01	58 (85.3)	>0.05
No	54 (69.2)		40 (88.9)	
Difficulty in breathing				
Yes	50 (82)	>0.05	—	—
No	116 (84.7)		—	—
Runny nose				
Yes	102 (81.6)	>0.05	37 (86)	>0.05
No	77 (83.7)		61 (87)	
Caregiver or patient preference to be treated by an antibiotic				
Yes	109 (88.6)	<0.01	59 (95.2)	>0.05
No	34 (68)		21 (79.3)	

Question: How many adult patients were prescribed antibiotics at primary healthcare units?

Bounding Box: [55, 116, 539, 779]

Short Answer: 45

CoE Generation of LAT:

<think>

The question asks for the number of adult patients prescribed antibiotics at primary healthcare units. I need to find Table 3, which presents univariate factors associated with antibiotic prescribing, and look for the row corresponding to "Primary healthcare units" under the "Adult Patients Prescribed Antibiotics" column.

1. Locate Table 3. It is clearly labeled and positioned in the middle of the page.

["bbox_2d": [50, 78, 540, 777], "Image_index": 1]

2. Find the row for "Primary healthcare units". This is the second row under the "Characteristics" column.

["bbox_2d": [58, 208, 285, 224], "Image_index": 1]

3. Look at the value in the "Adult Patients Prescribed Antibiotics" column for this row. The value is 45.

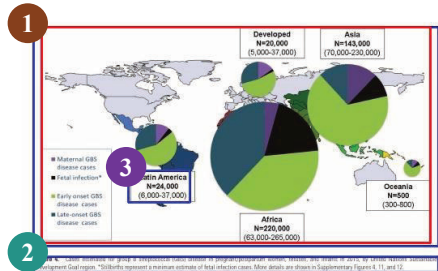
["bbox_2d": [362, 117, 533, 224], "Image_index": 1] </think>

<answer> 45.

["bbox_2d": [50, 78, 540, 777], "Image_index": 1]

</answer>

Figure 12: The Chain-of-Evidence reasoning process of LAT on the Paper-VISA. Blue boxes represent model predictions, while the red box indicates the ground truth bounding box from the dataset. (Case 1)



(UK, 8000–75000) stillbirths, and 67 000 (UK, 12 000–123 000) infant deaths (Figure 7). Serotype III is the most dominant serotype and colonizes 28% of newborn worldwide. It causes 40% of EICGBS, 70% LIGBS, and 24% of maternal GBS disease (Figure 8). A pentavalent vaccine (a/b/c/d/e) would cover 90% of worldwide colonizing maternal and infant GBS disease.

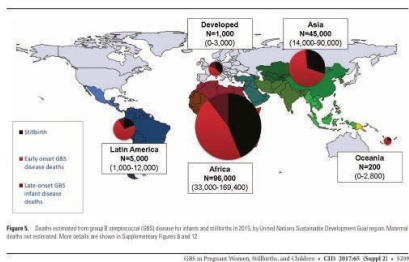


Figure 8. Deaths estimated from group B streptococcal (GBS) disease for infants and stillbirths in 2015, by United Nations Sustainable Development Goal region. Maternal deaths not estimated. More details are shown in Supplementary Figures 8 and 12.

Question: What is the estimated number of GBS disease cases in Latin America in 2015?

Bounding Box: [57, 39, 537, 305]

Short Answer: 24,000

CoE Generation of LAT:

<think>

The question asks about the estimated number of GBS disease cases in Latin America in 2015. I need to find the relevant figure or table that provides this information.

1. Scan the document for figures and tables related to GBS disease cases and Latin America.

Figure 4 seems relevant.

["bbox_2d": [46, 38, 546, 308], "Image_index": 1]

2. Examine Figure 4's caption and content. The caption mentions "Cases estimated for group B Streptococcal (GBS) disease in pregnant/postpartum women, fetuses, and infants in 2015".

["bbox_2d": [46, 321, 546, 340], "Image_index": 1]

3. Look for specific data related to Latin America. The pie chart section labeled "Latin America" shows a value of 24,000 cases.

["bbox_2d": [163, 215, 239, 254], "Image_index": 1]

Therefore, the estimated number of GBS disease cases in Latin America in 2015 is 24,000.

</think>

<answer> 24,000.

["bbox_2d": [46, 38, 546, 308], "Image_index": 1]

</answer>

Figure 13: The Chain-of-Evidence reasoning process of LAT on the Paper-VISA. Blue boxes represent model predictions, while the red box indicates the ground truth bounding box from the dataset. (Case 2)

Figure 15: The Chain-of-Evidence reasoning process of LAT on the Wiki-VISA. Blue boxes represent model predictions, while the red box indicates the ground truth bounding box from the dataset. (Case 4)



Question: How long can the isolation of corneal endothelial cell culture range from?

Bounding Box: -1

Short Answer: No Answer

CoE Generation of LAT:

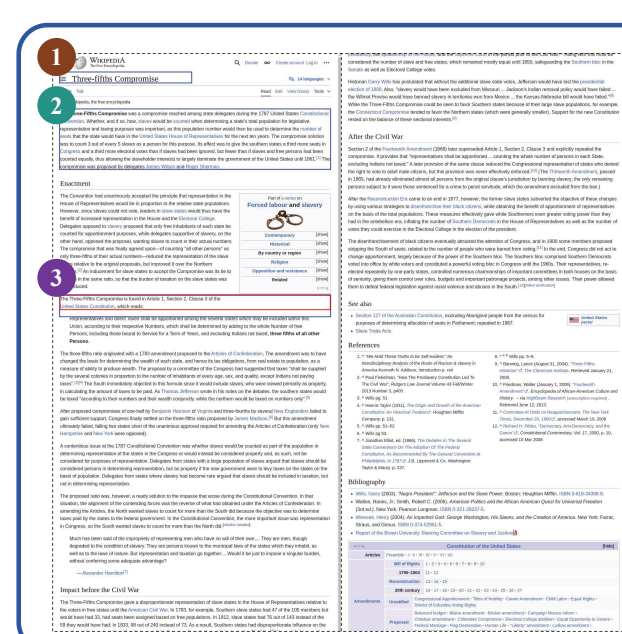
<think> The question asks about the time range for isolating corneal endothelial cell cultures. I need to find the section in the provided images that discusses this process and identify the time frame mentioned.

1. Scan the images for keywords related to "corneal endothelial cells" and "isolation". Image 3 seems most relevant as it mentions culturing of endothelial cells on microcarriers.
2. Examine Image 3 more closely. The section titled "Culturing of endothelial cells on microcarriers" describes the process.
3. Look for information on the time range. The text states that the cells were cultured in DMEM media containing 10 % heat-inactivated fetal bovine serum (GIBCO, Gaithersburg, USA), 100 U/ml penicillin, 100 U/ml streptomycin, and 2 mM glutamine. It also mentions that incubation was conducted at 37 °C in a 5 % CO₂ humidified atmosphere. However, there is no explicit mention of how long the isolation process takes. The text focuses on the culture conditions after isolation, not the duration of the isolation itself.

<answer> No answer.

<answer>

Figure 16: (Multi-image) The Chain-of-Evidence reasoning process of LAT on the Paper-VISA. This is a case in an unanswerable scenario. (Case 5)



Question: Where was the 3/5 compromise originally stated in the constitution ?

Bounding Box: [24, 820, 941, 872]

Short Answer: Article 1, Section 2, Clause 3

CoE Generation of LAT:

<think> The question asks about the original location of the Three-Fifths Compromise in the Constitution. I need to find the section of the document that discusses the Three-Fifths Compromise and identify where it is mentioned in the Constitution.

1. Scan the documents for the term "Three-Fifths Compromise". I can see this term in Document Image 3.
2. Read the introduction of the Three-Fifths Compromise section in Document Image 3. It states that the Three-Fifths Compromise was reached during the 1787 United States Constitutional Convention.
3. Locate the specific clause in the Constitution that mentions the Three-Fifths Compromise. The text says it's found in Article 1, Section 2, Clause 3.

<answer>

Article 1, Section 2, Clause 3 of the United States Constitution.

<answer>

Figure 17: (Multi-image) The Chain-of-Evidence reasoning process of LAT on the Wiki-VISA. Blue boxes represent model predictions, while the red box indicates the ground truth bounding box from the dataset. (Case 6)