A DOCCI-CRITIQUE AUTORATER LEADERBOARDS

This appendix presents the complete leaderboards detailing the performance of various Vision-Language Models (VLMs) as automated rankers (AutoRaters) on the DOCCI-Critique benchmark. These tables supplement the summary correlation metrics (Spearman's ρ and Kendall's τ) found in Section 4.2, Table 3. Our goal was to assess how well automated methods, including our VNLI-Critique, rank caption-generating VLMs by factual accuracy against human-derived Ground Truth rankings.

Three leaderboards are provided, each for a distinct factuality criterion:

- Response-Level Correctness: Percentage of entirely factually accurate paragraphs (Table 7).
- Correct Sentences Overall: Total percentage of correct sentences across all descriptions (Table 9).
- Correct Sentences per Description: Average percentage of correct sentences per description (Table 8).

In each leaderboard (Tables 7, 9, and 8), rows list the 14 caption-generating VLMs from DOCCI-Critique (details in Table 2). Columns denote automated ranking methods (e.g., 'Ours (VNLI-Critique)', 'GPT-4o'). Cells show the rank assigned by the column's method to the row's VLM, with the superscript indicating the raw metric score. The final two rows report Spearman's ρ (with p-value superscript) and Kendall's τ correlations against the Ground Truth for that criterion, offering a nuanced view of each AutoRater's performance.

Table 7: VLM AutoRater rankings for Response-Level Correctness on DOCCI-Critique. Each cell shows $Rank^{Metric-Score}$. Final two rows: Spearman's $\rho^{p-value}$ and Kendall's $\tau^{p-value}$ correlation against Human.

| Ranking Method Captioner VLM | Human | Ours | Gemini 2.0-Flash | GPT-40 | InstructBLIP | LLaVa-OV | Janus-Pro-7B | Qwen2.5-VL | mPLUG-Owl3-7B | Emu3-Chat |
|------------------------------|-------------|----------------------|---------------------|---------------|----------------|---------------|---------------|---------------|---------------|----------------|
| MiniGPT-4 | 140.04 | 140.06 | $14^{0.09}$ | $14^{0.12}$ | 70.96 | $14^{0.42}$ | 130.35 | 40.00 | 130.11 | 30.93 |
| MPlugOwl-2 | 130.11 | 120.12 | 130.28 | $12^{0.18}$ | 20.98 | 120.56 | $10^{0.56}$ | 40.00 | 120.17 | 90.55 |
| LLaVA | 110.19 | 110.17 | 90.40 | $11^{0.26}$ | 30.97 | $10^{0.71}$ | 50.68 | 40.00 | 80.21 | 120.48 |
| PALI-5B | 120.11 | 130.09 | 120.31 | $13^{0.17}$ | 90.96 | $11^{0.58}$ | 140.19 | 40.00 | 140.07 | 20.99 |
| VILA | $10^{0.21}$ | 100.20 | 100.39 | 90.33 | 130.91 | 90.85 | 10.79 | $4^{0.00}$ | 80.21 | 40.78 |
| InstructBLIP | 90.26 | 90.21 | 110.38 | $10^{0.31}$ | 90.96 | 130.54 | 90.60 | $2^{0.01}$ | $6^{0.26}$ | $10^{0.54}$ |
| Molmo-7B-D | $8^{0.31}$ | 80.27 | 70.68 | 80.58 | 10.99 | 50.93 | 40.70 | 40.00 | 80.21 | 40.78 |
| LLaVA-OV-7B-Chat | 70.36 | $6^{0.35}$ | 80.63 | 50.64 | 120.93 | 80.90 | 80.66 | 40.00 | 40.30 | 80.58 |
| Owen2-VL-7B-Instruct | 50.41 | 40.45 | 50.75 | $4^{0.65}$ | 30.97 | $6^{0.91}$ | $2^{0.78}$ | 40.00 | 50.29 | 110.49 |
| LLaVA-OV-7B | 50.41 | 70.34 | $6^{0.72}$ | 70.63 | 140.90 | $2^{0.97}$ | $6^{0.67}$ | $2^{0.01}$ | 20.41 | 70.61 |
| Gemini-1.5-Pro | 40.64 | 50.43 | 40.84 | 30.67 | 30.97 | $6^{0.91}$ | 120.49 | 40.00 | 70.24 | 140.31 |
| Gemini-1.5-Flash | 30.68 | $3^{0.52}$ | 30.88 | 50.64 | 110.94 | 40.94 | $11^{0.51}$ | 40.00 | 110.18 | 130.40 |
| mPLUG-Owl3-7B | $2^{0.71}$ | 20.69 | $2^{0.93}$ | $2^{0.77}$ | 70.96 | 10.98 | 30.73 | 10.02 | 10.75 | 11.00 |
| GPT-4o[2024-08-06] | $1^{0.83}$ | $1^{0.73}$ | $1^{0.94}$ | $1^{0.89}$ | 30.97 | $2^{0.97}$ | 60.67 | 40.00 | 30.37 | 60.74 |
| Spearman's Rank ρ | - | 0.98^{5e-10} | 0.97^{6e-9} | 0.92^{3e-6} | -0.06^{8e-1} | 0.88^{2e-5} | 0.30^{3e-1} | 0.29^{3e-1} | 0.73^{3e-3} | -0.2^{5e-1} |
| Kendall Tau τ | - | 0.93 ^{4e-6} | 0.89^{1e-5} | 0.82^{5e-5} | -0.05^{8e-1} | 0.76^{2e-4} | 0.21^{3e-1} | 0.25^{3e-1} | 0.27^{5e-3} | -0.17^{4e-1} |

Table 8: VLM AutoRater rankings for Average Percentage of Correct Sentences on DOCCI-Critique. Each cell shows $Rank^{Metric-Score}$. Final two rows: Spearman's $\rho^{p-value}$ and Kendall's $\tau^{p-value}$ correlation against Human.

| Ranking Method Captioner VLM | Human | Ours | Gemini 2.0-Flash | GPT-40 | InstructBLIP | LLaVa-OV | Janus-Pro-7B | Qwen2.5-VL | mPLUG-Owl3-7B | Emu3-Chat |
|---------------------------------|--------|---------------|---------------------|----------------|---------------|---------------|---------------|---------------|---------------|----------------|
| MiniGPT-4 | 140.46 | 140.48 | $14^{0.53}$ | $14^{0.42}$ | 80.99 | 140.84 | 130.77 | 130.05 | 140.51 | 30.99 |
| MPlugOwl-2 | 130.53 | 130.54 | 130.66 | 130.56 | 70.99 | 120.89 | 120.86 | 110.05 | 130.53 | 120.86 |
| LLaVA | 120.60 | $11^{0.59}$ | $11^{0.75}$ | $12^{0.59}$ | 90.99 | 110.91 | 90.90 | 90.08 | 110.58 | 140.83 |
| InstructBLIP | 110.61 | $12^{0.58}$ | 120.71 | $11^{0.62}$ | 120.98 | 130.86 | 100.89 | $10^{0.05}$ | 120.56 | 130.84 |
| PALI-5B | 100.67 | 100.68 | 100.79 | $10^{0.67}$ | 31.00 | 100.91 | 140.73 | 120.05 | 100.61 | 21.00 |
| VILA | 90.78 | 80.78 | 90.86 | 90.81 | 110.99 | 90.98 | 10.97 | 40.22 | 80.76 | 40.96 |
| mPLUG-Owl3-7B | 80.80 | 60.80 | 40.98 | 80.87 | 130.98 | 50.99 | 110.87 | 140.05 | 20.85 | 11.00 |
| LLaVA-OV-7B | 70.82 | 60.80 | 70.94 | $6^{0.91}$ | 140.97 | 11.00 | 50.94 | 60.20 | 70.81 | 80.92 |
| Molmo-7B-D | 60.83 | 90.78 | 60.95 | 70.90 | 11.00 | 70.99 | 40.94 | 80.10 | 90.73 | 50.95 |
| LLaVA-OV-7B-Chat | 50.86 | 50.85 | 80.94 | 40.94 | 100.99 | 80.99 | 30.94 | 10.28 | 10.85 | 70.93 |
| Owen2-VL-7B-Instruct | 40.88 | 40.88 | 50.97 | 50.93 | 41.00 | 60.99 | 20.97 | 50.22 | 50.83 | 90.91 |
| Gemini-1.5-Pro | 30.95 | 30.93 | 30.99 | 20.97 | 21.00 | 40.99 | 80.93 | 30.25 | 40.83 | 110.89 |
| Gemini-1.5-Flash | 20.96 | 20.94 | 20.99 | 30.95 | 50.99 | 30.99 | 70.93 | 20.27 | 30.84 | 100.90 |
| GPT-4o[2024-08-06] | 10.97 | 10.95 | 10.99 | 10.98 | 60.99 | 21.00 | 60.93 | 70.17 | 60.82 | 60.95 |
| Spearman's Rank ρ | - | 0.97^{1e-8} | 0.96^{9e-8} | 0.99^{7e-11} | 0.35^{2e-1} | 0.85^{1e-4} | 0.58^{3e-2} | 0.70^{5e-3} | 0.80^{6e-4} | 0.00^{1e-0} |
| Kendall Tau τ | - | 0.91^{7e-6} | 0.90^{2e-7} | 0.93^{1e-8} | 0.14^{5e-1} | 0.74^{7e-5} | 0.40^{4e-2} | 0.50^{1e-2} | 0.65^{7e-4} | -0.05^{8e-1} |

Table 9: VLM AutoRater rankings for Percentage of Correct Sentences Overall on DOCCI-Critique. Each cell shows $Rank^{Metric-Score}$. Final two rows: Spearman's $\rho^{p-value}$ and Kendall's $\tau^{p-value}$ correlation against Human.

| Ranking Method | Human | Ours | Gemini | GPT-40 | InstructBLIP | LLaVa-OV | Janus-Pro-7B | Owen2.5-VL | mPLUG-Owl3-7B | Emu3-Chat |
|----------------------|-------------|----------------------|---------------|---------------|-------------------|----------------|---------------|-------------------|--------------------|---------------|
| Captioner VLM | ruillali | Ours | 2.0-Flash | GF 1-40 | HISHUCIBLIF | LLa va-Ov | Janus-F10-7B | Qweii2.3-VL | IIIFLUG-UWI3-7B | Elliu5-Cliat |
| MiniGPT-4 | 140.48 | 140.49 | $14^{0.55}$ | $14^{0.44}$ | 70.99 | 140.83 | 130.77 | 110.06 | 130.52 | 30.99 |
| MPlugOwl-2 | $13^{0.52}$ | 130.53 | $13^{0.64}$ | $13^{0.56}$ | 7 ^{0.99} | 120.87 | $11^{0.85}$ | $13^{0.05}$ | 14 ^{0.51} | 120.86 |
| InstructBLIP | $12^{0.57}$ | 120.57 | $12^{0.68}$ | $11^{0.59}$ | 110.99 | 130.84 | 100.87 | 120.05 | 120.54 | $14^{0.82}$ |
| LLaVA | $11^{0.59}$ | 110.59 | $10^{0.74}$ | 120.59 | 70.99 | 100.90 | 90.89 | 90.08 | $11^{0.58}$ | 130.83 |
| PALI-5B | $10^{0.67}$ | 100.66 | $10^{0.74}$ | 100.62 | 41.00 | 110.88 | 140.73 | 100.07 | $10^{0.59}$ | $2^{1.00}$ |
| VILA | 90.79 | 80.79 | 90.86 | 90.81 | 110.99 | 90.98 | 10.97 | $4^{0.22}$ | 80.77 | 40.96 |
| LLaVA-OV-7B | 80.82 | 70.81 | $8^{0.95}$ | $6^{0.91}$ | 130.98 | 11.00 | 50.93 | $6^{0.21}$ | 20.85 | $9^{0.91}$ |
| mPLUG-Owl3-7B | 70.82 | $6^{0.82}$ | 40.96 | 80.84 | 140.98 | 80.99 | 120.81 | 140.05 | 30.84 | 11.00 |
| Molmo-7B-D | $6^{0.83}$ | 90.79 | $6^{0.95}$ | 70.90 | 11.00 | $6^{0.99}$ | 40.94 | 80.10 | 90.72 | 50.96 |
| Owen2-VL-7B-Instruct | 50.87 | 40.88 | 40.96 | 50.93 | 21.00 | 50.99 | 20.97 | 50.22 | 60.83 | 80.91 |
| LLaVA-OV-7B-Chat | 40.87 | 50.88 | $6^{0.95}$ | 40.94 | 100.99 | 70.99 | 30.96 | 10.31 | 10.87 | $7^{0.92}$ |
| Gemini-1.5-Pro | $3^{0.95}$ | 30.93 | 30.99 | $2^{0.97}$ | 21.00 | $4^{0.99}$ | 80.93 | $3^{0.25}$ | 50.84 | 110.89 |
| Gemini-1.5-Flash | 20.96 | $2^{0.94}$ | 10.99 | $3^{0.95}$ | 51.00 | 11.00 | 50.93 | 20.27 | 40.84 | 90.91 |
| GPT-4o[2024-08-06] | $1^{0.97}$ | $1^{0.95}$ | $1^{0.99}$ | $1^{0.98}$ | 5 ^{1.00} | 11.00 | 50.93 | 7 ^{0.16} | 70.81 | $6^{0.95}$ |
| Spearman's Rank ρ | - | 0.98 ^{1e-9} | 0.98^{2e-9} | 0.97^{2e-9} | 0.37^{2e-1} | 0.86^{10e-5} | 0.52^{6e-2} | 0.69^{6e-3} | 0.74^{2e-3} | 0.06^{8e-1} |
| Kendall Tau τ | - | 0.91 ^{5e-8} | 0.91^{5e-8} | 0.91^{5e-8} | 0.19^{4e-1} | 0.76^{4e-5} | 0.34^{1e-1} | 0.52^{10e-3} | 0.58^{3e-3} | 0.01^{1e+0} |

B VNLI-CRITIQUE MODEL DEVELOPMENT DETAILS

This section provides further details on the architecture, fine-tuning process, and computational resources utilized for the development of our VNLI-Critique model, as introduced in Section 4 of the main paper.

B.1 MODEL ARCHITECTURE

VNLI-Critique is developed by fine-tuning the PaliGemma 10B architecture Steiner et al. (2024). This architecture integrates a Gemma2-9B Large Language Model (LLM) Rivière et al. (2024) as its textual backbone and a SigLIP model Zhai et al. (2023) as its visual encoder. For visual processing, input images are standardized to a resolution of $448px^2$ pixels. At this resolution, the SigLIP visual encoder processes each image into a sequence of 1024 visual tokens, which are subsequently fed into the LLM component for multimodal understanding and generation tasks.

B.2 Fine-tuning Procedure

We performed full fine-tuning of the PaliGemma 10B model to develop VNLI-Critique. The fine-tuning process was conducted for 5 epochs. A batch size of 128 was used, with a dropout rate of 0.1 applied to aid regularization. No weight decay was utilized during training. The Adam optimizer Kingma & Ba (2015) was employed with its default hyperparameters, and a constant learning rate of 1×10^{-6} was maintained throughout the fine-tuning process.

B.3 Computational Resources

The training of the VNLI-Critique model was executed on Google Cloud TPUv5e Google Cloud (20xx) accelerators. Specifically, a configuration of 128 TPUv5e chips was utilized for the fine-tuning task. The total training time for the 5 epochs was approximately 1 hour and 30 minutes. Based on an estimated cost of \$1.20 per chip-hour, the total computational cost for training VNLI-Critique was approximately \$230.40.

C HUMAN ANNOTATION DETAILS

The creation of the DOCCI-Critique benchmark and the evaluation of our models' outputs, including critique generation and the Critic-and-Revise pipeline, relied on comprehensive human annotations. We engaged third-party human annotators sourced through Prolific². Each data entry subject to human evaluation, whether for sentence-level factuality in DOCCI-Critique or for the quality assessment of generated critiques, was independently assessed by five different annotators. This

²https://www.prolific.com/

multi-annotator approach helps ensure robustness and mitigate individual biases in the collected judgments. Annotators were compensated at a rate of \$20 per hour for their work.

This same rigorous 5-annotator protocol was used for annotating both the DOCCI-Critique benchmark (comprising 10,216 sentence-level judgments and the training set for VNLI-Critique (comprising 75,363 sentence-level annotations). To quantify annotation quality for the benchmark, we computed Fleiss' Kappa for factual correctness and achieved a score of 0.48. We interpret this as moderate agreement, reflecting the highly nuanced nature of fine-grained factual assessment.

The following subsections provide an illustrative overview of the annotation interfaces designed for the two primary human evaluation tasks: assessing the factuality of VLM-generated description sentences (Section C.1) and evaluating the quality of generated critiques (Section C.2). We also provide a detailed analysis of the error categories found in our dataset (Section C.3).

C.1 DESCRIPTION SENTENCES ANNOTATION INTERFACE

For the task of annotating sentence-level factuality within VLM-generated paragraph descriptions (as detailed in Section 3 for the DOCCI-Critique benchmark), annotators were presented with an interface displaying the source image, the full paragraph context, and the specific sentence under evaluation. Figure 3 illustrates a representative example of this annotation interface. Annotators were asked to judge whether the sentence accurately described the image content, providing labels such as 'Entailment', 'Neutral', or 'Contradiction', and to supply textual rationales for any non-entailed judgments.

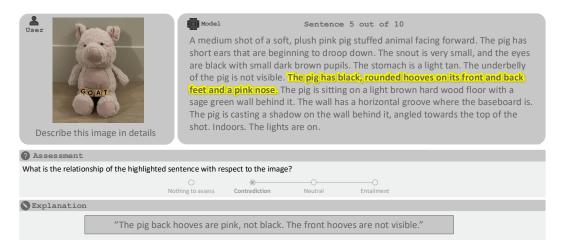


Figure 3: Example of the Description Sentences Annotation Interface. Annotators are shown the image, the full VLM-generated paragraph, and a highlighted sentence. They assess its factuality by selecting a label (here, 'Contradiction') and providing a textual explanation for any inaccuracies observed.

C.2 CRITIQUE ANNOTATION INTERFACE

To evaluate the quality of critiques generated by VNLI-Critique and other baseline models (as described in Section 4.3), a different interface was employed. This interface presented human annotators with the original image, the factually incorrect sentence that was critiqued, and the critique generated by the model under evaluation. Figure 4 shows an example of this interface. Annotators were tasked with judging whether the provided critique accurately and relevantly identified the factual error(s) present in the original sentence when compared against the visual evidence in the image.

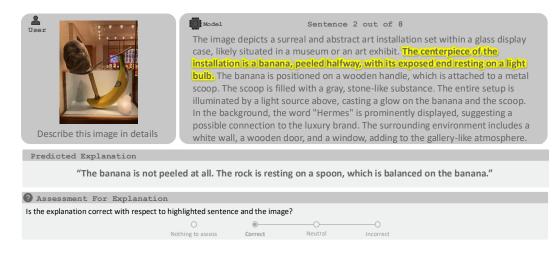


Figure 4: Example of the Critique Annotation Interface. Annotators assess if the 'Predicted Explanation' correctly identifies the error in the VLM's highlighted sentence relative to the image.

C.3 ERROR CATEGORIZATION ANALYSIS

To better understand the common failure modes of VLMs in detailed captioning, we conducted an analysis of the error types identified by the human-provided rationales in the DOCCI-Critique benchmark. We established a set of error categories from recurring patterns in the rationales and then used an LLM to systematically classify each explanation. Table 10 presents the distribution of these error categories, offering insights into the challenges of detailed captioning.

Table 10: Breakdown of error categories in DOCCI-Critique, based on human rationales

| Error Category | Description | Percentage |
|-------------------------------------|---|------------|
| Object Presence & Existence | The description includes objects that are not in the image or omits objects that are. | 22.99% |
| Spatial Relationship Error | An object's position or orientation is described incorrectly (e.g., "left" instead of "right"). | 16.34% |
| Attribute Error: Color & Appearance | An object's visual properties like color, texture, or shape are wrong. | 15.8% |
| Unverifiable Detail | A claim is made about something that is impossible to see or confirm from the image. | 13.25% |
| Incorrect Object Identification | An object is misidentified as something else (e.g., a toy is called a real animal). | 9.34% |
| Action & State Error | The action or state of a subject is wrong (e.g., "sitting" instead of "standing"). | 6.82% |
| Attribute Error: Count & Quantity | The number of objects described is incorrect. | 6.57% |
| Subjective or Unsupported Inference | The description includes a non-factual judgment, mood, or intent. | 4.89% |
| Attribute Error: Text & Numerals | Text or numbers within the image are misread. | 3.37% |
| Other | This category covers miscellaneous errors not fitting the above definitions. | 0.63% |

D ABLATION STUDIES

We include two ablation studies to further validate our methodological choices. The first study compares our two-step Critic-and-Revise pipeline against a unified end-to-end correction model. The second study investigates the pipeline's dependency on large LLMs by testing smaller, open-source models for the revision task.

D.1 END-TO-END VS. TWO-STEP PIPELINE COMPARISON

To validate our modular, two-step pipeline design, we compared it against a unified end-to-end (E2E) correction model. We used Gemini-2.0-Flash as the base model for both approaches to ensure a fair comparison. The E2E model was prompted to directly output a corrected sentence (or "YES" if the sentence was already correct) using the following prompt:

"Your task is to fix the target sentence if required to. If the target text is correct, reply only "YES". If the target text does not align with the image, fix it to be aligned. Given the image and the prompt prefix <PREFIX>Claim-Prefix</PREFIX>, does the following text align with the image: <TARGET>Target-Claim</TARGET>? Remember, answer ONLY, with YES or the fixed sentence."

The E2E setup yielded a Macro-F1 score of 0.54, only marginally better than a random classifier (0.45). This represents a substantial performance drop compared to the 0.74 Macro-F1 our two-step classification approach achieved with the same model (as shown in Table 4). This experiment confirms that simultaneously detecting and fixing sentences is a significantly more challenging task. Our two-step pipeline, by decoupling these actions, achieves higher accuracy and interpretability.

D.2 REVISION MODEL ACCESSIBILITY AND COST

To analyze the pipeline's dependency on large proprietary models, we conducted an ablation study on the revision step using smaller, open-source LLMs. We tested Gemma3-4B and Llama3.1-8B against the proprietary Gemini-2.0-Flash model. For a fair comparison, each LLM was provided the *exact same critique* from VNLI-Critique for a sample of 1,000 sentences. The factuality of each resulting revision was then assessed using our human annotation pipeline.

The results, shown in Table 11, demonstrate that smaller, open-source models like Llama3.1-8B perform comparably to the proprietary Gemini model. This confirms that our Critic-and-Revise pipeline is not dependent on large, costly models to be effective.

Table 11: Ablation study on revision LLM. Factual accuracy of revised sentences from a 1,000-sentence sample, using identical critiques from VNLI-Critique. Note: The Gemini score differs slightly from Table 6, which used a different sample size.

| Revision LLM | Factual Revised Sentences |
|------------------|----------------------------------|
| Gemma3-4B | 53.55% |
| Llama3.1-8B | 59.39% |
| Gemini-2.0-Flash | 61.93% |

E QUALITATIVE EXAMPLES

To further illustrate the core components and outputs of our work, this section provides additional qualitative examples, complementing the discussions and aggregated results presented in the main paper.

Table 12 showcases another detailed entry from the DOCCI-Critique benchmark. This example highlights the fine-grained nature of our sentence-level annotations, including the multi-rater judgments on whether a sentence makes a claim about the image, its factual correctness against the visual evidence, and the diverse human-written rationales provided by annotators for any identified inaccuracies. Such examples underscore the richness of the benchmark for evaluating nuanced understanding and error analysis.

Furthermore, Table 13 provides a step-by-step walkthrough of our Critic-and-Revise pipeline operating on an image description sourced from the PixelProse Singla et al. (2024) dataset. The example demonstrates: (1) the original VLM-generated description containing factual errors, (2) the specific unfactual sentences detected by VNLI-Critique, (3) the corresponding critiques generated by VNLI-Critique, (4) the individual sentence revisions made by the LLM based on these critiques, and (5) the final, more factually accurate revised description. This illustrates the practical application of our pipeline in automatically correcting errors in detailed image captions.

Table 12: Additional DOCCI-Critique benchmark annotation example (5 raters per assessment). Details sentence-level claims, factuality, and diverse human rationales for errors, showing varied perspectives.

Image



| Description Sentence | " Looking closely, we can see eight flamingos lined up" | " They are standing in a body of water, their reflection is seen in the water, and there are trees in the background" | " Flamingos primarily eat brine shrimp, blue-green algae, small insects, mollusks, and crustaceans" |
|--|---|---|--|
| Does the sentence include a claim about the image? (Answers from 5 raters) | ♥, ♥, ♥, ♥, ♥ | ▽ , ▽ , ▽ , ▽ , ▽ | X, X, X, X, X |
| Is the sentence factual? (Answers from 5 raters) | X, X, X, X, X | X, V , V , X, V | ♥, ♥, ♥, ♥, |
| Rationales | The count of eight flamingos is incorrect; I can see at least ten. There are 11 flamingos in the image. Incorrect number of flamingos stated, there appear to be more. I see eleven flamingos, not eight. | I don't see any prominent trees in the background, mostly just distant, blurry foliage or land. The background appears to be more of a distant shoreline or low vegetation, not distinct trees. | - |

1361

1362

1363

1364

1365

1366

13721373

1374 1375 1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393 1394 1395

Table 13: Table 11: Step-by-step illustration of the **Critic-and-Revise** pipeline in action on a sample from the PixelProse dataset. The 'Original Description' contains several inaccuracies. 'Detected Unfactual Sentences by VNLI-Critique' highlights these errors (e.g., regarding hand position, light source, text location). 'Predicted Critiques by VNLI-Critique' provides explanations for these errors. The 'Critic-and-Revise output' shows the LLM's corrected individual sentences, guided by the critiques. Finally, the 'Revised Description' integrates these corrections into a more factually grounded paragraph.

Image that is the sinful crown I shall adorn. A young man with short brown hair and dark brown eyes. He is wearing a black jacket and a white shirt. He has a serious expression on his face. He is looking at the viewer with his left hand on his chin and the **Original Description** other holding his jacket. There is a dark background with some light coming from the left side of the image. There is text at the top of the image that says "The right to use my friends as a weapon, that is the sinful crown I shall adorn - Shu Ouma". The text is in a white font. The image is in an anime style. There is text at the top of the image that says "The right to use He is looking at the viewer with There is a dark background with Detected Unfactual Sentences some light coming from the left my friends as a weapon, that is by VNLI-Critique other holding his jacket. side of the image. the sinful crown I shall adorn -Shu Ouma" The light is coming from the Predicted Critiques by He is looking at the viewer but The text is at the bottom of the right side of the image, not the VNLI-Critique his hands are not visible image and not the top. left. There is text at the bottom of the There is a dark background with image that says "The right to use He is looking at the viewer. Critic-and-Revise output some light coming from the my friends as a weapon, that is right side of the image. the sinful crown I shall adorn -Shu Ouma". A young man with short brown hair and dark brown eyes. He is wearing a black jacket and a white shirt. He has a serious expression on his face. He is looking at the viewer. There is a dark background with some light coming from the right side of the image. There is text at the bottom of the image that says **Revised Description** "The right to use my friends as a weapon, that is the sinful crown I shall adorn - Shu Ouma". The text is in a white font. The image is in an anime style.