# ImProvShow: Multimodal Fusion for Image Provenance Summarization

Alexander Black[1]
alex.black@surrey.ac.uk

Jing Shi[2]
jingshi@adobe.com

Yifei Fan[2]
yifan@adobe.com

John Collomosse[1,2]
collomos@adobe.com

[1] Centre for the Decentralized Digital Economy (DECaDE)
University of Surrey
Guildford, UK.

[2] Adobe Research
345 Park Avenue,
San Jose, CA. USA.

### Abstract

We present ImProvShow; a novel approach to summarizing the multi-stage edit history (or 'provenance') of an image. ImProvShow fuses visual and textual cues to succinctly summarize multiple manipulations applied to an image in a sequence; a novel extension of the classical image difference captioning (IDC) problem. ImProvShow takes as input several intermediate thumbnails of the image editing sequence, as well as any coarse human or machine-generated annotations of the individual manipulations at each stage, if available. We demonstrate that the presence of intermediate images and/or auxiliary textual information improves the model's edit captioning performance. To train ImProvShow, we introduce METS (Multiple Edits and Textual Summaries) – a new open dataset of image editing sequences, with textual machine annotations of each editorial step and human edit summarization captions after the 5th, 10th and 15th manipulation.

## 1 Introduction

With recent advancements in Generative AI, image manipulation becomes easier to perform and harder to notice, motivating new techniques for auditing the edit history (or '*provenance*') of an image. Often, multiple edits are applied in sequence by one or multiple editors, forming a provenance chain containing multiple versions of the image at different stages of the editing process. To mitigate the spread of disinformation, it is important to succinctly communicate the history of these changes to enable informed trust decisions [17].

Image difference captioning (IDC) usually aims to generate a difference caption given two images, the original and the edited one, regardless of the number of manipulations applied to the image. In this work, we explore image difference captioning with multiple inputs (IDC-MI), assuming access to multiple snapshots of the image editing sequence and/or auxiliary information about each individual edit. This commonly arises during a creative supply chain where multiple editors contribute to a final image. For example, emerging metadata standards for media provenance, such as the Coalition for Content Provenance and Authenticity (C2PA) [2] collect rich information on this edit process in a provenance 'manifest'.

This data structure contains multiple versions (thumbnails) of the image at different stages of the editing process, and optionally textual short descriptions of changes made (actions). Many image editing tools e.g. Adobe Photoshop now write C2PA manifests into editted images. ImProvShow considers the use case of IDC-MI to aggregate this multi-modal context (thumbnails and actions) and summarize it in a short textual description.
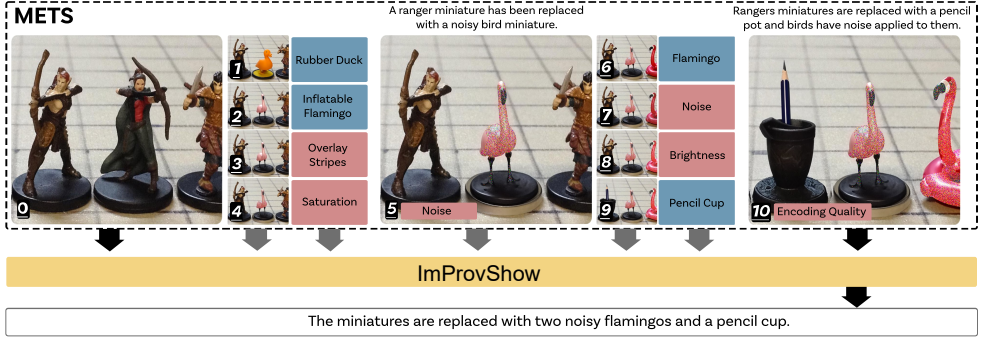


Figure 1: ImProvShow is capable of processing sequences of images, optionally accompanied by coarse edit annotations, to produce a succinct and informative summary of the differences. We train it with METS – a novel dataset of long image editing sequences paired with machine annotations and human-written summaries at multiple steps. The presence of visual and/or text information at any edit stage is optional, as denoted with grey arrows.

The first challenge in edit sequence captioning is the limited availability of training data. Most datasets for image difference captioning focus on image pairs rather than longer sequences. While the Magic Brush [53] dataset does provide multi-turn editing sequences, they are limited to three steps at most. Furthermore, all of the edits are applied to different non-overlapping objects, meaning that the final summary of all the manipulations could be constructed from a concatenation of the description of the individual steps. However, in real scenarios, the edits can be applied to the same area, potentially in a destructive or mutually exclusive manner, and the final summary should only describe the salient, still visible changes. For example, suppose the first manipulation changes the color of a bicycle, and the second one replaces the bicycle with a car. In that case, the final summary should not mention the color change as it is irrelevant to the final result. The second challenge lies in developing a methodology capable of handling interleaved multi-modal inputs. Many existing image difference captioning architectures are designed with exactly two image inputs in mind and would not be able to scale beyond that, either due to architectural constraints or memory limitations. The contributions of this paper are twofold:

1. First, we introduce METS (Multiple Edits and Textual Summaries) – a dataset of image editing sequences, with textual machine annotations of each editorial step and human edit summarization captions after the 5th, 10th, and 15th manipulation.

2. We train ImProvShow – a multi-modal LLM trained to fuse visual and textual cues to produce multi-edit summaries. We provide a comprehensive evaluation of the benefits of both additional visual and textual inputs at the IDC-MI task.

We demonstrate that the presence of intermediate images and/or auxiliary textual information improves the model's captioning performance. Note that whilst the proposed method

Figure 2: Illustration of the different types of manipulations performed. **(left)** Inpainting is done by using the word *background* as the prompt. **(middle)** property change is done by prompting GPT3.5 to output a likely change in color, material, texture or other applicable property of the object. **(right)** replacement is done by prompting GPT3.5 to output a replacement object of close match to the original shape but different semantically.

benefits from auxiliary textual data, it does not require it. Additionally, we demonstrate that fine-tuning a model trained on other synthetic data with METS helps to bridge the domain gap and improves zero-shot performance on real-world images. The illustration of ImProvShow and METS is shown in Fig. 1.

## 2 Related Work

Image difference captioning (IDC) is closely related to image captioning and visual question answering, both requiring a visual understanding system to model images and a language understanding system capable of generating syntactically correct captions. The revolution of IDC in recent years depends heavily on the advent of visual and text modeling approaches, together with cross-domain learning techniques that bridge the representation gap.

Initial methodologies for modeling visual content involve incorporating overarching CNN features such as VGG [10], and ResNet [41] into text generation models. Some methods [2, 18, 23, 51], partition images into discrete patches, extracting CNN features from each. Conversely, certain methodologies opt to utilize the outputs from an early ResNet layer, effectively capturing spatial attributes in a gridded format. In contrast, [2, 8, 23] employ Region Proposal Network (RPN) to extract features from potential object candidates. Other avenues of exploration include graph-based [55] and tree-based networks [57], aiming to capture object relations across varying levels of granularity.

Traditionally, RNN/LSTM architectures [16] have dominated text modeling. Variants like single-layer RNN [33, 53] or double-layer LSTM [2, 10, 57] are commonly utilized with diverse methods to embed image features into the recurrent process, such as additive attention [44]. During inference, captions are generated in a step-by-step manner, where the prediction of each word depends on all preceding words. Although this enhances linguistic coherence, RNN/LSTM-based approaches face challenges in modeling lengthy captions. Recent transformer-based methods employing full-attention [8, 32, 54], have alleviated this issue. Others such as BERT [9], GPT [5], and LLaMA [46] have demonstrated success across diverse visual-language tasks [15, 22, 30, 35, 59].

The objective of visual language modeling is to establish connections between image/video and text representations, catering to specific tasks like joint embedding (e.g., CLIP [40] and LIMoE [36] for cross-domain retrieval), text-to-image tasks (e.g., Stable Diffusion [42], InstructPix2Pix [4]), and image-to-text tasks (e.g., visual question answering [1, 54], or instructions [12, 15]). Image captioning, strategies for mapping images to text can be classified into two main approaches. The first involves the early fusion of image and text features to enhance alignment between image objects and textual descriptions [30, 35, 48, 54]. These methods employ BERT-like training strategies, where a pair of images and a masked caption are inputted, replacing the masked words during inference. The second approach centers on learning a direct conversion from image to text embedding. Initial CNN-based methods incorporate image features as the hidden states of LSTM text modules [10, 27, 41, 53, 57], whereas later transformer-based techniques favor cross-attention mechanisms [8, 32]. Notably, recent trends in both approaches involve harnessing powerful pretrained large language and vision models to establish a straightforward mapping between the two domains [6, 13, 29, 34, 35, 48].

Image difference captioning represents a specialized form of image captioning, aiming to disregard common objects across images and instead accentuate subtle alterations between them. Spot-the-Diff [25] introduces potential change clusters, employing an LSTM-based network to model them. However, their approach relies on pixel-level differences between input images, rendering it sensitive to noise and geometric transformations. In contrast, DUDA [39] computes image differences at the semantic level using CNNs, enhancing robustness against minor global alterations. Several approaches extend the foundation laid by DUDA. SRDRL+AVS [50] assesses the correlation between the subtracted difference and image pairs to ascertain the occurrence of the change, incorporating part-of-speech information. M-VAM [43] and VACC [28] propose a viewpoint encoder to mitigate viewpoint disparities, while VARD [51] suggests a viewpoint invariant representation network to explicitly capture changes. Additionally, [45] integrates bidirectional encoding to refine change localization, and NCT [52] utilizes a transformer to aggregate neighboring features. These methodologies concentrate on the image modality, exploiting benchmark-specific characteristics such as nearly identical views in Spot-the-Diff [25] or synthetic scenes with limited objects and change types in CLEVR [39]. More recently, IDC-PCL [56] and CLIP4IDC [19] have adopted BERT-like training approaches to model difference captioning language.

# 3 Methodology

We describe the method for generating the METS (Multiple Edits and Textual Summaries) dataset and model training for the multi-input image difference captioning (IDC-MI) task.

## 3.1 Data generation

We generate a dataset of image editing sequences, with textual machine annotations of each editorial step and human edit summarization captions after the 5th, 10th, and 15th manipulation, as shown in Fig. 3. Binary masks of the manipulation regions at each step are also included. Our dataset covers a wide variety of pixel-level and generative manipulations. The prompt for each manipulation is generated using GPT-3.5 to ensure diverse manipulations.

### 3.1.1 Individual Edits

We identify two main categories of edits: pixel-level and generative manipulations. Pixel-level edits are simple manipulations such as changing the brightness of an image or applying

| | Original | Edit 5 | Edit 10 | Edit 15 |
|---|---|---|---|---|

Machine Annotations

| Edit 5 | Edit 10 | Edit 15 |
|---|---|---|
| 1: Duck, replacement: background<br>2: Object was removed, nothing applied<br>3: Duck, random_noise, variance: 0.1<br>4: Duck, replacement: background<br>5: Object was removed, nothing applied | 1 ... 5<br>6: Goose, replacement: pink flamingo<br>7: Pink flamingo, sharpness, decreased severely<br>8: Pink flamingo, sharpness, increased moderately<br>9: Pink flamingo, saturation, increased moderately<br>10: Goose, replacement: rubber duck | 1 ... 10<br>11: duck, sharpness, decreased severely<br>12: duck, contrast, increased severely<br>13: Rubber duck, contrast, increased slightly<br>14: Duck, replacement: swan<br>15: Swan, saturation, increased moderately |

*Human Annotations* — '*Two geese are removed.*' — '*Two birds are removed, one is slightly changed, and one is replaced with a flamingo.*' — '*Two Canada gooses are missing, and one is replaced with a swan.*'
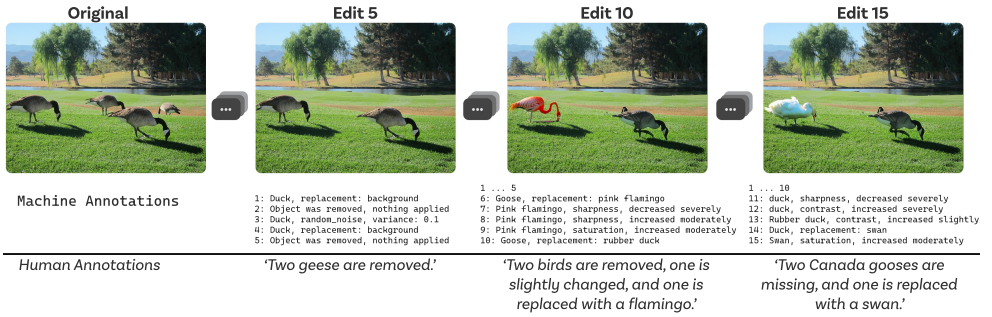
Figure 3: An example of a sequence of manipulations in METS. The original image is shown in the first column, followed by the manipulated images. The binary masks of the manipulated regions are superimposed on the images. The machine annotations generated during the sequence creation are shown in orange, while the human annotations are shown in blue.

a blur filter. Generative manipulations change the semantic content of the image.

The image, its localized narrative, object class name, and segmentation mask are sampled from the OpenImages dataset. The localized narrative and class name are used to construct a prompt for GPT3.5, which outputs a likely replacement candidate object or a property change. The prompt templates are manipulation-type specific and can be seen in suppmat. In the case of inpainting, the GPT3.5 block is omitted, and the prompt is simply *background*. The pre-processing of the segmentation mask ensures that no part of the object remains outside of the mask. The generative manipulation is then conditioned on the image, the mask, and the prompt and applied using Firefly Generative Fill.

Pixel level manipulations are performed using the Augly[58] library, with random augmentations including brightness, contrast, saturation, and encoding quality; blur, noise and sharpness filters; and overlaid color stripes. We further divide generative manipulations into three categories: **inpainting** where an object is removed from the image, **replacement** where an object is replaced with another object, and **property change** where the object's material properties are altered. We illustrate different types of manipulations in Fig. 2.

Generative manipulations are applied using the Adobe Firefly Generative Fill[1] tool, which is a language-guided inpainting GenAI model. In addition to the image itself, the model is provided with a segmentation mask and a text prompt. We generate a convex hull of the segmentation mask and apply dilation to it to ensure that no part of the object remains outside of the mask. The origin of the text prompt depends on the type of manipulation. For **inpainting** we use the word *background*, which was shown to perform on par with inpainting-specific models. For **replacement**, we use GPT3.5 in a few-shot learning manner, prompting with a localized narrative for the whole image, a bounding box of the mask, and the class label of the mask to come up with a probable replacement candidate object that would be a close match to the shape of the original object. We use a similar strategy for **property change**, but prompting GPT3.5 to output a likely property change.

### 3.1.2 Sequence Generation

We sample OpenImages, using images with at least 5 non-overlapping segmentation masks. We then follow a procedure illustrated in Fig. 4 to apply a sequence of edits to the image.
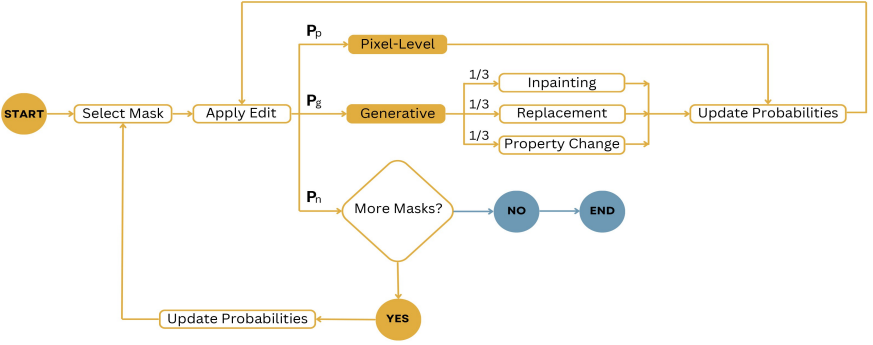
---

[1]https://firefly.adobe.com/upload/inpaint

Figure 4: The diagram of the sequence generation process. For each image, we first go through up to 15 segmentation masks and apply edits, chosen randomly, where the probabilities of choices depend on the number of edits already applied to the mask.

At each iteration step, we pick a segmentation mask and either apply a generative or a pixel-level manipulation to that area of the image or move on to the next mask. The probability of switching to the next mask is proportional to the number of manipulations already applied.

Formally, we define the probabilities of applying a generative manipulation $P_g$, a pixel-level manipulation $P_p$ and moving on to the next mask $P_n$ as follows: $P_g = g - \frac{n}{2}$; $P_p = (1 - g) - \frac{n}{2}$; $P_n = 1 - P_g - P_p$, where $g = 0.9$ if no generative manipulations have been applied to the mask yet and $g = 0.1$ otherwise. The value of $n$ is proportional to the number of manipulations already applied to the mask, defined as follows:

$$n = \max(0, \frac{40 \times (i - i_{min})}{100}), \tag{1}$$

where $i$ is the current step and $i_{min}$ is the minimum number of steps required to move on to the next mask. We set $i_{min} = 5$.

After each manipulation step, we record the type of manipulation, the parameters of the manipulation, and the binary mask used to apply the manipulation. This information is saved in a text format. For pixel-level manipulations, the text format is as follows: `Object: obj_name, manipulation: edit_name, intensity: intensity;` where `obj_name` is the name of the object annotated within OpenImages, `edit_name` is the manipulation type and `intensity` is chosen at random from a set of predefined parameters, individual for each manipulation type.

For generative manipulations, the format is: `Object: obj_name, replacement: prompt ;` where `prompt` is either *background* for inpainting or the output of GPT3.5 for replacement and property change manipulations. Examples of the template-generated text can be seen from Fig. 3, marked as *machine annotation*.

As a result, for each input image, we obtain a sequence of manipulated versions applied on top of each other and a list of annotations describing each manipulation step type, parameters, and location. We generate 1000 such sequences averaging 21.4 steps/sequence.

### 3.1.3   Labelling

We collect human annotations for difference summarization at the 5th, 10th, and 15th step of the manipulation sequence. In each task, the users are presented with the input image $I$ and
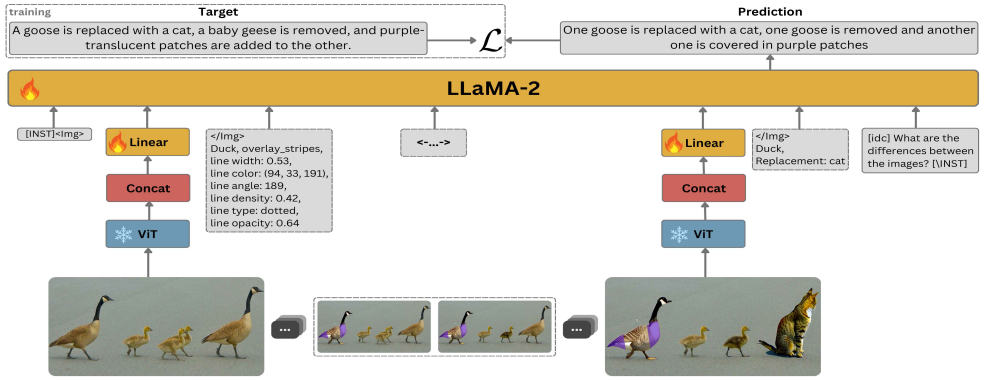
Figure 5: Architecture of the ImProvShow model. The LLaMA-2 language model is conditioned using the multi-modal instruction template, which includes at least two image features and optional auxiliary textual information. All optional content is placed within dashed boxes. The image features extracted from the ViT image encoder are concatenated in groups of 4 and projected to the LLM embedding space.

an output image $I'_n, n \in [5, 10, 15]$ and are asked to provide a short one-sentence summary of all of the differences they see between the two images (Fig. 3).

## 3.2 Architecture

Our architecture is illustrated in Fig. 5. Our setup consists of a Vision Transformer (ViT) [□] image encoder and the open-sourced LLaMA2-chat (7B) large language model [□□]. The visual tokens are concatenated in groups of 4 and projected to the language model's embedding space with a linear projection layer. During training, the visual encoder weights are frozen, and only the language model and the projection layer are trained.

Uniquely, we use multiple images as input to the model and train it for the task of image difference captioning. We note that this approach is capable of handling an arbitrary number of input images, which allows us to input several snapshots of editing sequence at once.

Optionally, we provide the model with auxiliary textual information in the form of machine annotations, described in Section 3.1. The annotations for each manipulation are interleaved with the image features and are used to guide the model's attention. We follow the multi-modal instructional template from [□] and adjust it to our task:

*[INST] <Img><ImageFeature></Img> T ... <Img><ImageFeature></Img> T [idc] ins [/INST]*

where the image feature tags are repeated for each input image in the sequence, T is the optional auxiliary textual information, [idc] is the task identifier for image difference captioning and ins is the instruction that is chosen at random from a set of predefined instructions, all synonymous with *describe the differences between the images.*.

Training minimizes the captioning loss $\mathcal{L} = -\sum_{i=1}^{m} l(s^v, s_1^t, \ldots, s_i^t)$, where $m$ is a variable token length and $l$ is next-token log-probability conditioned on the previous sequence elements: $l(s^v, s_1^t, \ldots, s_i^t) = \log p(t_i|x, t_1, \ldots, t_{i-1})$.

### 3.2.1 Training

All models are trained on a single A100 GPU with 80GB of memory for 300 epochs with 1000 steps per epoch and batch size 6. We use AdamW with a cosine learning rate scheduler

with an initial learning rate of $10^{-5}$ and a warmup learning rate of $10^{-6}$ for 1000 steps. The input image size is $448 \times 448$, and maximum token length is 1024.

# 4    Experiments

## 4.1    Datasets

In addition to our own **METS** dataset, we train and evaluate our model on a number of other datasets used in the image difference captioning literature (sup. mat shows examples).

**CLEVR-Change** [26] consists of 67,660, 3,976, 7,970 training, validation, and test image pairs, respectively. The images are generated using the CLEVR engine and contain renders of primitive 3D shapes. The types of edits include changes in shape, color, material, size, and position of the objects. This dataset serves as a good benchmark due to its large volume and precise annotations. However, the synthetic nature of the images creates a large domain gap, making it difficult to generalize to real-world images.

**Spot-the-Diff** [25] is a dataset of 13,192 well-aligned image pairs from CCTV cameras. There are no viewpoint changes, and the edits are limited to object addition, deletion, or movement. We follow the official dataset split of 80%, 10% and 10%.

**PSBattles** [20] is a dataset of real-world image pairs collected from the Photoshop Battles subreddit. The difference captions for a subset of the dataset were collected by [3] in a user study. We use this dataset to evaluate generalization to real-world images.

**InstructPix2Pix** [4] is a dataset of $\sim$1M image pairs generated with prompt-to-prompt [21] approach. The difference captions are later generated by [3] using chatGPT-3. We use this dataset for pre-training of the model during the evaluation in the PSBattles dataset to assess the benefits of fine-tuning on the METS dataset for domain adaptation.

**MagicBrush** [58] contains sequences of edited images generated in a manner similar to ours, but with human supervision. Due to the need for human supervision, the maximum length of the sequences is limited to 3 steps. Of 878 training sequences, only 304 have a length of 4 (including the original image), and 547 have a length of 3. We use this dataset to evaluate the model's performance in the IDC-MI setting, using only the samples that have a length of 4. The target annotation is a concatenation of the instructions for each step. As input, we use either the first and the last image in the sequence or all four images in the sequence.

## 4.2    Evaluation

We evaluate the performance of our model on the standard IDC setting on the CLEVR-Change, InstructPix2Pix, and PSBattles datasets. We evaluate the performance of our model in the IDC-MI setting on the MagicBrush and our proposed METS datasets. In both cases, we use the standard n-gram based metrics BLEU-4 (B4), CIDEr (C), METEOR (M), ROUGE-L (R) and SPICE (S) to evaluate the performance of our model. Additionally, we use LLM-as-judge metric to assess the semantic similarity of the captions that n-gram based metrics struggle to capture. We use GPT4 to score the semantic similarity of each text pair as 'low', 'medium' or 'high' and report the percentage of medium and high scores.

### 4.2.1    Evaluating IDC with Multiple Inputs

For IDC-MI, we evaluate the model's performance while varying the number of input images and the presence of auxiliary textual information. The intermediate images are sampled to be equally spaced in the sequence, and the textual information is provided in the form of machine annotations described in Section 3.1. We baseline the performance of our model with

Table 1: Performance evaluation in the IDC-MI setting shows BLEU-4 (**B4**), CIDEr (**C**), METEOR (**M**), ROUGE-L (**R**) and LLM as judge medium (**L (M)**) and high (**L (H)**) scores. We report the performance of our model and compare it with GPT3.5 and GPT4-V, varying the number of input images and the presence of auxiliary textual information.

| Model | Images | Text | B4 | C | M | R | L (M) | L (H) |
|---|---|---|---|---|---|---|---|---|
| METS | | | | | | | | |
| GPT3.5 [5] | 0 | yes | 1.6 | 8.6 | 10.4 | 15.1 | 16.2 | 0.6 |
| GPT4-V [14] | 2 | no | 4.0 | 18.6 | **14.0** | 20.3 | 22.2 | 2.6 |
| GPT4-V[14] | 2 | yes | 1.3 | 0.3 | 11.5 | 13.5 | 19.7 | 0.9 |
| GPT4-V[14] | 4 | no | 3.0 | 15.1 | <u>13.4</u> | 19.9 | <u>26.9</u> | 1.9 |
| GPT4-V[14] | 4 | yes | 1.4 | 0.4 | 11.6 | 12.9 | 24.1 | 1.2 |
| ImProvShow-2 (ours) | 2 | no | 5.8 | 20.7 | 11.4 | 23.1 | 22.6 | 9.4 |
| ImProvShow-2T (ours) | 2 | yes | <u>7.8</u> | <u>25.8</u> | 13.0 | <u>26.0</u> | 24.3 | <u>11.0</u> |
| ImProvShow-4 (ours) | 4 | no | 6.6 | 23.5 | 12.3 | 24.3 | 22.6 | 9.6 |
| ImProvShow-4T (ours) | 4 | yes | **8.2** | **25.9** | <u>13.4</u> | **26.3** | **30.1** | **12.4** |
| MagicBrush | | | | | | | | |
| ImProvShow-2 (ours) | 2 | no | 4.9 | 29.4 | 13.3 | 28.1 | - | - |
| ImProvShow-4 (ours) | 4 | no | **6.8** | **44.5** | **15.6** | **31.2** | - | - |

GPT4-V, which has multi-modal capabilities and is capable of taking multiple images and/or text as input. We compare with GPT3.5, which serves as a text-only baseline (Table 1). Our method is able to take advantage of the additional inputs, achieving the best performance when both intermediate images and auxiliary textual information are present.

Compared to the base case of just two-image input, the addition of text to our model improves the performance by an average of 18.9% across all metrics, and intermediate images improve the performance by an average of 10.1% across all metrics. The combination of intermediate images and textual information shows improvement of 22.4% across all metrics. On the other hand, the performance of GPT4-V suffers from the addition of intermediate images, decreasing with the addition of extra images and text.

### 4.2.2 Evaluating IDC with Two Inputs

In the IDC setting, shown in Table 2, the model achieves competitive performance on the CLEVR-Change dataset, outperforming the previous state-of-the-art VARD on the CIDEr and ROUGE-L metrics. On the InstructPix2Pix dataset, the model outperforms VIXEN only on the METEOR metric. However, it shows a better capability to generalize to real-world images, outperforming VIXEN on the PSBattles dataset. Fine-tuning the model on the METS dataset further improves its performance on PSBattles, showing the dataset's ability to bridge the domain gap between synthetic and real-world images.

## 5   Conclusion

We introduced ImProvShow, a novel multimodal approach for summarizing multi-stage image edit histories (provenance), extending the conventional image difference captioning

Table 2: Image difference captioning performance evaluation on CLEVR-Change and PS-Battles. We compare our model with the state-of-the-art models and report BLEU-4 (**B4**), CIDEr (**C**), METEOR (**M**) and ROUGE-L (**R**) scores.

| MODEL | TRAINING DATA | B4 | C | M | R | S |
|---|---|---|---|---|---|---|
| CLEVR CHANGE | | | | | | |
| DUDA [59] | CLEVR | 47.3 | 112.3 | 33.9 | - | - |
| IFDC [24] | CLEVR | 49.2 | 118.7 | 32.5 | 69.1 | - |
| $R^3$NET+SSP [49] | CLEVR | 54.7 | 123.0 | 39.8 | 73.1 | - |
| SGCC [57] | CLEVR | 51.1 | 121.8 | 40.6 | 73.9 | - |
| NCT [52] | CLEVR | 55.1 | 124.1 | 40.2 | 73.8 | - |
| SRDL+AVS [50] | CLEVR | 54.9 | 122.2 | 40.2 | 73.3 | - |
| VARD [51] | CLEVR | **55.2** | 124.1 | **40.8** | 74.1 | - |
| IMPROVSHOW-2 (OURS) | CLEVR | 54.7 | **151.8** | 40.0 | **77.1** | - |
| SPOT-THE-DIFF | | | | | | |
| SRDL+AVS [50] | SPOT-DIFF | - | 35.3 | 13.0 | 31.0 | 18.0 |
| $R^3$NET+SSP [49] | SPOT-DIFF | - | 36.6 | 13.1 | **32.6** | 18.8 |
| VARD-LSTM [51] | SPOT-DIFF | - | 39.3 | 13.1 | 33.1 | 17.5 |
| VARD-TRANSFORMER [51] | SPOT-DIFF | - | 30.3 | 12.5 | 29.3 | 17.3 |
| IMPROVSHOW-2 (OURS) | SPOT-DIFF | - | **45.5** | **13.7** | 28.7 | **19.3** |
| PSBATTLES | | | | | | |
| VIXEN-C [3] | IP2P | 4.5 | 7.7 | 9.5 | 20.5 | - |
| IMPROVSHOW-2 (OURS) | IP2P | 5.3 | 10.3 | 10.8 | 22. | - |
| IMPROVSHOW-2 (OURS) | IP2P + METS | **5.5** | **14.2** | **11.2** | **22.6** | - |

(IDC) task to image difference captioning with multiple inputs (IDC-MI). Our method effectively fuses visual and textual information to generate concise summaries of complex editing sequences. ImProvShow is motivated by the real-world task of summarizing provenance metadata, to communicate the change history (provenance) of an image. Such metadata (*e.g.* C2PA [7]) includes the original and final editted image with intermediate image and text data, and is now output by many tools (*e.g.* Adobe Photoshop, Adobe Lightroom).

To support this research, we introduced METS—a dataset of multi-step image editing sequences containing both machine-generated annotations and human-authored summarization captions. We showed that incorporating intermediate images and auxiliary text significantly improves performance: ImProvShow achieves a 30.1% LLM similarity score with four images and text, outperforming GPT-4V's 26.9%, and reaches a CIDEr score of 25.9, surpassing GPT-4V's 15.1. In the two-input IDC setting, ImProvShow achieves state-of-the-art results, scoring 151.8 CIDEr on CLEVR-Change and 45.5 CIDEr on Spot-the-Diff. Fine-tuning on METS further improves generalization, boosting CIDEr on PSBattles from 10.3 to 14.2. Future work could explore evaluation on more recent vision-language models and self-verification mechanisms to improve reliability.

# Acknowledgements

# References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *NeurIPS*, 35:23716–23736, 2022.

[2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proc. CVPR*, pages 6077–6086, 2018.

[3] Alexander Black, Jing Shi, Yifei Fai, Tu Bui, and John Collomosse. Vixen: Visual text comparison network for image difference captioning, 2024.

[4] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. *arXiv preprint arXiv:2211.09800*, 2022.

[5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 33:1877–1901, 2020.

[6] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning, 2023.

[7] Coalition for Content Provenance and Authenticity. Technical specification 1.3. Technical report, C2PA, 2023. URL https://c2pa.org/specifications/specifications/1.3/specs/_attachments/C2PA_Specification.pdf.

[8] Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. Meshed-memory transformer for image captioning. In *Proc. CVPR*, pages 10578–10587, 2020.

[9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[10] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proc. CVPR*, pages 2625–2634, 2015.

[11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021.

[12] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model. In *arXiv preprint arXiv:2303.03378*, 2023.

[13] Constantin Eichenberg, Sidney Black, Samuel Weinbach, Letitia Parcalabescu, and Anette Frank. Magma–multimodal augmentation of generative models through adapter-based finetuning. *arXiv preprint arXiv:2112.05253*, 2021.

[14] Josh Achiam *et al.* Gpt-4: OpenAI technical report, 2024.

[15] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023.

[16] Alex Graves and Alex Graves. Long short-term memory. *Supervised sequence labelling with recurrent neural networks*, pages 37–45, 2012.

[17] S. Gregory. Ticks or it didn't happen. https://lab.witness.org/ticks-or-it-didnt-happen/, 2019. Accessed: 2024-01-20.

[18] Jiuxiang Gu, Jianfei Cai, Gang Wang, and Tsuhan Chen. Stack-captioning: Coarse-to-fine learning for image captioning. In *Proc. AAAI*, volume 32, 2018.

[19] Zixin Guo, Tzu-Jui Wang, and Jorma Laaksonen. Clip4idc: Clip for image difference captioning. In *Proc. Conf. Asia-Pacific Chapter Assoc. Comp. Linguistics and Int. Joint Conf. NLP*, pages 33–42, 2022.

[20] S. Heller, L. Rossetto, and H. Schuldt. The PS-Battles Dataset – an Image Collection for Image Manipulation Detection. *CoRR*, abs/1804.04866, 2018. URL http://arxiv.org/abs/1804.04866.

[21] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.

[22] Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. Scaling up vision-language pre-training for image captioning. In *Proc. CVPR*, pages 17980–17989, 2022.

[23] Lun Huang, Wenmin Wang, Yaxian Xia, and Jie Chen. Adaptively aligned image captioning via adaptive attention time. *NeurIPS*, 32, 2019.

[24] Qingbao Huang, Yu Liang, Jielong Wei, Yi Cai, Hanyu Liang, Ho-fung Leung, and Qing Li. Image difference captioning with instance-level fine-grained feature representation. *IEEE Transactions on Multimedia*, 24:2004–2017, 2022. doi: 10.1109/TMM. 2021.3074803.

[25] Harsh Jhamtani and Taylor Berg-Kirkpatrick. Learning to describe differences between pairs of similar images. In *Proc. Conf. Empirical Methods NLP*, pages 4024–4034, 2018.

[26] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910, 2017.

[27] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proc. CVPR*, pages 3128–3137, 2015.

[28] Hoeseong Kim, Jongseok Kim, Hyungseok Lee, Hyunsung Park, and Gunhee Kim. Agnostic change captioning with cycle consistency. In *Proc. ICCV*, pages 2095–2104, 2021.

[29] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.

[30] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Proc. ECCV*, pages 121–137. Springer, 2020.

[31] Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *Proc. CVPR*, pages 375–383, 2017.

[32] Yunpeng Luo, Jiayi Ji, Xiaoshuai Sun, Liujuan Cao, Yongjian Wu, Feiyue Huang, Chia-Wen Lin, and Rongrong Ji. Dual-level collaborative transformer for image captioning. In *Proc. AAAI*, volume 35, pages 2286–2293, 2021.

[33] Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, Zhiheng Huang, and Alan Yuille. Deep captioning with multimodal recurrent neural networks (m-rnn). In *Proc. ICLR*, 2015.

[34] Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. *arXiv preprint arXiv:2209.15162*, 2022.

[35] Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021.

[36] Basil Mustafa, Carlos Riquelme Ruiz, Joan Puigcerver, Rodolphe Jenatton, and Neil Houlsby. Multimodal contrastive learning with limoe: the language-image mixture of experts. In *NeurIPS*, 2022.

[37] Ariyo Oluwasanmi, Enoch Frimpong, Muhammad Umar Aftab, Edward Yellakuor Baagyere, Zhiguang Qin, and Kifayat Ullah. Fully convolutional captionnet: Siamese difference captioning attention model. *IEEE Access*, 7:175929–175939, 2019. URL https://api.semanticscholar.org/CorpusID:209382557.

[38] Zoë Papakipos and Joanna Bitton. Augly: Data augmentations for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 156–163, 2022.

[39] Dong Huk Park, Trevor Darrell, and Anna Rohrbach. Robust change captioning. In *Proc. ICCV*, pages 4624–4633, 2019.

[40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. ICML*, pages 8748–8763. PMLR, 2021. URL https://github.com/OpenAI/CLIP.

[41] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical sequence training for image captioning. In *Proc. CVPR*, pages 7008–7024, 2017.

[42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. CVPR*, pages 10684–10695, 2022.

[43] Xiangxi Shi, Xu Yang, Jiuxiang Gu, Shafiq Joty, and Jianfei Cai. Finding it at another side: A viewpoint-adapted matching encoder for change captioning. In *Proc. ECCV*, pages 574–590. Springer, 2020.

[44] Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. From show to tell: A survey on deep learning-based image captioning. *IEEE TPAMI*, 45(1):539–559, 2022.

[45] Yaoqi Sun, Liang Li, Tingting Yao, Tongyv Lu, Bolun Zheng, Chenggang Yan, Hua Zhang, Yongjun Bao, Guiguang Ding, and Gregory Slabaugh. Bidirectional difference locating and semantic consistency reasoning for change captioning. *IJIS*, 37(5):2969–2987, 2022.

[46] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[47] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[48] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *NeurIPS*, 34: 200–212, 2021.

[49] Yunbin Tu, Liang Li, Chenggang Yan, Shengxiang Gao, and Zhengtao Yu. R^3Net:relation-embedded representation reconstruction network for change captioning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9319–9329, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021. emnlp-main.735. URL https://aclanthology.org/2021.emnlp-main.735.

[50] Yunbin Tu, Tingting Yao, Liang Li, Jiedong Lou, Shengxiang Gao, Zhengtao Yu, and Chenggang Yan. Semantic relation-aware difference representation learning for change captioning. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 63–73, 2021.

[51] Yunbin Tu, Liang Li, Li Su, Junping Du, Ke Lu, and Qingming Huang. Viewpoint-adaptive representation disentanglement network for change captioning. *IEEE Transactions on Image Processing*, 32:2620–2635, 2023. doi: 10.1109/TIP.2023.3268004.

[52] Yunbin Tu, Liang Li, Li Su, Ke Lu, and Qingming Huang. Neighborhood contrastive transformer for change captioning, 2023.

[53] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *Proc. CVPR*, pages 3156–3164, 2015.

[54] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision. In *Proc. ICLR*, 2021.

[55] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *Proc. CVPR*, pages 10685–10694, 2019.

[56] Linli Yao, Weiying Wang, and Qin Jin. Image difference captioning with pre-training and contrastive learning. In *Proc. AAAI*, volume 36, pages 3108–3116, 2022.

[57] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Hierarchy parsing for image captioning. In *Proc. ICCV*, pages 2621–2629, 2019.

[58] Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36, 2024.

[59] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proc. CVPR*, pages 5579–5588, 2021.