

Structured Over Scale: Learning Spatial Reasoning from Educational Video

Bishoy Galoaa^{*1} Xiangyu Bai^{*1} Sarah Ostadabbas¹

Abstract

Vision-language models (VLMs) demonstrate impressive performance on standard video understanding benchmarks yet fail systematically on simple reasoning tasks that preschool children can solve, including counting, spatial reasoning, and compositional understanding. We hypothesize that the pedagogically-structured content of educational videos provides an ideal training signal for improving these capabilities. We introduce DoraVQA, a dataset of 5,344 question-answer pairs automatically extracted from 8 seasons of Dora the Explorer with precise timestamp alignment. Each episode follows a consistent *context-question-pause-answer* structure that creates a self-contained learning environment analogous to interactive tutoring. We fine-tune both Qwen2 and Qwen3 using Group Relative Policy Optimization (GRPO), leveraging the clear correctness signals and structured reasoning traces inherent in educational content. Despite training exclusively on 38 hours of children’s educational videos, our approach achieves improvements of 8-14 points on DoraVQA and state-of-the-art 86.16% on CVBench, with strong transfer to Video-MME and NExT-QA, demonstrating effective generalization from narrow pedagogical content to broad multimodal understanding. Through cross-domain benchmarks, we show that VLMs can perform tasks that require robust reasoning learned from structured educational content, suggesting that content structure matters as much as content scale. Code and data are available at <https://github.com/ostadabbas/DORA-Learning-Spatial-Reasoning>. Please refer to our appendix and supplemental materials for more dataset samples and qualitative results.

¹Northeastern University, Boston, MA, USA. Correspondence to: Bishoy Galoaa <galoaa.b@northeastern.edu>, Xiangyu Bai <bai.xiang@northeastern.edu>, Sarah Ostadabbas <s.ostadabbas@northeastern.edu>.

1. Introduction

Vision language models (VLMs) have made rapid progress on standard video question-answering benchmarks (Xu et al., 2017; Xiao et al., 2021), yet they continue to fail on elementary reasoning tasks that even young children solve reliably. On diagnostic benchmarks targeting counting, spatial relations, and directional understanding, state-of-the-art models such as GPT-4V achieve only 50–60% accuracy (Fu et al., 2024), far below the 95%+ performance observed in humans. These failures persist even in simple scenes with fewer than ten objects (Tong et al., 2024), suggesting that gains in video understanding often arise from scaling data volume and statistical pattern matching rather than grounded spatial reasoning.

Recent research indicates that this limitation is not from architectural decisions, but data-driven: while these models exhibit superficial reasoning ability, they largely reflect the learned data distribution from large-scale web videos and generic video-text pairs that *display* spatial relations, rather than genuinely learning the rules beyond simple appearance. Effective learning of spatial concepts requires repeated, explicit supervision that links language to visual evidence through instruction, feedback, and temporal scaffolding. In contrast, most large-scale video datasets lack clear correctness signals and offer little incentive for models to ground linguistic concepts such as “over,” “behind,” or “more than” in perceptual structure.

Children’s educational television offer a uniquely structured alternative. Programs such as *Dora the Explorer* and *Micky Mouse Clubhouse* are explicitly designed around a pedagogical loop that repeatedly teaches spatial reasoning through a consistent *context-question-pause-answer* format. Each episode establishes context through extended visual scenes and previous conversations, poses direct spatial questions to the viewer (e.g., “How do we get to the blue tree?”, “Which path should we take?”), pauses for several seconds while visually emphasizing relevant objects via gestures, zooms, or highlights, and then reveals an unambiguous answer with verbal explanation. This structure appears systematically throughout every episode, with an average of 15–20 explicit spatial questions and clear correctness signals per episode. Importantly, the pause segment functions as an implicit reasoning trace: it visually isolates the evidence required to

answer the question, creating a self-supervised environment for grounding language in perception.

This pedagogical design has been shown to produce measurable learning gains in children. Linebarger and Walker (Linebarger & Walker, 2005) report significantly higher vocabulary acquisition rates for viewers of educational television, while Angin (Angin, 2017) demonstrates a 44% improvement in spatial concept learning following structured educational viewing. Most noticeably, the importance of pause-based visual context is supported by prior work (Pozzulo et al., 2012), which reports that children’s accuracy on *Dora the Explorer* questions decreases 25% when the visual frames are not given during the pause. We hypothesize that the same structure that supports human learning can serve as an efficient and underexplored training signal for vision language models to enable spatial reasoning.

We argue that the pedagogical structure adopted from children’s educational programs enables VLMs to improve spatial reasoning skills that generalize across diverse tasks and datasets. To support our hypothesis, we introduce **DoraVQA**, a dataset extracted from 8 seasons (96 episodes) of *Dora the Explorer*, comprising 5,344 question–answer pairs with precise temporal alignment to video frames. Each example preserves the show’s pedagogical structure, including extended context segments, explicit questions, visually salient pause intervals, and ground-truth answers revealed in the transcripts. Using a large language model (LLM) agent, we process transcribed dialogue to form cohesive question-answer pairs with corresponding video timestamps and extract the pause-based visuals with surrounding textual context. DoraVQA spans multiple reasoning modes, including immediate spatial grounding, sequential navigation, compositional constraints, and recall of knowledge (Figure 1).

To further illustrate that structured supervision improves spatial reasoning, we fine-tune Qwen2-VL (Wang et al., 2024a) and Qwen3-VL (Yang et al., 2025) using Group Relative Policy Optimization (GRPO) (Shao et al., 2024). Compared with supervised fine-tuning (SFT) and other reinforcement learning (RL) approaches, our goal is to stabilize and enable spatial reasoning with existing knowledge, rather than adapting the model to specific content. GRPO aligns with our task: for educational programs, answers are objectively correct or incorrect, and the show’s explanations provide implicit reasoning cues without requiring manual annotation. Additionally, to showcase the knowledge transfer across tasks and datasets, we train on open-ended answer generation from DoraVQA but evaluate using additional multiple-choice reasoning benchmarks, including the DoraVQA test split, Video-MME (Fu et al., 2025), CVBench (Zhu et al., 2025), and NExT-QA (Xiao et al., 2021). The deliberate mismatch tests whether improvements reflect transferable spatial reasoning rather than answer-format memorization.

Despite training on approximately 38 hours of narrowly scoped educational video, our approach yields consistent gains over the Qwen2-VL and Qwen3-VL baselines. These results suggest that pedagogical structure can compensate for scale, enabling models to better reason and learn effectively from fewer, simpler, but better-organized examples. At the same time, our analysis reveals persistent limitations—particularly in counting—highlighting that structure alone does not solve all perceptual deficits and underscoring the need for stronger visual grounding in future models.

Our contribution can be listed as follows:

- We formulate a novel learning structure inspired by children’s educational content that strengthens VLMs’ spatial reasoning abilities without large-scale data corpora, yielding cross-benchmark improvements.
- We show that the feature-rich visual and textual pause context embedded in pedagogically structured question-pause-answer pairs serves as an implicit reasoning supervision, substantially improving a VLM’s correspondence between questions and answers.
- We introduce DoraVQA, the first video question-answering dataset extracted from children’s educational television, with temporal alignment between visual frames, questions, pause context, and answers.
- We present an automated pipeline that extracts self-supervised training data from structured video content, reducing reliance on manual annotation while preserving supervision signals through the content’s inherent reasoning structure.

2. Related Work

Vision Language Models and Temporal Understanding.

Recent work has scaled VLMs to capture temporal dynamics unavailable in static images. Video-LLaVA (Lin et al., 2024) demonstrated that joint training on videos and images improves both modalities through unified visual representations; InternVideo2 (Wang et al., 2024b) expanded video encoders using progressive training and masked video modeling; LLaVA-Video (Zhang et al., 2024) introduced dense frame sampling to capture fine-grained temporal dynamics, improving temporal reasoning. However, these approaches train predominantly on web videos (Miech et al., 2019; Abu-El-Haija et al., 2016; Kay et al., 2017), which contain diverse activities but lack explicit pedagogical structure. While these videos display spatial relationships and temporal sequences, they rarely verbalize spatial concepts explicitly or provide structured question-answer pairs that ground visual patterns in linguistic descriptions. Our work differs by leveraging educational content where spatial reasoning

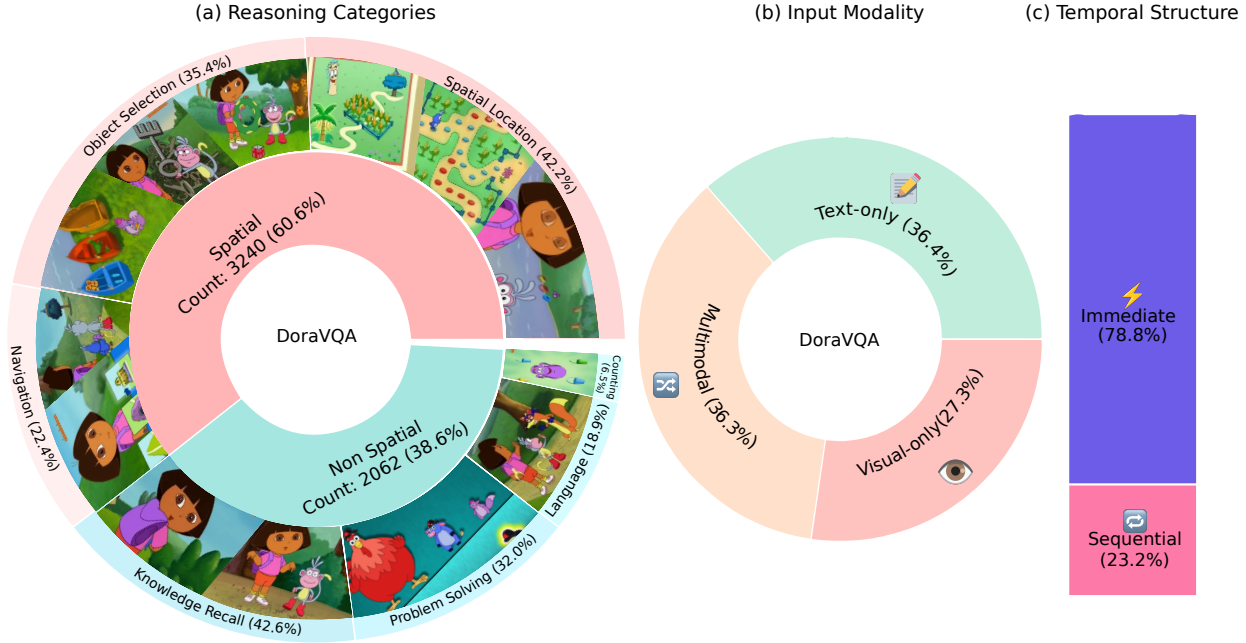


Figure 1. Our DoraVQA dataset composition across three key dimensions. (a) **Reasoning categories** divide into spatial tasks (60.6%), including object selection (35.4%), spatial location (42.2%), and navigation (22.4%), and non-spatial tasks (38.6%), encompassing language (18.9%), counting (6.5%), knowledge recall (42.6%), and problem solving (32.0%). (b) **Input modality** shows balanced distribution across text-only (36.4%), visual-only (27.3%), and multimodal questions (36.3%). (c) **Temporal structure** reveals that most questions require immediate reasoning (78.8%), while sequential reasoning across multiple frames accounts for 23.2%.

is explicitly taught through scaffolded question-answer sequences with clear correctness signals.

Spatial Reasoning in VLMs. Despite achieving strong performance on existing video question-answering benchmarks (Xiao et al., 2021; Xu et al., 2017), VLMs exhibit systematic failures on spatial reasoning tasks. Fu et al. (Fu et al., 2024) introduced BLINK, a benchmark for video spatial understanding, revealing VLMs achieve only half of human accuracy with noticeable failures on spatial reasoning tasks. Tong et al. (Tong et al., 2024) demonstrated that current VLMs achieve only 24-30% on tasks requiring precise spatial understanding, barely above random chance. These failures extend to fundamental spatial primitives: models struggle with counting objects beyond small numbers (Paiss et al., 2023), fail at compositional spatial relationships (Thrush et al., 2022), and fail to reason about object positions across video frames (Li et al., 2024). Importantly, these deficiencies persist even in simple visual scenes with minimal occlusion and clear object boundaries, indicating that the limitation is not perceptual but conceptual: models lack the structured spatial reasoning that children acquire through explicit instruction.

Educational and Instructional Video as Training Data. Instructional videos have been explored as training data for procedure understanding and action recognition. Miech et al. (Miech et al., 2019) demonstrated that instructional

videos provide natural language supervision for video-text embeddings. However, the narrations describe actions chain rather than teaching concepts; The CrossTask (Zhukov et al., 2019) and COIN (Tang et al., 2019) datasets extract procedural steps from instructional videos but focus on temporal ordering rather than spatial reasoning, yet these datasets capture implicit learning through observation rather than explicit instruction through question-answer structure. Closest to our work, Anderson and Burns (Anderson et al., 2000) analyzed the pedagogical structure of children’s educational television, documenting how the question-pause-answer format creates structured learning opportunities. We build on this insight by extracting the explicit question-pause-answer pairs from educational content as direct supervision signals, leveraging the pedagogical design for model training rather than treating video as passive observation.

Reinforcement Learning for VLM Fine-tuning. Group Relative Policy Optimization (GRPO) (Shao et al., 2024) has emerged as an efficient alternative to Proximal policy optimization (PPO) (Schulman et al., 2017) for reasoning tasks, eliminating the value network while maintaining strong performance through group-relative advantage estimation. DeepSeek-R1 (DeepSeek-AI, 2025) demonstrated that GRPO with rule-based rewards achieves reasoning capabilities competitive with proprietary models on mathematical reasoning benchmarks. Recent work has successfully applied GRPO to vision-language tasks: VLM-R1 (Shen

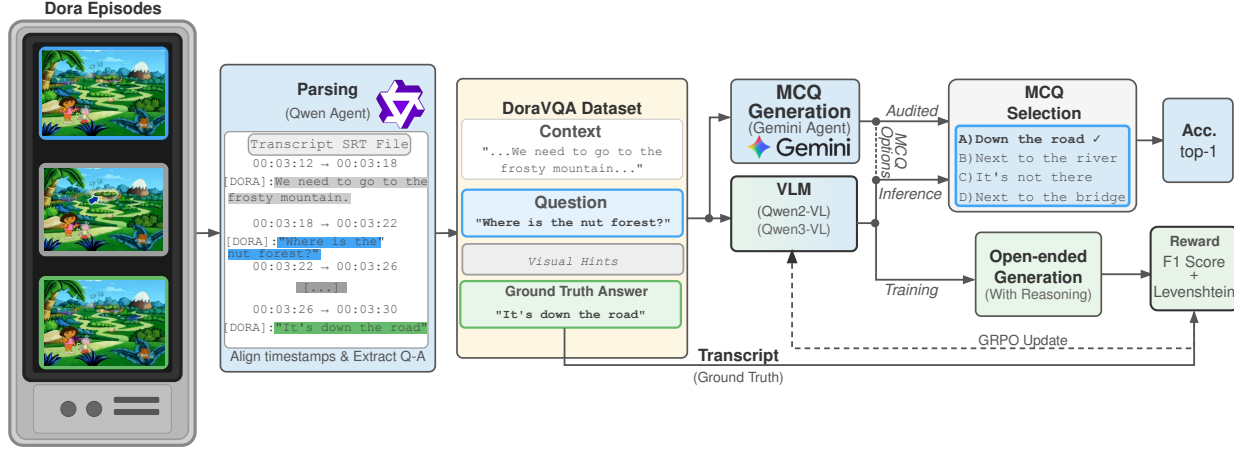


Figure 2. DoraVQA pipeline overview. We extract question-answer pairs from Dora the Explorer episodes by parsing SRT transcript files with a Qwen agent, aligning timestamps to identify the show’s pedagogical *context-question-pause-answer* structure. Each detected question is paired with its surrounding context window and the ground truth answer that follows. During **training**, we fine-tune Qwen2-VL and Qwen3-VL using GRPO on open-ended generation with reasoning; rewards are computed from F1 score and normalized Levenshtein distance against transcript ground truth. A Gemini agent generates multiple-choice distractors from the ground truth answers, which are human-audited for quality. During **inference**, the MCQ options are provided to the fine-tuned VLM for selection, creating a deliberate train-test format mismatch (open-ended → MCQ) that evaluates transferable reasoning. We report top-1 accuracy as the evaluation metric.

et al., 2025) showed that GRPO fine-tuning achieves superior out-of-distribution generalization compared to SFT on visual reasoning tasks, while CrowdVLM-R1 (Liu et al., 2025) developed specialized reward functions for counting tasks in crowded scenes. These works primarily design reward models for mathematical reasoning or general-purpose instruction following. In contrast, we apply GRPO to pedagogically-structured video content where the reward signal is implicit in the educational design, where ground-truth answers are directly in its structure. This eliminates the need for external reward model while providing clear correctness signals, creating a self-supervised RL training environment where the pedagogical content naturally defines the learning objectives.

Our work is positioned at the intersection of these research directions: we address video-language model spatial reasoning failures by training on educational content that explicitly teaches spatial concepts through structured question-answer pairs. Unlike general instructional videos that show procedures, educational television is designed to teach reasoning through direct instruction, repetition, and scaffolded examples. We leverage this pedagogical structure with GRPO fine-tuning, using the natural correctness signals from educational content as supervision without requiring manual annotation or reward engineering.

3. Method

Given a vision-language model pre-trained on large-scale multimodal data, we seek to improve its spatial reasoning capabilities by fine-tuning on pedagogically-structured

video content. The core challenge is to leverage the natural question-answer structure of educational videos as a self-supervised training signal, where the model must observe visual context, process explicit questions, and generate answers that align with the ground-truth responses provided by the educational content itself.

3.1. DoraVQA: Structure and Formatting

The DoraVQA dataset capitalizes on the inherent *context-question-pause-answer* structure of educational television (Figure 2). Each episode follows a deliberate pedagogical pattern: context is presented through both video frames (visual) and previous dialog (textual), a question is explicitly posed to the viewer, a pause allows for processing time, and finally a clear answer is provided with explanation. This structure creates a learning structure where every question exists within complete narrative and context. The answers are also definite, either directly contained within the context or can be deduced through simple steps.

To better format and adapt this structure for fine-tuning purposes, we format each example from DoraVQA dataset \mathcal{D} as a multimodal input $x = \{I, T, Q, a^*\} \in \mathcal{D}$ where I represents the visual frames sampled before and at the question timestamp, T is the transcript providing narrative context within a fixed time window before the question, Q is the explicit question and a^* is ground truth answer. The policy π_θ receives this input as a structured message combining images and text: the visual frames are embedded alongside a text prompt that concatenates the system instruction, transcript context, and question. The model processes this multimodal input through its vision-language encoder, gen-

erating a probability distribution over answer tokens. a^* is extracted from the content transcript, and serves as the target for reward computation, creating a self-supervised learning signal where the model learns to generate responses that match the show’s intended answers.

3.2. Group Relative Policy Optimization

We fine-tune the model using Group Relative Policy Optimization (GRPO), which enables efficient reinforcement learning without requiring a separate value network. GRPO is particularly well-suited for this task because educational content provides clear correctness signals (answers are objectively right or wrong) eliminating the need for complex reward modeling while the pedagogical structure offers implicit reasoning traces through the show’s explanations. For each input x , the policy generates K candidate answers $\{a_1, \dots, a_K\}$ through sampling. Each generated answer a_i is evaluated against the ground-truth a^* using a reward function $r(a_i, a^*)$ that computes token-level similarities between the extracted final answer and the ground truth:

$$r(a_i, a^*) = \alpha \text{F1}(a_i, a^*) + \beta \left(1 - \frac{\text{lev}(a_i, a^*)}{\max(|a_i|, |a^*|)} \right), \quad (1)$$

where the F1 score captures token-level semantic overlap while remaining insensitive to minor surface variations, and the normalized Levenshtein distance (lev) complements it by measuring character-level edits, penalizing incorrect re-ordering and near-miss errors. The weights α and β are fixed to 0.3 and 0.7 based on empirical validation. Because target answers are short and objectively defined, and optimization is performed via GRPO’s group-relative advantage, this reward discourages degenerate copying while providing a stable continuous signal: exact matches score 1, semantically correct answers with minor variations receive high scores (e.g., 0.65–0.95), and incorrect answers receive near-zero reward. Strong transfer to external multiple-choice benchmarks further suggests that learning extends beyond surface-level matching.

GRPO computes the advantage for each generation relative to the group average rather than requiring an absolute value estimate. For a group of K generations with rewards $\{r(a_1, a^*), \dots, r(a_K, a^*)\}$, the advantage for generation i is computed as:

$$A(a_i, a^*) = r(a_i, a^*) - \frac{1}{K} \sum_{j=1}^K r(a_j, a^*). \quad (2)$$

This group-relative advantage eliminates the need for a separate value network while providing effective learning signals. The policy is updated to maximize the expected advantage, increasing the probability of generating answers with higher relative rewards. The training objective optimizes:

$$\mathcal{L}(\theta) = -\mathbb{E}_{x \sim \mathcal{D}, a \sim \pi_\theta(\cdot|x)} [A(a, a^*) \log \pi_\theta(a|x)] \quad (3)$$

Algorithm 1 GRPO Training with Context-Question-Pause-Answer Structure

Require: Dataset \mathcal{D} with examples (I, T, Q, a^*) ; policy π_θ ; K generations per prompt

- 1: **for** each batch $\mathcal{B} \sim \mathcal{D}$ **do**
- 2: **for** each $(I, T, Q, a^*) \in \mathcal{B}$ **do**
- 3: Format input: $x = \{\text{images} : I, \text{context} : T, \text{question} : Q\}$
- 4: Generate K answers: $\{a_j\}_{j=1}^K \sim \pi_\theta(\cdot|x)$ with sampling
- 5: Compute rewards: $r_j = \text{F1}(\text{extract}(a_j), a^*)$ for $j = 1, \dots, K$
- 6: Compute advantages: $A_j = r_j - \frac{1}{K} \sum_{k=1}^K r_k$ for $j = 1, \dots, K$
- 7: **end for**
- 8: Update: $\theta \leftarrow \theta - \alpha \nabla_\theta \sum_{i,j} A_{i,j} \log \pi_\theta(a_{i,j}|x_i)$
- 9: **end for**

where $A(a, a^*) = \{A(a_1, a^*), \dots, A(a_K, a^*)\}$ is the group-relative advantage computed from the reward function. Through this process, the model learns to generate answers that better match the ground-truth responses from the educational content, improving its reasoning capabilities on visual understanding and compositional reasoning tasks.

An important aspect of our training setup is the distinction between training and evaluation: during training, the model learns to generate open-ended answers to questions from the educational content, where the reward signal incentivizes close word- and semantic-match between generated text and ground-truth free-form answers. However, for evaluation on benchmarks like Video-MME, we evaluate the model’s performance on multiple-choice questions (MCQs), where the model is prompted to only select from a set of candidate options. This training-evaluation format mismatch tests the model’s ability to transfer learned reasoning capabilities from open-ended generation to structured selection tasks, demonstrating generalization beyond the specific answer format seen during training.

3.3. Training

The GRPO training procedure for the context-question-pause-answer structure is presented in Algorithm 1 and illustrated in Figure 2. For each example x , π_θ generates K candidate answers (temperature=1.0, top-p=0.9). Each answer is evaluated using reward function $r(a_i, a^*)$, and group-relative advantages are computed. The policy is updated to maximize expected advantage. Training hyperparameters (learning rate=1e-4, KL coefficient=0.01, reward scaling=2.0, group size=8) are selected via sequential sweeps on validation data (Appendix A).

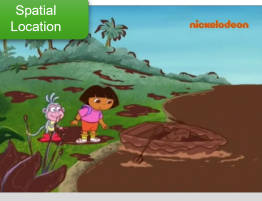

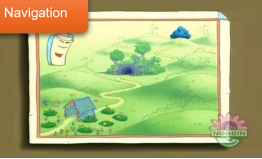

	Baseline	+GRPO (Ours)		Baseline	+GRPO (Ours)
Spatial Location  <p>Question: Do you see the chocolate boat? Is this the chocolate boat?</p> <p>Ground Truth Answer: That's the chocolate boat!</p>	<p>Qwen2-VL-2B ×No, we don't see the chocolate boat in this scene.</p> <p>Qwen3-VL-8B ×No, this is not the chocolate boat.</p>	<p>Qwen2-VL-2B + GRPO ✓Yes, the Chocolate Boat is the Chocolate Boat.</p> <p>Qwen3-VL-8B + GRPO ✓Yes, this is the chocolate boat.</p>	Object Selection  <p>Question: Which boat should we choose? Blue?</p> <p>Ground Truth Answer: The blue boat.</p>	<p>Qwen2-VL-2B ~We should choose the boat on the right side of the screen, as it has no holes and has a life jacket attached to it.</p> <p>Qwen3-VL-8B ×No, you should not choose the blue boat.</p>	<p>Qwen2-VL-2B + GRPO ✓The blue boat.</p> <p>Qwen3-VL-8B + GRPO ✓The blue boat.</p>
Navigation  <p>Question: Where do we go next? Garden... Web... Big Blue Bush.</p> <p>Ground Truth Answer: Spider Web.</p>	<p>Qwen2-VL-2B ×We go on to the Big Blue Bush next.</p> <p>Qwen3-VL-8B ×Based on the visual evidence and the context provided, the characters have just escaped the Flowery Garden by using the word "abre" to open a flower and rescue Mami Bugga Bugga.</p>	<p>Qwen2-VL-2B + GRPO ✓The web.</p> <p>Qwen3-VL-8B + GRPO ✓Web.</p>	Counting  <p>Question: Did we get past all four crocodiles?</p> <p>Ground Truth Answer: We got past them all.</p>	<p>Qwen2-VL-2B ×No, we did not get past all four crocodiles.</p> <p>Qwen3-VL-8B ×No, we did not get past all four crocodiles yet.</p>	<p>Qwen2-VL-2B + GRPO ✓Yes, we got past all four crocodiles.</p> <p>Qwen3-VL-8B + GRPO ✓Yes, we got past all four.</p>

Figure 3. Qualitative comparison on challenging spatial reasoning tasks. Our GRPO-finetuned models demonstrate superior performance on: (1) *Spatial location* where the chocolate boat is camouflaged against the chocolate river, (2) *Object selection* requiring identification of distant boats in the background, where baselines correctly identify spatial position but hallucinate object properties, (3) *Navigation* requiring sequential recall of past events to determine the next destination, and (4) *Counting* of visually similar crocodiles positioned close together, requiring frame-by-frame enumeration. Baseline models (orange boxes) fail or hallucinate, while our models (blue boxes) provide accurate, concise answers. See Appendix B and supplementary material for more examples.

4. Experiments and Results

In this section, we first detail the composition of the DoraVQA dataset and compare with other similar VQA datasets. Through cross benchmark transfer, we then show that large-scale data corpora are not strictly necessary to improve VLMs’ performance; carefully constructed pedagogical content can substantially enhance the reasoning and deduction capabilities of VLMs. Finally, we conduct ablation studies to prove the effectiveness of our modifications of training data structure and strategies. Notably, with only 5.3K QA pairs, our approach achieves state-of-the-art performance on CVBench (86.16%) and competitive results on NExT-QA (81.80%), outperforming commercial models trained on significantly larger datasets.

4.1. Datasets

We evaluate on the following VQA datasets: DoraVQA, Video-MME (Fu et al., 2025), CVBench (Zhu et al., 2025), NExT-QA (Xiao et al., 2021). For evaluation metrics, we report top-1 accuracy as percentage for all experiments.

DoraVQA: We introduce DoraVQA, a dataset of 5,344 question-answer pairs automatically extracted from transcripts of 8 seasons (96 episodes) of *Dora the Explorer* collected from publicly available online sources, with precise timestamp alignment. The dataset follows the show’s consistent *context-question-pause-answer* pedagogical structure, providing clear supervision. As illustrated in Figure 1, the

benchmark consists of 60.6% spatial tasks—comprising spatial location (42.2%), object selection (35.4%), and navigation (22.4%)—and 38.6% non-spatial tasks including knowledge recall, problem-solving, and counting. The dataset is balanced across text-only (36.4%), visual-only (27.3%), and multimodal (36.3%) input modalities, with 78.8% of questions requiring immediate reasoning and 23.2% requiring sequential reasoning across frames. While we plan to open-source our models and the DoraVQA annotations, **we do not redistribute any video frames or audio**; instead, we release only episode identifiers, temporal timestamps, question-answer annotations, and the annotation pipeline code, and users must independently access the original content under its existing licensing terms.

Video-MME is a comprehensive benchmark for evaluating multimodal LLMs across short-, medium-, and long-duration videos using visual, subtitle, and audio cues; in this work, we report results on the short partition only.

CVBench is designed to evaluate cross-video synergies, specifically testing a model’s proficiency in complex multimodal understanding and reasoning across different video segments.

NExT-QA focuses on high-level temporal action explanation, challenging models to perform causal and temporal reasoning to answer “why” and “how” questions beyond basic description.

Table 1. Performance comparison on DoraVQA test set. We evaluate our GRPO-finetuned models against state-of-the-art video-language models across different reasoning categories. Gray rows indicate proprietary models; underlined scores indicate results surpassed by our best model (Qwen3-VL-8B + GRPO). **Green** indicates improvement over baseline, **red** indicates degradation.

Model	Overall	Spatial	Counting	Navigation	Knowledge
Gemini-3.0-Flash	76.10	67.83	81.25	65.38	88.52
Gemini-3.0-Flash (low thinking)	75.84	66.81	87.50	69.23	86.89
GPT-4V	<u>67.79</u>	<u>60.34</u>	50.00	69.23	<u>68.85</u>
Gemini-2.5-Pro	<u>64.41</u>	62.62	50.00	<u>60.78</u>	<u>71.93</u>
LLaVA-Video-7B	<u>55.41</u>	<u>54.02</u>	<u>31.25</u>	66.00	<u>67.31</u>
Video-LLaVA-7B	<u>37.82</u>	<u>47.49</u>	<u>41.67</u>	<u>46.34</u>	<u>32.43</u>
InternVideo2.5-8B	<u>57.68</u>	<u>46.55</u>	<u>31.25</u>	<u>59.62</u>	<u>68.85</u>
Qwen2-VL-2B (baseline)	41.36	38.77	37.50	38.46	45.90
Qwen2-VL-2B + GRPO	55.11 (+13.75)	55.45 (+16.68)	28.57 (-8.93)	50.00 (+11.54)	50.67 (+5.77)
Qwen2-VL-7B (baseline)	56.74	53.02	43.75	51.92	60.66
Qwen2-VL-7B + GRPO	62.38 (+5.64)	59.74 (+6.72)	43.75 (±0.00)	64.58 (+12.66)	63.93 (+3.27)
Qwen3-VL-8B (baseline)	58.08	52.16	43.75	41.18	65.57
Qwen3-VL-8B + GRPO	67.98 (+9.90)	62.50 (+10.34)	37.50 (-6.25)	59.62 (+18.44)	72.13 (+6.56)

4.2. Main Results on DoraVQA

Table 1 presents our main results on the DoraVQA test set. GRPO fine-tuning on pedagogically-structured content yields substantial improvements across all model sizes, with gains inversely correlated with model capacity; demonstrating that structured data provides particularly strong training signals for smaller models.

The Qwen2-VL-2B model shows the largest improvement, jumping from 41.36% to 55.11% (+13.75 points) overall accuracy, with spatial reasoning improving by 16.68 points. The Qwen2-VL-7B model improves from 56.74% to 62.38% (+5.64 points), while Qwen3-VL-8B advances from 58.08% to 67.98% (+9.90 points). Notably, our Qwen3-VL-8B + GRPO model surpasses both Gemini-2.5-Pro (64.41%) and GPT-4V (67.79%), despite training exclusively on children’s educational content.

The improvements are particularly pronounced on spatial reasoning tasks: navigation accuracy increases by 12.66 points for Qwen2-VL-7B and 18.44 points for Qwen3-VL-8B, demonstrating that the pedagogical structure effectively teaches spatial concepts. Figure 3 illustrates these improvements qualitatively; our models correctly identify camouflaged objects (chocolate boat), avoid hallucinating object properties in background scenes (boat selection), successfully perform sequential reasoning (navigation), and accurately count visually similar objects (crocodiles), while baseline models fail or generate unwanted details. However, counting tasks demand explicit visual grounding and cannot be solved through linguistic pedagogy alone; our pedagogical structure can teach reasoning patterns, but it cannot replace visual perception in counting tasks. Specifically, we observe performance degradation in both Qwen2-VL-2B and Qwen3-VL-8B, suggesting that structured pedagogy alone is insufficient for this task. Knowledge recall shows

consistent gains across all models (3.27-6.56 points), indicating effective transfer of the show’s explicit teaching of facts and problem-solving strategies.

4.3. Cross-Benchmark Transfer

Table 2 demonstrates that training on pedagogically-structured content transfers effectively to diverse video understanding benchmarks, rather than being constrained by the training data distribution. All models show consistent improvements across external benchmarks despite zero-shot evaluation and using only 5.3K QA pairs compared to millions of videos used by competing approaches.

The largest gains are observed on CVBench, which explicitly requires cross-video integration and knowledge transfer, where structured knowledge recall proves particularly effective. Qwen2-VL-2B improves by 28.93 points, Qwen2-VL-7B by 24.51 points, and Qwen3-VL-8B by 40.36 points, reaching state-of-the-art performance of 86.16%. NExT-QA, which emphasizes causal and temporal reasoning, exhibits similarly strong transfer, with improvements of 20.26, 12.46, and 19.70 points, respectively. Video-MME shows more modest but consistent gains (1.62–12.01 points) across all models. Together, these results demonstrate that structured pedagogical training induces reasoning capabilities that generalize beyond the narrow domain of children’s educational content, providing empirical support for our structure-over-scale hypothesis.

4.4. Ablation Studies

We conduct ablation studies on training strategies and input modalities using a Qwen2-VL-2B model fine-tuned for 100 epochs on DoraVQA, evaluating performance on 10% randomly sampled test subsets of Video-MME, CVBench, and NExT-QA.

Table 2. Cross-benchmark transfer evaluation highlighting **structure vs. scale**. Despite training on only 5.3K QA pairs (38 hours), our models achieve competitive or state-of-the-art results against models trained on orders of magnitude more data. Gray rows indicate proprietary models; yellow highlight indicates our state-of-the-art result; blue cells highlight our minimal fine-tuning scale. Green indicates improvement over baseline. † indicates zero-shot evaluation; * indicates training set seen during pre-training. Baseline “Training Scale” refers to pre-training data; our GRPO models use the same base but fine-tune on only 5.3K pedagogically-structured QA pairs (38 hours of educational video).

Model	Params	Training Scale	DoraVQA	Video-MME	CVBench	NExT-QA
GPT-4V [†]	~1.8T	~10T tokens	67.79	59.9	52.4	68.2
Gemini-2.5-Pro [†]	—	—	64.41	85.2	62.4	74.6
Gemini-3.0-Flash [†]	—	—	76.10	86.9	67.2	80.4
InternVideo2.5-8B	8B	16M clips	57.68	65.4	57.3 [†]	71.5*
LLaVA-Video-7B	7B	1.3M videos	55.41	63.3	52.6 [†]	83.2*
Video-LLaVA-7B	7B	760K videos	37.82	45.3	28.1 [†]	52.1*
Qwen2-VL-2B (baseline)	2B	1.2T tokens	41.36	50.10	31.38	52.60
Qwen2-VL-2B + GRPO [†]	2B	5.3K QA	55.11 (+13.75)	62.11 (+12.01)	60.31 (+28.93)	72.86 (+20.26)
Qwen2-VL-7B (baseline)	7B	1.2T tokens	56.74	67.5	50.7	67.0
Qwen2-VL-7B + GRPO [†]	7B	5.3K QA	62.38 (+5.64)	69.12 (+1.62)	75.21 (+24.51)	79.46 (+12.46)
Qwen3-VL-8B (baseline)	8B	1.0T tokens	58.08	71.4	45.8	62.1
Qwen3-VL-8B + GRPO [†]	8B	5.3K QA	67.98 (+9.90)	76.78 (+5.38)	86.16 (+40.36)	81.80[†] (+19.70)

Table 3. Ablation study on training method. We compare supervised fine-tuning (SFT) against our GRPO approach on in-domain (DoraVQA) and out-of-domain benchmarks.

Method	Dora	Video-MME	CVB	NExT
No finetuning	41.36	50.10	31.38	52.60
+ SFT	3.56	53.33	43.56	66.63
+ GRPO	55.11	55.40	60.31	72.86

Training Method. Our goal is to enhance the VLM’s reasoning capabilities in general domains rather than tailoring it to a specific dataset or task. As shown in Table 3, we demonstrate that the RL-based GRPO policy enables Qwen2-VL to transfer the knowledge acquired from generating free-form answers to making discrete choices. In contrast, SFT substantially reduces the model’s flexibility due to overfitting: specifically for DoraVQA, after SFT, the model can no longer make choices and instead produces only free-form responses. Additionally, GRPO showed significant performance gain across the board compared with SFT, indicating better generalization ability.

Context Modalities. In Table 4, we evaluate how different input modalities affect model performance. We observed that the model learns more effectively from transcripts than from visual frames, suggesting that its text-processing component is stronger than its visual counterpart, yielding better alignment to language patterns than visual grounding. This finding is also consistent with the observations in the Qwen2 technical report (Wang et al., 2024a), revealing a limitation of current VLMs. Combining the two modalities yields a modest performance improvement across all benchmarks, indicating the effectiveness of our proposed

Table 4. Ablation study on context modalities. We evaluate the impact of transcript context, visual frames, and their combination on spatial reasoning performance.

Context	Dora	Video-MME	CVB	NExT
Visual (V) only	43.42	60.74	47.35	67.91
Transcript (T) only	50.95	60.00	50.38	70.48
V+T, no Context (C)	54.72	61.11	53.41	70.83
Full (V+T+C)	55.11	62.20	60.31	72.86

structured learning format and DoraVQA dataset. Incorporating text context before the question further enhances performance – CVBench, in particular, gains 7%, reflecting its synergistic nature and its reliance on more connected contextual information to perform well.

5. Conclusion

We demonstrated that pedagogically-structured educational videos provide effective training signals for vision-language model spatial reasoning. By fine-tuning Qwen2-VL and Qwen3-VL on just 5.3K QA pairs from DoraVQA using GRPO, we achieved state-of-the-art performance on CVBench (86.16%) and consistent improvements across Video-MME and NExT-QA—outperforming models trained on orders of magnitude more data. These results validate that content structure can compensate for content scale. Future work should incorporate visual reward modeling to capture motion-based teaching signals during pause, and expand beyond a single show to a comprehensive Pedagogical Interactive Structure (PIS) dataset spanning diverse educational content (Blue’s Clues, Khan Academy, PhET simulations) across multiple levels and reasoning domains.

6. Impact Statements

This paper presents DoraVQA and our effort to fine-tune open-source VLMs using this dataset and subsequently evaluate them across existing VLMs and benchmarks. Our goal is to advance the field of computer vision and machine learning. Although there are many potential societal consequences of our work, we would specifically highlight the following. First, our dataset does not redistribute any video frames, audio, or raw episode files. All visual and audio content remains hosted by the original rights holders, and users are required to independently obtain access to the corresponding episodes under the platforms’ existing terms of service. We release only episode identifiers, temporal timestamps, processed question–answer annotations, and manual transcript span indices, along with the code to align these annotations with externally hosted content. As the source material consists of publicly broadcast children’s educational programming, the dataset does not introduce new privacy risks. Nevertheless, DoraVQA is intended strictly for research and evaluation purposes, and we discourage any use outside this scope. Second, there are no human subjects involved in any of our experiments.

References

- Abu-El-Haija, S., Kothari, N., Lee, J., Natsev, P., Toderici, G., Varadarajan, B., and Vijayanarasimhan, S. Youtube-8m: A large-scale video classification benchmark. In *arXiv preprint arXiv:1609.08675*, 2016.
- Anderson, D. R., Bryant, J., Wilder, A., Santomero, A., Williams, M., and Crawley, A. M. Researching blue’s clues: Viewing behavior and impact. *Media Psychology*, 2(2):179–194, 2000. doi: 10.1207/S1532785XMEP0202.4.
- Angin, D. E. The effects of Dora the Explorer on preschool children’s spatial concept acquisition and spatial ability. *European Scientific Journal*, 13(1):39–53, 2017. doi: 10.19044/esj.2017.v13n1p39.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Fu, C., Dai, Y., Luo, Y., Li, L., Ren, S., Zhang, R., Wang, Z., Zhou, C., Shen, Y., Zhang, M., et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 24108–24118, 2025.
- Fu, X., Hu, Y., Li, B., Feng, Y., Wang, H., Lin, X., Roth, D., Smith, N. A., Ma, W.-C., and Krishna, R. Blink: Multi-modal large language models can see but not perceive. In *ECCV*, 2024.
- Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- Li, Z., Xu, Y., Wang, Z., Hu, Y., Xie, H., Liu, M., Chen, T., Huang, F., Huang, Z., and Yang, H. Tomato: Assessing visual temporal reasoning capabilities in multimodal foundation models. *arXiv preprint arXiv:2311.14468*, 2024.
- Lin, B., Zhu, B., Ye, Y., Ning, M., Jin, P., and Yuan, L. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2024.
- Linebarger, D. L. and Walker, D. Infants’ and toddlers’ television viewing and language outcomes. *American Behavioral Scientist*, 48(5):624–645, 2005. doi: 10.1177/0002764204271505.
- Liu, Y., Chen, X., Wang, H., and Zhang, M. Crowdvlm-r1: Advancing vision-language models for crowd counting via reinforcement learning. *arXiv preprint arXiv:2503.03724*, 2025.
- Miech, A., Zhukov, D., Alayrac, J.-B., Tapaswi, M., Laptev, I., and Sivic, J. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2630–2640, 2019.
- Paiss, R., Ephrat, A., Tschannen, O., Zhai, X., Gilad-Bachrach, R., Amir, I., Grattarola, D., Rubenstein, M., Baldrige, J., Elharrar, M., et al. Teaching clip to count to ten. *arXiv preprint arXiv:2302.12066*, 2023.
- Pozzulo, J. D., Dempsey, J., Bruer, K., and Sheahan, C. The culprit in target-absent lineups: Understanding young children’s false positive responding. *Journal of Police and Criminal Psychology*, 27(1):55–62, 2012.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Zhang, M., Li, Y., Wu, Y., and Guo, D. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Shen, Y., Zhang, T., Liu, X., Wang, C., and Zhang, L. Vlm-r1: Enhancing vision-language models with structured reasoning via reinforcement learning. *arXiv preprint arXiv:2504.07615*, 2025.
- Tang, Y., Ding, D., Rao, Y., Zheng, Y., Zhang, D., Zhao, L., Lu, J., and Zhou, J. Coin: A large-scale dataset for

- comprehensive instructional video analysis. In *CVPR*, 2019.
- Thrush, T., Jiang, R., Bartolo, M., Singh, A., Williams, A., Kiela, D., and Ross, C. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5238–5248, 2022.
- Tong, S., Liu, Z., Zhai, Y., Ma, Y., LeCun, Y., and Xie, S. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *CVPR*, 2024.
- Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024a.
- Wang, Y., Li, K., Li, X., Yu, J., He, Y., Cheng, G., Chen, B., Ouyang, T., Dou, X., Liu, J., et al. Internvideo2: Scaling video foundation models for multimodal video understanding. In *European Conference on Computer Vision (ECCV)*, pp. 543–563, 2024b.
- Xiao, J., Shang, X., Yao, A., and Chua, T.-S. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9777–9786, 2021.
- Xu, D., Zhao, Z., Xiao, J., Wu, F., Zhang, H., He, X., and Zhuang, Y. Video question answering via gradually refined attention over appearance and motion. In *ACM Multimedia*, 2017.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Zhang, Y., Li, B., Zhang, H., Li, L., Gao, D., Zhang, R., and Liu, Z. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024.
- Zhu, N., Dong, Y., Wang, T., Li, X., Deng, S., Wang, Y., Hong, Z., Geng, T., Niu, G., Huang, H., et al. Cvbench: Evaluating cross-video synergies for complex multimodal understanding and reasoning. *arXiv preprint arXiv:2508.19542*, 2025.
- Zhukov, D., Alayrac, J.-B., Cinbis, R. G., Fouhey, D., Laptev, I., and Sivic, J. Crosstask: A dataset for cross-task learning from instructional videos. In *CVPR*, 2019.

A. Hyperparameter Sweeps

We conduct sequential hyperparameter sweeps to identify optimal GRPO training configurations. Each sweep fixes previously optimized hyperparameters while varying a single dimension.

Table 5. Sequential hyperparameter sweeps for GRPO training. Each sweep fixes previously optimized parameters. Bold indicates selected values used in main experiments.

Parameter	Values Tested	Selected	DoraVQA
KL Coefficient	{0.0, 0.001, 0.01 , 0.1}	0.01	{42.96, 40.71, 46.38 , 45.19}
Learning Rate	{1e-6, 1e-5, 1e-4 , 1e-3}	1e-4	{43.69, 45.19, 50.09 , 38.89}
Reward Scaling	{0.5, 1.0, 2.0 , 5.0}	2.0	{45.92, 45.19, 45.69 , 43.93}

Table 5 presents our hyperparameter optimization results. The KL coefficient controls the divergence penalty between the policy and reference model; we find 0.01 provides the best balance, with both higher (0.1) and lower (0.001) values degrading performance. For learning rate, 1e-4 achieves optimal performance (50.09%), while higher rates (1e-3) cause training instability (38.89%) and lower rates (1e-6) result in insufficient learning. Reward scaling at 2.0 provides marginal improvements over the default 1.0, amplifying the learning signal from correct answers without overwhelming the policy gradient.

Table 6 examines training duration and overfitting characteristics. Qwen2-VL-7B peaks at 150 steps (62.96%) before degrading sharply at 250 steps (0.56%), indicating catastrophic overfitting to the narrow educational domain. Qwen3-VL-2B shows more robust training, peaking later at 250 steps (67.97%) with gradual performance decline. These results highlight the importance of early stopping when fine-tuning on domain-specific content to preserve generalization capabilities.

Table 6. Ablation study on training steps for both models. Performance peaks at 150 steps then degrades due to overfitting on the narrow educational domain.

Model	50	100	150	200	250	300
Qwen2-VL-7B	57.49	60.29	62.96	60.98	0.56	60.84
Qwen3-VL-2B	63.67	66.67	65.91	65.54	67.97	66.60

Based on these sequential sweeps, we use KL coefficient = 0.01, learning rate = 1e-4, and reward scaling = 2.0 for all main experiments reported in this paper. Training duration is set at 150 steps for Qwen2-VL-7B and 250 steps for Qwen3-VL-2B to avoid overfitting while maximizing in-domain performance.

B. Additional Qualitative Results

Figure 4 presents additional challenging examples from DoraVQA that highlight common failure modes of baseline models. The spatial location example (left) demonstrates partial occlusion challenges: Swiper the fox is only partially visible behind a blue wall, requiring the model to recognize objects from incomplete visual information. The baseline Qwen2-VL-2B fails to detect the occluded character, while the GRPO-finetuned version correctly identifies the location. The counting example (right) illustrates fine-grained enumeration difficulties: the key contains 8 identical points positioned closely together, requiring systematic frame-by-frame counting. Both baseline models overestimate the count (9-10), while GRPO models produce the correct answer (8). These examples demonstrate that pedagogical training improves both spatial perception under occlusion and precise counting of visually similar objects. For more examples and interactive demonstrations, refer to the interactive supplementary material.

C. Data Access and Licensing

DoraVQA does not redistribute any video frames, audio, or raw episode files. All visual and audio content remains hosted by the original rights holders, and users are required to independently obtain access to the corresponding episodes under the platforms’ existing terms of service. We release only episode identifiers, temporal timestamps, question–answer annotations, and transcript span indices, along with the code to align these annotations with externally hosted content. As the source

material consists of publicly broadcast children’s educational programming, the dataset does not introduce new privacy risks; nevertheless, DoraVQA is intended strictly for research and evaluation purposes, and we discourage any use outside this scope.

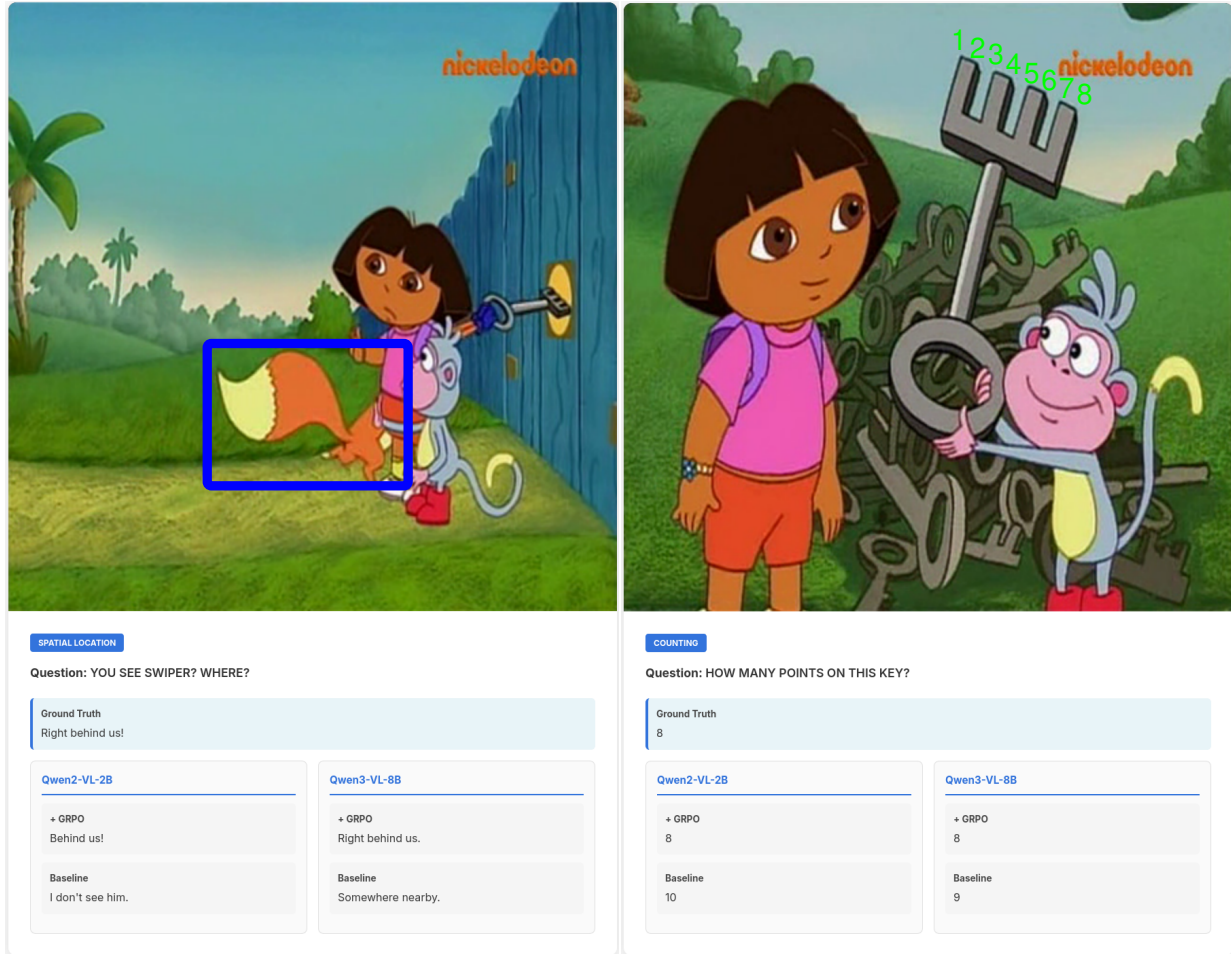


Figure 4. **Additional challenging examples from DoraVQA.** Left: Spatial location task where Swiper the fox is partially occluded behind Dora, requiring detection of partially visible objects. Right: Counting task requiring enumeration of 8 identical points on a key positioned closely together. Baseline models fail on both tasks, while our GRPO-finetuned models provide correct answers.