

# ALDEN: Reinforcement Learning for Active Navigation and Evidence Gathering in Long Documents

Anonymous ACL submission

## Abstract

While Vision-language models (VLMs) interpret text-rich images effectively, they struggle with reasoning across long, multi-page documents. We present Active Long-DocumEnt Navigation (ALDEN), a multi-turn reinforcement learning framework that fine-tunes VLMs as interactive agents capable of actively navigating long, visually rich documents rather than passive readers. ALDEN features a novel fetch action that allows direct page indexing, complementing the classic search action and better exploiting document structure. To ensure training efficiency and stability, we introduce a rule-based cross-level reward for dense supervision and a visual-semantic anchoring mechanism utilizing dual-path KL-divergence constraints. We train ALDEN on a curated corpus built from open-source datasets where trivial samples are filtered, and queries are rewritten to incentivize multi-turn navigation and fetch usage. Empirically, ALDEN achieves state-of-the-art results on five long-document benchmarks, offering a more accurate and efficient path for long-document understanding.

## 1 Introduction

Visually rich documents (VRDs) are fundamental to real-world knowledge storage, as they combine text, tables, and figures within complex, semantically structured layouts. Unlike plain text, understanding VRDs requires joint reasoning over both textual and visual content and structural organization. This has given rise to the task of visually rich document understanding (VRDU) (Wang et al., 2023; Ding et al., 2022) which aims to automate analysis and question answering across these multi-modal formats, underpinning critical tasks such as scientific multi-modal question answering (Liang et al., 2024) and key information extraction from business documents (Rombach and Fettke, 2024).

While vision-language models (VLMs) excel on single-page or short documents analysis (Xie et al.,

2024; Lv et al., 2023; Hu et al., 2024), they struggle with long documents where full-context processing is computationally prohibitive and noisy (Cho et al., 2024). Current solutions typically adopt Retrieval-Augmented Generation (RAG) pipelines (Cho et al., 2024; Chen et al., 2025a), utilizing retrievers to select query-relevant pages (Faysse et al., 2025) and prompting VLMs to perform fixed subtasks like query reformulation, retrieved content summarization, or final answer synthesis (Han et al., 2025; Wang et al., 2025b). While effective, these systems rely on static reasoning patterns and rigid workflows, limiting their ability to generalize or adapt strategies to diverse user queries. This motivates a shift toward the **Agentic VRDU** (A-VRDU) task, which requires the model to act as an agent that can actively navigate and reason over long documents to deliver accurate and adaptive question answering beyond fixed RAG pipelines.

Recent studies (Chen et al., 2025b; Jin et al., 2025; Song et al., 2025) show that modeling search as an action and optimizing the workflow with outcome-based reinforcement learning (RL) yields more generalizable agents capable of active information gathering. While promising for A-VRDU, adapting this framework to VLMs presents unique challenges. Standard semantic retrieval lacks the precision for queries referencing specific pages or requiring reasoning across consecutive pages. Moreover, document-level information gathering typically demands multi-turn interaction, where sparse and delayed outcome-based rewards fail to reinforce helpful intermediate steps or discourage redundant actions (Li et al.). A further challenge arises from the high-dimensional visual inputs. We empirically observe that fully masking the visual tokens when computing the policy gradient, as done in existing approaches (Jin et al., 2025; Song et al., 2025), leads to unstable training dynamics and can even cause collapse.

These limitations motivate our framework,

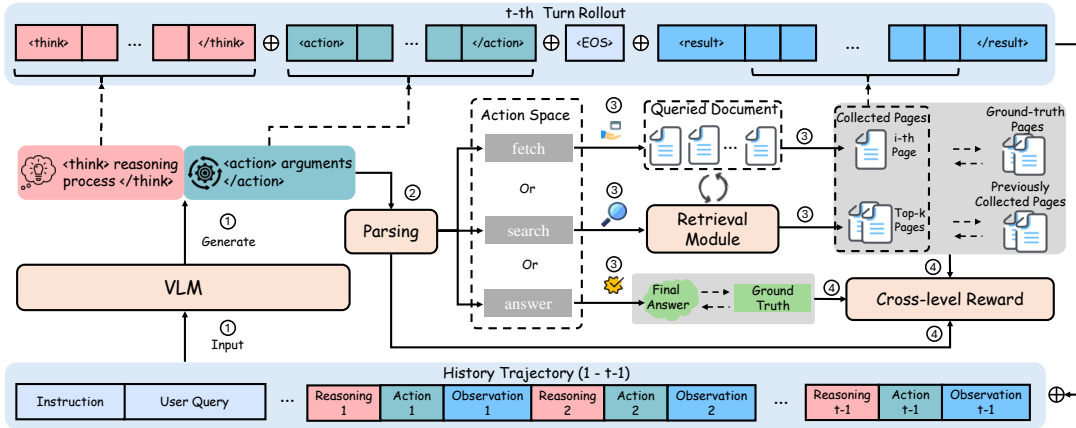


Figure 1: Overview of the rollout process. At each turn: (1) the VLM generates a response conditioned on the dialogue history; (2) the response is parsed into an action (search, fetch, or answer); (3) the action is executed, where search or fetch collect document pages and answer terminates the process; and (4) the cross-level reward function assigns rewards based on execution outcomes and parsing results.  $\oplus$  denotes the concatenation operation.

084 **Active Long-DocumEnt Navigation (ALDEN)**, a  
 085 multi-turn RL framework that trains VLMs as inter-  
 086 active agents for navigation in long, visually-rich  
 087 documents. The overall reasoning-action rollout  
 088 of ALDEN is illustrated in Fig. 1. We expand the  
 089 action space with the *fetch* action, enabling direct  
 090 page access to complement search-based retrieval  
 091 and efficiently handle diverse queries. To overcome  
 092 sparse rewards, we incorporate a *cross-level reward*  
 093 *function*, which integrates rule-based turn-level su-  
 094 pervision with a token-level repetition penalty to  
 095 provide fine-grained process supervision, encourag-  
 096 ing informative evidence collection while discourag-  
 097 ing redundant action invoking. Finally, we incor-  
 098 porate a *visual semantic anchoring* mechanism,  
 099 which constrains the hidden states of generated and  
 100 visual tokens separately during training to preserve  
 101 the grounding of visual-token representations and  
 102 improve overall training robustness.

103 To train ALDEN, we curate a corpus of 30k  
 104 samples from DUDE (Van Landeghem et al.,  
 105 2023), MPDocVQA (Tito et al., 2023a), and Slide-  
 106 VQA (Tanaka et al., 2023a), where we filter trivial  
 107 documents and rewrite the user queries through  
 108 LLMs to incentivize multi-turn navigation and  
 109 the use of fetch actions. Experimental results  
 110 across five benchmarks demonstrate that ALDEN  
 111 achieves state-of-the-art performance, significantly  
 112 outperforming strong baselines. Overall, A-VRDU  
 113 marks a departure from passive processing toward  
 114 autonomous navigation, and ALDEN’s perfor-  
 115 mance validates this scalable framework for robust  
 116 document understanding.

117 Overall, our main contribution can be summa-

rized as follows:

- We propose the agentic visually-rich document  
understanding (A-VRDU) task that aims to de-  
velop agents that can actively navigate and reason  
over long visually rich documents.
- To perform the A-VRDU task, we introduce  
**ALDEN**, a multi-turn RL framework with three  
key components: an expanded action space fea-  
turing a novel *fetch* action, a cross-level reward  
function, and a visual semantic anchoring mech-  
anism, which together enable efficient and robust  
training.
- We construct a training corpus for training the  
A-VRDU agent and conduct extensive experi-  
ments on five commonly used VRDU bench-  
marks, showing that ALDEN significantly out-  
performs the strongest baseline, improving the  
answer accuracy by 9.14% on average.

## 2 Related Work

### 2.1 Visually-rich documents understanding

Existing VLMs (Hu et al., 2024; Xie et al., 2024; Feng et al., 2024; Liu et al., 2024b) achieve high performance on single-page documents (Mathew et al., 2021; Masry et al., 2022) but face scalability issues with long, multi-page contexts (Deng et al., 2024; Ma et al., 2024b). While semantic retrieval mitigates the computational cost of full-context processing (Tito et al., 2023a; Hu et al., 2024), current methods are limited to passive, prompting-based workflows (Han et al., 2025; Wang et al., 2025b). We advance this paradigm by treating

VRDU as an agentic task (A-VRDU), employing RL to train agents that actively navigate and reason over document structures rather than relying on static retrieval pipelines.

## 2.2 RL Training for LLMs/VLMs

The success of RLHF (Ziegler et al., 2019; Ouyang et al., 2022) based on the Proximal Policy Optimization (PPO) algorithm (Schulman et al., 2017) and RL with Verifiable Rewards (RLVR) (Shao et al., 2024) based on the Group Relative Policy Optimization (GRPO) algorithm has inspired a new class of "active" RAG agents that optimize retrieval workflows via reinforcement learning (Jin et al., 2025; Song et al., 2025). While concurrent work VRAG-RL (Wang et al., 2025c) applies RL to enhance visual RAG, it focuses on intra-image perception via cropping and zooming. Consequently, fine-tuning VLMs for A-VRDU remains largely unexplored. Unlike open-domain retrieval, A-VRDU requires exploiting explicit document structure (e.g., page indices), denser supervision to guide multi-turn navigation, and stability against the large number of unconstrained visual tokens introduced by high-resolution document pages, which motivates the development of new RL frameworks tailored for this task.

## 3 Preliminaries

### 3.1 Problem Formulation

We define A-VRDU as an interactive task where an agent answers a query  $q_u$  by navigating a document  $\mathcal{D} = (p_1, p_2, \dots, p_{|\mathcal{D}|})$ , composed of a sequence of pages  $p_i$ . Since direct access to the full document is restricted, the agent must plan a trajectory of navigational or answering steps to generate a final answer  $y'$ . We model this as a Hierarchical Markov Decision Process (MDP) to handle the dual granularity of the task. A **High-Level** MDP manages the discrete action choices at each turn, while a low-level MDP handles the token-by-token realization of those actions. Formally, the high-level MDP operates at the turn level: given a state  $s_t$  encapsulating the interaction history, the agent selects a textual action  $a_t$ , receives pages as an observation  $o_t$ , and a reward  $r_t$ . The transition is a deterministic concatenation:  $s_{t+1} = [s_t, a_t, o_t]$ . The generation of  $a_t$  is handled by the **Low-Level** MDP. Here, the state  $s_t^i = (s_t, a_t^{1:i-1})$  includes the high-level context and tokens generated so far. The agent selects a token  $a_t^i$  from the vocabulary  $\mathcal{V}$ , receives a reward

$r_t^i$  and transitions to  $s_t^{i+1} = [s_t^i, a_t^i]$ .

### 3.2 Joint Policy Optimization

The previously defined hierarchical MDP reduces to a standard RLHF formulation when the horizon is collapsed to a single turn. In this setting, the distinction between high-level planning and low-level generation dissolves, and a unified policy  $\pi_\theta$  generates the complete token sequence directly from the initial state. This policy is typically optimized using PPO algorithm, which maximizes the expected reward via a clipped surrogate objective:

$$\mathcal{L}_{\text{PPO}} = \mathbb{E}_i \left[ \min \left( \rho_i A_1^i, \text{clip}(\rho_i, 1 - \epsilon, 1 + \epsilon) A_1^i \right) \right] \quad (1)$$

where  $\rho_i$  is the probability ratio between the current and old policies. To balance bias and variance in credit assignment, the advantage term  $A_1^i$  is computed using Generalized Advantage Estimation (GAE) (Schulman et al., 2015). This involves estimating the temporal-difference (TD) error  $\delta_i = r_1^i + \gamma_{\text{token}} V_\phi(s_1^{i+1}) - V_\phi(s_1^i)$  using a learned value function  $V_\phi$ , and aggregating these errors over the trajectory:  $A_1^i = \sum_{k=0}^{L-i-1} (\gamma_{\text{token}} \lambda_{\text{token}})^k \delta_{i+k}$ , where  $\gamma_{\text{token}}$  is the discount factor,  $\lambda_{\text{token}}$  controls the trade-off between bias and variance and  $L$  denotes the length of the response. To extend this optimization paradigm from single-turn generation to the multi-turn interaction required for A-VRDU, we propose the ALDEN framework, which is detailed in the following section

## 4 Methodology

We present **Active Long-DocumEnt Navigation** (ALDEN), a framework for training interactive VLM agents to navigate VRDs via a multi-turn reasoning-action loop. To this end, ALDEN introduces three key components. **(i) Expanded action space:** the agent is equipped with both a semantic search action for retrieving pages and a novel fetch action for direct page access, enabling flexible exploitation of document structure (§4.1). **(ii) Cross-level reward function:** supervision is provided jointly at the turn level and the token level, guiding the agent toward effective evidence collection and accurate answer generation (§4.2). **(iii) Visual semantic anchoring:** to stabilize RL training, ALDEN constrains the hidden-state evolution of generated and visual tokens separately, mitigating drift and preserving semantic grounding during optimization (§4.3). The overall RL training pipeline of ALDEN is illustrated in Fig. 2 and Alg. 1.

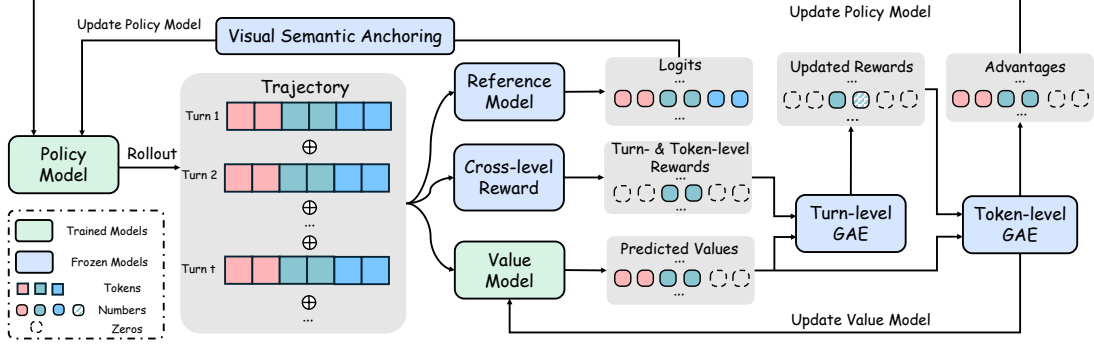


Figure 2: Overview of RL training in ALDEN. The policy model generates multi-turn trajectories, which are scored by a **cross-level reward function** and a **value model**. **Turn-level GAE** integrates future rewards to update the cross-level reward, and **token-level GAE** produces advantages for policy updates. A **reference model** supplies logits for both generated and visual tokens, which the **visual semantic anchoring** mechanism uses to constrain hidden-state evolution during optimization.

#### 4.1 Expanded Action Space

In A-VRDU, agents must navigate using both semantic and structural cues. Standard semantic retrieval struggles with structural dependencies, such as specific page indices (e.g., “see page 12”) or sequential reading. ALDEN bridges this gap by introducing a fetch action for direct page access, complementing the classic search operation. The action space thus consists of three options, each expressed in a structured format that combines free-form reasoning with executable commands as shown in Fig. 1:

- **Search.** The agent generates a reasoning trace within `<think>` tags, followed by a semantic query enclosed in `<search>` tags. An external retrieval module then returns a ranked list of pages based on semantic similarity.
- **Fetch.** Distinct from search, the agent specifies a target page index within `<fetch>` tags to access a page directly, bypassing semantic matching.
- **Answer.** The agent outputs a reasoning trace followed by the final response within `<answer>` tags. This action terminates the rollout and provides the final output for the user query.

Once the action is parsed, the document returns the corresponding page images enclosed within the `<result>` tag. For the search action, the associated page numbers are also returned to provide cues of document structure.

#### 4.2 Cross-level Reward Modeling

Since sparse outcome rewards are insufficient for guiding complex navigation, ALDEN introduces a cross-level reward function to provide dense process supervision. This mechanism operates on two

levels: turn-level rewards for assessing the strategic utility of actions, and token-level rewards for shaping local generation dynamics.

**Turn-level Reward.** The immediate turn-level reward  $r_t$  is defined as  $r_t = f_t + u_t$ , comprising a format constraint  $f_t$  and a utility score  $u_t$ . The format reward  $f_t$  is given by:

$$f_t = \begin{cases} 0, & \text{if the format is correct} \\ -1, & \text{otherwise} \end{cases} \quad (2)$$

Thus, only well-formed responses avoid penalty, enforcing structural validity. The utility reward  $u_t$  evaluates the action outcome. Let  $\mathcal{H}_{t-1}$  denote the history of visited pages,  $\mathcal{G}$  denote the ground-truth evidence pages. For a fetch action targeting page  $p_i$ , the current page set is  $\mathcal{C}_t = \{p_i\}$ . For a search action, the current page set is defined as the collection of the top- $K$  retrieved pages, i.e.,  $\mathcal{C}_t = \{p_1, \dots, p_K\}$ . The reward is defined as:

$$u_t = \begin{cases} \alpha \cdot \text{F1}(y, y') & \text{if } a_t = \text{answer} \\ f_{\text{prox}}(p_i, \mathcal{G}) - \cdot f_{\text{rep}}(\mathcal{C}_t, \mathcal{H}_{t-1}) & \text{if } a_t = \text{fetch} \\ \frac{|\mathcal{C}_t \cap \mathcal{G}|}{|\mathcal{C}_t|} & \text{if } a_t = \text{search} \end{cases} \quad (3)$$

where  $f_{\text{prox}}(p, \mathcal{G}) = \exp\left(-\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} |p - g|\right)$  rewards the geometric proximity of page  $p_i$  to the evidence set  $\mathcal{G}$ , the character-level F1 score evaluates the character-level overlap between the generated answer and the ground-truth answer. The fetch action is subject to a repetition penalty  $f_{\text{rep}} = \frac{|\mathcal{C}_t \cap \mathcal{R}|}{|\mathcal{C}_t|}$ , penalizing the re-acquisition of known information. We scale the reward of answer action by  $\alpha > 1$  to ensure answer quality dominates the cumulative intermediate rewards.

While standard RLHF assigns rewards to the final token of a single response, this approach is

insufficient for multi-turn tasks where actions have delayed consequences. We therefore replace  $r_t$  with a value estimate  $\hat{V}_t$  derived via GAE (Wang et al., 2025a). We first define the turn-level TD error as  $\delta_k = r_k + \gamma_{\text{turn}} V_\phi(s_{k+1}^L) - V_\phi(s_k^L)$ . The effective reward signal is then computed as:  $\hat{V}_t = V_\phi(s_t^L) + \sum_{k=0}^{T-1-t} (\gamma_{\text{turn}} \lambda_{\text{turn}})^k \delta_{t+k}$ , where  $T$  denotes the total number of turns. By replacing  $r_t$  with  $\hat{V}_t$ , we inject long-horizon strategic information into the token-level optimization.

**Token-level Reward.** Unlike the atomic fetch argument (a single page index), the search action generates a multi-token query. Coarse turn-level penalties often fail to identify specific redundant phrases within these queries, leading to inefficient loops. To address this, we introduce a token-level repetition penalty applied specifically to the search query span. For any search action after the first, we quantify redundancy by computing the maximum Jaccard similarity between the n-grams of the current query  $q_t$  and those of all historical queries  $\{q_j\}_{j<t}$ :

$$\text{overlap}_t = \max_{j<t} \frac{|Q_n(q_t) \cap Q_n(q_j)|}{|Q_n(q_t) \cup Q_n(q_j)|} \quad (4)$$

where  $Q_n(q)$  denotes the set of n-grams in query  $q$ . We distribute this penalty to individual tokens to precisely penalize repeated segments. For each token  $u$  in the query span  $a_t^{\text{query}}$ , we assign a weight  $w_u = \frac{c_u}{\sum_{v \in a_t^{\text{query}}} c_v}$ , where  $c_u$  counts the number of overlapping n-grams that contain token  $u$ .

Finally, the reward  $r_t^i$  assigned to each generated token  $a_t^i$  within turn  $t$  is defined by combining turn-level and token-level signals:

$$r_t^i = \begin{cases} \hat{V}_t & \text{if } i = L \\ -w_i \cdot \text{overlap}_t & \text{if } t > 1 \wedge a_t = \text{search} \\ & \wedge a_t^i \in a_t^{\text{query}} \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

This yields a unified cross-level signal that penalizes local redundancy without disrupting global credit assignment. Advantages are subsequently derived by applying token-level GAE to this reward stream.

### 4.3 Visual Semantic Anchoring

While existing studies generally mask observation tokens to isolate the generative action space (Jin et al., 2025; Song et al., 2025), this approach proves insufficient for A-VRDU. Document pages

introduce a massive volume of visual tokens, creating a fundamental vulnerability, i.e., an under-constrained visual manifold. Since the mapping from the visual latent space to the textual action space is non-injective, relying solely on action-space KL divergence leaves a high-dimensional null space in the optimization landscape (verified formally in Appx. B). Consequently, aggressive reward-driven gradients can distort visual representations without triggering the trust region penalty, leading to severe training instability and rapid entropy collapse (see Fig. 3).

To remedy this, we propose Visual Semantic Anchoring, a stabilization mechanism that introduces dual-path KL regularization. Beyond the standard textual constraint, we impose a secondary KL penalty on the visual hidden states. This enforces a smoothness constraint on the visual encoder, preserving Lipschitz continuity and maintaining the robustness of pre-trained features against drastic policy updates. Formally, we define:

$$\begin{aligned} \mathcal{L}_{\text{policy}} = \mathbb{E}_t [\mathbb{E}_i [\min & \left[ \rho_t^i A_t^i, \text{clip} \left( \rho_t^i, 1 - \epsilon, 1 + \epsilon \right) A_t^i \right] \\ & + \beta_{\text{gen}} \text{KL}(\pi_\theta(a_t^i | s_t^i) || \pi_{\text{ref}}(a_t^i | s_t^i))] \\ & + \mathbb{E}_j [\beta_{\text{obs}} \text{KL}(\pi_\theta(o_t^j | o_t^{<j}, a_t, s_t) || \pi_{\text{ref}}(o_t^j | o_t^{<j}, a_t, s_t))] \end{aligned} \quad (6)$$

where  $\beta_{\text{gen}}$  and  $\beta_{\text{obs}}$  are independent coefficients. In practice, we set  $\beta_{\text{obs}} > \beta_{\text{gen}}$  to tightly regularize the much larger observation-token set while allowing more flexibility for generated tokens to adapt to the task.

## 5 Experiments

We conduct experiments on long VRDU benchmarks to (i) compare ALDEN with strong baselines and (ii) assess the contribution of its key components, including expanded action space, cross-level reward, and visual semantic anchoring, to navigation accuracy, answer quality, and training stability. We first outline datasets, baselines, implementation details, and evaluation metrics (§5.1), then present main results (§5.2), followed by ablations (§5.3) and detailed component analyses (§5.4).

### 5.1 Experimental Setup

**Datasets.** We build a challenging training corpus by filtering DUDE (Van Landeghem et al., 2023), MPDocVQA (Tito et al., 2023b), and SlideVQA (Tanaka et al., 2023b) for documents exceeding 3 pages. To enrich query diversity, we use GPT-4o (Hurst et al., 2024) to rewrite part of

MPDocVQA, increasing the proportion of page-index-referenced queries in the final training corpus. The evaluation is conducted mainly on the following VRDU benchmarks: **MMLongBench** (Ma et al., 2024b), **LongDocURL** (Deng et al., 2024), **PaperTab** (Hui et al., 2024), **PaperText** (Hui et al., 2024), and **FetaTab** (Hui et al., 2024). To validate the fetch mechanism, we introduce DUDE-sub, a balanced validation set comprising 960 queries that contain both general and structure-dependent questions (e.g., sequential cues). More details about the dataset can be seen in Appx. A.

**Baselines.** To validate ALDEN’s effectiveness, we compare it with three categories of baselines. (1) **Full-Document Input:** SoTA VLMs prompted with the entire document as context to answer user queries. (2) **Visual RAG:** methods that retrieve the most relevant document pages using the user query, including M3DocRAG (Cho et al., 2024), and an ALDEN variant trained with GRPO adapted from a fully textual method ReSearch (Chen et al., 2025b). (3) **Hybrid RAG:** approaches that augment page images with OCR-extracted text for retrieval and reasoning, including MDocAgent (Han et al., 2025), VidoRAG (Wang et al., 2025b). More details are presented in Appx. C.

**Implementation Details.** Both the policy and value models are initialized from Qwen2.5-VL-7B-Instruct (Bai et al., 2025), and all Visual RAG and Hybrid RAG baselines use the same backbone for fairness. During training, we adopt the single-vector retriever vdr-2b-v1 (Ma et al., 2024a) for images and e5-large-v2 (Wang et al., 2022) for text. For evaluation, we also report results with the multi-vector retrievers ColQwen2-v1.0 (ColQwen) (Faysse et al., 2025) for images and ColBERT-v2.0 (ColBERT) (Santhanam et al., 2021) for text. Unless otherwise noted, each search action retrieves the top-1 candidate page, with a maximum of  $T = 6$  reasoning-action turns. On average, ALDEN collects 1.87 unique pages per query; hence, single-turn RAG baselines are set to retrieve the top-2 pages for a fair comparison. Further implementation details are provided in Appx. D.

**Evaluation Metrics.** The primary evaluation metric is GPT-4o-judged answer accuracy (**Acc**) on each benchmark. For finer-grained analysis of ALDEN’s components, we further assess navigation quality using trajectory-level retrieval recall (**Rec**), precision (**Pre**), F1-score (**F1**), and the number of unique collected pages (**#UP**). Detailed defi-

nitions of these metrics are provided in Appx. E.

## 5.2 Main Results

Table 1 reports answer accuracy across all baselines. Directly prompting large VLMs with the entire document performs poorly ( $\text{Acc} < 0.30$ ), confirming the difficulty of long-document reasoning where irrelevant content overwhelms true evidence. Retrieval-based methods achieve substantially better results. Among Visual RAG approaches, ALDEN with ColQwen attains the highest average accuracy (0.410), surpassing M3DocRAG by 3.2 points. In Hybrid RAG, baselines such as VidoRAG and MDocAgent benefit from textual signals but are limited by fixed reasoning pipelines. ALDEN with hybrid retrievers achieves the best overall performance, exceeding the strongest hybrid baseline by +7.47% relative improvement. These results highlight ALDEN’s ability to generalize across benchmarks by actively collecting and reasoning over evidence, though modest performance on scientific-paper datasets (PaperText, PaperTab) suggests domain knowledge remains a limiting factor. The notably larger gain over GRPO underscores the limitations the limitation of outcome-based rewards for training multimodal agents in multi-turn, long-horizon settings from base VLMs, which is one of the key motivations for this work. Moreover, ALDEN achieves higher accuracy with a multi-vector retriever at inference despite being trained with a single-vector retriever, indicating that strategies learned with a weaker retriever generalize to stronger ones and suggesting a path to more efficient training. A case study demonstrating learned action patterns is provided in Appx. F.

## 5.3 Ablation Study

To understand the contribution of each component in ALDEN, we further conduct ablation studies on the five benchmarks. Table 2 reports the Acc metric results for the full model and three variants: (i) *w/o Fetch*, which removes the index-based fetch action and relies solely on semantic retrieval; (ii) *w/o Cross-level Reward*, which uses only outcome-level supervision without our designed turn- and token-level reward shaping; and (iii) *w/o Visual Semantic Anchoring*, which omits the constraint on visual hidden states during optimization. Removing any component consistently lowers accuracy, with the largest drop from omitting fetch, underscoring the value of direct page-index access. Excluding the cross-level reward also substantially

Table 1: Answer accuracy comparison on five VRDU benchmarks. † indicates the strongest non-ALDEN baseline used to compute the relative improvement (%). **Bold** indicates the best result per dataset.

Method	MMLongBench	LongDocUrl	PaperTab	PaperText	FetaTab	Avg
<i>Full Document Input</i>						
SmolVLM-Instruct (Marafioti et al.)	0.072	0.165	0.065	0.142	0.148	0.118
Phi-3.5-Vision-Instruct (Abdin et al.)	0.141	0.285	0.068	0.174	0.232	0.180
mPLUG-DocOwl2 (Hu et al.)	0.159	0.273	0.072	0.162	0.288	0.191
Qwen2-VL-7B-Instruct (Wang et al.)	0.177	0.280	0.077	0.146	0.339	0.203
LEOPARD (Jia et al.)	0.196	0.313	0.112	0.189	0.341	0.230
Qwen2.5-VL-7B-Instruct (Bai et al.)	0.221	0.375	0.131	0.265	0.336	0.265
InternVL3.5-8B-Instruct (Wang et al.)	0.219	0.381	0.130	0.271	0.348	0.270
<i>Visual RAG methods</i>						
GRPO (ColQwen)	0.274	0.384	0.150	0.295	0.406	0.302
M3DocRAG (ColQwen)†	0.330	0.464	0.201	0.350	0.547	0.378
ALDEN (vdr-2b-v1)	0.335	0.513	0.201	0.342	0.542	0.386
ALDEN (ColQwen)	0.367	0.526	0.211	0.345	0.603	0.410
Relative Improvement (%)	11.21	13.36	4.98	-1.43	10.23	10.81
<i>Hybrid RAG methods</i>						
ViDoRAG (ColQwen + ColBERT)	0.215	0.323	0.158	0.264	0.358	0.264
MDocAgent (ColQwen + ColBERT)†	0.347	0.494	0.221	0.408	0.607	0.415
ALDEN (vdr-2b-v1 + e5-large-v2)	0.385	0.542	0.228	0.416	0.611	0.436
ALDEN (ColQwen + ColBERT)	<b>0.392</b>	<b>0.551</b>	<b>0.245</b>	<b>0.421</b>	<b>0.623</b>	<b>0.446</b>
Relative Improvement (%)	12.97	11.54	10.86	3.18	2.63	7.47

Table 2: Answer accuracy for different ablations of ALDEN on five VRDU benchmarks. **Bold** indicates the best result per dataset.

Method	MMLongBench	LongDocUrl	PaperTab	PaperText	FetaTab	Avg
Full ALDEN	<b>0.335</b>	<b>0.513</b>	<b>0.201</b>	<b>0.342</b>	<b>0.542</b>	<b>0.386</b>
w/o Fetch	0.301	0.469	0.140	0.258	0.443	0.322
w/o Cross-level Reward	0.329	0.483	0.148	0.301	0.518	0.356
w/o Visual Semantic Anchoring	0.326	0.502	0.181	0.328	0.529	0.373

hurts performance, confirming the importance of fine-grained reward shaping. In contrast, removing visual-semantic anchoring causes milder yet consistent degradation. Building on these results, we next provide a detailed component analysis to understand the specific roles of each key design choice in ALDEN.

#### 5.4 Component Analysis

**Fetch vs. Search** To assess the effect of the proposed fetch action, we compare the full ALDEN agent with a *search-only* variant that disables direct page-index access and relies solely on semantic retrieval. Evaluation on the DUDE-sub dataset, which contains explicit page references and structured navigation queries, shows clear benefits of fetch (Table 3). Acc improves from 0.545 to 0.653 and Rec from 0.471 to 0.598, while Pre and F1 also increase, indicating more accurate evidence retrieval. The number of unique pages rises from 1.03 to 1.19, reflecting broader coverage. These results confirm that combining index-based fetch with semantic search enables more flexible and efficient navigation, especially for queries that ref-

Table 3: Comparison between search-only and full ALDEN on the DUDE-sub dataset.

Method	Acc	Rec	Pre	F1	#UP
Search-only	0.545	0.471	0.841	0.531	1.03
Full ALDEN	<b>0.653</b>	<b>0.598</b>	<b>0.874</b>	<b>0.628</b>	<b>1.19</b>

Table 4: Effect of reward design of outcome-based, turn-level and outcome-based, and full ALDEN on Long-DocURL.

Method	Acc	Rec	Pre	F1	#UP
Outcome-based Only	0.483	0.483	0.612	0.520	1.27
Turn-level + Outcome	0.509	0.497	0.608	0.522	1.22
Full ALDEN	<b>0.513</b>	<b>0.506</b>	<b>0.612</b>	<b>0.526</b>	<b>1.39</b>

erence specific pages or require traversal across consecutive pages.

**Effect of Reward Design.** We evaluate how different reward schemes affect ALDEN’s retrieval and reasoning (Table 4). (i) Outcome-based Only assigns a single scalar reward for final answer correctness. (ii) Turn-level + Outcome adds rule-based turn-level supervision, improving Acc from 0.483 to 0.509 and Rec from 0.483 to 0.497, showing that denser feedback aids evidence localization. (iii)

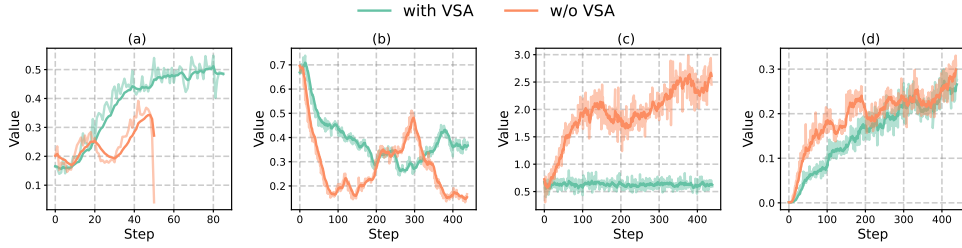


Figure 3: Training dynamics of ALDEN with and without Visual Semantic Anchoring (VSA). Panel (a) shows the turn-level reward of the answer action, panel (b) shows token-level entropy, panels (c) and (d) plot the KL divergence of visual tokens and generated tokens, respectively.

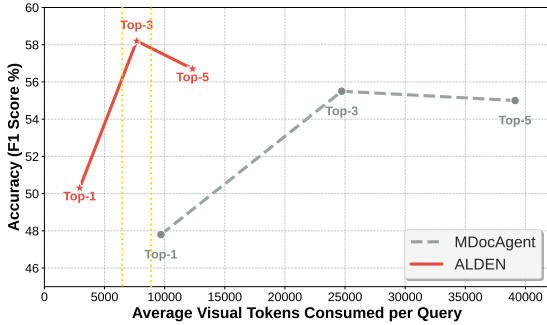


Figure 4: Efficiency and accuracy comparison between ALDEN and the MDocAgent baseline.

Full ALDEN further introduces token-level shaping, yielding a smaller but consistent gain (Acc 0.513, Rec 0.506) and increasing unique pages from 1.22 to 1.39, indicating reduced query repetition and broader exploration. Overall, the cross-level reward design fosters richer query reformulation and more thorough evidence gathering, enhancing both navigation and answer quality.

**Effect of Visual Semantic Anchoring.** We evaluate the effect of Visual Semantic Anchoring (VSA) on training stability and representation drift, as shown in Figure 3. With a larger batch size (512) than in the main experiments (128), the VSA-enabled model achieves steadily increasing answer rewards, while the non-VSA variant fluctuates and collapses (a). VSA also maintains higher policy entropy, supporting healthier exploration (b). Besides, KL divergence of visual tokens grows unchecked without VSA, indicating hidden-state drift, whereas VSA constrains these values while allowing moderate growth for action tokens (c,d). Overall, VSA achieves stabilizing RL training and preventing drift in visual representations.

## 5.5 Efficiency Analysis

We investigate the token efficiency of ALDEN on the LongDocURL dataset. Specifically, ALDEN is trained using the top-3 retrieved results as observations. During evaluation, we broaden the scope

to compare ALDEN against the strongest baseline, MDocAgent, across top-1, 3, and 5 retrieval settings. For a fair comparison, both methods utilize only query-image retrieval. As illustrated in Fig. 4, the baseline relies on expanding the retrieval scope to Top-5 for performance gains, incurring a linear cost penalty. In contrast, ALDEN reaches saturation at Top-3 by effectively identifying signals within noisy contexts. This results in a Pareto improvement, surpassing the baseline’s best configuration (+3.2% accuracy) with an approximate 3× reduction in token usage. This confirms that equipping the model with agentic reasoning is a fundamentally more efficient strategy than the passive reading, even when the latter is enhanced by multi-agent mechanisms.

## 6 Conclusions

We introduced the **Agentic VRDU** task and proposed **ALDEN**, a reinforcement-learning framework that trains VLMs as autonomous agents capable of multi-turn navigation and evidence gathering. The framework integrates a fetch action for direct page access, a cross-level reward for fine-grained reward modeling, and a visual semantic anchoring mechanism for stable training. Extensive experiments on multiple long-document benchmarks demonstrate that ALDEN achieves state-of-the-art accuracy while maintaining significantly lower token consumption. Furthermore, ablation studies validate the critical role of the fetch action in evidence localization and the effectiveness of visual semantic anchoring in stabilizing the training process, offering broader insights for multi-turn RL in multimodal agents. Ultimately, the A-VRDU paradigm marks a fundamental shift from passive document reading to autonomous navigation across vast information landscapes. The robust performance of ALDEN highlights the potential of such agents to deliver scalable, adaptive, and accurate understanding of complex, visually rich documents.

## 601 Limitations

602 Despite the promising results, ALDEN has several  
603 limitations. First, the agent still faces challenges in  
604 optimally balancing exploration and exploitation  
605 within large document spaces; it may occasionally  
606 terminate the search prematurely or navigate re-  
607 dundantly when the target information is buried  
608 deep in the document. Second, identifying true  
609 evidence pages remains non-trivial, as the agent  
610 can sometimes be misled by pages that are visually  
611 or semantically similar to the target but lack the  
612 precise answer.

613 Future work could address these issues by con-  
614 structing larger-scale datasets with high-quality tra-  
615 jectory annotations to improve sample efficiency.  
616 Additionally, we plan to leverage trajectories from  
617 stronger, closed-source models to guide training,  
618 integrating validation and reflection mechanisms to  
619 reduce hallucinations. Finally, adopting curriculum  
620 learning—starting from shorter documents and pro-  
621 gressively moving to complex ones—could help  
622 the agent better generalize across tasks of varying  
623 difficulty.

## 624 LLM Usage Statement

625 Large Language Models (LLMs) were used as  
626 general-purpose writing and editing aids. Specifi-  
627 cally, OpenAI’s ChatGPT (GPT-5) assisted in pol-  
628 ishing grammar, improving clarity, and suggesting  
629 alternative phrasings. All research ideas, exper-  
630 imental design, data processing, model develop-  
631 ment, and analysis were conceived and executed  
632 solely by the authors. The LLM provided no novel  
633 research insights or substantive scientific contribu-  
634 tions.

## 635 Reproducibility Statement

636 We are committed to ensuring the reproducibility  
637 of our results. To this end, we will release:

- 638 • All source code for training, evaluation,  
639 and data preprocessing, including scripts for  
640 dataset construction, reward computation, and  
641 reinforcement-learning training with ALDEN.
- 642 • The processed training corpus derived from  
643 DUDE, MPDocVQA, and SlideVQA, along with  
644 instructions to regenerate it from the original pub-  
645 lic datasets.
- 646 • Detailed configuration files specifying model hy-  
647 perparameters, random seeds, and hardware set-

tings. 648

- Checkpoints for both the policy and value mod- 649  
els, and prompts used for GPT-4o evaluation. 650

Our experiments were run on NVIDIA A100 GPUs 651  
(80GB) with PyTorch 2.4 and HuggingFace Trans- 652  
formers 4.49; exact package versions will be pro- 653  
vided in the released code. These resources will 654  
allow other researchers to fully reproduce our train- 655  
ing, evaluation, and analysis results. 656

## References 657

- 658 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien  
659 Bubeck, Ronen Eldan, Suriya Gunasekar, Michael  
660 Harrison, Russell J Hewett, Mojan Javaheripi, Piero  
661 Kauffmann, and 1 others. 2024. Phi-4 technical re-  
662 port. *arXiv preprint arXiv:2412.08905*.
- 663 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-  
664 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie  
665 Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl  
666 technical report. *arXiv preprint arXiv:2502.13923*.
- 667 Jian Chen, Ruiyi Zhang, Yufan Zhou, Tong Yu, Franck  
668 Dernoncourt, Jiuxiang Gu, Ryan A. Rossi, Changyou  
669 Chen, and Tong Sun. 2025a. [SV-RAG: LoRA-  
670 contextualizing adaptation of MLLMs for long docu-  
671 ment understanding](#). In *The Thirteenth International  
672 Conference on Learning Representations*.
- 673 Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou,  
674 Chenzheng Zhu, Haofen Wang, Jeff Z Pan, Wen  
675 Zhang, Huajun Chen, Fan Yang, and 1 others. 2025b.  
676 Learning to reason with search for llms via reinforce-  
677 ment learning. *arXiv preprint arXiv:2503.19470*.
- 678 Jaemin Cho, Debanjan Mahata, Ozan Irsoy, Yujie  
679 He, and Mohit Bansal. 2024. M3docrag: Multi-  
680 modal retrieval is what you need for multi-page  
681 multi-document understanding. *arXiv preprint  
682 arXiv:2411.04952*.
- 683 Chao Deng, Jiale Yuan, Pi Bu, Peijie Wang, Zhong-  
684 zhi Li, Jian Xu, Xiao-Hui Li, Yuan Gao, Jun Song,  
685 Bo Zheng, and 1 others. 2024. Longdocurl: a com-  
686 prehensive multimodal long document benchmark  
687 integrating understanding, reasoning, and locating.  
688 *arXiv preprint arXiv:2412.18424*.
- 689 Yihao Ding, Zhe Huang, Runlin Wang, Yanhang Zhang,  
690 Xianru Chen, Yuzhong Ma, Hyunsuk Chung, and  
691 Soyeon Caren Han. 2022. V-doc: Visual questions  
692 answers with documents. In *Proceedings of the  
693 IEEE/CVF conference on computer vision and pat-  
694 tern recognition*, pages 21492–21498.
- 695 Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Om-  
696 rani, Gautier Viaud, CELINE HUDELLOT, and Pierre  
697 Colombo. 2025. [Colpali: Efficient document re-  
698 trieval with vision language models](#). In *The Thir-  
699 teenth International Conference on Learning Repre-  
700 sentations*.



814	Alexander Michael Rombach and Peter Fettke. 2024.	Kangrui Wang, Pingyue Zhang, Zihan Wang, Yaning	871
815	Deep learning based key information extraction from	Gao*, Linjie Li, Qineng Wang, Hanyang Chen, Chi	872
816	business documents: Systematic literature review.	Wan, Yiping Lu, Zhengyuan Yang, Lijuan Wang,	873
817	<i>ACM Computing Surveys</i> .	Ranjay Krishna, Jiajun Wu, Li Fei-Fei, Yejin Choi,	874
		and Manling Li. 2025a. Reinforcing visual state	875
818	Keshav Santhanam, Omar Khattab, Jon Saad-Falcon,	reasoning for multi-turn vlm agents.	876
819	Christopher Potts, and Matei Zaharia. 2021. Col-		
820	berty2: Effective and efficient retrieval via	Liang Wang, Nan Yang, Xiaolong Huang, Binxing	877
821	lightweight late interaction. <i>arXiv preprint</i>	Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,	878
822	<i>arXiv:2112.01488</i> .	and Furu Wei. 2022. Text embeddings by weakly-	879
		supervised contrastive pre-training. <i>arXiv preprint</i>	880
823	John Schulman, Philipp Moritz, Sergey Levine, Michael	<i>arXiv:2212.03533</i> .	881
824	Jordan, and Pieter Abbeel. 2015. High-dimensional		
825	continuous control using generalized advantage esti-	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-	882
826	mation. <i>arXiv preprint arXiv:1506.02438</i> .	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin	883
		Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei	884
827	John Schulman, Filip Wolski, Prafulla Dhariwal,	Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang	885
828	Alec Radford, and Oleg Klimov. 2017. Proxi-	Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-	886
829	mal policy optimization algorithms. <i>arXiv preprint</i>	vl: Enhancing vision-language model’s perception	887
830	<i>arXiv:1707.06347</i> .	of the world at any resolution. <i>arXiv preprint</i>	888
		<i>arXiv:2409.12191</i> .	889
831	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	Qiuchen Wang, Ruixue Ding, Zehui Chen, Weiqi Wu,	890
832	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	Shihang Wang, Pengjun Xie, and Feng Zhao. 2025b.	891
833	Zhang, YK Li, Y Wu, and 1 others. 2024. Deepseek-	Vidorag: Visual document retrieval-augmented gen-	892
834	math: Pushing the limits of mathematical reason-	eration via dynamic iterative reasoning agents. <i>arXiv</i>	893
835	ing in open language models. <i>arXiv preprint</i>	<i>preprint arXiv:2502.18017</i> .	894
836	<i>arXiv:2402.03300</i> .		
837	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen,	Qiuchen Wang, Ruixue Ding, Yu Zeng, Zehui Chen,	895
838	Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-	Lin Chen, Shihang Wang, Pengjun Xie, Fei Huang,	896
839	Rong Wen. 2025. R1-searcher: Incentivizing the	and Feng Zhao. 2025c. VRAG-RL: Empower vision-	897
840	search capability in llms via reinforcement learning.	perception-based RAG for visually rich information	898
841	<i>arXiv preprint arXiv:2503.05592</i> .	understanding via iterative reasoning with reinforce-	899
		ment learning. In <i>The Thirty-ninth Annual Confer-</i>	900
842	Ryota Tanaka, Kyosuke Nishida, Kosuke Nishida, Taku	<i>ence on Neural Information Processing Systems</i> .	901
843	Hasegawa, Itsumi Saito, and Kuniko Saito. 2023a.		
844	Slidevqa: A dataset for document visual question	Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu,	902
845	answering on multiple images. In <i>Proceedings of</i>	Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin	903
846	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	Jing, Shenglong Ye, Jie Shao, and 1 others. 2025d. In-	904
847	ume 37, pages 13636–13645.	ternvl3. 5: Advancing open-source multimodal mod-	905
		els in versatility, reasoning, and efficiency. <i>arXiv</i>	906
848	Ryota Tanaka, Kyosuke Nishida, Kosuke Nishida, Taku	<i>preprint arXiv:2508.18265</i> .	907
849	Hasegawa, Itsumi Saito, and Kuniko Saito. 2023b.		
850	Slidevqa: a dataset for document visual question	Zilong Wang, Yichao Zhou, Wei Wei, Chen-Yu Lee, and	908
851	answering on multiple images. In <i>Proceedings of</i>	Sandeep Tata. 2023. Vrdu: A benchmark for visually-	909
852	<i>the Thirty-Seventh AAAI Conference on Artificial</i>	rich document understanding. In <i>Proceedings of the</i>	910
853	<i>Intelligence and Thirty-Fifth Conference on Inno-</i>	<i>29th ACM SIGKDD Conference on Knowledge Dis-</i>	911
854	<i>vative Applications of Artificial Intelligence and</i>	<i>covery and Data Mining</i> , pages 5184–5193.	912
855	<i>Thirteenth Symposium on Educational Advances in</i>		
856	<i>Artificial Intelligence</i> , AAAI’23/IAAI’23/EAAI’23.	Xudong Xie, Hao Yan, Liang Yin, Yang Liu, Jing Ding,	913
857	AAAI Press.	Minghui Liao, Yuliang Liu, Wei Chen, and Xiang	914
		Bai. 2024. Wukong: A large multimodal model for	915
858	Rubèn Tito, Dimosthenis Karatzas, and Ernest Valveny.	efficient long pdf reading with end-to-end sparse sam-	916
859	2023a. Hierarchical multimodal transformers for	pling. <i>arXiv preprint arXiv:2410.05970</i> .	917
860	multipage docvqa. <i>Pattern Recognition</i> , 144:109834.		
861	Rubèn Tito, Dimosthenis Karatzas, and Ernest Valveny.	Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B	918
862	2023b. Hierarchical multimodal transformers for	Brown, Alec Radford, Dario Amodei, Paul Chris-	919
863	multipage docvqa. <i>Pattern Recogn.</i> , 144(C).	tiano, and Geoffrey Irving. 2019. Fine-tuning lan-	920
		guage models from human preferences. <i>arXiv</i>	921
864	Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann,	<i>preprint arXiv:1909.08593</i> .	922
865	Michał Pietruszka, Pawel Joziak, Rafal Powalski,	<b>A Datasets</b>	923
866	Dawid Jurkiewicz, Mickaël Coustaty, Bertrand Anck-	<b>A.1 Training dataset</b>	924
867	aert, Ernest Valveny, and 1 others. 2023. Document	<b>Training.</b> We construct our training dataset	925
868	understanding dataset and evaluation (dude). In <i>Pro-</i>	by combining samples from three publicly	926
869	<i>ceedings of the IEEE/CVF International Conference</i>		
870	<i>on Computer Vision</i> , pages 19528–19540.		

Table 5: Statistics of the training dataset. #GQ and #PQ denote the numbers of general user queries and page-index-referenced queries, respectively.

Sub-dataset	DUDE	SlideVQA	MPDocVQA
#GQ	6,943	10,615	7,992
#PQ	1,011	2	4,165
Sum	7,954	10,617	12,157

available multi-page document understanding datasets: DUDE (Van Landeghem et al., 2023), MPDocVQA (Tito et al., 2023b), and SlideVQA (Tanaka et al., 2023b). These datasets provide diverse document layouts and question-answering formats, making them well-suited for training models on complex multi-turn document question answering tasks.

DUDE is a large-scale benchmark designed for multi-page, visually rich document understanding. It covers diverse domains such as scientific articles, financial and legal reports, technical manuals, and presentations. Each example consists of a full PDF document rendered into page images, paired with a natural-language query and a free-form textual answer, along with page-level ground-truth evidence annotations. SlideVQA contains questions grounded in slide decks, where understanding layout and inter-slide referencing is crucial. It contains slide decks from diverse topics such as education, business, and research talks, requiring models to reason across sequential pages that mix text, charts, and images. Each example provides a slide deck rendered as ordered page images, a natural-language question, and a free-form textual answer, with annotations of relevant slides for evidence grounding. MPDocVQA extends the traditional single-page VQA setting (originally based on DocVQA) by concatenating additional pages to the original single-page input, while retaining the same set of user questions. However, since many of these questions were authored under the assumption that only one page is visible (e.g., "What is the date?" or "Who is the author?"), they often lack sufficient context to guide document retrieval or navigation. To address this, we first use GPT-4o (Hurst et al., 2024) to automatically identify this kind of samples. Then we integrate the index of referred pages into the questions to get page-index-referenced questions, e.g., "In page 5, what is the date?". The prompt we used is shown below:

### Prompt for Filtering Queries

You are given a question from a multi-page document VQA dataset. Some questions are not suitable for training an agent to autonomously locate the target page, because they assume the agent already knows which page is relevant. These questions are often vague, layout-based, or refer to elements only visible on a known page (e.g., "What is the PVR no given in the approval sheet?", or "What is written at the top right?"). Your task is to assign a label to each question:

- 1 if the question belongs to this kind of problem, i.e., it assumes the correct page is known and cannot be answered without it.
- 0 if the question does not belong to this kind of problem, i.e., it can be answered after locating the page based on content in the question.

Respond with a JSON object containing only the field "label". Examples:

Question: What is the PVR no given in the approval sheet? Answer: { "label": 1 }

Question: What is the project name mentioned in the title block? Answer: { "label": 0 }

Question: What is the symposium organized by Division of Agricultural and Food Chemistry? Answer: { "label": 0 }

Question: What is written on the top right corner? Answer: { "label": 1 }

Question: What is the page number? Answer: { "label": 1 }

Question: What is the Date? Answer: { "label": 1 }

Now, label the following question:

Question: {question}

To ensure that our model is consistently exposed to multi-page reasoning scenarios, we additionally discard any documents with fewer than 10 pages from all three datasets. This helps avoid biasing the model toward short-context behavior and ensures a consistent level of document complexity.

After merging and filtering, we obtain a training set consisting of 30,728 samples, each comprising a user query and its corresponding multi-page document context, answer and the index of evidence pages. Finally, we proportionally sample 1,024 samples from the validation set of these three datasets as our validation set.

## A.2 Benchmarks

We evaluate our method on a diverse set of benchmarks: MMLongBench (Ma et al., 2024b), LongDocURL (Deng et al., 2024), PaperTab (Hui et al., 2024), PaperText (Hui et al., 2024), and FetaTab (Hui et al., 2024). These datasets span a wide range of scenarios, including both open-domain and closed-domain tasks, and include textual as well as visual content. The documents also vary in length and structure, ranging from short forms to complex, multi-page documents. This diversity ensures a comprehensive and fair evaluation of our model’s performance across real-world document understanding tasks.

- **MMLongBench-Doc** is a large-scale benchmark designed to evaluate how multimodal large language models handle long, visually rich documents. It contains over a thousand expert-annotated questions drawn from lengthy PDFs (averaging 50 pages and 20k tokens) that mix text, tables, charts, and images. Tasks require single-page, cross-page, and sometimes unanswerable reasoning, testing a model’s ability to retrieve and integrate evidence across multiple modalities and extended contexts.
- **LongDocURL** is a benchmark for evaluating large vision-language models on long, multimodal documents by combining three core task types: understanding, numerical reasoning, and element locating. It includes 2,325 high-quality question-answer pairs over 396 documents totaling over 33,000 pages, with an average of 85.6 pages per document. Tasks vary in their evidence requirements: some require single-page evidence, others multi-page, and many involve locating evidence across different layout elements (text, tables, figures, and layout).
- **PaperText** is a subset in the UDA benchmark made up of academic papers (in PDF form) used for retrieval-augmented generation / document question answering tasks. Each document comes with multiple question-answer pairs drawn from “Qasper” (an academic paper reading comprehension dataset), where questions may be extractive, yes/no, or free-form. The dataset preserves full documents to allow answering from context, rather than just small passages.
- **PaperTab** is another subset in UDA also based on academic papers, but the focus is on Q&A

pairs where evidence comes from or interacts with tables inside papers. Like PaperText, it retains full PDF documents so that models must locate and reason over tabular content, as well as textual content. The questions are similarly diverse (extractive, yes/no, free-form), and the average size is modest (10–11 pages per document).

- **FetaTab** is a subset of the UDA (Unstructured Document Analysis) benchmark that focuses on free-form question answering over Wikipedia tables in both HTML and PDF formats. It comprises 878 documents and 1,023 QA pairs, averaging about 14.9 pages per document. The questions are “free-form” (i.e. natural language answers, not limited to extractive spans or simple yes/no), which requires models to understand table content, context, and sometimes cross-format layout.

## B Derivation about the Visual Semantic Anchoring

Let the input to the VLMs head be a concatenation of visual tokens  $Z_v \in \mathbb{R}^{N_v \times d}$  and textual tokens  $Z_t \in \mathbb{R}^{N_t \times d}$ . The policy output (logits) is a function  $f : \mathbb{R}^{(N_v+N_t) \times d} \rightarrow \mathbb{R}^V$ . We examine the update of the visual encoder parameters  $\phi$ , which solely influence  $Z_v$ .

Our verification starts with the Block-Jacobian Decomposition. The linearization of the change in logits  $\Delta y$  with respect to the input features is given by the total derivative. Since the input is concatenated  $X = [Z_v, Z_t]$ , the Jacobian  $J_f$  can be decomposed into two block matrices:

$$J_f = \begin{bmatrix} \frac{\partial f}{\partial Z_v} & \mid & \frac{\partial f}{\partial Z_t} \end{bmatrix} = [J_v \mid J_t]$$

where  $J_v \in \mathbb{R}^{V \times (N_v \cdot d)}$  is the Partial Jacobian with respect to visual features.

During the PPO update, we are concerned with the gradients flowing into the visual encoder parameters  $\phi$ . These parameters only affect  $Z_v$ . The textual tokens  $Z_t$  are either fixed (from history) or updated by separate parameters (LLM weights). Thus, relative to the visual encoder update,  $\Delta Z_t = 0$ . The constraint imposed by the text-only KL divergence ( $\|\Delta y\| < \epsilon$ ) simplifies to:

$$\begin{aligned} \Delta y &\approx J_v \cdot \text{vec}(\Delta Z_v) + J_t \cdot \mathbf{0} \\ &\approx J_v \cdot \text{vec}(\Delta Z_v) \end{aligned} \quad (7)$$

The "freedom" available to the visual encoder to drift without triggering the KL penalty is determined solely by the null space of the Partial Jacobian  $J_v$ . The dimensionality of this subspace is:

$$\dim(\text{null}(J_v)) = (N_v \cdot d) - \text{rank}(J_v)$$

Even with the presence of  $Z_t$ , the rank of  $J_v$  is still upper-bounded by the bottleneck dimension of the model (or vocabulary size  $V$ ). It does not gain rank from the text tokens. Thus the available "drift space" scales linearly with the number of visual tokens:

$$\dim(\text{null}(J_v)) \propto N_v$$

The Null Space represents the "freedom" the visual encoder has to change its weights ( $\Delta Z$ ) without violating the text-KL constraint. As the number of visual tokens  $N$  increases, the dimensionality of this unconstrained subspace grows. Consequently, the optimizer has significantly more degrees of freedom to introduce aggressive, potentially destructive updates to the visual representations that satisfy the short-term reward while remaining "invisible" to the text-only KL penalty. This confirms that VLMs with higher visual token counts are mathematically more susceptible to representation drift when lacking visual-semantic anchoring.

## C Baselines

To evaluate the effectiveness of ALDEN, we compare it against three categories of methods:

- **Base VLMs supporting multi-image input.** These models directly take the entire multi-page document as context without retrieval, leveraging their built-in multi-page visual processing capabilities. For fairness, we select open-source VLMs of similar scale to Qwen2.5-VL-7B, including LLaVA-v1.6-Mistral-7B (Liu et al., 2024a), Phi-3.5-Vision-Instruct (Abdin et al., 2024), LLaVA-One-Vision-7B (Li et al., 2024), SmolVLM-Instruct (Marafioti et al., 2025), mPLUG-DocOwl2 (Hu et al., 2024), LEOPARD (Jia et al., 2024), InternVL3.5-8B-Instruct (Wang et al., 2025d).
- **Visual RAG methods.** These methods use the user query to retrieve the most relevant document pages and feed them into the model as context. We include M3DocRAG (Cho et al., 2024) as a strong baseline, as well as our proposed ALDEN.

To isolate the impact of our reward function design, we additionally evaluate a variant that trains the same backbone with GRPO using only outcome-based rewards (no turn-level shaping), mirroring common text-only RLHF setups as in ReSearch (Chen et al., 2025b). Specifically,

- M3DocRAG is a multi-modal document understanding framework designed for multi-page and multi-document question answering. It first encodes each page into joint visual-text embeddings using a multi-modal encoder, then retrieves the top-K relevant pages via a MaxSim-based retrieval mechanism, optionally accelerated with FAISS for large-scale documents. Finally, a multi-modal language model processes the retrieved pages to generate precise answers, effectively handling complex queries that require reasoning over both textual and visual content.
- ReSearch introduces a framework that trains large language models to integrate reasoning and search in a unified process. The model learns, via reinforcement learning, when and how to perform search actions during multi-step reasoning, using search results to guide subsequent reasoning steps. By treating search as part of the reasoning chain, ReSearch enables LLMs to solve complex multi-hop tasks, demonstrate self-correction and reflection, and generalize effectively across benchmarks, achieving significant performance gains over baseline models.
- **Hybrid RAG methods.** These approaches combine visual and textual retrieval by first applying an OCR tool to extract all text from the document. The query is then used to retrieve both the most relevant page image and the most relevant OCR-extracted text, which are jointly fed into the model. We evaluate MDocAgent (Han et al., 2025) and VidoRAG (Wang et al., 2025b) as a representative method in this category.
  - MDocAgent is a multi-modal, multi-agent framework for document understanding that combines Retrieval-Augmented Generation (RAG) with specialized agents to handle complex documents. The system employs a General Agent for multi-modal context retrieval, a Critical Agent for identifying

key information, a Text Agent for analyzing textual content, an Image Agent for interpreting visual elements, and a Summarizing Agent to synthesize results. By coordinating these agents, MDocAgent effectively integrates textual and visual reasoning, achieving significant improvements in accuracy and error reduction compared to existing large vision-language models and RAG-based methods. For all five agents in this framework, we consistently use the original LLaMA3.1-8B as the LLM for the text agent, while employing a consistent VLMs, i.e., Qwen2.5-VL-7B, for remaining agents.

- ViDoRAG is a multi-agent framework designed to enhance the understanding of visually rich documents. It employs a Gaussian Mixture Model (GMM)-based hybrid retrieval strategy to effectively handle multi-modal retrieval, integrating both textual and visual information. The framework incorporates a dynamic iterative reasoning process, utilizing agents such as Seeker, Inspector, and Answer to iteratively refine the understanding and generation of responses. This approach addresses challenges in traditional Retrieval-Augmented Generation (RAG) methods by improving retrieval accuracy and enabling complex reasoning over visual documents. We use Qwen2.5-VL-7B as backbone for all agents in this methods.

## D Implementation Details

Our implementation is based on the EasyR1<sup>1</sup> framework. We adopt the default optimization hyperparameters from the EasyR1 framework, specifically the learning rates and KL coefficients, to ensure training stability. Both the policy model and the value function are initialized from Qwen2.5-VL-7B-Instruct (Bai et al., 2025). We use a batch size of 128, with fixed learning rates of  $1 \times 10^{-6}$  for the policy model and  $1 \times 10^{-5}$  for the value function. The maximum number of interaction turns is set to  $T = 6$ . For visual inputs, we constrain the number of image pixels to lie between 261,070 and 2,508,800. Based on these settings, we set the maximum number of tokens in the trajectory as 19000. The KL coefficients for generated tokens and observation tokens are set to  $\beta_{\text{gen}} = 0.001$  and  $\beta_{\text{obs}} = 0.01$ , respectively. For the search actions,

<sup>1</sup><https://github.com/hiyouga/EasyR1>

we used only the top-1 retrieved pages. Besides, we set the scale coefficient  $\alpha = 5$ . The weight of repetition penalty is set as  $\eta = 0.5$ . For the calculation of GAE, we set  $\gamma_{\text{token}} = 1.0$ ,  $\gamma_{\text{turn}} = 0.9$  and  $\lambda_{\text{token}} = \lambda_{\text{turn}} = 1.0$ . During training, we adopt the single-vector retriever vdr-2b-v1 (Ma et al., 2024a) for images and e5-large-v2 (Wang et al., 2022) for text for training efficiency. For evaluation, we also report results with the multi-vector retrievers ColQwen2-v1.0 (ColQwen) (Faysse et al., 2025) for images and ColBERT-v2.0 (ColBERT) (Santhanam et al., 2021) for text. All experiments are conducted on 16 NVIDIA A100-80Gb GPUs.

The system prompt that we used during training of Visual RAG variant of ALDEN is shown in Fig. 5.

The system prompt that we used during training of Hybrid RAG variant of ALDEN is shown in Fig. 5.

## E Evaluation Metrics

We evaluate models using both answer quality and intermediate navigation metrics.

**Model-based Accuracy (Acc).** Answer quality is assessed with an LLM-as-judge protocol. Given a predicted answer and the ground-truth reference, GPT-4o is prompted to classify the prediction as *Correct*, *Incorrect*, or *Tie/Unclear*. We compute accuracy for each benchmark as the percentage of responses judged *Correct* over all responses:

$$\text{Acc} = \frac{\#\text{Correct}}{N}, \quad (8)$$

where  $N$  is the number of test instances.

**Trajectory-level Recall (Rec).** Let  $\mathcal{G}$  denote the set of ground-truth evidence pages for a given query, and let  $\mathcal{T}$  denote the set of pages collected by the agent along a trajectory. The trajectory-level recall is defined as:

$$\text{Rec} = \frac{|\mathcal{T} \cap \mathcal{G}|}{|\mathcal{G}|}. \quad (9)$$

This metric measures the fraction of ground-truth pages successfully retrieved by the agent over the course of a trajectory, providing an indicator of how effectively the agent gathers relevant information.

**Trajectory-level Precision (Pre).** Let  $\mathcal{G}$  denote the set of ground-truth evidence pages for a given query, and let  $\mathcal{T}$  denote the set of pages collected by the agent along a trajectory. The trajectory-level precision is defined as:

## System prompt of ALDEN with Visual RAG

You are a helpful assistant designed to answer user questions based on a user-provided multi-page document. The document can not be input directly with the question, you must reason step by step to determine how to obtain evidence document pages by optimally utilizing tools and analyze the relevant content in the obtained document pages to precisely answer user’s question. Your reasoning process MUST BE enclosed within `<think>` `</think>` tags. Your answer MUST BE enclosed within `<answer>` `</answer>` tags. In the last part of the answer, the final exact answer is enclosed within `\boxed{ }` with latex format. The available tool is a **search tool**. After reasoning, you can invoke the search tool by generating `<search>` your search query here `</search>` to retrieve document pages most relevant to your search query. For example, your response could be in the format of `<think>` your reasoning process `</think>` `<search>` search query `</search>`, or `<think>` your reasoning process `</think>` `<answer>` your answer here. The final answer is `\[ \boxed{ { answer here} } \]` `</answer>`. After invoking a tool, the user will return obtained document pages inside `<result>` `</result>` tags to you. Besides, the user will additionally provide the page number of the obtained page.

**\*\*Important constraints\*\*:**

- Only if you get all the potential evidence pages and find that there is no evidenced answer or the document content is irrelevant to the user query, you can respond with `<think>` your reasoning process `</think>` `<answer>` The final answer is `\[ \boxed{The problem is not answerable} \]` `</answer>`.
- If multiple valid answers are found, return them separated by semicolons.
- You may not get the true evidence page in one-shot, carefully check whether the obtained pages are the true evidence page. If not, try different rewritings of your query or try different tool usage strategy several times.

Figure 5: System prompt of ALDEN with Visual RAG

$$\text{Pre} = \frac{|\mathcal{T} \cap \mathcal{G}|}{|\mathcal{T}|}. \quad (10)$$

This metric measures the fraction of pages collected by the agent that are actually relevant, providing an indicator of how accurately the agent identifies evidence pages during a trajectory.

**F1 Score (F1).** Based on the trajectory-level precision and recall, the trajectory-level F1 score is defined as the harmonic mean of the two:

$$\text{F1} = 2 \cdot \frac{\text{Pre} \cdot \text{Rec}}{\text{Pre} + \text{Rec}}. \quad (11)$$

This metric provides a balanced measure of the agent’s performance, accounting for both its ability to collect relevant pages (recall) and to avoid collecting irrelevant ones (precision) over a trajectory.

**Number of uniquely collected pages (#UP).**

## F Case Study

In this section, we present typical examples from the LongDocURL benchmark to analyze the behavior of the ALDEN model. As illustrated in the

successful cases in Tab. 6 and Tab. 7, the agent effectively learns to combine the search and fetch actions to gather evidence and correctly answer the question. In particular, the example in Tab. 7 shows the agent using fetch to complement search when the latter fails to retrieve the true evidence page in a single step, highlighting the benefit of ALDEN’s expanded action space. Alongside these successes, we also present failure cases in Tab. 8 and Tab. 9. These examples reveal that the model still struggles to verify evidence pages and explore sufficiently when the search action returns results far from the ground-truth pages—a challenge we leave for future work.

1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280

<b>User Question</b>	What's name of the figure at the page which contains a table whose name is "Table 18: Quality control related uses by type of use"?
<b>Ground-truth Pages</b>	40
<b>Generated Trajectory</b>	<p><b>assistant</b>  &lt;think&gt; The question asks for the name of the figure at the page that contains a table with the name "Table 18: Quality control related uses by type of use". To answer this, I should use the search tool to look for the table with the specified name in the document. The fetch tool can be used to obtain the specific page where this table is located. &lt;/think&gt; &lt;search&gt; name of the figure at the page which contains a table whose name is "Table 18: Quality control related uses by type of use" &lt;/search&gt;</p> <p><b>user</b>  &lt;result&gt; Page 39: Image tokens &lt;/result&gt;</p> <p><b>assistant</b>  &lt;think&gt; The table mentioned in the search query is not found on the provided page. The table in the document is "Table 17: Regulatory uses by main types of uses". The search query asks for the figure at the page which contains the table "Table 18: Quality control related uses by type of use". I should use the fetch tool to obtain the page with the table "Table 18: Quality control related uses by type of use". &lt;/think&gt; &lt;fetch&gt; 40 &lt;/fetch&gt;</p> <p><b>user</b>  &lt;result&gt; Image tokens &lt;/result&gt;</p> <p><b>assistant</b>  &lt;think&gt; The table mentioned in the search query is found on the provided page. The table in the document is "Table 18: Quality control related uses by type of use". The search query asks for the figure at the page which contains the table "Table 18: Quality control related uses by type of use". The figure mentioned in the document is "Figure 20: Evolution of total numbers of uses of animals for pyrogenicity testing between 2015 and 2018". The figure name is "Figure 20". &lt;/think&gt; &lt;answer&gt; The final answer is \[ \boxed{Figure 20} \]&lt;/answer&gt;</p>

Table 6: Example of a generated case showing the user question and ALDEN's full reasoning-action trajectory.

### System prompt of ALDEN with Hybrid RAG

You are a helpful assistant designed to answer user questions based on a user-provided multi-page document. Each page exists in two modalities: the original image and an OCR text extraction. You cannot access the full document directly; instead, you must reason step by step to determine how to obtain evidence document pages by optimally utilizing tools and analyze the relevant content in the obtained document pages to precisely answer user's question. Your reasoning process **MUST BE** enclosed within `<think>` `</think>` tags. Your answer **MUST BE** enclosed within `<answer>` `</answer>` tags. In the last part of the answer, the final exact answer should be enclosed within `\boxed{{}}` with latex format. The available tools include a **\*\*search tool\*\*** and a **\*\*fetch tool\*\***. After reasoning, you can invoke either the search tool by generating `<search>` your search query here `</search>` to retrieve relevant document pages in both modalities or the fetch tool by generating `<fetch>` modal, page number `</fetch>` to obtain a specific document page in the specified modal, where the modal should be 'image' or 'text' and the page number should be a integrity number chosen from the user specified page number range. For example, your response could be in the format of '`<think>` your reasoning process `</think>` `<search>` search query `</search>`', or '`<think>` your reasoning process `</think>` `<fetch>` image, page number `</fetch>`', or '`<think>` your reasoning process `</think>` `<fetch>` text, page number `</fetch>`', or '`<think>` your reasoning process `</think>` `<answer>` your answer here. The final answer is `\[ \boxed{{ answer here }} \]` `</answer>`'. After invoking a tool, the user will return obtained document pages inside `<result>` `</result>` tags to you. For the search tool, the user will return both the relevant image pages and the relevant OCR text pages and attach them with corresponding page numbers. For the fetch tool, the user will only return either the image page or the OCR text page according to your input arguments.

**\*\*Important constraints\*\*:**

- Only if you get all the potential evidence pages and find that there is no evidenced answer or the document content is irrelevant to the user query, you can respond with '`<think>` your reasoning process `</think>` `<answer>` The final answer is `\[ \boxed{The problem is not answerable} \]` `</answer>`'.
- If multiple valid answers are found, return them separated by semicolons.
- Only one page can be fetched at a time using the fetch tool.
- You may not get the true evidence page in one-shot, carefully check whether the obtained pages are the true evidence page. If not, try different rewritings of your query or try different tool usage strategy several times.
- Page numbers shown in the document pages may not be consistent with user specified page number range. In case of any discrepancy, the user defined page number range shall prevail.
- You need to invoke the tools at least once and can invoke up to 5 times. When you output the answer, the interaction stops.

Figure 6: System prompt of ALDEN with Hybrid RAG.

<b>User Question</b>	What is the percentage decrease in hospitalizations between May 14 and May 20, 2020? (A). -6.5(B). +6.5(C). -65(D). +65Choose the letter name in front of the right option from A, B, C, D.
<b>Ground-truth Pages</b>	21
<b>Generated Trajectory</b>	<p><b>assistant</b></p> <p>&lt;think&gt; The question asks for the percentage decrease in hospitalizations between May 14 and May 20, 2020. To answer this, I need to find information in the document that discusses hospitalizations for these specific dates and the percentage change between them. I will use the search tool to look for the relevant information in the document. &lt;/think&gt; &lt;search&gt; percentage decrease in hospitalizations between May 14 and May 20, 2020 &lt;/search&gt;</p> <p><b>user</b></p> <p>&lt;result&gt;Page 21: Image tokens &lt;/result&gt;</p> <p><b>assistant</b></p> <p>&lt;think&gt; The image contains the information needed to answer the question. According to the document, there has been a -6.5% decrease in hospitalized patients between May 14 and May 20, 2020. The correct answer is (A). &lt;/think&gt;</p> <p>&lt;answer&gt; The final answer is <math>\boxed{A}</math> &lt;/answer&gt;</p>

Table 7: Example of a generated case showing the user question and ALDEN’s full reasoning–action trajectory.

<b>User Question</b>	What is the interval between specific lunar/solar events in days?
<b>Ground-truth Pages</b>	46, 47, 48, 55, 57, 58
<b>Generated Trajectory</b>	<p><b>assistant</b></p> <p>&lt;think&gt; The question asks for the interval between specific lunar/solar events in days. To answer this, I should look for information in the document that discusses lunar and solar events and their intervals. The search tool seems the best option to find this information. I will use the search tool first. &lt;/think&gt;</p> <p>&lt;search&gt; interval between lunar/solar events in days &lt;/search&gt;</p> <p><b>user</b></p> <p>&lt;result&gt;Page 25: Image tokens &lt;/result&gt;</p> <p><b>assistant</b></p> <p>&lt;think&gt; The image contains the information needed to answer the question. The document states that the interval between the lunar eclipse on April 15, 2014, and the lunar eclipse on September 28, 2015, is exactly 266 days. This matches the information in the question. The interval between specific lunar/solar events in days is 266 days. The final answer is <math>\boxed{266}</math></p> <p>&lt;/think&gt; &lt;answer&gt; The final answer is <math>\boxed{266}</math> &lt;/answer&gt;</p>

Table 8: Example of a generated case showing the user question and ALDEN’s full reasoning–action trajectory.

<b>User Question</b>	What penalties are specified for possession of category 1 or category 2 items?
<b>Ground-truth Pages</b>	21
<b>Generated Trajectory</b>	<p><b>assistant</b></p> <p>&lt;think&gt; The question asks for the penalties specified for possession of category 1 or category 2 items. To answer this, I should look for information in the document that discusses the penalties for these categories. The search tool seems the best option to start with as it can help locate the relevant information in the document. I will use the search tool first. &lt;/think&gt; &lt;search&gt; penalties for possession of category 1 or category 2 items &lt;/search&gt;</p> <p><b>user</b></p> <p>&lt;result&gt;Page 21: Image tokens &lt;/result&gt;</p> <p><b>assistant</b></p> <p>&lt;think&gt; The image contains the information needed to answer the question. According to the document, the penalties for possession of a category 1 item or a category 2 item are specified in the Misuse of Drugs Act 1981. For a simple offence, the penalty is \$12,000 or imprisonment for 3 years or both. The relevant information is found in the document. The final answer is \[ \boxed{{\\$12,000 or imprisonment for 3 years or both}} \]&lt;/think&gt; &lt;answer&gt; The final answer is \[ \boxed{{\\$12,000 or imprisonment for 3 years or both}} \]&lt;/answer&gt;</p>

Table 9: Example of a generated case showing the user question and ALDEN’s full reasoning–action trajectory.

---

**Algorithm 1** PPO with Dual KL Regularization for Multi-Turn VRDU Agents

---

**Require:** Actor  $\pi_\theta$ , Critic  $V_\phi$ , Reference model  $\pi_{\text{ref}}$ , KL weights  $\beta_{\text{gen}}, \beta_{\text{obs}}$ , discount factors  $\gamma_{\text{token}}, \gamma_{\text{turn}}$ , GAE parameters  $\lambda_{\text{token}}, \lambda_{\text{turn}}$ , replay buffer  $\mathcal{B}$

- 1: Initialize replay buffer  $\mathcal{B}$
- 2: **for** iteration = 1, 2, ... **do**
- 3:     Sample  $|\mathcal{B}|$  queries from the dataset
- 4:     **for** each query **do**
- 5:         Reset: query  $q$ , empty retrieval history,  $t \leftarrow 1$
- 6:         **while**  $t < T$  **and**  $a_{t-1} \neq \text{answer}$  **do**
- 7:              $\pi_\theta$  generates a token sequence  $a_t \sim \pi_\theta(\cdot | s_t)$
- 8:             Parse the discrete action (search, fetch, or answer) from  $a_t$
- 9:             Execute action  $\rightarrow$  obtain new state  $s_{t+1}$  and turn reward  $r_t$
- 10:             Store  $\{a_t, s_{t+1}, r_t\}$  in  $\mathcal{B}$
- 11:              $t \leftarrow t + 1$
- 12:         **Turn-level value estimation:**
- 13:         **for** each episode in  $\mathcal{B}$  **do**
- 14:             Estimate  $V_\phi(s_t)$  at final token of each turn
- 15:             Compute target turn value  $\hat{V}_t$  via turn-level GAE
- 16:             Assign token-level reward  $\tilde{r}_t \leftarrow \hat{V}_t$
- 17:         **Dual KL penalty computation:**
- 18:         **for** each token in  $\mathcal{B}$  **do**
- 19:             **if** token is generated **then**
- 20:                 Compute  $A_t^i$  via token-level GAE using  $\tilde{r}_t$
- 21:                 Compute  $\text{KL}(\pi_\theta(\cdot | s) \parallel \pi_{\text{ref}}(\cdot | s))$  with weight  $\beta_{\text{gen}}$
- 22:                 **else if** token is observation **then**
- 23:                 Compute  $\text{KL}(\pi_\theta(\cdot | s) \parallel \pi_{\text{ref}}(\cdot | s))$  with weight  $\beta_{\text{obs}}$
- 24:         **PPO update:**
- 25:         Update  $\theta$  by maximizing policy loss  $\mathcal{L}_{\text{policy}}$
- 26:         Update  $\phi$  by minimizing value loss  $\mathcal{L}_{\text{value}}$

---