FUSING VISUAL AND TEXTUAL CUES FOR SEQUENTIAL IMAGE DIFFERENCE CAPTIONING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

Paper under double-blind review

Abstract

We present Fusing Visual and Textual Cues (FVTC); a novel technique for image difference captioning that is able to benefit from additional visual and/or textual inputs. FVTC is able to succinctly summarize multiple manipulations that were applied to an image in a sequence. Optionally, it can take several intermediate thumbnails of the image editing sequence as input, as well as coarse machine-generated annotations of the individual manipulations. We demonstrate that the presence of intermediate images and/or auxiliary textual information improves the model's captioning performance. To train FVTC, 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

1 INTRODUCTION

With recent advancements in Generative AI, image manipulation becomes increasingly easier to perform and harder to notice, motivating new techniques for auditing the provenance and edits made to an image. Often, multiple edits are applied in sequence by one or multiple editors, forming a provenance graph containing multiple versions of the same image at different stages of the editing process. To avoid the spread of misinformation, it is important to be able to communicate the history of these changes to the end user succinctly to enable them to make informed trust decisions Gregory (2019).

032 Image difference captioning (IDC) usually aims to generate a difference caption given two images, 033 the original and the edited one, regardless of the number of manipulations applied to the image. In 034 this work, we explore image difference captioning with multiple inputs (IDC-MI), assuming access 035 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 037 for Content Provenance and Authenticity (C2PA) Coalition for Content Provenance and Authenticity (2023) collect rich information on this edit process as a provenance graph. This data structure contains multiple versions (thumbnails) of the image at different stages of the editing process, and 040 optionally textual short descriptions of changes made. One use case for IDC-MI is aggregate this 041 multi-modal context and summarize it in a short textual description. 042

The first challenge in edit sequence captioning is the limited availability of training data. Most 043 datasets for image difference captioning focus on image pairs rather than longer sequences. While 044 the Magic Brush Zhang et al. (2024) 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 046 objects, meaning that the final summary of all the manipulations could be constructed from a con-047 catenation of the description of the individual steps. However, in real scenarios, the edits can be 048 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 051 case, the final summary should not mention the color change as it is irrelevant to the final result. The 052 second challenge lies in developing a methodology capable of handling interleaved multi-modal in-

¹The METS dataset will be released for open access upon acceptance.



Figure 1: FVTC 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. Optional image and text inputs are denoted with gray arrows.

puts. 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 074 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 FVTC (Fusing Visual and Textual Cues) a multi-modal LLM with multiple visual inputs and provide a comprehensive evaluation of the benefits of both additional visual and textual inputs.

We demonstrate that the presence of intermediate images and/or auxiliary textual information improves the model's captioning performance. 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-life images. The illustration of FVTC and METS is shown in Fig. 1.

2 **RELATED WORK**

090 Image difference captioning (IDC) is closely related to image captioning and visual question answer-091 ing, both requiring a visual understanding system to model images and a language understanding 092 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 between them. 094

Initial methodologies for modeling visual content involve incorporating overarching CNN features 096 such as VGG Donahue et al. (2015), and ResNet Rennie et al. (2017) into text generation models. This integration capitalizes on the dense and meaningful representations these models provide. 098 To enhance the representation of multiple objects and their interrelations, various techniques have 099 emerged. Some methods Lu et al. (2017); Gu et al. (2018); Anderson et al. (2018); Huang et al. (2019), partition images into discrete patches, extracting CNN features from each. Conversely, cer-100 tain methodologies opt to utilize the outputs from an early ResNet layer, effectively capturing spatial 101 attributes in a gridded format. In contrast, Cornia et al. (2020); Anderson et al. (2018); Huang et al. 102 (2019) employ Region Proposal Network (RPN) to extract features from potential object candi-103 dates, thereby improving alignment with the semantic entities referenced in paired captions. Other 104 avenues of exploration include graph-based Yang et al. (2019) and tree-based networks Yao et al. 105 (2019), aiming to capture object relations across varying levels of granularity. 106

Traditionally, RNN/LSTM architectures Graves & Graves (2012) have dominated text modeling 107 owing to their intrinsic sequential nature. Variants like single-layer RNN Vinyals et al. (2015);

067

068

069

070

071 072

073

075 076

077

078

079

080

081

082

084

085

087

089

054

055

108 Mao et al. (2015) or double-layer LSTM Donahue et al. (2015); Anderson et al. (2018); Yao et al. 109 (2019) are commonly utilized, often coupled with diverse methods to embed image features more 110 deeply into the recurrent process, such as additive attention Stefanini et al. (2022). During inference, 111 captions are generated in a step-by-step manner, where the prediction of each word depends on 112 all preceding words. Although this enhances linguistic coherence, RNN/LSTM-based approaches face challenges in modeling lengthy captions. Recent transformer-based methodologies, like those 113 employing a full-attention mechanism Luo et al. (2021); Wang et al. (2021); Cornia et al. (2020), 114 have alleviated this issue. Advanced transformer-based models such as BERT Devlin et al. (2018), 115 GPT Brown et al. (2020), and LLaMA Touvron et al. (2023a) have demonstrated success across 116 diverse visual-language tasks Hu et al. (2022); Mokady et al. (2021); Gao et al. (2023); Zhang et al. 117 (2021); Li et al. (2020). 118



133 134

119

Figure 2: Illustration of the different types of manipulations performed using Firefly Generative Fill. (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 likely replacement candidate object that would be a close match to the shape of the original, but different semantically.

141

142 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 Radford et al. (2021) and 143 LIMoE Mustafa et al. (2022) for cross-domain retrieval), text-to-image tasks (e.g., Stable Diffusion 144 Rombach et al. (2022) for text-based image generation, InstructPix2Pix Brooks et al. (2022) for 145 image editing), and image-to-text tasks (e.g., visual question answering Alayrac et al. (2022); Wang 146 et al. (2021), visual instructions Gao et al. (2023); Driess et al. (2023)). In the realm of image 147 captioning, strategies for mapping images to text can be classified into two main approaches. The 148 first approach involves the early fusion of image and text features to enhance alignment between 149 image objects and textual descriptions Tsimpoukelli et al. (2021); Mokady et al. (2021); Wang et al. 150 (2021); Li et al. (2020). These methods employ BERT-like training strategies, where a pair of 151 images and a masked caption are inputted, replacing the masked words during inference with either 152 a start token or a prefixed phrase like 'A picture of'. The second approach centers on learning a direct conversion from image to text embedding. Initial CNN-based methods incorporate image 153 features as the hidden states of LSTM text modules Donahue et al. (2015); Vinyals et al. (2015); 154 Yao et al. (2019); Karpathy & Fei-Fei (2015); Rennie et al. (2017), whereas later transformer-based 155 techniques favor cross-attention mechanisms Luo et al. (2021); Cornia et al. (2020). Notably, recent 156 trends in both approaches involve harnessing powerful pretrained large language and vision models 157 to establish a straightforward mapping between the two domains Merullo et al. (2022); Eichenberg 158 et al. (2021); Li et al. (2023); Tsimpoukelli et al. (2021); Mokady et al. (2021); Chen et al. (2023). 159

Image difference captioning represents a specialized form of image captioning, aiming to disregard
 common objects across images and instead accentuate subtle alterations between them. Pioneering
 this domain, Spot-the-Diff Jhamtani & Berg-Kirkpatrick (2018) introduces potential change clusters,

162 employing an LSTM-based network to model them. However, their approach relies on pixel-level 163 differences between input images, rendering it sensitive to noise and geometric transformations. In 164 contrast, DUDA Park et al. (2019) computes image differences at the semantic level using CNNs, en-165 hancing robustness against minor global alterations. Several approaches extend the foundation laid 166 by DUDA. For example, SRDRL+AVS Tu et al. (2021b) initially assesses the correlation between the subtracted difference and image pairs to ascertain the occurrence of the change. Subsequently, 167 it incorporates part-of-speech information to dynamically leverage visual data. M-VAM Shi et al. 168 (2020) and VACC Kim et al. (2021) propose a viewpoint encoder to mitigate viewpoint disparities, while VARD Tu et al. (2023a) suggests a viewpoint invariant representation network to explicitly 170 capture changes. Additionally, Sun et al. (2022) integrates bidirectional encoding to refine change 171 localization, and NCT Tu et al. (2023b) utilizes a transformer to aggregate neighboring features. 172 These methodologies concentrate on the image modality, exploiting benchmark-specific characteris-173 tics such as nearly identical views in Spot-the-Diff Jhamtani & Berg-Kirkpatrick (2018) or synthetic 174 scenes with limited objects and change types in CLEVR Park et al. (2019). More recently, IDC-PCL 175 Yao et al. (2022) and CLIP4IDC Guo et al. (2022) have adopted BERT-like training approaches to 176 model difference captioning language, achieving state-of-the-art performance.

177 178

179 180

181

182

183

184 185

187

191

192 193

194

195

196

197

198 199

200

201

202

203 204

205 206 207

3 METHODOLOGY

In this section, we describe the methodology behind the dataset generation as well as IDC model training. Subsection 3.1 describes the data generation process used to create the METS dataset. Subsection 3.2 describes the architecture of the model used to train on the METS dataset to perform the multi-input image difference captioning task.

3.1 DATA GENERATION

We generate a dataset of image editing sequences, with textual machine annotations of each editorial 188 step and human edit summarization captions after the 5th, 10th, and 15th manipulation, as shown in 189 Fig. 4. Binary masks of the manipulation regions at each step are also included. Our dataset covers 190 a wide variety of pixel-level and generative manipulations. The prompt for each manipulation is generated using GPT-3.5 to ensure plausible and 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 a blur filter. Generative manipulations change the semantic content of the image, altering the story that the image tells.



208 Figure 3: METS image generation pipeline for generative manipulations. The image, its localized 209 narrative, object class name, and segmentation mask are sampled from the OpenImages dataset. The 210 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 211 specific and can be seen in suppmat. In the case of inpainting, the GPT3.5 block is omitted, and 212 the prompt is simply *background*. The pre-processing of the segmentation mask ensures that no part 213 of the object remains outside of the mask. It involves generating a convex hull of the mask and 214 applying dilation to it. The generative manipulation is then conditioned on the image, the mask, and 215 the prompt and applied using Firefly Generative Fill.

225

231

247

248

250

251

252 253 254

255 256

257

258

259

260

261 262

263

264 265 266

267

268 269

216 Pixel level manipulations are performed using the AuglyPapakipos & Bitton (2022) image augmen-217 tation library, with a random choice of augmentation type and parameters. Augmentation types 218 include brightness, contrast, saturation, and encoding quality changes; blur, noise and sharpness 219 filters; and overlaying random stripes of the color of different widths.

220 We further divide generative manipulations into three categories: inpainting where an object is re-221 moved from the image, **replacement** where an object is replaced with another object, and **property** 222 change where the object's material properties are altered. We illustrate different types of manipula-223 tions in Fig. 2. 224

Generative manipulations are applied using Firefly Generative Fill², which is a language-guided inpainting model. In addition to the image itself, the model is provided with a segmentation mask 226 and a text prompt. We generate a convex hull of the segmentation mask and apply dilation to it to 227 ensure that no part of the object remains outside of the mask. The origin of the text prompt depends 228 on the type of manipulation. For **inpainting** we use the word *background*, which was shown to 229 perform on par with inpainting-specific models. For **replacement**, we further illustrate the image 230 editing pipeline in Fig. 3, where 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 232 to come up with a probable replacement candidate object that would be a close match to the shape of 233 the original object. We use a similar strategy for **property change**, but prompting GPT3.5 to output a likely property change. 234



Figure 4: 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 249 superimposed on the images. The machine annotations generated during the sequence creation are shown in orange, while the human annotations are shown in blue. Note that only edit steps 5, 10, and 15 are shown, as these are the steps for which human annotations were collected. All other data types are available for all steps.

3.1.2 SEQUENCE GENERATION

We sample the images from the OpenImages dataset, making use of the provided segmentation masks and localized narratives. We choose the images with at least 5 non-overlapping segmentation masks. We then follow a procedure illustrated in Fig. 5 to apply a sequence of edits to the image. 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 to the mask.

Formally, we define the probabilities of applying a generative manipulation P_q , a pixel-level manipulation P_p and moving on to the next mask P_n as follows:

$$P_g = g - \frac{n}{2}, \quad P_p = (1 - g) - \frac{n}{2}, \quad P_n = 1 - P_g - P_p,$$
 (1)

where q = 0.9 if no generative manipulations have been applied to the mask yet and q = 0.1otherwise. The value of n is proportional to the number of manipulations already applied to the

²https://firefly.adobe.com/upload/inpaint

Figure 5: 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. The probability of applying a generative manipulation is greatly lowered if a generative manipulation has already been applied. This lowers, but does not eliminate, the chance of making destructive or mutually exclusive manipulations.

mask, defined as follows:

$$n = \max(0, \frac{40 \times (i - i_{min})}{100}), \tag{2}$$

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 as annotated within the OpenImages dataset, edit_name is the manipulation type and intensity is chosen at random from a set of predefined parameters, individual for each manipulation type.

³⁰² For generative manipulations, the text format is as follows:

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. 4, 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 with an average of 21.4 steps per sequence.

312

281

282

283

284

285

286 287 288

289 290

297

298

303 304

305

313 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 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. Examples of such summaries can be seen in Fig. 4, marked as *human annotations*.

319

321

320 3.2 ARCHITECTURE

Our architecture is illustrated in Fig. 6. Our setup consists of a Vision Transformer (ViT) Dosovitskiy et al. (2021) image encoder and the open-sourced LLaMA2-chat (7B) large language model Touvron et al. (2023b). The visual tokens are concatenated in groups of 4 and projected to the language





Figure 6: Architecture diagram of the model. The LLaMA-2 language model is conditioned using the multi-modal instruction template, which includes at least two image features and optional aux-342 iliary textual information. All optional content is placed within dashed boxes. The image features 343 extracted from the ViT image encoder are concatenated in groups of 4 and projected to the LLM embedding space with a linear projection layer. The visual encoder weights are frozen, and only the language model and the projection layer are trained. 345

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 350 captioning. We note that this approach is capable of handling an arbitrary number of input images, 351 which allows us to input several snapshots of the image editing sequence at once. 352

353 Optionally, we provide the model with auxiliary textual information in the form of machine anno-354 tations, described in Section 3.1. The annotations for each manipulation are interleaved with the 355 image features and are used to guide the model's attention to the relevant parts of the image.

We follow the multi-modal instructional template from Chen et al. (2023) and adjust it to our task:

361 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 362 is the instruction that is chosen at random from a set of predefined instructions, all synonymous with describe the defferences between the images. 364

365 The model is trained to minimize the captioning loss, which is defined as

$$\mathcal{L} = -\sum_{i=1}^{m} l(s^{v}, s_{1}^{t}, \dots, s_{i}^{t}),$$
(3)

where m is a variable token length and l is next-token log-probability conditioned on the previous sequence elements

371 372 373

374

376

370

340

341

344

346 347

348

349

356

$$l(s^{v}, s_{1}^{t}, \dots, s_{i}^{t}) = \log p(t_{i} | x, t_{1}, \dots, t_{i-1}).$$
(4)

375 3.2.1 TRAINING

All of the models are trained on a single A100 GPU with 80GB of memory for 300 epochs with 377 1000 steps per epoch and a batch size of 6. We use AdamW optimizer with a cosine learning rate scheduler with an initial learning rate of 10^{-5} and a warmup learning rate of 10^{-6} for a warmup period of 1000 steps. The input image size is 448×448 , and the maximum token length is 1024.

4 EXPERIMENTS

4.1 DATASETS



Figure 7: Examples of images and annotations from the CLEVR-Change, Spot-the-Diff, MagicBrush and PSBattles datasets.

In addition to our own dataset, we train and evaluate our model on a number of other datasets used in the image difference captioning literature illustrated in Fig. 7.

CLEVR-Change Johnson et al. (2017) 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.

413 Spot-the-Diff Jhamtani & Berg-Kirkpatrick (2018) is a dataset of 13,192 well-aligned image pairs
414 from CCTV cameras. There are no viewpoint changes, and the edits are limited to object addition,
415 deletion, or movement. The dataset is split into training, validation, and test sets following the
416 official split of 80%, 10% and 10%.

PSBattles Heller et al. (2018) is a dataset of real-world image pairs collected from the Reddit Photoshop Battles subreddit. The difference captions for a subset of the dataset were collected by Black et al. (2024) in a user study. We use this dataset for the evaluation of the model's generalization capability to real-world images.

InstructPix2Pix Brooks et al. (2022) is a dataset of ~1M image pairs generated with prompt-toprompt Hertz et al. (2022) approach. The difference captions are later generated by Black et al.
(2024) 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 Zhang et al. (2024) 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 in two different settings: standard IDC with two images as input and 'image difference captioning with multiple inputs' (IDC-MI). The former setting is the most common in the literature, while the latter is a novel setting that we introduce in this work.

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

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	С	М	R	L (M)	L (H)
METS								
GPT3.5 Brown et al. (2020)	0	yes	1.6	8.6	10.4	15.1	16.2	0.6
GPT4-V et al. (2024)	2	no	4.0	18.6	14.0	20.3	22.2	2.6
GPT4-Vet al. (2024)	2	yes	1.3	0.3	11.5	13.5	19.7	0.9
GPT4-Vet al. (2024)	4	no	3.0	15.1	13.4	19.9	26.9	1.9
GPT4-Vet al. (2024)	4	yes	1.4	0.4	11.6	12.9	24.1	1.2
FVTC-2 (ours)	2	no	5.8	20.7	11.4	23.1	22.6	9.4
FVTC-2T (ours)	2	yes	7.8	<u>25.8</u>	13.0	<u>26.0</u>	24.3	<u>11.0</u>
FVTC-4 (ours)	4	no	6.6	23.5	12.3	24.3	22.6	9.6
FVTC-4T (ours)	4	yes	8.2	25.9	<u>13.4</u>	26.3	30.1	12.4
MagicBrush								
FVTC-2 (ours)	2	no	4.9	29.4	13.3	28.1	-	-
FVTC-4 (ours)	4	no	6.8	44.5	15.6	31.2	-	-

For the IDC-MI setting, 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 compare the performance of our model with GPT4-V, which has multi-modal capabilities and is capable of taking multiple images and/or text as input. Additionally, we compare with GPT3.5, which serves as a text-only baseline, taking as input only the auxiliary text and no images.

The results of the IDC-MI setting are shown in Table 1. We demonstrate that 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. On the other hand, GPT4-V suffers from the addition of intermediate images and text, showing a decrease in performance in both cases.

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 both intermediate images and textual information shows an average improvement of 22.4% across all metrics.

On the other hand, the performance of GPT4-V suffers from the addition of intermediate images, decreasing in performance with the addition of both extra images and text.

Table 2: Image difference captioning performance evaluation on the CLEVR-Change and PSBattles 487 datasets. We report the performance of our model and compare it with the state-of-the-art models 488 and report BLEU-4 (B4), CIDEr (C), METEOR (M) and ROUGE-L (R) scores. 489

MODEL	TRAINING DATA		С	Μ	R	S
CLEV	/R CHANGE					
DUDA Park et al. (2019)	CLEVR	47.3	112.3	33.9	-	-
IFDC HUANG ET AL. (2022)	CLEVR	49.2	118.7	32.5	69.1	-
R^3 Net+SSP Tu et al. (2021a)	CLEVR	54.7	123.0	39.8	73.1	-
SGCC OLUWASANMI ET AL. (2019)	CLEVR	<u>51.1</u>	121.8	<u>40.6</u>	73.9	-
NCT TU ET AL. (2023B)	CLEVR	<u>55.1</u>	124.1	40.2	73.8	-
SRDL+AVS TU ET AL. (2021b)	CLEVR	54.9	122.2	40.2	73.3	-
VARD TU ET AL. (2023A)	CLEVR	55.2	124.1	40.8	<u>74.1</u>	-
FVTC-2 (OURS)	CLEVR	54.7	151.8	40.0	77.1	-
Spot-the-Diff						
SRDL+AVS TU ET AL. (2021B)	SPOT-DIFF	-	35.3	13.0	31.0	18.0
R^3 Net+SSP Tu et al. (2021a)	Spot-Diff	-	36.6	13.1	32.6	18.8
VARD-LSTM TU ET AL. (2023A)	Spot-Diff	-	39.3	13.1	33.1	17.5
VARD-TRANSFORMER TU ET AL. (2023A)	Spot-Diff	-	30.3	12.5	29.3	17.3
FVTC-2 (OURS)	SPOT-DIFF	-	45.5	13.7	28.7	19.3
PSBATTLES						
VIXEN-C BLACK ET AL. (2024)	IP2P	4.5	7.7	9.5	20.5	-
FVTC-2 (OURS)	IP2P	5.3	10.3	10.8	22.	-
FVTC-2 (OURS)	IP2P + METS	5.5	14.2	11.2	22.6	-

510 511 512

513

4.2.2 EVALUATING IDC WITH TWO INPUTS

514 We observe that in the IDC setting, shown in Table 2, the model achieves competitive performance 515 on the CLEVR-Change dataset, outperforming the previous state-of-the-art model VARD on the 516 CIDEr and ROUGE-L metrics. On the InstructPix2Pix dataset, the model outperforms VIXEN only 517 on the METEOR metric. However, it shows a better capability to generalize to real-world images, 518 outperforming VIXEN on the PSBattles dataset for all metrics. Additionally, fine-tuning the model 519 on the METS dataset further improves its performance on PSBattles, showing the dataset's ability 520 to bridge the domain gap between synthetic and real-world images. 521

4.3 LIMITATIONS

As with most LLMs, FVTC can occasionally hallucinate details that are not present in the input 525 images. When errors are made, they are most commonly occurrences of miscounting. For example, 526 the model may state that multiple occurrences of an object have been replaced with another instead of a single occurence or vice versa (i.e. use of plural versus singular). The remaining category of failure cases observed involves cases where level of detail may be too succinct, for example stating 528 an object is replaced but not stating with what.

529 530 531

532

527

522

523 524

5 CONCLUSION

We have introduced a novel task of image difference captioning with multiple inputs and demon-534 strated that the presence of additional visual and/or textual inputs improves the model's captioning 535 performance. We have introduced METS – a new dataset of long image editing sequences paired 536 with machine annotations and human edit summarization captions. We have trained a multi-modal 537 LLM with multiple visual inputs and provided a comprehensive evaluation of the benefits of both additional visual and textual inputs. Additionally, we have demonstrated that fine-tuning a model 538 that is trained on other synthetic data with METS helps to bridge the domain gap and improves zero-shot performance on real-life images.

490 491

540 REFERENCES 541

+ 1 12 13 14	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. <i>NeurIPS</i> , 35:23716–23736, 2022.
15 16 17	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 <i>Proc. CVPR</i> , pp. 6077–6086, 2018.
18 19 50	Alexander Black, Jing Shi, Yifei Fai, Tu Bui, and John Collomosse. Vixen: Visual text comparison network for image difference captioning, 2024.
51 52	Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. <i>arXiv preprint arXiv:2211.09800</i> , 2022.
53 54 55 56	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. <i>NeurIPS</i> , 33:1877–1901, 2020.
57 58 59	Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Kr- ishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large lan- guage model as a unified interface for vision-language multi-task learning, 2023.
50 51 52	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.
5 5 5	Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. Meshed-memory trans- former for image captioning. In <i>Proc. CVPR</i> , pp. 10578–10587, 2020.
6 7	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> , 2018.
o 9 0 1	Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venu- gopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In <i>Proc. CVPR</i> , pp. 2625–2634, 2015.
2 3 1 5	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.
; ; ; ; ;	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 <i>arXiv preprint arXiv:2303.03378</i> , 2023.
	Constantin Eichenberg, Sidney Black, Samuel Weinbach, Letitia Parcalabescu, and Anette Frank. Magma–multimodal augmentation of generative models through adapter-based finetuning. <i>arXiv</i> preprint arXiv:2112.05253, 2021.
	Josh Achiam et al. Gpt-4: OpenAI technical report, 2024.
	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. <i>arXiv preprint arXiv:2304.15010</i> , 2023.
0 1	Alex Graves and Alex Graves. Long short-term memory. <i>Supervised sequence labelling with recurrent neural networks</i> , pp. 37–45, 2012.
5	S. Gregory. Ticks or it didn't happen. https://lab.witness.org/ ticks-or-it-didnt-happen/, 2019. Accessed: 2024-01-20.

610

616

623

628

- Jiuxiang Gu, Jianfei Cai, Gang Wang, and Tsuhan Chen. Stack-captioning: Coarse-to-fine learning for image captioning. In *Proc. AAAI*, volume 32, 2018.
- 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, pp. 33–42, 2022.
- 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.
- 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.
- Kiaowei 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*, pp. 17980–17989, 2022.
- Lun Huang, Wenmin Wang, Yaxian Xia, and Jie Chen. Adaptively aligned image captioning via adaptive attention time. *NeurIPS*, 32, 2019.
- 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.
- Harsh Jhamtani and Taylor Berg-Kirkpatrick. Learning to describe differences between pairs of
 similar images. In *Proc. Conf. Empirical Methods NLP*, pp. 4024–4034, 2018.
- 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*, pp. 2901–2910, 2017.
- Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proc. CVPR*, pp. 3128–3137, 2015.
- Hoeseong Kim, Jongseok Kim, Hyungseok Lee, Hyunsung Park, and Gunhee Kim. Agnostic change
 captioning with cycle consistency. In *Proc. ICCV*, pp. 2095–2104, 2021.
- 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.
- Kiujun 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*, pp. 121–137. Springer, 2020.
- 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*, pp. 375–383, 2017.
- 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, pp. 2286–2293, 2021.
- 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.
- Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. *arXiv preprint arXiv:2209.15162*, 2022.
- 647 Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021.

656

657

658

659

675

686

687

688

- 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.
- 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.
 - Zoë Papakipos and Joanna Bitton. Augly: Data augmentations for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 156–163, 2022.
- Dong Huk Park, Trevor Darrell, and Anna Rohrbach. Robust change captioning. In *Proc. ICCV*, pp. 4624–4633, 2019.
- 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*, pp. 8748–8763. PMLR, 2021. URL https://github.com/OpenAI/CLIP.
- Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical
 sequence training for image captioning. In *Proc. CVPR*, pp. 7008–7024, 2017.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. CVPR*, pp. 10684–10695, 2022.
- Kiangxi 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*, pp. 574–590.
 Springer, 2020.
- 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.
- 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.
- 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*, 2023a.
 - 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*, 2023b.
- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Mul timodal few-shot learning with frozen language models. *NeurIPS*, 34:200–212, 2021.
- Yunbin Tu, Liang Li, Chenggang Yan, Shengxiang Gao, and Zhengtao Yu. R³Net:relationembedded representation reconstruction network for change captioning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 9319–9329, Online and Punta Cana, Dominican Republic, November 2021a. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.735. URL https://aclanthology.org/2021. emnlp-main.735.
- 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 *Find*ings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pp. 63–73, 2021b.

702 703 704	Yunbin Tu, Liang Li, Li Su, Junping Du, Ke Lu, and Qingming Huang. Viewpoint-adaptive represen- tation disentanglement network for change captioning. <i>IEEE Transactions on Image Processing</i> , 32:2620–2635, 2023a. doi: 10.1109/TIP.2023.3268004.
705 706 707	Yunbin Tu, Liang Li, Li Su, Ke Lu, and Qingming Huang. Neighborhood contrastive transformer for change captioning, 2023b.
708 709 710	Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In <i>Proc. CVPR</i> , pp. 3156–3164, 2015.
710 711 712	Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision. In <i>Proc. ICLR</i> , 2021.
713 714	Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In <i>Proc. CVPR</i> , pp. 10685–10694, 2019.
715 716 717	Linli Yao, Weiying Wang, and Qin Jin. Image difference captioning with pre-training and contrastive learning. In <i>Proc. AAAI</i> , volume 36, pp. 3108–3116, 2022.
718 719	Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Hierarchy parsing for image captioning. In <i>Proc. ICCV</i> , pp. 2621–2629, 2019.
720 721 722 723	Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
724 725 726 727	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 <i>Proc. CVPR</i> , pp. 5579–5588, 2021.
728 729 730	
731 732	
733 734	
735 736	
737 738 739	
740 741	
742 743	
744 745	
746 747	
748 749	
750 751	
752 753	
754 755	