DiffTell: A High-Quality Dataset for Describing Image Manipulation Changes

Zonglin Di¹*, Jing Shi², Yifei Fan², Hao Tan², Alexander Black³, John Collomosse^{2,3}, Yang Liu¹ University of California, Santa Cruz, ²Adobe Research, ³CVSSP, University of Surrey

{zdi, yangliu}@ucsc.edu, {jingshi, yifan, hatan, collomos}@adobe.com, {alex.black, j.collomosse}@surrey.ac.uk

Abstract

The image difference captioning (IDC) task is to describe the distinctions between two images. However, existing datasets do not offer comprehensive coverage across all image-difference categories. In this work, we introduce a high-quality dataset, DiffTell, with various types of image manipulations, including global image alterations, objectlevel changes, and text manipulations. Data quality is controlled through careful human review and filtering. Additionally, to scale up the data collection without prohibitive human labor costs, we explore the possibility of automatically filtering for quality control. We demonstrate that both traditional methods and recent multimodal large language models (MLLMs) exhibit performance improvements on the IDC task after training on the DiffTell dataset. Through extensive ablation studies, we provide a detailed analysis of the performance gains attributed to DiffTell. Experiments show DiffTell significantly enhances the availability of resources for IDC research, offering a more comprehensive foundation and benchmark for future investigations. The dataset is available at https://huggingface.co/ datasets/zodi1121/DiffTell.

1. Introduction

Given the tremendous progress in image generation [51, 53], disseminating AI-modified fake images can lead to widespread misinformation, erosion of public trust, and manipulation of public opinion on critical issues. Emerging open standards, such as C2PA (Coalition for Content Provenance and Authenticity, 2023), outline provenance frameworks that utilize perceptual hashing techniques to link images found in the public domain with a federated database of original content [7, 47]. Upon retrieving the source image, image difference captioning (IDC) models can describe the discrepancies between the circulated manipulated image and its original, enabling individuals to make more informed and nuanced trust assessments. The overview of

the entire system is given in Appendix A. IDC has been researched with various algorithms [17, 48, 59–61, 65, 73]. However, the image domain and the types of visual differences in the current IDC dataset are limited or small-scale, as summarized in Table 1. This makes the generalization ability of the current model unsatisfactory. Although some datasets are larger, such as OneDiff [20], they tend to be noisy and lack human verification. Thus, a comprehensive, clean IDC dataset on a large scale is needed.

The IDC dataset consists of the data triplet, including one image pair (the original and the manipulated) and one language caption describing the difference between them. The formal definition is given in Section 3.1. As shown in Table 1, existing datasets focus either on domain-specific images, such as Spot-the-diff [22], which uses frames of the surveillance videos, or rendered scenes with limited geometric objects and change types (color, texture, add, drop, remove) in CLEVR [46]. Even though image editing request (IER) has various types of editing on the real natural images, it is limited in volume ($\sim 4K$) since manual human editing is costly and time-consuming, making it harder to scale up [59]. Given the development of generative AI and image editing technologies, language-guided AI-manipulated image data have been created using data triplets: the before-edited image, the after-edited image, and the language editing request. InstructPix2Pix [9] leverages GPT-3 [10] to scale up possible editing commands and resort to prompt2prompt [18] for automatic editing. However, we find that it has a high error rate exceeding 60%. MagicBrush [74] provides 10K manually annotated real image editing triplets with careful quality control, but only contains local edits. It has showcased the importance of highquality data for language-guided image editing. Therefore, we identify the need for an IDC dataset that is varied in manipulation types and maintains high quality on a large scale.

To better support research in image difference captioning, we introduce the *DiffTell* dataset, specifically created to encompass a broader range of editing types, including both real and synthetic image pairs, while maintaining careful quality control. We categorize image differences into four categories: background change, local object change,

^{*}This work was done when ZD interned at Adobe Research

Dataset	Size	Real	Syn.	Human Annotation	Categories	Domain
CLEVR-Change [46]	70K	×	/	X	Local object	primitive shapes
Birds-to-Words [15]	4.8K	1	X	✓	Local object	Birds
Spot-the-Diff [23]	13K	1	X	✓	Local object	top-down street view
IER [59]	4K	1	X	✓	Comprehensive	varied natural images
PSBattle [8]	100	1	X	✓	Comprehensive	varied natural images
One-Diff [20]	316K	1	1	×	Comprehensive	varied natural images & genAI
DiffTell (Ours)	70K	1	✓	✓	Comprehensive	varied natural images & genAI

Table 1. The comparison involves *DiffTell* and currently available datasets designed for the image difference captioning (IDC) task. "Real" and "Syn." signify the presence of real and synthetic images in the datasets, respectively. The term "comprehensive" category denotes that the dataset can encompass all the categories outlined in Section 3.2. A more detailed existing dataset description is given in Appendix B.

text manipulation, and image style change from various data sources. Examples of the DiffTell dataset are illustrated in Figure 1. We first include two accessible languageguided image editing datasets InstructPix2Pix [9] and MagicBrush [74]. We manually filtered out the noisy, lowquality data in InstructPix2Pix. As text manipulation is critical in creating fake news, we enriched the text addition and removal data by inpainting the text in MARIO-10M images [12]. In addition, we extended the object addition and removal by inpainting the COCO [33] dataset. All AIgenerated editing outcomes have passed the quality filtering process. Moreover, since the labor cost of manual quality filtering can be expensive when scaled up, we further develop an automatic data filtering model to reduce the cost and observed the benefit of such an auto filtering process according to model captioning performance.

Multimodal large language model (MLLMs) have become increasingly popular in the research community due to their strong general-purpose capability. By linking large language models (LLMs) with visual conditioning [38, 77], MLLMs have shown impressive results in natural instruction-following and visual reasoning capabilities. Meanwhile, the *DiffTell* dataset can serve as a visual instruct finetuning [38] step upon the multiple MLLM models. We demonstrate the general improvement of IDC performance using the *DiffTell* dataset on various baselines, indicating its value and benefits. In summary, our contributions are

- Proposing the *DiffTell* dataset that includes various kinds of changes with high-quality samples on a larger scale than previous datasets;
- Proving *DiffTell* can boost the IDC on various baselines on both IER and PSBattle datasets on various models;
- A detailed analysis of how the *DiffTell* dataset enhances IDC in different editing categories;
- Probing the model-based data filtering given the fixed amount of human-filtered data, allowing potential data scale-up.

2. Related Work

2.1. Multimodal Large Language Models

With the development of visual encoder and its combination to large language models (LLMs), multimodal large language models (MLLMs) [36, 37, 39, 77] show promising capability to understand images, accept text inputs, and generate natural-language responses. Increasing the model capacity and dataset size can generally improve the capability of MLLMs [4, 14, 75]. Visual encoders [30, 31, 49] are applied to encode visual information into visual tokens, providing input for the LLMs. Other strategies like expanding the instruction-tuning dataset [34] and increasing the visual resolution [4, 35, 69] can also improve the performance of the MLLMs. Recently, MLLMs have been used to understand fine-grained images, such as in local region understanding [13, 40]. Image difference captioning is closely related to fine-grained image understanding with multipleimage input.

2.2. Image Difference Captioning

As mentioned above, MLLMs are used to understand the local region. Image difference captioning (IDC) is more challenging because the model needs to not only understand each image correctly but also capture and identify the difference between two images correctly and express it precisely in language. In IDC, the caption aims to describe the differences between the images while ignoring their commonalities. The first work on IDC, Spot-the-Diff [23], categorizes different types of changes and uses an LSTM-based network to model them. DUDA [45] enhances the robustness against slight global changes by analyzing image differences at the CNN semantic level instead. Viewpoint invariant encoders have been proposed in M-VAM [57], VACC [27], and VARD [63] to mitigate potential viewpoint differences, while [58] uses bidirectional encoding to improve change localization and NCT [64] aggregates neighboring features with a transformer. IDC-PCL [72] and CLIP4IDC [16] adopt BERTlike training strategies to model the difference-captioning

language. SCORER [66] applies a self-supervised crossview representation reconstruction technique for difference captioning. Recently, with the advancement of MLLMs, more datasets have integrated the existing IDC dataset to train powerful MLLMs with diverse capabilities. For instance, LLaVA-OneVision [29] includes the CLEVR dataset, and Mantis-Instruct [24] incorporates the Spot-the-Diff dataset. [20] proposes OneDiff, a large-scale dataset incorperating several existing IDC dataset as well as the synthetic data using ChatGPT.

2.3. Image Editing

One of the biggest challenges in IDC is the shortage of high-quality, comprehensive datasets of paired images. The development of the diffusion model [19] significantly improves the quality and controllability of the generated images. By controlling the cross-attention, diffusion models can transform the image globally [53, 55]. Local editing depends on the fine-grained predicted or user-provided mask, such as inpainting [2, 41, 43]. Different from the image transformation and local editing, the input of the instruction-guided image editing is in the command format rather than the detailed description and mask [9]. DiffTell significantly benefits from the progress in image generation models [54], especially the local editing model, leveraging their capabilities to enhance the quality and diversity of the dataset. HQ-Edit introduces a high-quality instructionbased image editing dataset with around 200,000 edit and a scalable data collection pipeline leveraging advanced foundation models [21].

3. Problem Formulation and Dataset Construction

3.1. Problem Definition

For IDC problem, when presented with two similar images, denoted as I_1 and I_2 , our objective is to employ a vision-language (VL) model, f_{θ} , to articulate the distinctions between I_1 and I_2 in natural language. This can be represented as: $T_{I_1,I_2} = f_{\theta}(I_1,I_2)$, where T_{I_1,I_2} represents the descriptive caption text provided by the model regarding the dissimilarities between the images, and θ signifies the model parameters within the VL model. The elements I_1 , I_2 , and T_{I_1,I_2} collectively form the constituents of each sample within the IDC dataset.

3.2. IDC Categories

Considering that our main motivation is to alleviate the misinformation and spreading of doctored images, we focus on the image pairs created by manipulation or editing and exclude the pairs without any correlation or that cannot be easily obtained by human/AI editing. To further concretize the research problem, we categorize four image dif-

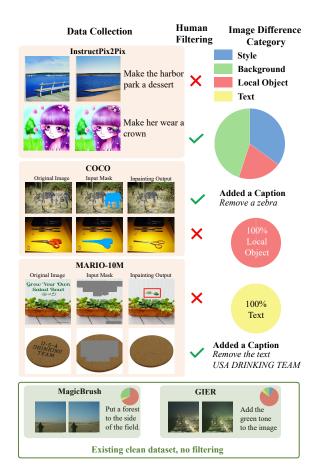


Figure 1. The data collection pipeline and the data distribution regarding the image difference categories. The data collection involves two steps. Initially, data is gathered from COCO, MARIO-10M, InstructPix2Pix, MagicBrush, and GIER. For COCO and MARIO-10M, an in-painting process is applied to the images with the help of masks, and the labeling team subsequently filters out unsuccessful cases. The three images are the original image, the input mask, and the output from Firefly Generative Fill, from left to right. In the second (lower) COCO example, where the scissors remain unaltered, the labeling team excludes this case from the dataset. Similarly, for the first (upper) MARIO-10M example, although the text in green is removed, the generation model introduces an additional element outlined in the red box, leading to the exclusion of this example as well. In the case of InstructPix2Pix, the labeling team verifies the alignment between image pairs and language instructions. Instances with unsuccessful modification (e.g., the dessert modification in the top example) are removed from the dataset. For the MagicBrush and GIER datasets, there is no need to filter the image as they have already undergone manual filtering. The final stage involves compiling the filtered data, resulting in the creation of the DiffTell dataset.

ference types as *background change*, *local object modification*, *style change* and *text manipulation*. **Background change** is alterations related to the background, such as removing, adding, or changing the background of an image. **Text manipulation** involves addition, removal, or modi-

fication of text within the original image. **Local object change** is about object re-colorization, appearance editing, object removal, insertion, or translation. **Style change** is the artistic style change, such as realistic photo to painting, and photo-realistic style change, such as adjusting the brightness or tone. Existing datasets, such as IER, mainly include the first three categories but lack text manipulation. However, text manipulation is crucial in our scope since some text changes can flip the message of an image, leading to fake news and forged messages. For example, the message of a smiling face image can be changed from happiness to sarcasm by adding the sentence "absolutely thrilled to be overworked and underpaid." Therefore, we put additional effort into text manipulation data collection. The detailed elaboration of each difference category is as follows.

3.3. Dataset Collection Pipeline

Based on the definition in Section 3.1, the triplet (I_1,I_2,T_{I_1,I_2}) reflecting the four categories given above is the fundamental element to build an image difference captioning (IDC) dataset. As the mirrored task of IDC, the instruction-guided image editing dataset is considered, which provides (I_1,I_2,T_{I_1,I_2}) exactly. We select InstructPix2Pix [9], GIER [56], and MagicBrush [74] as the subset of our dataset due to the editing types, dataset sizes/qualities. The difference categories of those three datasets are given in Table 2.

Most existing vision datasets only provide I_1 and its corresponding annotations, like the object segmentation mask or the object's name. Empowering the generative model [52, 71], we can remove an object from the image to generate I_2 , although a quality check step is necessary due to the limitation of the generative model. The difference caption T_{I_1,I_2} can be generated based on the editing operation from the generative model. For datasets only providing I_1 , such as COCO and MARIO-10M, we mainly focus on object change and text manipulation. For the generation of I_2 , we apply the inpainting model Firefly Generative Fill and the details of how to generate images are given in Appendix I. T_{I_1,I_2} is based on the template Add <Text> / <Object> or Remove <Text> / <Object> depending on the order of I_1 and I_2 , which is determined by a random number generator whose probability is 0.5. For the datasets providing I_1, I_2 and T_{I_1,I_2} without manually filtering like InstructPix2Pix, we ask the labeling team to filter them. We provide the details of each subset and the annotation details below.

InstructPix2Pix [9] provides I_1 , I_2 and T_{I_1,I_2} , where (I_1,I_2) are generated by StableDiffusion [53] in combina-

Datasets	Syn. Image	F. Rate (%)	Dataset Size
InstructPix2Pix	/	35.13	17,592
GIER	×	100.00	6,179
MagicBrush	×	100.00	8,807
MARIO-10M	×	26.86	30,903
COCO	×	30.93	3,986
DiffTell	✓	-	67,589

Table 2. Summary of the source datasets from which we derived our dataset. "Syn. Image" indicates whether the image domain contains synthetic images, while the "F. rate" denotes the ratio of images retained after manual filtering by our labeling team if needed, which is equal to (100% - Rejection Rate).

tion with Prompt-to-Prompt, and T_{I_1,I_2} is produced by a finetuned GPT-3 [10]. It is a large (450K+) dataset with various image-difference categories thanks to the automated process. However, the automated process occasionally mismatches the image pair and its corresponding instruction. We present such a noisy sample in Figure 1. The instruction "Make the harbor park a dessert" does not describe the difference between the image pair. To mitigate this, our labeling team meticulously reviews a subset to retain clear and accurate samples. After reviewing 50,012 selected triplets from the InstructPix2Pix dataset, we obtain 17,592 image pairs covering background, style, and local object change.

GIER [56] also provides the (I_1, I_2, T_{I_1,I_2}) triplet, presenting 6,179 image pairs. IER and GIER are both from the same source and complementary to each other. More specifically, they are both from the human Photoshop-edited images based on the language editing instructions. GIER is mostly characterized by its global tone and lighting editing. We employ these pairs along with expert annotations as I_1 , I_2 , and T_{I_1,I_2} respectively, while standardizing the language style by removing unnecessary politeness indicators like "Please."

MagicBrush [74] constitutes a high-quality dataset for multi-turn image editing, meticulously curated through manual filtering, providing (I_1, I_2, T_{I_1, I_2}) triplets in high quality, which can be used directly in IDC task. To adapt this multi-turn editing to fit our framework, we segmented it into several single-turn edits and randomized their order. As a result, we incorporate 8,807 image pairs from Magicbrush into *DiffTell*.

MARIO-10M [12]: Text manipulation data is gathered based on MARIO-10M, a dataset offering rough segmentation masks and optical character recognition (OCR) results for text within images. The dataset only provides I_1 , and we use FireFly Generative Fill to remove the masked text from the images to generate I_2 with the input of I_1 and its corresponding mask. We apply mask dilation, enlarging the original mask by 5 pixels to make the region of interest (ROI) covered by the mask as much as possible. Our labeling team

¹As a type of artificial intelligence that can translate text and other inputs into extraordinary results, Firefly Generative Fill model can generate the image according to the image or text input and be accessed at https://firefly.adobe.com.

carefully verifies the resulting images to ensure that the text is fully removed and there is no additional element added in I_2 , leading to the retention of 30,903 image pairs out of 115,059 in our dataset. For filtered image pairs (I_1,I_2) , the language templates T_{I_1,I_2} , we use are "add text" or "remove text," depending on the order of the image pair. We also add the OCR results to the caption, with examples given in Figure 1.

COCO [33]: Similar to MARIO-10M dataset, COCO dataset only provides I_1 and we need to generate I_2 and T_{I_1,I_2} . We initially generated masks for each instance from the annotations in the training set. Different from MARIO-10M, the mask cannot be used directly because some object masks are tiny, while others occupy almost the entire image, despite the object being the same. To ensure proper object sizes, a mask filtering technique is applied, selecting objects within a specific size range based on the distribution of mask sizes within each class. For each class, we select the images with the masks whose area is 50%-75% of the largest area to ensure that the change within the image pairs is obvious and meaningful while not occupying the full image. This process results in a selection of 128,969 images from an initial pool of 860,001. Similar to the MARIO-10M approach, mask dilation is applied in case of potential detail loss in polygon masks. Objects are in-painted using FireFly Generative Fill, and the resulting images are scrutinized by our labeling team, resulting in a final selection of 3,986 image pairs out of 12,886 for our dataset. After getting the image pairs with and without the object from inpainting, we follow the language template in MARIO-10M, which is add <object> or remove <object> as shown in Figure 1. The COCO subset in DiffTell focuses on local object change.

Quality Check Statistics We use LabelBox² as our crowd-sourcing platform. Each sample added to *DiffTell* is initially labeled by an annotator and then reviewed by a high-performing annotator selected by us. To identify high-performing annotators, we have each annotator label 500 images to assess their understanding of the task, and we manually evaluate their accuracy. The top 30% of annotators are selected as high-performing and assist with the review process on a larger scale. On average, the labeling time is 56.73 seconds, while the reviewing time averages 72.44 seconds.

Rationality of Data Construction with Generative Model Considering the circulated deceptive doctored images are usually edited by humans or AI, we also create image pairs with human or AI manipulation. Instruct-Pix2Pix, MagicBrush, MARIO-10M, and COCO are AI-edited, while GIER is human-Photoshopped. And we can control the type of difference in the dataset based on the editing we applied, allowing future balancing and debiasing

of various IDC categories.

Our work specifically focuses on image manipulation difference captioning. In the existing image difference captioning datasets, the image pairs typically consist of one real-world image and another that has been manipulated using software. Using only real-world images for pairs would significantly constrain the scale and diversity of the dataset.

3.4. Dataset Analysis

Following the dataset collection, we conduct a statistical analysis of the DiffTell dataset based on the four categories in Section 3.2. The contribution to each editing category within each subset of DiffTell is presented in Figure 1. Background and image style changes are from GIER and InstructPix2Pix. MARIO-10M is for text manipulation. Local object change is from all the subsets except MARIO-10M. Over 72.9% images' resolution is 512×512 . The largest image is 1024×1024 , which is over 10%. The ratio of the images of other resolutions is less than 1.5%. The average length of the difference description is 9.72 words. The longest description contains 66 words, while the shortest has 33 words. The mode of the description length is 9 words. The description length distribution, addional dataset illustration, and how the labeling team works to filter the data are provided in the Appendix I.

4. Experiments

4.1. Experiment Setup

Benchmark Datasets and Evaluation Metrics We conduct experiments on the IER dataset [59] and the PSBattle dataset [8], which encompass a wide range of image editing differences. The PSBattle dataset is sourced from the PSRequest channel on Reddit³, comprising 100 pairs of images, each associated with at least three captions depicting image modifications [8]. For the out-of-distribution evaluation using the existing IDC dataset, we employ Spot-the-Diff [23] and CLEVR-DC [26], and the details are given in Appendix G.

In the case of IER, we evaluate the performance on the testing set by comparing models trained exclusively on the IER training set and those trained on a combination of the IER training set and the *DiffTell* dataset. There is overlap between the GIER and IER datasets, and we exclude the samples in GIER that also appear in the IER testing set. For the PSBattle dataset, we adopt it as an out-of-domain dataset to test the zero-shot capability of our model. Aligned with prior captioning research, we employ BLEU@4 [44] (B@4), METEOR [5] (M), CIDEr [67] (C), and ROUGE-L [32] (R-L), SPICE [1] (S) and BertScore- F_1 [76] (B- F_1) as the evaluation metrics. In addition, we

²https://labelbox.com

³https://www.reddit.com/r/photoshopbattles/

Testing Set		CLIP	4IDC	OpenFlar	ningo-3B	Fuy	u-8B	LLaVA-Ir	nterleave-7B	Idefic	s3-8B	Qwen2	-VL-7B
	w/ DiffTell	Х	✓	X	✓	X	1	×	✓	×	1	Х	✓
	BLEU@4	5.65	8.84	4.45	6.49	4.85	9.59	6.09	11.06	11.45	17.29	17.15	17.76
IER	METEOR	10.23	13.54	14.87	16.68	11.84	16.52	14.05	17.35	19.28	21.64	19.84	22.02
IEK	CIDEr	22.52	28.14	15.80	21.04	23.67	41.05	29.69	44.79	53.06	64.35	64.61	67.71
	ROUGE-L	28.95	36.84	29.79	31.36	28.10	35.44	32.67	37.21	37.37	41.86	44.58	45.35
	BLEU@4	0.00	3.08	2.35e-4	2.12	1.38	2.15	2.60	4.13	5.37	5.87	5.06	5.33
PSBattle	METEOR	3.08	6.25	2.33	6.60	4.79	7.57	8.88	9.39	9.96	12.93	10.68	12.12
PSBattle	CIDEr	1.59	3.63	4.79	7.75	4.05	4.19	7.86	8.55	11.35	12.97	14.29	16.66
	ROUGE-L	13.83	21.22	16.24	19.10	12.23	13.73	18.01	21.09	20.51	22.11	22.34	22.91

Table 3. Comparison of the methods fine-tuned on IER training set with and without *DiffTell*. The testing sets are the IER testing set and the PSBattle dataset. The results of SPICE and BertScore are given in Table 11 in Appendix F.

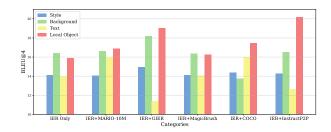


Figure 2. Category-wise BLEU@4 comparison on IER testing set using Qwen2-VL-7B trained with different subsets in *DiffTell*.

also conduct human evaluation. The results of human evaluation [25] are given in Appendix F.

Baselines and Implementation Details We implement several baseline methods for IDC to comprehensively illustrate the benefits of the *DiffTell* dataset, including both IDC-specific and MLLM methods. For IDC-specific methods, we use CLIP4IDC [17]. For MLLM methods, we report OpenFlamingo-3B [3], Fuyu-8B [6], Llave-Interleave-8B [36], Idefics3-8B [28] and Qwen2-VL-7B [68] here. We follow the instruction tuning methods to train the MLLMs. Without further clarification, the prompt we use across all the experiments is "What is the difference between two images?". We also try the diverse instruction prompts, and the results are given in Appendix E, but the difference is not significant. The implementation details and the results of more baselines [62, 64] are given in Appendix C.

4.2. Main Result

Quantitative Result We report the experiment results on the IER testing set and PSBattle dataset with and without DiffTell in Table 3. Results demonstrate DiffTell's ability to enhance performance across all evaluation metrics and for all baseline methods, underscoring the contribution of the DiffTell dataset on IDC. Notice that OpenFlamingo-3B with an LLM backbone is less capable than CLIP4IDC with a much smaller model size. We suspect that the Flamingo model does not have direct modeling of the in-

teraction between the two images because each image feature is cross-attentioned by language tokens, and then the language tokens will interact via causal attention. In contrast, in CLIP4IDC, the two image patch features extracted by CLIP are fused using a transformer, which is a direct information interaction among image tokens, serving as a strong condition to guide the transformer decoder to generate the language that describes the visual difference. There is no image encoder in the Fuyu model, and the image is patched linearly to the transformer. Thus, Fuyu can accept an image of arbitrary size, improving its capability to detect tiny differences and small objects. This can be the reason why Fuyu improves greatly after fine-tuning. For Llava-Interleave-7B, the pre-trained interleaved dataset provides a good knowledge base for the model to understand the context with multiple image inputs. Thus, it outperforms the IDC-specific model without *DiffTell* and can perform best among all the baselines. Qwen2-VL-7B achieves the best performance over all the comparison methods.

Qualitative Study We compare the prediction from Qwen2-VL-7B models trained with and without *DiffTell*. The visualization examples of IER and PSBattle testing set are shown in Figures 3 and 4, respectively.

As depicted in Figure 3, the model demonstrates enhanced proficiency in describing local object changes, text detection and recognition, background alterations, and image style changes. Notably, in the local object change example, the model accurately identifies the addition of a tiger which is tiny. In addition, the model can describe the spatial relationship in the text output. In the second example of background change. Compared with the model without DiffTell, the model can describe the general change and the object unchanged, showing that the caption is generated by not just learning the image manipulation, but recognizing and understanding the image. Furthermore, in the third example depicting an image style change, the model with DiffTell recognizes the color name and describes it. In the text manipulation example, the model exhibits OCR capabilities without relying on existing OCR techniques. Moreover, the model can recognize subtle and long text. In the

zero-shot testing scenario of PSBattle, it is crucial to observe the spatial capability is boost a lot by using DiffTell dataset. Moreover, the model can observe the more changes in addition to the most obvious one.

By analyzing the qualitative results, we find the model trained with DiffTell have a obvious capability boost in spatial relationship, tiny object recognition, complex text OCR, geometric operations and complex change identification which includes several changes between the two images. More examples are given in Appendix L.2.

<Local object change>









w/o DiffTell: Add a football to the mans hand.

GT: Add a tiger, Insert the tiger behind the man

<Image style change>





w/ DiffTell: Change the color of the

w/o DiffTell: Change the color of the sky to blue

GT: Change the color of the sky from orange to blue, color the sky blue

<Background change>





w/ DiffTell: Remove the background and leave just the cat. w/o DiffTell: Remove the background.

GT: Remove all background except for the cats face

<Text manipulation>





w/ DiffTell: Add the text Pool Party.

w/o DiffTell: Add text. GT: add text pool party

Figure 3. Visual comparison that illustrates the impact of utilizing the DiffTell dataset on Qwen2-VL-7B's performance across four distinct categories in the IER testing set. Our dataset demonstrates its effectiveness in enhancing performance, especially in local object description, text detection and recognition, spatial recognition, and image style description. The text in green shows an obviously precise expression over the text in red.









w/ DiffTell: Change the bird on the left to a hamster w/o DiffTell: Replace the bird with a hamster GT: The bird on the left has been replaced by a small

w/ DiffTell: Rem Idiers. Add a cartoor ter, Add big eyes to the helicopter w/o DiffTell: Add eyes to the helicopter, Add a cartoon character in front of the helicopter.

GT: The soldiers entering the helicopter were noved and replaced with a yellow creature

Figure 4. The visual comparison illustrates the impact of utilizing the DiffTell dataset on the Qwen2-VL-7B model's performance in the PSBattle dataset.

4.3. Ablation Study

Since *DiffTell* is a dataset with several subsets contributing to different image difference categories, it is necessary to

	B@4	M	С	R-L	S	$B-F_1$
IER	17.15	19.84	64.61	44.58	18.73	90.17
+ InstructP2P	16.52	21.11	66.23	46.00	19.90	90.28
+ OCR	16.02	21.09	62.00	46.00	20.63	90.27
+ MagicBrush	16.24	20.68	59.31	43.47	18.83	90.13
+ COCO	15.66	20.66	59.92	43.42	18.70	90.15
+ GIER	17.19	20.61	66.97	45.62	19.92	90.41
+ DiffTell	17.76	22.02	67.71	46.35	21.86	90.48

Table 4. Results of IER testing set from Qwen2-VL-7B model finetuned on different datasets.

study the contribution of each subset to the IDC performance. We consider two parts: the contribution of each subset to the general performance and the contribution of each subset to each category. We show the performance on the IER testing set from the Qwen2-VL-7B model finetuned with the IER training set and each subset in DiffTell in Table 4. Almost every subset can improve the performance, and in summary, the DiffTell can boost the performance further.

We show another ablation study on the category-wise contribution. To better study the performance of each category, we compute the statistics of the IER testing set based on the category given in Section 3.2. The statistics are given in Table 6 in the Appendix. Figure 2 provides an overview of the contributions based on the IER testing set and Qwen2-VL-7B of each subset in DiffTell to each category, regarding BLEU@4. Compared to the model trained exclusively on IER, the model trained on our subset derived from MARIO-10M shows a notable performance improvement, benefiting from the versatility of words in various real-life scenarios. Our subset derived from GIER contributes positively to overall performance, except for text manipulation, where no such data exists in the GIER dataset. The absence of background change data in the MagicBrush dataset leads to a performance decrease in the background change category. COCO, designed for local object changes, enhances performance in this category. In the InstructPix2Pix dataset, the lack of text data results in a performance decrease. In summary, the subset belonging to the specific categories can generally contribute to the corresponding categories in the IER testing set.

4.4. Automatic Data Filtering

The cost of manual data filtering can become a bottleneck when scaling up this dataset. To address this, we propose an alternative automatic data filtering pipeline, as shown in Figure 5. Using a dataset previously reviewed by humans, we compile both accepted and rejected samples as the training set for a binary classifier. The classifier's input consists of features extracted by the Owen2-VL-7B model, which has been fine-tuned on the IDC task. This classifier can assist annotators in more efficiently filtering the data.

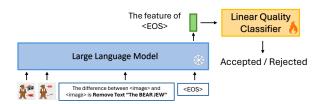


Figure 5. The framework of the automatic data filtering pipeline. The image pair and difference caption will be passed to the Open-Flamingo model, and the output feature of <EOS> token will be used for the classification of acceptance or rejection

Training Set	B@4	M	С	R	S	$B-F_1$
IER	17.15	19.84	64.61	44.58	18.73	90.17
IER + 10K(R)	17.13	19.99	64.31	44.60	18.85	90.17
IER + 10K (C)						
IER + 10K (H)	17.27	21.10	69.47	46.15	19.63	90.46

Table 5. The results of performance on IER testing set using the MARIO-10M data with automatic classifier or not.

To validate the effectiveness of our pipeline, we train a quality classifier on an annotator-validated subset of MARIO-10M, comprising 10K accepted and 10K rejected samples. We use an SVM as the classifier, splitting 16K samples for training and 4K for testing, achieving an accuracy of 89.47%. The classifier is then applied to unseen data from MARIO-10M, filtering 10K accepted samples. This unseen data is newly in-painted using FireFly Generative Fill, as explained in Section 3.3, and generation stops once 10K accepted samples are collected through the classifier. We compare the performance on IER dataset of the IDC model (Qwen2-VL-7B) trained on three subsets from MARIO-10M: 10K auto-filtered samples by classifier (C), 10K randomly selected samples (R), and 10K filtered samples by human (H). The randomly selected data is taken directly from the in-painted model without quality control, while the manually filtered data is a subset of MARIO-10M used in DiffTell. The results in Table 5 demonstrate that the auto-filtered training data can achieve much better performance than unfiltered data (random data), and be comparable to human filtered training data. Such a result shows the necessity of the filtering step in our designed pipeline and highlights the classifier's effectiveness and the potential for scaling data collection using this auto-filtering pipeline. More experiments and analysis are given in Appendix J.

4.5. Failure Cases

Although the model gains performance improvement in IDC, there are still some cases where the model fails to predict correctly. We illustrate the failure cases in Figure 6. The model may sometimes limit its predictions to local changes rather than providing a comprehensive description. In the first example shown in Figure 6, the model









w/ DiffTell: Swap their mouths.
w/o DiffTell: Make the man in the back smile.
GT: The faces of the two basketball players have been swapped.

w/ DiffTell: Add a hat and a gun to the bird on the right.

w/o DiffTell: Add a hat and a gun.
GT: Added Head hair in *left eagle* and cap and gun in the right one.

Figure 6. Illustration of the failure cases from the model trained with *DiffTell*. The examples are from Qwen2-VL-7B on PSBattle.

exclusively identifies the difference in the mouth from the face, neglecting the other facial elements and the relationship between the two faces. Although the model recognizes the change in the second example, it produces an inaccurate description. These shortcomings may result from the limited diversity in the dataset. A predominant portion of the images in *DiffTell* originates from real-life scenarios. The model struggles to capture surreal or fantastical compositions, such as facial swaps, as the training data may not adequately represent those instances. Following our methodology in creating *DiffTell*, incorporating more data sources covering a wider range of fine-grained domains may help the model to establish connections between objects and accurately identify specific object categories.

5. Conclusion and Limitation

In this study, we introduce DiffTell, an extensive and highquality dataset for image difference captioning (IDC). This dataset addresses the gaps in diversity and scale that were previously present in the IDC task. Through comprehensive experiments conducted on diverse testing sets and employing various baseline methods, we demonstrate the efficacy of our dataset in enhancing performance. Additionally, we analyze to understand the improvement contributed by each component of DiffTell to different image difference categories. We aspire that DiffTell will play a significant role in advancing the development of more sophisticated multi-modality models for IDC and language-guided image editing in the future. With DiffTell, the model capability in spatial relationship, tiny object recognition, long text OCR and complex change identifaction is enhanced. As for future work, we find the current model cannot identify the difference beyond the real images, like cartoons, science fiction, etc., due to the lack of training data. We hope to utilize the human-filtered data (acceptance and rejection) for preference optimization [42, 50] to boost the performance. We also plan to augment Chain-of-Thought [70] data into DiffTell, possibly by training a model with RL, enabling the thinking process for image difference prediction.

Acknowledgments

Zonglin Di is partially supported by the National Science Foundation (NSF) under grants IIS-2007951 and IIS-2143895.

References

- [1] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *European conference on computer vision*, pages 382–398. Springer, 2016. 5
- [2] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 18208–18218, 2022. 3
- [3] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An opensource framework for training large autoregressive visionlanguage models. arXiv preprint arXiv:2308.01390, 2023.
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv* preprint arXiv:2308.12966, 2023. 2
- [5] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization, pages 65–72, 2005. 5
- [6] Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. Introducing our multimodal models, 2023. 6
- [7] Alexander Black, Tu Bui, Simon Jenni, Viswanathan Swaminathan, and John Collomosse. Vpn: Video provenance network for robust content attribution. In *Proceedings of the 18th ACM SIGGRAPH European Conference on Visual Media Production*, pages 1–10, 2021. 1
- [8] Alexander Black, Jing Shi, Yifei Fan, Tu Bui, and John Collomosse. Vixen: Visual text comparison network for image difference captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 846–854, 2024. 2, 5, 16
- [9] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. 1, 2, 3, 4
- [10] 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. Advances in neural information processing systems, 33:1877–1901, 2020. 1, 4
- [11] Alexander Bukharin and Tuo Zhao. Data diversity matters for robust instruction tuning. arXiv preprint arXiv:2311.14736, 2023. 13

- [12] Jingye Chen, Yupan Huang, Tengchao Lv, Lei Cui, Qifeng Chen, and Furu Wei. Textdiffuser: Diffusion models as text painters. arXiv preprint arXiv:2305.10855, 2023. 2, 4
- [13] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. arXiv preprint arXiv:2306.15195, 2023. 2
- [14] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. arXiv preprint arXiv:2311.12793, 2023. 2
- [15] Maxwell Forbes, Christine Kaeser-Chen, Piyush Sharma, and Serge Belongie. Neural naturalist: Generating fine-grained image comparisons. arXiv preprint arXiv:1909.04101, 2019.
- [16] 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. 2
- [17] Zixin Guo, Tzu-Jui Julius Wang, and Jorma Laaksonen. Clip4idc: Clip for image difference captioning. *arXiv* preprint arXiv:2206.00629, 2022. 1, 6
- [18] 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. 1
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020.
- [20] Erdong Hu, Longteng Guo, Tongtian Yue, Zijia Zhao, Shuning Xue, and Jing Liu. Onediff: A generalist model for image difference captioning. In *Proceedings of the Asian Conference on Computer Vision*, pages 2439–2455, 2024. 1, 2, 3
- [21] Mude Hui, Siwei Yang, Bingchen Zhao, Yichun Shi, Heng Wang, Peng Wang, Yuyin Zhou, and Cihang Xie. Hq-edit: A high-quality dataset for instruction-based image editing. arXiv preprint arXiv:2404.09990, 2024. 3
- [22] Harsh Jhamtani and Taylor Berg-Kirkpatrick. Learning to describe differences between pairs of similar images. arXiv preprint arXiv:1808.10584, 2018. 1
- [23] 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. 2, 5, 12
- [24] Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max W.F. Ku, Qian Liu, and Wenhu Chen. Mantis: Interleaved multiimage instruction tuning. arXiv2405.01483, 2024. 3
- [25] Jungo Kasai, Keisuke Sakaguchi, Lavinia Dunagan, Jacob Morrison, Ronan Le Bras, Yejin Choi, and Noah A Smith. Transparent human evaluation for image captioning. arXiv preprint arXiv:2111.08940, 2021. 6
- [26] Hoeseong Kim, Jongseok Kim, Hyungseok Lee, Hyunsung Park, and Gunhee Kim. Viewpoint-agnostic change captioning with cycle consistency. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2095–2104, 2021. 5

- [27] Hoeseong Kim, Jongseok Kim, Hyungseok Lee, Hyunsung Park, and Gunhee Kim. Agnostic change captioning with cycle consistency. In *Proc. ICCV*, pages 2095–2104, 2021.
- [28] Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better understanding visionlanguage models: insights and future directions. In Workshop on Responsibly Building the Next Generation of Multimodal Foundational Models, 2024. 6
- [29] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. arXiv preprint arXiv:2408.03326, 2024. 3
- [30] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 2
- [31] 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. 2
- [32] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004. 5
- [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pages 740–755. Springer, 2014. 2, 5
- [34] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 2
- [35] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. arXiv preprint arXiv:2310.03744, 2023. 2
- [36] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 2,
- [37] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 2
- [38] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. arXiv preprint arXiv:2304.08485, 2023. 2
- [39] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 2
- [40] Zhaoyang Liu, Yinan He, Wenhai Wang, Weiyun Wang, Yi Wang, Shoufa Chen, Qinglong Zhang, Yang Yang, Qingyun Li, Jiashuo Yu, et al. Internchat: Solving vision-centric tasks by interacting with chatbots beyond language. arXiv preprint arXiv:2305.05662, 2023.
- [41] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting

- using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 3
- [42] Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *arXiv* preprint arXiv:2405.14734, 2024. 8
- [43] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 3
- [44] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the* Association for Computational Linguistics, pages 311–318, 2002. 5
- [45] Dong Huk Park, Trevor Darrell, and Anna Rohrbach. Robust change captioning. In *Proc. ICCV*, pages 4624–4633, 2019.
- [46] Dong Huk Park, Trevor Darrell, and Anna Rohrbach. Robust change captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4624–4633, 2019. 1, 2, 12
- [47] Ed Pizzi, Sreya Dutta Roy, Sugosh Nagavara Ravindra, Priya Goyal, and Matthijs Douze. A self-supervised descriptor for image copy detection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 14532–14542, 2022.
- [48] Yue Qiu, Shintaro Yamamoto, Kodai Nakashima, Ryota Suzuki, Kenji Iwata, Hirokatsu Kataoka, and Yutaka Satoh. Describing and localizing multiple changes with transformers. In *Proceedings of the IEEE/CVF International Confer*ence on Computer Vision, pages 1971–1980, 2021. 1
- [49] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [50] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. Advances in Neural Information Processing Systems, 36, 2024. 8
- [51] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [52] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 4
- [53] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 3, 4

- [54] 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. 3
- [55] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. Advances in Neural Information Processing Systems, 35:36479–36494, 2022. 3
- [56] Jing Shi, Ning Xu, Trung Bui, Franck Dernoncourt, Zheng Wen, and Chenliang Xu. A benchmark and baseline for language-driven image editing. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 4
- [57] 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. 2
- [58] 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. 2
- [59] Hao Tan, Franck Dernoncourt, Zhe Lin, Trung Bui, and Mohit Bansal. Expressing visual relationships via language. arXiv preprint arXiv:1906.07689, 2019. 1, 2, 5, 12
- [60] 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.
- [61] Yunbin Tu, Liang Li, Li Su, Junping Du, Ke Lu, and Qingming Huang. Adaptive representation disentanglement network for change captioning. *IEEE Transactions on Image Processing*, 2023.
- [62] 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. 6, 12
- [63] 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. 2
- [64] Yunbin Tu, Liang Li, Li Su, Ke Lu, and Qingming Huang. Neighborhood contrastive transformer for change captioning, 2023. 2, 6
- [65] Yunbin Tu, Liang Li, Li Su, Ke Lu, and Qingming Huang. Neighborhood contrastive transformer for change captioning. *IEEE Transactions on Multimedia*, pages 1–12, 2023. 1, 12
- [66] Yunbin Tu, Liang Li, Li Su, Zheng-Jun Zha, Chenggang Yan, and Qingming Huang. Self-supervised cross-view representation reconstruction for change captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2805–2815, 2023. 3
- [67] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evalua-

- tion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015. 5
- [68] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024. 6
- [69] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. arXiv preprint arXiv:2311.03079, 2023. 2
- [70] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022. 8
- [71] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. ACM Computing Surveys, 56(4): 1–39, 2023. 4
- [72] Linli Yao, Weiying Wang, and Qin Jin. Image difference captioning with pre-training and contrastive learning. In *Proc. AAAI*, pages 3108–3116, 2022.
- [73] Linli Yao, Weiying Wang, and Qin Jin. Image difference captioning with pre-training and contrastive learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3108–3116, 2022. 1
- [74] Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. arXiv preprint arXiv:2306.10012, 2023. 1, 2, 4
- [75] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068, 2022. 2
- [76] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675, 2019. 5
- [77] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592, 2023. 2

DiffTell: A High-Quality Dataset for Describing Image Manipulation Changes

Supplementary Material

Appendix

The supplementary material is composed as follows.

- Appendix A illustrates the motivation for designing the different captioning models for content authenticity.
- Appendix B presents a detailed description of the existing datasets in IDC.
- Appendix C gives the implementation details.
- Appendix D presents more baselines that are not included in the main paper.
- Appendix E presents the details using various instruction prompts.
- Appendix F presents the more results of SPICE and BertScore- F_1 as well as the human evaluation.
- Appendix G presents the out-of-distribution (OOD) results using Spot-the-Diff and CLEVR-DC datasets.
- Appendix H presents the zero-shot or few-shot performance on LLMs without being finetuned on IER testing set.
- Appendix I discusses more about the dataset collection.
- Appendix J presents more experiments and analysis about the automatic data filtering pipeline.
- Appendix K provides more details about PSBattle testing set
- Appendix L provides more visual results, including how we filter the data, more successful cases, etc.

A. Motivation from Content Provenance and Authenticity

Figure 7 shows the motivation to design the different captioning models to ensure the authenticity of the content, which helps users decide where the image is from and what is modified over its original version.

B. Existing Datasets

The most commonly used datasets in the IDC task are CLEVR change [46], Spot-the-Diff [23], and Image Editing Request (IER) [59]. CLEVR change constitutes a sizable synthetic dataset characterized by moderate viewpoint variations. Spot-the-Diff is composed of pairs of frames extracted from video surveillance footage and the corresponding textual descriptions of visual changes. IER is crawled from the practical image editing requests from the Reddit channel, consisting of 3,939 pairs of real images, accompanied by 5,695 editing instructions. Each image pair in the training set is associated with one instruction. In contrast, each image pair is linked to three instructions for a more objective evaluation in the validation and testing sets.

Because IER is collected from a real-world scenario, it covers more image difference categories, such as background change, text manipulation, and local object change. The definition of the image difference categories can be found in Section 3.2. Due to the single domain in CLEVR and Spotthe-Diff datasets, we mainly use IER in this work as the testing set, which aligns our scope to have a comprehensive, diverse, and practical dataset. For these two datasets, which are mainly about a single domain, we use them as out-of-distribution evaluation, which is given in Appendix G.

Category	Background	Text	Local object	Image style
Number of Images	117	53	277	223

Table 6. Statistics of each image difference category in the IER testing set.

C. Implementation Details

C.1. Training Details

For CLIP4IDC, We adopt the official implementation of CLIP4IDC. However, as it lacks the training script and the pretrained weights for IER, we reproduce the CLIP4IDC⁴ model trained on IER exactly following its provided training hyper-parameter settings of the CLEVR dataset. For VARD-LSTM⁵ and NCT⁶, there is still no official implementation for IER and we reproduce them using the settings in CLEVR dataset. The pre-trained Biaffine Parse in NCT we use is from Diaparser⁷. For OpenFlamingo-3B, the vision encoder and language encoder are ViT-L-14 and anas-awadalla/mpt-1b-redpajama-200b.

The cross attention interval is 1. For LLaVA-Interleave-7B, the language model we use is meta-llama/Meta-Llama-3-8B-Instruct.

For Fuyu-8B, we use adept/fuyu-8b. The training platform we use is 8 NVIDIA A100s with the 80GB GPU memory. The training epochs is 2 for the MLLMs. For the other hyper-parameters like learning rate, weight decay, batch size, please refer to Table 7.

D. The Performance of More Baselines

Besides the methods in the main text, we test more baselines including NCT [65] and VARD-LSTM [62] given in Table 8

 $^{^{\}bf 4} \texttt{https://github.com/sushizixin/CLIP4IDC}$

⁵https://github.com/tuyunbin/VARD

⁶https://github.com/tuyunbin/NCT

⁷https://github.com/Unipisa/diaparser

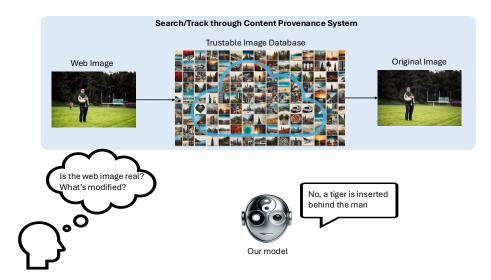


Figure 7. The overview of difference captioning for content provenance and authenticity. When the user asks if the web image is trustworthy, the content provenance system can search for its original version in the trustworthy image database if it exists. Our model is able to tell the user what is modified based on image difference captioning on the web image and the original one.

VL Model	Learning Rate	Epochs	Warmup Ratio	Weight Decay	Batch Sizes
OpenFlamingo-3B	5e-6	2	0.03	0.01	128
Fuyu-8B	1e-5	2	0.03	0.01	16
LLaVA-Interleave-7B	1e-5	2	0.03	0.0	16
Idefics3	5e-6	2	0.03	0.01	16
Qwen2-VL-7B-Instruct	5e-6	2	0.03	0.01	16

Table 7. The hyperparameters we use in this paper

Testing Set	Method	w/ DiffTell	BLEU@4	METEOR	CIDEr	ROUGE-L
IER	NCT	X ./	1.64 1.94	7.97 9.63	7.47 7.58	19.40 23.79
IER	VARD-LSTM	× ✓	1.60 1.71	8.06 8.54	5.49 6.02	18.87 20.08
PSBattle	NCT	X /	2.78e-08 1.65e-06	0.73 1.22	1.12 3.11	4.53 9.78
PSBattle	VARD-LSTM	X /	1.49e-08 7.46e-07	0.43 0.88	1.56 2.07	7.01 7.79

Table 8. The comparison of the methods fine-tuned on image editing request (IER) training set with and without *DiffTell* using more baselines.

E. The Experiments with Diverse Prompts

In instruction tuning, incorporating diverse prompts enhances model robustness, making them more adaptable and better at generating accurate responses across varying contexts [11]. Initially, we use a uniform prompt "What is the difference between two images?" across all datasets and ask the model to provide an answer. To ablate this, we expand the prompt into nine different variations and compare the

performance against the single-prompt approach, as shown in Table 9. The nine prompts we use are as follows. The model we use is OpenFlamingo-3B. As a complex vision-language task, it is more important for the model to understand two images, identify the difference and express the answer. Thus, to improve the vision encoder could be more useful.

• Please tell me the editing instruction of how to edit < | image | > to look like < | image | >.

Testing Set	w/ DiffTell	Diverse Prompt	BLEU@4	METEOR	CIDEr	ROUGE-L
	Х	×	4.45	14.87	15.80	29.79
IER	1	×	6.49	16.68	21.04	31.36
	✓	✓	6.32	16.59	23.88	30.34
	X	X	2.35e-04	2.33	7.71	19.24
PSBattle	✓	X	2.12	6.60	4.02	16.10
	✓	✓	1.77	6.45	4.48	16.46

Table 9. The results of performance on the IER testing set using the diverse prompts. The model we use is OpenFlamingo-3B.

Model		Idefics3-8B		Qwen2-VL-7B			
w/ DiffTell	Fluency (†)	Correctness (†)	Relevance (†)	Fluency (†)	Correctness (†)	Relevance (†)	
×	3.49	3.59	3.54	3.71	3.68	3.55	
	3.53	3.66	3.67	3.84	3.79	3.60	

Table 10. The results of human evaluation on Idefics3-8B and Qwen2-VL-7B. The data we use is the 50 pairs of images in the IER testing set. Each pair of images is labeled by 5 persons.

- Identify the transformations applied to <|image|> to achieve the appearance of <|image|>.
- Outline the steps required to edit < | image | > so that it matches the look of | image | >.
- Explain the edits necessary to convert < | image | > into < | image | >.
- What alterations were made to < | image | > to create < | image | >?
- Detail the changes from < | image | > to < | image | >.
- < | image | > is image1, < | image | > is image2, tell me what the change is between these two images.
- < | image | > is image1, < | image | > is image2, tell me what the change is from image1 to image2.

F. The Performance of Additional Evaluation Metrics

As mentioned in Section 4.1, we also evaluate the performance using SPICE and BertScore- F_1 . The results are presented in Table 11.

For human evaluation, we randomly select 10% of the IER testing set (50 pairs of images) and let humans evaluate the output of Idefics3-8B and Qwen2-VL-7B trained with and without *DiffTell*. We use Amazon Mechanical Turk (MTurk) for this user study with each caption evaluated by 5 persons, ranging from 1 (worst) to 5 (best), resulting in 500 samples. The results shown in Table 10 demonstrate the effectiveness of *DiffTell*.

G. Out-of-Distribution Results

To evaluate the generalization capability of *DiffTell*, we test the model trained with and without *DiffTell* dataset on Spotthe-Diff and CLEVR-DC dataset without training on these two datasets. The results of Spot-the-Diff and CLEVR-DC using Qwen2-VL-7B are given in Table 12, showing that *DiffTell* can boost the performance of out-of-distribution (OOD) data, which is a good proof of its comprehensiveness

H. Zero-shot/Few-shot Prompt Results

Investigating the potential of zero-shot learning is essential for methods utilizing LLM. For few-shot prompt testing, we randomly choose three examples from the IER training set. Performance results on the PSBattle dataset are not reported due to the absence of training data in that specific dataset. The detailed results can be found in Table 13. The few-shot prompt example is shown in Figure 8. The results show that image difference caption (IDC) is a hard task for the current LLMs, although they are trained on a huge amount of data. Even with few-shot prompt, the results are still not satisfying.

I. Dataset Collection Details

Image In-painting We use FireFly Generative Fill to inpaint the image. The inputs we can provide are the original image and the prompt for the generative model. There is no need for us to select the parameters. The illustration is given in Figure 11. We generate I_2 for COCO and MARIO-10M subsets in *DiffTell*.

Data Filtering The illustration of how the annotators filter the data is given in Figure 12, 13, and 14, which are for InstructPix2Pix, COCO, and MARIO-10M subsets, respectively. For InstructPix2Pix, the annotators filter whether the T_{I_1,I_2} matches (I_1,I_2) or whether the change reflects on I_1 and I_2 because (I_1,I_2,T_{I_1,I_2}) has already been provided.

Testing Set		IER Testing Set				PSBattle			
w/ DiffTell	Х	✓	×	✓	×	✓	X	✓	
Models	SPICE		BertSc	BertScore- F_1		SPICE		BertScore- F_1	
OpenFlamingo-3B	9.31	10.07	87.13	87.80	3.77	5.15	85.50	87.51	
Fuyu-8B	11.34	16.07	87.57	88.95	5.25	6.60	84.61	87.56	
Llava-Interleave-7B	12.73	16.86	88.51	89.20	8.79	8.88	86.62	87.67	
Idefics-8B	19.49	22.64	89.54	90.49	11.72	13.37	86.35	87.45	
Qwen2-VL-7B	18.73	21.86	90.17	90.48	11.22	11.78	87.13	87.64	

Table 11. The results of SPICE and BertScore- F_1 corresponding to Table 3 (main experiments) in the main paper.

Testing Set	w/ DiffTell	BLEU@4	METEOR	CIDEr	ROUGE-L	SPICE	BertScore- F_1
Spot-the-Diff	X	2.51	4.88	4.12	11.45	6.51	85.34
	✓	5.86	5.60	9.24	16.03	6.94	86.01
CLEVR-DC	X	0.99	2.53	2.27	6.71	17.07	85.71
	✓	6.50	4.72	7.66	17.27	21.94	86.26

Table 12. The results of out-of-distribution evaluation on Spot-the-Diff and CLEVR-DC datasets.

For COCO and MARIO-10M only providing I_1 , the annotators filter whether the object or the text is successfully inpainted from I_1 .

J. Extensive Experiments and Analysis on the Automatic Filtering

We further evaluate our automatic filter pipeline in the InstructPix2Pix dataset, using Qwen2-VL-7B as the feature extractor, in the same manner described in Section 4.4 in

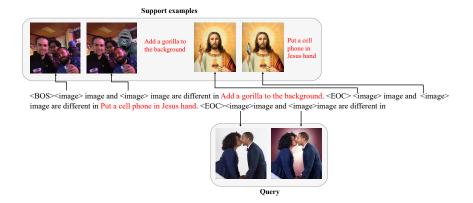


Figure 8. The example of how we construct the few-shot prompt.

Method	Few-shot	BLEU@4	METEOR	CIDEr	ROUGE-L
OpenFlamingo-3B	X	1.18	8.07	8.72	16.63
	✓	0.84	7.64	4.09	17.54
OpenFlamingo-9B	×	1.15	8.26	6.04	19.00
	✓	1.99	9.18	5.01	20.93

Table 13. The results of zero-shot or few-shot prompt results on the IER testing set. The few-shot prompt is the composition of 3 training examples from the training set.

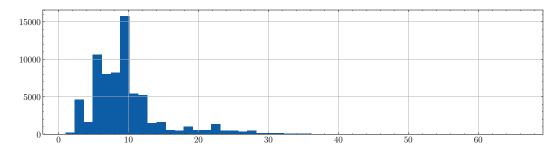


Figure 9. The difference description length distribution in DiffTell

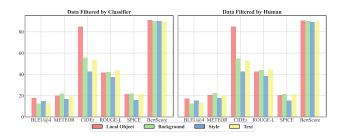


Figure 10. Comparison between data filtered by the classifier and human in terms of all the metrics.

the main paper. The accuracy of the data filter classifiers for InstructPix2Pix is 87.35%, respectively, indicating that the classifier has a satisfactory accuracy to keep clean data from a noisy dataset. The filtered data is used to train the Qwen2-VL-7B IDC model in Table 14. Due to the time limit, we compared the results with and without the automatically filtered 10K data from InstructPix2Pix in the last column of Table 14, indicating the **effectiveness** of the autofiltered InstructPix2Pix data and the generalization of our automatic data filtering pipeline.

We further categorize the data into different types of differences and analyze the data filter accuracy on each type of difference, so we can analyze potential bias. We also conduct the **trade-off analysis** between the human and automatic data filtering using InstructPix2Pix data, with the result displayed in Figure 10. However, we do not observe an obvious performance trade-off in each difference type.

Training Set	B@4	M	С	R-L	S	$B-F_1$
IER IER + 10K (R) IER + 10K (C)	17.15	19.84	64.61	44.58	18.73	90.17
IER + 10K(R)	17.20	19.66	63.63	45.02	19.53	89.95
IER + 10K(C)	17.53	20.41	66.90	45.15	19.69	90.34
IER + 10K (H)	17.97	20.67	68.02	45.25	19.70	90.46

Table 14. The results of data filtering on the InstructPix2Pix dataset. R, C, and H refer to random data, data filtered by the classifier, and human filtering.

K. PSBattle Dataset

The PSBattle dataset is another practical dataset used in [8] that consists of images edited in Adobe PhotoshopTM and is curated from the "Photoshopbattles" subreddit. We include this dataset only for the evaluation of out-of-domain data to test the generalizability of the models. This dataset comprises over 10,000 images, each paired with several modified variants generated according to editing instructions provided by users. In total, there are 102,208 variants created by 31,000 different artists. For our study, we randomly selected 100 image pairs, each accompanied by three captions obtained through crowd-sourced annotation on MTurk. The illustration of PSBattle dataset is shown in Figure 15.

L. More Visual Results

L.1. Failure Cases in Data Filtering

As mentioned in Section 3.3, we present the importance of the data filtering by showing more cases in InstructP2P, COCO, and MARIO-10M datasets in Figure 16, 17, and 18, respectively.

L.2. More Successful Cases

To better illustrate the improvement from *DiffTell*, we select another two prediction results in the IER testing set from the four categories, respectively, shown in Figure 19. The model we use is Qwen2-VL-7B.







(b) The image after in-painting

Figure 11. We in-paint the image using Firefly Generative Fill in Photoshop. For each image, we provide the original image (I_1) and the corresponding mask. The mask is used to identify the selected area shown with the red arrow. We use a prompt to ask Firefly to in-paint the image and fit the background. Normally, the Firefly will return 3 to 4 in-painted images.

Instructions: Given an input image, the output image and the editting instruction. The meanings of the terms are as follows:

- Input Image: The original image we want to edit.
- Output Image: The image generated by AI model based on the input image.
- · Editting Instruction: The instruction used to guide the AI model to generate the output image from the input image.

You are supposed to evaluate whether the ouput image is matched with the input image and the editting instruction. After carefully check the images and the instruction, you should select the quality score for the output image. Please check the Yes for the successful editting while No for the unacceptable editting.



input image

output image

Editting Instruction: make the creek dry

Figure 12. The labeling illustration of InstructPix2Pix subsets. The two images are I_1 and I_2 . T_{I_1,I_2} is given in **Editing Instruction**. The annotator is asked to identify whether the T_{I_1,I_2} matches (I_1,I_2) or whether the change reflects on I_1 and I_2 and give the answer "Yes" or "No". We keep those which are identified as "Yes".

Instructions: Given an input image, the input mask and a object-free image. The meanings of these 3 images are as follows:

- · Input Image: The original image we want to remove the object.
- Input Mask: The region of the object generated by AI model. Ideally the mask should cover the object we want to remove.
- <object>-free Image: The image processed by Al model. Ideally, there should not exist the <object> covered the mask and no extra element should be added. The <object> here is a placeholder which be will replaced by a specific object word.

You are supposed to evaluate the object-free image, whether the object is fully removed without changing the original image content. After carefully compare the object-free image and the input image, you should select the quality score for how well the object is removed. We set 2 levels regarding the quality of the object-free image which are:

- · Acceptable
- Unacceptable

The detailed criterion for the 2 categories and the corresponding example are given in the instruction document.



input image input mask airplane-free image

Figure 13. The labeling illustration of COCO subsets. From left to right, the first, second, and third images are the original image (I_1) , the input mask, and the in-painted image. We provide the input mask and object name to remind the annotator which area to focus on. The annotator selects "Acceptable" and "Unacceptable". We keep those which are identified as "Acceptable".

Instructions: Given an input image, the input mask and a text-free image. The meanings of these 3 images are as follows:

- Input Image: The original image we want to remove the text.
- Input Mask: The region of the text generated by AI model. Ideally the mask should cover all the text.
- · Text-free Image: The image processed by AI model. Ideally, there should not exist text and no extra element should be added.

You are supposed to evaluate the text-free image, whether the text is fully removed without changing the original image content. After carefully compare the mask-free image and the input image, you should select the quality score for how well the text is removed. We set 2 levels regarding the quality of the object-free image which are:

- Acceptable
- Unacceptable

The detailed criterion for the 2 categories and the corresponding example are given in the instruction document.



input image input mask text-free image

Figure 14. The labeling illustration of MARIO-10M subsets. From left to right, the first, second, and third images are the original image (I_1) , the input mask, and the in-painted image. We provide the input mask and object name to remind the annotator which area to focus on. The annotator selects "Acceptable" and "Unacceptable". We keep those which are identified as "Acceptable".





- · hanging person added
- The right image has a person hanging off the end of the track with a horrified expression on his face.
- On the right, a man is clinging to the bomb bay door, about to fall. He is not there at all on the left.





- A new face has been given to batman. I think it is the face of Will Ferral.
- The mask only covers part of the face and the man wears glasses now.
- Batman has been given a bushy head of hair and a large pair of glasses.





- In the right picture the gun is visible
- Added Head hair in left eagle and cap and gun in the left one.
- Hawks are fighting each others in second one Hawk kept machine gun.





- The hippo is wearing a cross and holding a bible.
- The hippo is now carrying a bible and a crucifix necklace.
- The hippo is holding a bible and a crucifix in one of its hooves.

Figure 15. Four examples in PSBattle dataset.



















Add a river. as a carto

Figure 16. Within the InstructP2P dataset, we have identified four sets of images, each composed of the original image, the altered image, and the corresponding instruction. All four of these image sets represent instances of failure. In the first pair of images, not only is the mural altered as per the instruction, but there are also changes to the face of the person in white and the text on the wall. The second pair exhibits subtle changes that are unrelated to the provided instruction. For the third pair, the images undergo significant alterations, including the addition of a river, surpassing the intended modifications. In the fourth pair, the changes between the two images fail to accurately reflect the given instruction. The InstructP2P dataset is characterized by a high noise ratio, leading to a low acceptance rate of 35.13% during manual filtering.



Figure 17. We choose four sets of images in the COCO dataset, each comprising the original image, the dilated mask, and the in-painted image. The initial two sets depict instances of failure, whereas the latter two sets showcase successful outcomes. The initial failure occurs when the mask fails to adequately cover the object, and the second failure is attributed to the inadvertent addition of another object despite the mask effectively covering the intended object. The labeling team is instructed to exclude images falling into *DiffTell*.



Figure 18. We select four image sets from the MAIRO-10M dataset, each including the original image, the dilated mask, and the in-painted image. All four of these cases have been deemed failures and subsequently excluded by the labeling team. The mask in the MARIO-10M dataset is not very precise. All of these 4 image sets show this issue. In the first image set, the text is not very clear, either. Besides the inadequate mask and addition objects which exist in the COCO dataset, another issue of MARIO-10M dataset is the existence of low-quality images.

Local Object







w/ DiffTell: Remove all the writing from the image.
W/o DiffTell: Remove all the text .
GT: Remove the background with all the signatures

w/ DiffTell: Add masks to all three people.W/o DiffTell: Add masks.GT: Add mask to their faces.

Background









w/ DiffTell: Lighten the entire image.W/o DiffTell: Remove the collar.GT: lighten backgound

 $\ensuremath{\mathbf{w}}\xspace/\ensuremath{\mathbf{DiffTell:}}$ Crop the photo and make look clear.

W/o DiffTell: Remove the background . GT: Make picture clear'

Image Style







w/ DiffTell: Rotate the image 90 degrees crop the image.
W/o DiffTell: Make the drawing more detailed and darker.

GT: Rotate image counter clockwise

w/ DiffTell: Crop the image .W/o DiffTell: Remove the background.GT: Crop this image and make it lighter

Text Editing









w/ DiffTell: change the background to black add the text IM SORRY CANT HEAR YOU OVER THE SOUND OF HOW AWESOME AM

 $\ensuremath{\text{w/o}}$ DiffTell: change the background to black add a yellow frame add text.

GT: add a black background to the whole image, crop the image, add the light on camera, add the whole image to the frame, add the text under the frame. w/ DiffTell: remove the text URGENT MESSAGE WANTED APPROACH WITH CAUTION GENDER MALE WEAPON SIGNAL JAMMER..

W/o DiffTell: remove the wanted poster.

GT: Crop out the man with the camera and add him into a different background

Figure 19. More examples from the IER testing dataset regarding the four categories from Qwen2-VL-7B.