

Learning Temporal Relations for Evaluating Instruction-Guided Image Editing

Pia Donabauer^{1,2}, Udo Kruschwitz¹, Alexander Tack²

1. Introduction

Instruction-Guided Image Editing enables automated image modification based on natural language instructions, using trained AI models.



(a) Example of an image edit
The child's right hand made a scissors gesture

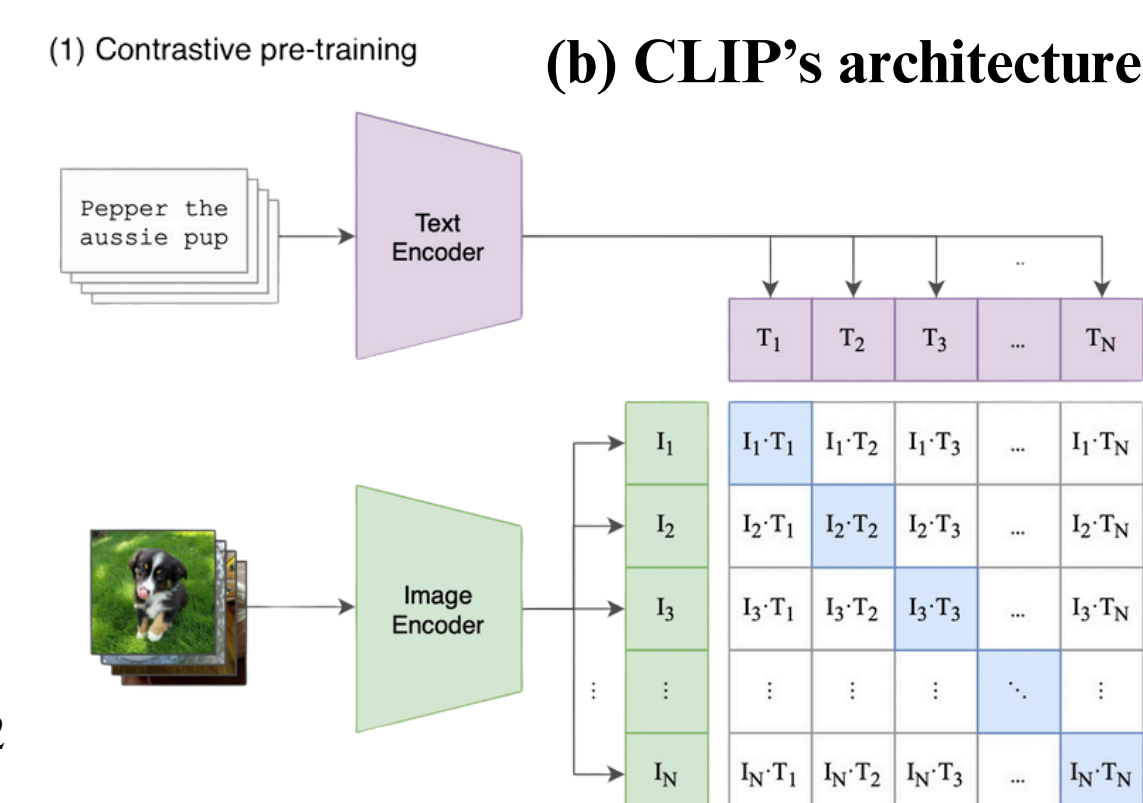
Problem: Lack of an established evaluation framework tailored to text-guided image editing.^{1,2,3}
Existing metrics are often inadequate, risking misrepresentation of results and slowing research progress.

Limitations:

- **Automated metrics** are efficient,² but measure image quality rather than instruction fidelity.³ Weak correlation with human judgement.⁴
- **Human evaluation** is more reliable, but costly, time-consuming and difficult to scale.^{5,6}
- **Visual Language Models (VLMs)**-as-a-judge show promise,⁷ but may hallucinate, be biased or misunderstand visuals.^{8,9}

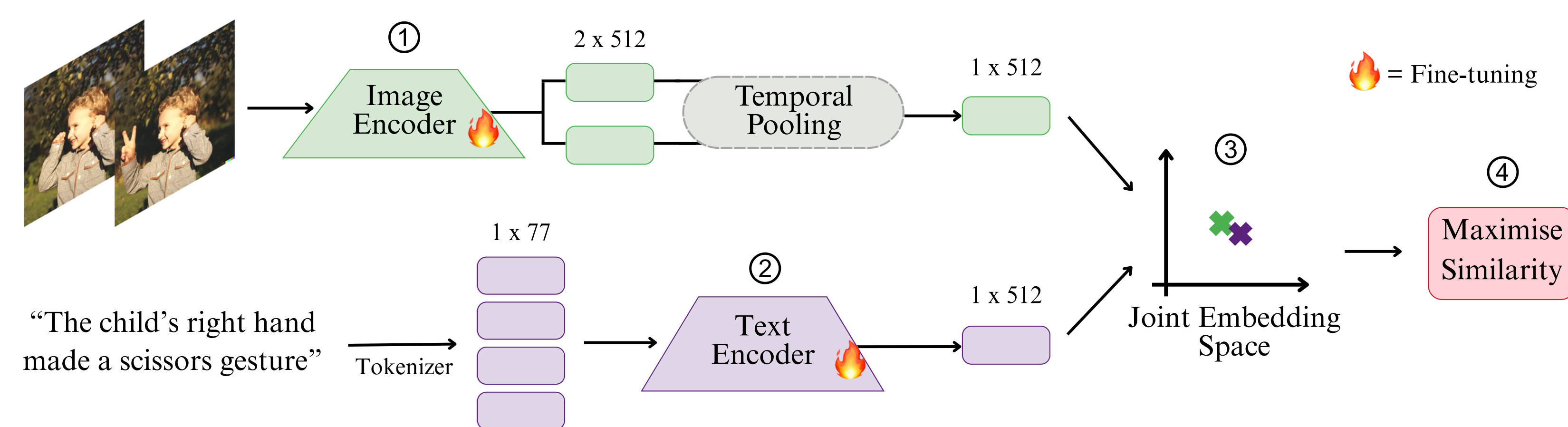
2. Approach

- Advances in **video understanding** using *Contrastive Language-Image Pre-Training*¹⁰ (**CLIP**) suggest its potential to comprehend **temporal dependencies**.^{11,12}



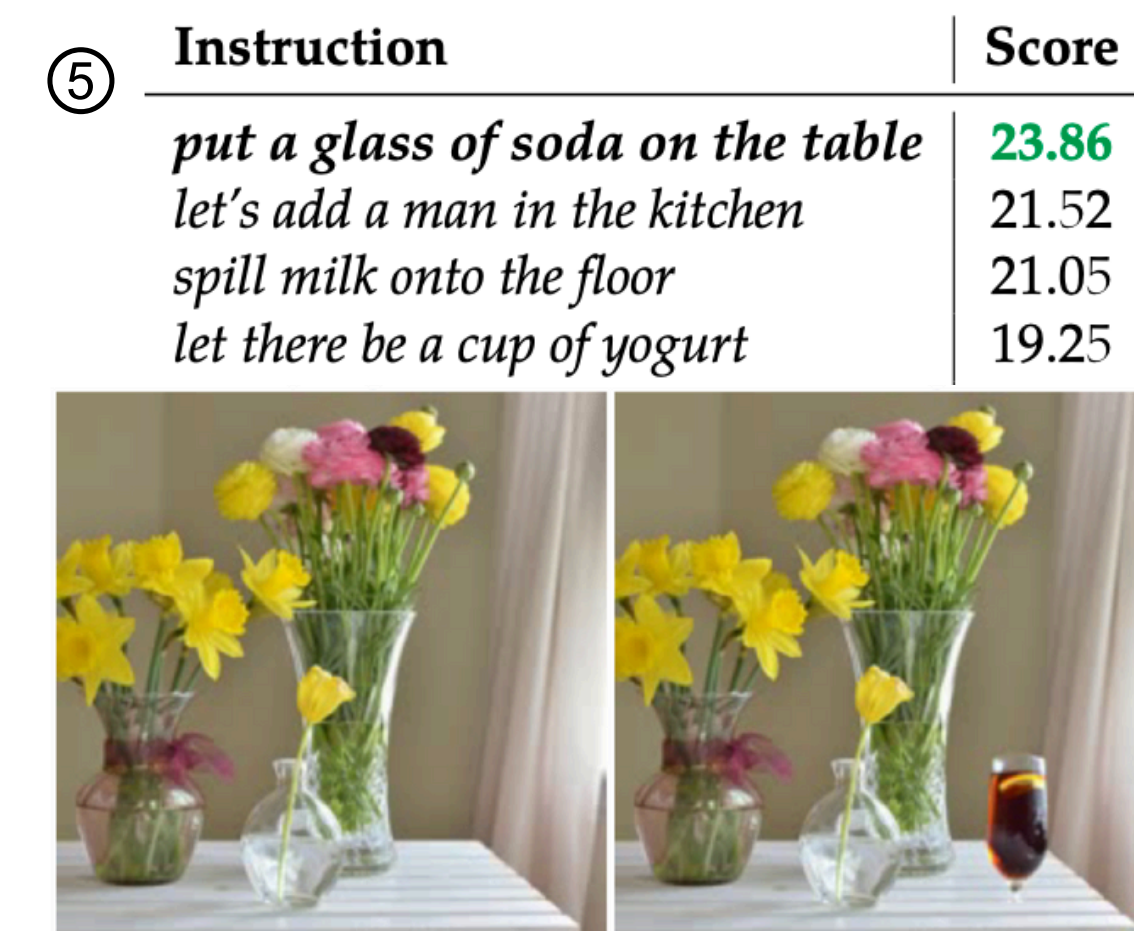
- An image edit can be interpreted as a transformation from an original to an edited state.
→ Use similar modelling techniques to evaluate image edits.
→ Learn **spatio-temporal relationships** of video sequences to better capture changes between image edit pairs and their instructions.

3. Methodology

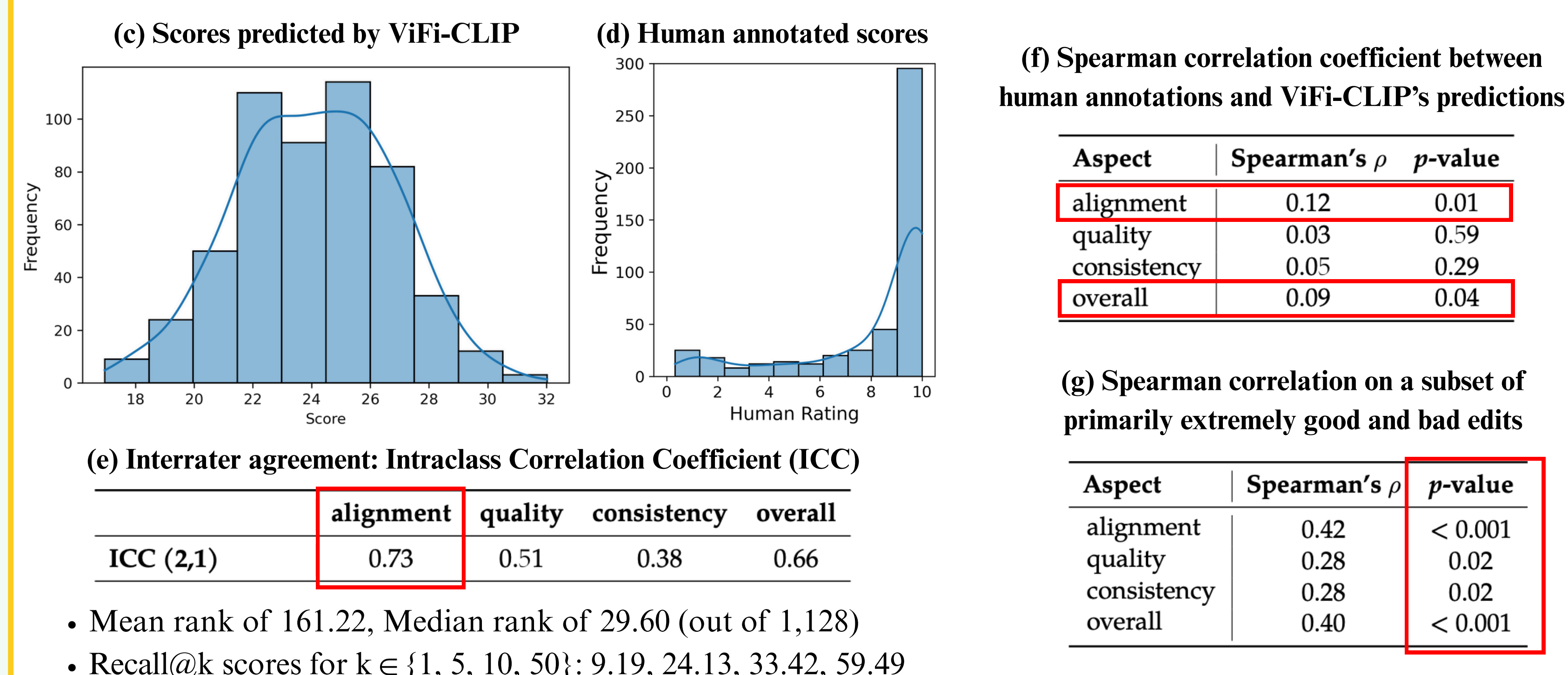


- ① Image edit pairs are encoded as **2-frame videos**: [original image → edited image].
- ② Instruction is encoded.
- ③ Both the embeddings are projected into a **joint embedding space**.
- ④ **Training:** *Video-Finetuned CLIP*¹³ (**ViFi-CLIP**) is fine-tuned on a **video-to-text retrieval task**, ranking edit-caption pairs by maximising similarity.
 - Dataset: **HumanEdit**¹⁴ (5,751 samples, human-curated)
- ⑤ **Inference:** Assign each image pair and instruction a **similarity score**.

Instruction	Score
put a glass of soda on the table	23.86
let's add a man in the kitchen	21.52
spill milk onto the floor	21.05
let there be a cup of yogurt	19.25
- ⑥ **Evaluation:** Conduct **human annotation study** to assess human alignment.
 - Dataset: **MagicBrush**¹⁵ (10,388 samples, varying quality)
- ⑦ Classify edits using **GPT-4o** for a more detailed analysis.

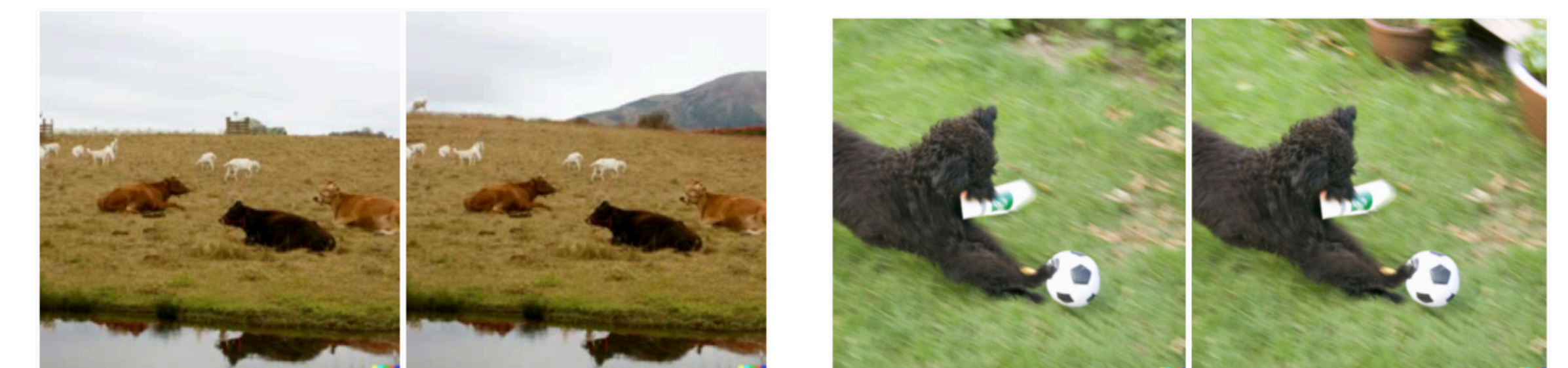


4. Quantitative Results



5. Qualitative Results

(h) High human ratings (0.90* to 1.00) and low ViFi-CLIP scores (0.07 to 0.13)



Can we have mountains on the background?

Let there be potted plant

(i) Low human ratings (0.04 and 0.14) and medium ViFi-CLIP scores (0.40 and 0.59)



Add a cotton candy machine

Add a shark next to the surfboard



(j) High disagreement among human raters (9-point difference on a 10-point Likert scale).

Let the umbrella be striped

*all values min-max normalised

6. Discussion

- **Highly subjective** nature of evaluating visual content.
- ViFi-CLIP shows **moderate retrieval performance**, struggling to capture fine-grained distinctions.
- Performance **improves** significantly for:
 - Extremely well or poorly executed image edits.
 - Specific edit types (e.g. object removals, food, **drastic** changes).
 - Larger high-quality training data available.
- **Lower** performance observed for attribute changes, **subtle** changes, edits involving people, backgrounds and objects.

Why is performance **limited**? (1) **Average pooling** may oversimplify complex instructions. (2) **CLIP's limitations** in fine-grained perception (pre-trained on 224x224 and prominent objects). (3) **Discrepancy** between training and eval data distribution.