# AutoSplice: A Text-prompt Manipulated Image Dataset for Media Forensics

Shan Jia　　　Mingzhen Huang　　　Zhou Zhou　　　Yan Ju　　　Jialing Cai　　　Siwei Lyu

University at Buffalo, State University of New York, NY, USA

`{shanjia, mhuang3, zzhou38, yanju, jialingc, siweilyu}@buffalo.edu`

## Abstract

*Recent advancements in language-image models have led to the development of highly realistic images that can be generated from textual descriptions. However, the increased visual quality of these generated images poses a potential threat to the field of media forensics. This paper aims to investigate the level of challenge that language-image generation models pose to media forensics. To achieve this, we propose a new approach that leverages the DALL-E2 language-image model to automatically generate and splice masked regions guided by a text prompt. To ensure the creation of realistic manipulations, we have designed an annotation platform with human checking to verify reasonable text prompts. This approach has resulted in the creation of a new image dataset called* AutoSplice, *containing 5,894 manipulated and authentic images. Specifically, we have generated a total of 3,621 images by locally or globally manipulating real-world image-caption pairs, which we believe will provide a valuable resource for developing generalized detection methods in this area [1]. The dataset is evaluated under two media forensic tasks: forgery detection and localization. Our extensive experiments show that most media forensic models struggle to detect the AutoSplice dataset as an unseen manipulation. However, when fine-tuned models are used, they exhibit improved performance in both tasks.*

## 1. Introduction

The proliferation of digital media and AI technology have made it easier to manipulate and fabricate digital content. In recent years, the rapid development of powerful deep generative models, such as Variational Autoencoders (VAEs) [14, 35], Generative Adversarial Networks (GAN) [32, 33], diffusion-based models [16, 56], and the latest large-scale language-image (LLI) models [22, 53–55, 73], have brought new challenges to the authentication of digital media. The generated images have become increas-

---

[1] The AutoSplice dataset is available from https://github.com/shanface33/AutoSplice_Dataset
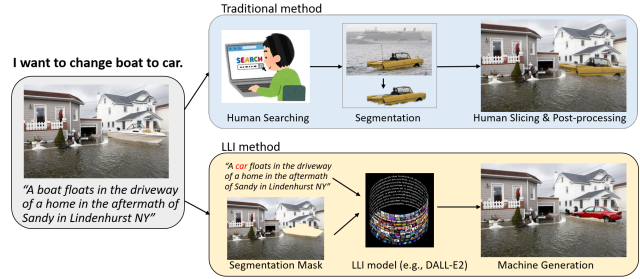


Figure 1. Comparison of our text-prompt-based image manipulation pipeline and traditional manual pipeline.

ingly realistic and convincing so that it can be difficult for human eyes to discern as artificial [59, 67].

Many efforts [29, 63, 66, 67] have investigated the challenges of GAN-generated images to media forensics due to the surge of GAN generation models, such as ProGAN [30], BigGAN [3], StyleGAN [32], etc. In addition, synthetic images from diffusion models can be identified with high accuracy using similar diffusion models in training, as demonstrated by [9]. The majority of previous studies have focused on entire synthesis using generative models. However, with the emergence of large-scale language-image models, local image manipulations guided by text prompts have become more accessible. Local manipulation of image regions, regardless of size, tends to be more realistic and challenging to detect than entire synthesis. In contrast to local manipulation techniques such as manual copy-move, slicing, or inpainting as described in [64], the latest language and vision models such as DALL-E [53, 54] allow for fully automatic and realistic local edits to images guided by a text prompt while preserving semantic information and stylistic elements, as shown in Figure 1. These models offer improved generation efficiency, image quality, and content flexibility compared to traditional local manipulation methods, and potentially revolutionize image manipulation.

A key question we investigate in this paper is: how much threat does the state-of-the-art language-image model pose to the current media forensic techniques? To explore this question, we first create a novel text-prompt manipulated

image dataset using DALL-E2 [2]. To create high-quality manipulations, we designed a semi-automatic annotation platform with manual checking. To ensure a diverse set of real-world media data, we utilized caption-image pairs from the Visual News dataset [43] as source data. Using automatic media analysis tools and human annotations, we extracted potential object regions to be manipulated from the images and parsed and replaced the corresponding text prompts in the captions. These captions were then taken as input to the DALL-E2 model for local image manipulation. For each caption-image pair input, we generated manipulated images (three by DALL-E2) and the corresponding manipulation mask. After data cleaning, our AutoSplice dataset contains 3,621 manipulated images and 2,273 authentic images.

Our dataset has several advantages over existing relevant media forensics datasets or methods, including high flexibility in content generation, high diversity in manipulation region, and good and reasonable generation quality. We use the large-scale language-image model, DALL-E2, for automatic local manipulation to create realistic forgery images, unlike previous local image manipulation methods based on manual and random object copy-move, slicing, or inpainting. Additionally, we only partially manipulate the image guided by the input region mask instead of generating the entire image, unlike recent semantic editing tools that use diffusion-based LLI models such as DiffEdit [10], Prompt-to-prompt [25], and Imagic [34]. Thanks to the powerful DALL-E2 generation model and human annotations, our dataset contains highly diverse and realistic synthesized images. We evaluate two media forensics tasks, namely, image forgery detection and image manipulation localization, on the AutoSplice dataset with lossless and lossy compression. Results show that pre-trained methods have limited generalization ability and unreliable prediction in detecting AutoSplice images. Models with fine-tuning on the dataset achieve improved performance, but also show vulnerability to compression.

## 2. Related Work

### 2.1. LLI Synthesis Models

Recent advancements in attention-based transformer and diffusion models have significantly improved text-to-image generation in the past two years. Several large-scale language-image models have been developed. The DALL-E model [54], proposed by OpenAI in 2021, uses an autoregressive transformer to achieve high-quality image generation on the MS-COCO dataset [40] without using any training labels. Other models, such as CogView [17], Parti [70], and Make-A-Scene [22], have also trained autoregressive transformer models on text and image tokens for text-to-image generation. In 2022, an updated version of DALL-E,

DALL-E2 [53], was developed using a diffusion model with CLIP image embeddings, making it computationally more efficient and able to produce higher-quality and more diverse samples. Other models, such as GLIDE [50], Stable-Diffusion [55], and Imagen [57], have also used diffusion models to improve text-to-image synthesis. Inspired by these powerful LLI synthesis models, several studies have developed text-guided image editing models, including DiffEdit [10], Prompt-to-prompt [25], Null-text Inversion [47], Imagic [34], and Muse [5]. These models apply local semantic editing to an image given a text input (with the desired edit) and an optional scene layout (segmentation map). However, their optimization tends to maximize the similarity to the original image while maintaining the ability to perform meaningful editing on local regions. This kind of entire synthesis can be easily identified if seen in training data [9]. To create more challenging fake media, we utilize the DALL-E2 model with high-quality local image editing techniques, which can generate text-guided pixels only in erased image regions.

### 2.2. Image Forensic Datasets

Two types of fake media datasets are relevant to our work: AI-synthesized image datasets and local image manipulation datasets. Several large-scale AI-synthesized image datasets have been collected from various GAN and VAE models, including DFFD [13] with GAN-based face attribute manipulations and entire face synthesis, CNNDetection [66] created using 11 CNN-based image generators, $DF^3$ [28] with entire face generation from six generation models (i.e., StyleGAN2 [33], StyleGAN3 [31], 3DGAN [4], Taming Transformers [20], LSGM [62], and Stable Diffusion [55]), and DMimageDetection [9] with diverse images from different GAN and diffusion models. The recent diffusion-based generation models are a class of likelihood-based models [51], which perturb data through successive addition of Gaussian noise and learn to recover the data by reversing this noising process. Although diffusion-based models achieve superior generation quality to GAN models [48, 60], a study [9] showed that current GAN image detectors can achieve near-perfect detection on similar diffusion models when trained with images generated with the diffusion models.

For local image manipulation, current research primarily focuses on techniques such as image slicing and copy-move. These methods involve copying and pasting specific regions of an image onto another part of the same image or a different one. Several widely used datasets, including Columbia [26, 49], CASIA [19], NIST16 [1], Coverage [68], Realistic Tampering [37], and IMD2020 [52], offer various types of locally manipulated images that are created either by manual operations or random slicing. However, these datasets have certain limitations, such as small

Table 1. Summary of previous image manipulation datasets and our work.

| Dataset | Year | # Forged Image | # Authentic Image | Image Size | Format | Manipulation Method |
|---|---|---|---|---|---|---|
| Columbia [49] | 2004 | 912 | 933 | $128 \times 128$ | BMP | Random |
| Columbia [26] | 2006 | 180 | 183 | $757 \times 568$ - $1152 \times 768$ | TIF | Random |
| NIST16 [1] | 2016 | 564 | 875 | $500 \times 500$ - $5616 \times 3744$ | JPEG | Manual |
| CASIA v1 [19] | 2013 | 921 | 800 | $374 \times 256$ | JPEG | Manual |
| CASIA v2 [19] | 2013 | 5,123 | 7,200 | $320 \times 240$ - $800 \times 600$ | JPEG, BMP, TIF | Manual |
| Coverage [68] | 2016 | 100 | 100 | $400 \times 486^*$ | TIF | Manual |
| Realistic Tampering [37] | 2016 | 220 | 220 | $1920 \times 1080$ | TIF | Manual |
| IMD2020 [52] | 2020 | 2,010 | 414 | $1062 \times 866^*$ | JPEG, PNG | Collected from Internet |
| **AutoSplice (ours)** | **2023** | **3,621** | **2,273** | **$256 \times 256$ - $4232 \times 4232$** | **JPEG** | **LLI model** |

\* Using the average image size.

data size (e.g., Columbia, Coverage, and NIST16), low authenticity level (e.g., Columbia's random region copy-move), or low flexibility and efficiency in generation due to careful and manual operations (e.g., CASIA and Realistic Tampering).

Recent advances in large-scale language-image models have demonstrated remarkable abilities in text-guided image manipulation and generation. Leveraging the power of these models, we introduce AutoSplice, a text-prompt guided image manipulation dataset that is built using the DALL-E2 model for automatic image editing. Table 1 provides further details regarding the existing image local manipulation datasets and our AutoSplice dataset.

## 2.3. Image Forgery Detection

The advancement of image forgery detection methods is a critical step in identifying manipulated/synthetic images for media forensics. Deep learning techniques have become increasingly popular for designing effective image forgery detectors. Most approaches consider forgery detection as a binary classification task and utilize well-designed deep neural networks to learn discriminative features automatically. Studies in this area can be divided into two categories: image-level forgery detection and pixel-level forgery detection (i.e., localization).

The former category concentrates on extracting global artifacts that synthesis models leave on the entire image, such as using augmented CNN models [66], frequency analysis [21], and re-synthesis residuals [24]. To improve the generalization ability for local image manipulation detection, recent studies [15, 28, 29, 41, 74] demonstrate the effectiveness of fusing local and global features in detecting different types of image forgery.

For manipulation localization, which aims to identify modifed image regions at the pixel level, existing methods mainly focus on identifying image tampering involving copy-move, splicing, removal, and faceswap. A subset of techniques formulates this task as a local anomaly detection problem and designs methods for capturing anomalies [6, 11, 69]. Several methods [2, 38, 39, 42, 45] uti-lize compression artifacts for forgery detection considering that the manipulation often involves double or more times of compressions. In addition, a branch of methods explores distinctive noise patterns in forged images, such as the RGB-N model [71] fusing RGB image content and image noise features, Noiseprint model with Siamese Network [12], [36], and ViT-VAE [6] to combine Noiseprint, High-pass filtering residuals, and Laplacian edge maps using Vision Transformer (ViT).

Given that the latest LLI models offer a powerful tool for realistic image generation and manipulation, presenting a potential threat of their misuse for spreading disinformation, our goal is to investigate the performance of current detectors on the newly created synthetic images. We will also pay particular attention to the detectors' ability in challenging social-network scenarios with image resizing and compression.

## 3. AutoSplice Dataset

To evaluate the difficulty of detecting media generated by recent LLI models using current forensic techniques, we introduce a new dataset called AutoSplice. In this section, we provide a comprehensive overview of the AutoSplice dataset creation process. The entire generation pipeline is illustrated in Figure 2. We begin by outlining the data pre-processing techniques and human annotation process. Next, we describe the data cleaning process, followed by a summary of the dataset and an analysis of its statistics.

## 3.1. DALL-E2

DALL-E2 [53] is a generative model that uses multiple modes to create synthetic images based on given text inputs. The model utilizes diffusion models to produce images based on CLIP image embeddings. Unlike traditional image synthesis approaches, DALL-E2 can also perform image inpainting by using both the input text and a region mask.
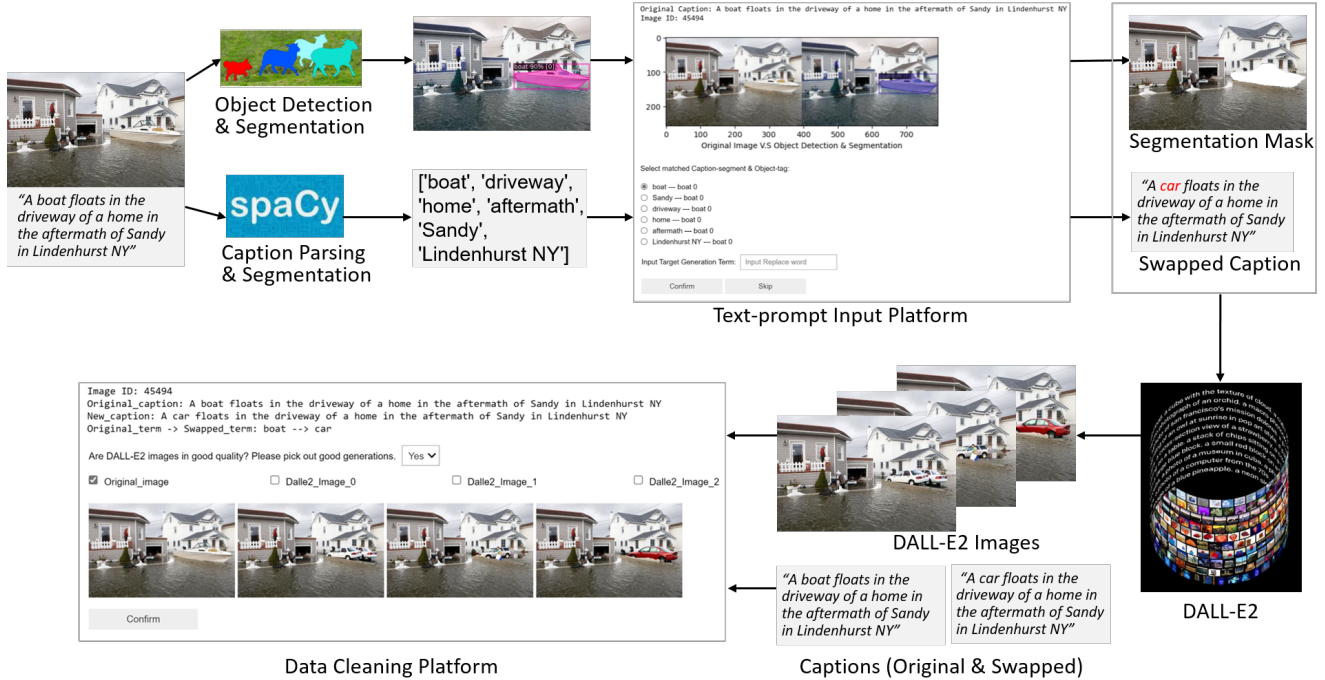
Figure 2. Pipeline of our text-prompt-based manipulation and annotation.

## 3.2. Dataset Construction

In this section, we introduce the details of the construction of the AutoSplice dataset.

### 3.2.1 Pre-processing

The pre-processing step for the DALL-E2 model involves providing a manipulation mask and a contextual text description to perform local image editing. In order to generate realistic manipulations, we propose using a specific object region with a modified caption where the corresponding text-prompt to the object is replaced with a target generation term. To achieve this, we use an object detection model to extract a list of object regions and a text parsing tool to segment text terms. We compare these terms with corresponding object-term pairs to facilitate further replacement and manipulation. Specifically, we use the Visual News dataset [43], which contains over one million news images along with their corresponding captions and metadata obtained from reputable real-world news sources (The Guardian, BBC, USA TODAY, and The Washington Post). This dataset has been utilized for various media forensic tasks, including the media manipulation detection [58] and text-image inconsistency detection [27]. For each sample in Visual News, we utilize the Detic model [72] to extract and segment object regions with detected object tags in the image. We also use spaCy ³ for sentence segmentation and noun term extraction, as shown in Figure 2. Human anno-

³ https://github.com/explosion/spaCy.

tations are then used to select the object with corresponding descriptions in the caption, and input target terms to replace the object description in the caption.

### 3.2.2 Human Annotations

Five annotators who are graduate and undergraduate students with professional backgrounds and have a clear understanding of the data annotation task for DALL-E2 input, strictly follow the steps outlined in Figure 2 during the data annotation process.

1. Select the matching object region tag in the real-world image and the corresponding text description term in the caption (if present).
2. Provide a target generation term that is similar but inconsistent with the original term and image.
3. Ensure that the modified caption has the correct syntax.

For each caption-image sample with a matched object-term pair, human annotations provide the two required inputs for the DALL-E2 model to perform local image generation: the segmented object region as the erased manipulation mask and the modified caption as the text prompt. The DALL-E2 model returns a group of three manipulation outputs for each generation.

### 3.2.3 Post-processing

To address the limitations of the DALL-E2 model in generating human, text, and abstract concepts [57], we conducted manual data cleaning to filter out images with visible visual

artifacts or caption-image pairs with undesirable consistencies. Given an original image-caption pair with the swapped caption and three DALL-E2 generated images, annotators were required to assess the visual quality of each generated image and identify images with good quality (i.e., no obvious artifacts). Since the definition of "good quality" is subjective and may vary among different annotators, each image was assessed multiple times by different annotators. We only retained images that received consistent labels from at least three annotators.

The high-quality DALL-E2 images were resized to match the dimensions of their corresponding authentic images. Despite being initially in PNG format, we compressed the generated DALL-E2 images using JPEG for two reasons. Firstly, their corresponding authentic images included in our dataset are in JPEG format. It is essential to eliminate format-level clues in the binary image forgery classification task. Secondly, JPEG is the most important and widely used image compression format [61], particularly on social media and websites, due to its simplicity and efficiency. Therefore, we chose both lossless (with a JPEG quality factor of 100) and gently lossy compression (with a quality factor of 90) formats to produce two variations of our DALL-E2 dataset.

### 3.3. Dataset Summary

Our AutoSplice dataset includes $3,621$ high-quality manipulated images and $2,273$ authentic images for each compression version. The data has been cleaned, and manipulation masks have been applied, allowing for further evaluation in both image forgery detection and image manipulation localization tasks. Figure 3 presents statistical information about the size of the manipulation region within the dataset, indicating that the dataset has a high diversity in the manipulation region. We further show some image examples in Figure 4.

## 4. Experimental Evaluation

This section outlines the evaluation experiments conducted on the AutoSplice dataset using state-of-the-art image-level and pixel-level forgery detection methods. The first experiment examines the generalization ability of existing pre-trained detectors to the LLI model manipulated images. Following this, we analyze the performance limits of these detectors in in-domain testing scenarios, where the models are fine-tuned on our AutoSplice dataset for both detection and localization tasks.

### 4.1. Evaluation Baselines

We evaluated five image forgery detection methods on our dataset: CNN-aug [66], ResNet50 Nodown [23], BeyondtheSpectrum [24], PSM [29], and GLFF [28]. These
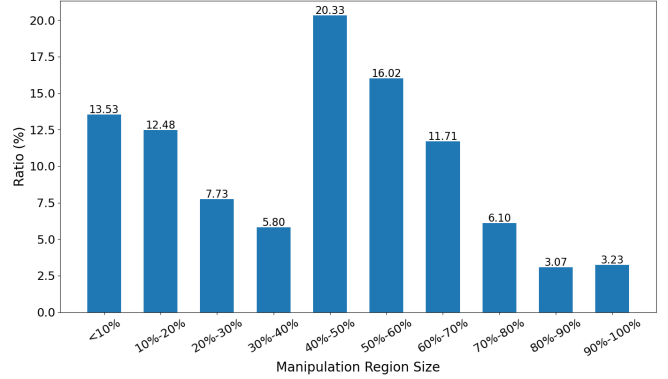


Figure 3. Statistical distribution of AutoSplice dataset in manipulation region size.



Figure 4. Examples in our AutoSplice dataset. The first column shows authentic images, while the second column displays forged images. The third column shows forgery masks.

models were chosen because they have demonstrated excellent performance in detecting image forgery, have been evaluated on both globally and locally manipulated images, and provide open-source codes and pre-trained models.

For pixel-level forgery detection and localization, we evaluated seven techniques that employ different feature learning strategies: Noiseprint [12], ManTra-Net [69], ForensicsGraph [46], CAT-Net [39], MVSS-Net [7], PSCC-Net [44], and ViT-VAE [6]. Table 2 presents detailed information on the feature design, training set, and software codes of these baseline methods.

### 4.2. Evaluation Metrics

To detect binary image forgery, we assign Positive (1) to the tampered image/pixel and Negative (0) to the authentic image/pixel. In image-level detection, we use standard metrics such as Area Under ROC Curve (AUC), True Positive

Table 2. Image forgery detection and localization baselines

| Reference | Year | Feature | Training set | Software Code |
|---|---|---|---|---|
| CNN-aug [66] | 2020 | Augmented CNN features | ProGAN [66] (720K images) | https://github.com/PeterWang512/CNNDetection |
| ResNet50 Nodown [23] | 2021 | No down-sampling CNN features | ProGAN [66] (720K images) | https://github.com/grip-unina/GANimageDetection |
| Beyondthe-Spectrum [24] | 2021 | Re-synthesis residuals | ProGAN [66] (720K images) | https://github.com/SSAW14/BeyondtheSpectrum |
| PSM [29] | 2022 | Global & local features | ProGAN [66] (720K images) | https://github.com/littlejuyan/FusingGlobalandLocal |
| GLFF [28] | 2022 | Multi-scale features | ProGAN [66] (720K images) | https://github.com/littlejuyan/GLFF |
| Noiseprint [12] | 2019 | Noise residuals | 4 datasets with 125 cameras | https://grip-unina.github.io/noiseprint/ |
| ManTra-Net [69] | 2019 | Anomalous features | 4 synthetic datasets | https://github.com/ISICV/ManTraNet |
| ForensicsGraph [46] | 2020 | Similarity graph | 4 million image patches from 80 cameras | https://gitlab.com/omayer/forensic-graph |
| CAT-Net [39] | 2021 | Compression artifacts | 4 synthetic datasets (960K images) | https://github.com/mjkwon2021/CAT-Net |
| MVSS-Net [7] | 2021 | Multi-view features | 1 dataset (CASIA v2) | https://github.com/dong03/MVSS-Net |
| PSCC-Net [44] | 2022 | Spatio-channel correlation | 0.38 million images | https://github.com/proteus1991/PSCC-Net |
| ViT-VAE [6] | 2023 | Multi-modal features | −* | https://github.com/media-sec-lab/ViT-VAE |

* Using run-time training where the ViT-VAE model requires an independent training phase for each test image.

Rate (TPR), and False Positive Rate (FPR). For pixel-level detection, we calculate F1 score, Intersection over Union (IoU), precision, and accuracy (ACC) by comparing the results with the binary ground-truth mask using a fixed threshold of 0.5, as done in previous studies [6–8, 18, 65]. The average scores for all testing images are reported.

**Experiment Settings**. To evaluate our dataset, we adhered to the parameter settings used in the baseline implementation for each task. During the fine-tuning process, we set a 6:4 split between training and testing data to ensure that there was no overlapping between the two sets.

## 4.3. Comparisons on Image Forgery Detection

### 4.3.1 Pre-trained Models

We started by using five pre-trained models for image forgery detection and tested their ability to identify the AutoSplice forgery. The results, including AUC, TPR, and FPR, are presented in Table 3. Our analysis shows that all models trained on the ProGAN dataset [66], which includes 720K images (360K real images and 360K fake images across 20 object categories), experienced a performance decrease in detecting image forgery in the AutoSplice dataset. Only the ResNet Nodown [23] model achieved the best AUC on two compression sets. All other models had an AUC lower than 0.600. The low TPR and FPR scores suggest that most LLI model forged images were incorrectly classified as authentic, whereas authentic images were correctly identified. Additionally, all models performed poorly on JPEG-90 images with mild compression, indicating that the compression process reduces the distinguishability of features.

### 4.3.2 Fine-tuned Models

We fine-tuned four models with training codes on the AutoSplice training dataset, after considering that data-driven classification methods tend to perform better when the domain discrepancy is alleviated [44]. To ensure that compression artifacts did not influence the binary forgery detection task, we compressed the original images produced

by the DALL-E2 model using the same JPEG compression quality factor (75 [4]) as the authentic images derived from Visual News dataset [43]. We evaluated the detection performance on two testing subsets (JPEG-100 and JPEG-90) and reported the results in Table 4. As expected, most methods showed a significant improvement in the detection AUC and TPR when evaluated in the in-domain testing scenario. The CNN-aug [66] achieved the best performance on both compression sets. However, after fine-tuning, most methods had a significant increase in FPR, indicating that a greater number of authentic images were incorrectly classified.

Table 3. Image forgery detection results on AutoSplice dataset using pre-trained models. Best results are shown in bold.

| Method | JPEG - 100 | | | JPEG - 90 | | |
|---|---|---|---|---|---|---|
| | AUC↑ | TPR↑ | FPR↓ | AUC↑ | TPR↑ | FPR↓ |
| CNN-aug [66] | 0.597 | 0.025 | 0.004 | 0.551 | 0.004 | 0.004 |
| ResNet50 Nodown [23] | **0.750** | 0.070 | **0.002** | **0.664** | 0.004 | **0.002** |
| Beyondthe-Spectrum [24] | 0.547 | **0.335** | 0.290 | 0.503 | **0.303** | 0.290 |
| PSM [29] | 0.586 | 0.038 | **0.002** | 0.535 | 0.005 | **0.002** |
| GLFF [28] | 0.572 | 0.055 | 0.013 | 0.526 | 0.016 | 0.013 |

Table 4. Image forgery detection results on AutoSplice dataset using fine-tuned models. Best results are shown in bold.

| Method | JPEG - 100 | | | JPEG - 90 | | |
|---|---|---|---|---|---|---|
| | AUC↑ | TPR↑ | FPR↓ | AUC↑ | TPR↑ | FPR↓ |
| CNN-aug [66] | **0.979** | **0.981** | 0.372 | **0.948** | **0.932** | 0.372 |
| Beyondthe-Spectrum [24] | 0.797 | 0.785 | 0.341 | 0.787 | 0.762 | 0.341 |
| PSM [29] | 0.880 | 0.847 | 0.257 | 0.882 | 0.841 | 0.257 |
| GLFF [28] | 0.926 | 0.776 | **0.077** | 0.908 | 0.735 | **0.077** |

## 4.4. Comparisons on Image Forgery Localization

### 4.4.1 Pre-trained Models

We evaluated seven baselines for image forgery localization, and reported their pixel-level metrics in Table 5. Re-

---

[4]The JPEG-75 compressed images are included in our AutoSplice dataset along with JPEG-100 and JPEG-90 versions at https://github.com/shanface33/AutoSplice_Dataset.

Table 5. Image forgery localization results on AutoSplice dataset using pre-trained models. Best results are shown in bold.

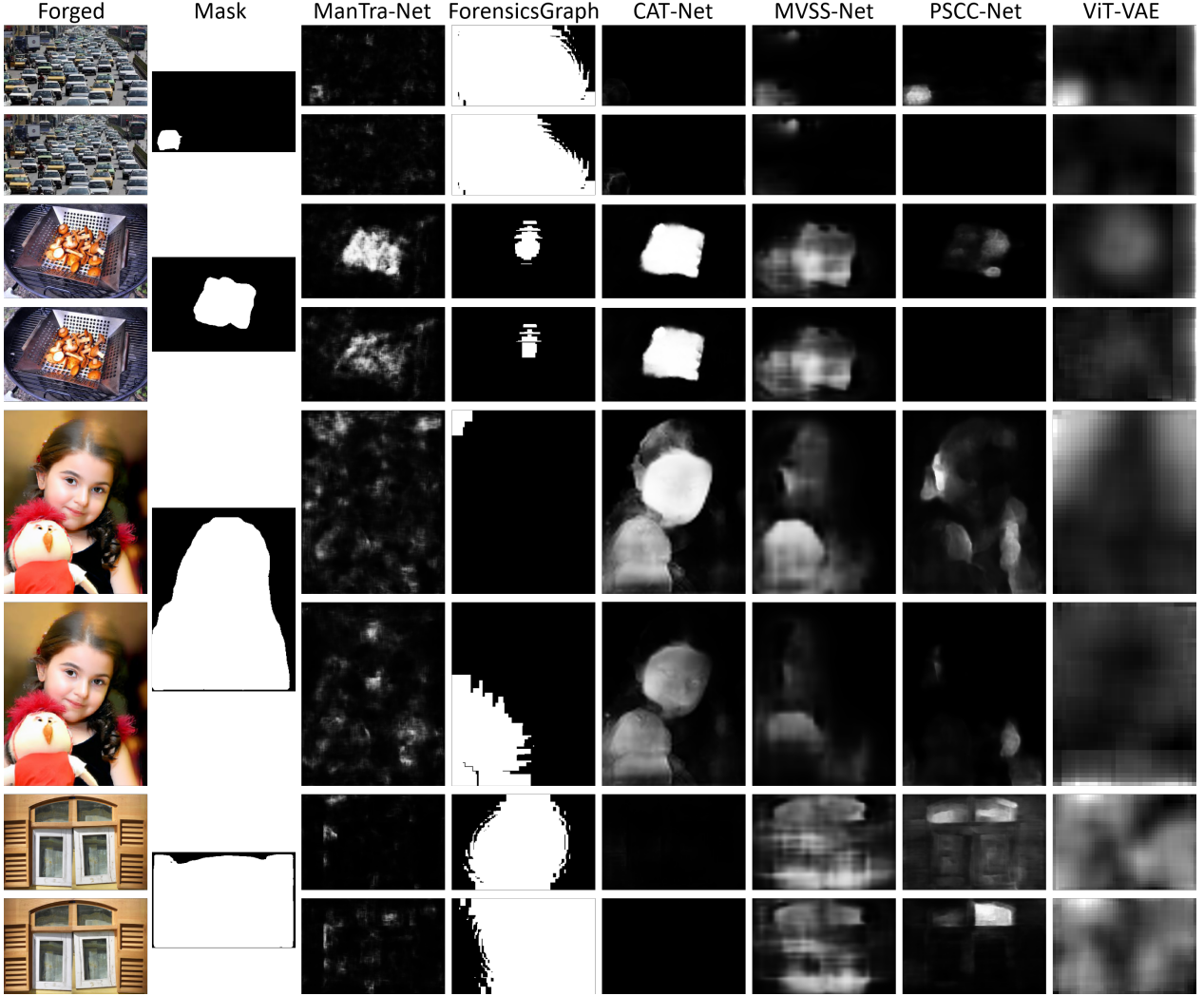| Method | Forged JPEG-100 | | | | Forged JPEG-90 | | | | Authentic |
|---|---|---|---|---|---|---|---|---|---|
| | F1↑ | IoU↑ | Precision↑ | ACC↑ | F1↑ | IoU↑ | Precision↑ | ACC↑ | ACC↑ |
| Noiseprint [12] | 0.333 | 0.217 | 0.390 | 0.480 | 0.316 | 0.205 | 0.373 | 0.467 | 0.594 |
| ManTra-Net [69] | 0.179 | 0.120 | 0.639 | 0.586 | 0.062 | 0.035 | 0.586 | 0.716 | 0.992 |
| ForensicsGraph [46] | 0.362 | 0.289 | 0.393 | 0.530 | 0.354 | 0.253 | 0.438 | 0.518 | 0.584 |
| CAT-Net [39] | **0.751** | **0.648** | 0.884 | **0.827** | **0.676** | **0.579** | **0.833** | **0.793** | 0.933 |
| MVSS-Net [7] | 0.330 | 0.238 | 0.734 | 0.677 | 0.141 | 0.093 | 0.516 | 0.612 | 0.991 |
| PSCC-Net [44] | 0.558 | 0.447 | **0.899** | 0.725 | 0.056 | 0.036 | 0.295 | 0.591 | **0.998** |
| ViT-VAE [6] | 0.156 | 0.115 | 0.245 | 0.560 | 0.244 | 0.183 | 0.275 | 0.534 | 0.835 |



Figure 5. Localization results of pre-trained models in detecting AutoSplice forged images with different manipulation regions. The images in odd rows are compressed using JPEG-100, while the images in even rows are compressed using JPEG-90.

sults varied significantly across models due to their different training data and forensics cues (detailed in Table 2). The CAT-Net [39] model performed the best thanks to its two-stream network that learns compression artifacts from both RGB and DCT domains, as well as its extensive and diverse training data. The PSCC-Net [44] also outperformed other methods on the JPEG-100 testing set, but its performance degraded significantly on the JPEG-90 set. Most models showed poor generalization ability on the AutoSplice dataset, with F1 and IoU lower than 0.37 and 0.29, respectively. For authentic images, where every pixel in the ground truth mask is negative, F1, IoU, and Precision met-

Table 6. Image forgery localization results on AutoSplice dataset using fine-tuned models. Best results are shown in bold.

| Method | JPEG - 100 | | | | JPEG - 90 | | | | Authentic |
|---|---|---|---|---|---|---|---|---|---|
| | F1↑ | IoU↑ | Precision↑ | ACC↑ | F1↑ | IoU↑ | Precision↑ | ACC↑ | ACC↑ |
| CAT-Net [39] | 0.762 | 0.658 | **0.882** | 0.837 | 0.693 | 0.594 | **0.844** | 0.805 | 0.927 |
| PSCC-Net [44] | **0.862** | **0.794** | 0.847 | **0.919** | **0.771** | **0.693** | 0.775 | **0.886** | **0.993** |

rics are not appropriate. We reported the ACC in Table 5, and most methods achieved high accuracy in detecting authentic images derived from real-world media data.

We want to note that the optimal threshold for each localization method may not be exactly 0.5, and it varies for different models and images. To eliminate the influence of the threshold, we further compared several examples with the predicted masks of these models without binarizing the map. Figure 5 shows the results on AutoSplice forged images with two compression versions and different manipulation regions. We observed the influence of even mild compression on the localization performance when comparing the results in the odd rows on JPEG-100 compressed images and even rows on JPEG-90 images. Moreover, the size of the manipulation region appeared to be another crucial factor affecting the accuracy of localization. The majority of models struggled to detect forgeries containing large tampered regions in the AutoSplice dataset.

### 4.4.2 Fine-tuned Models

We fine-tuned two localization methods, CAT-Net [38, 39] and PSCC-Net [44], on the AutoSplice training set using JPEG-75 compressed images (as previously described in Section 4.3.2) to evaluate their performance in localizing pixel-level forgery. As expected, both models outperformed pre-trained models on the JPEG-100 and JPEG-90 testing sets, with the PSCC-Net model demonstrating significant improvement. However, the performance on the JPEG-90 set decreases obviously in comparison to the results on the JPEG-100 set, which aligns with the findings from Table 5.

## 5. Conclusions

This paper investigates the challenge posed by language-image generation models to media forensics and proposes a new approach that utilizes the DALL-E2 language-image model to splice masked regions guided by a text prompt. To ensure the creation of realistic manipulations, we have developed an annotation platform with human verification to validate reasonable text prompts. The approach has resulted in the creation of a new image dataset called *AutoSplice*, containing 5,894 manipulated and authentic images, including $3,621$ images generated by locally or globally manipulating real-world image-caption pairs, which we believe will be a valuable resource for future research. We have evaluated the effectiveness of several state-of-the-art forgery detectors in various testing scenarios. However,

our experiments with pre-trained models revealed unsatisfactory generalization performance in forgery detection and localization. Including our dataset in training could enhance the performance of existing models during in-domain testing. This finding emphasizes that fine-tuning on datasets with homogeneous characteristics results in significant performance improvements in media forensics. Nevertheless, achieving balanced performance across different JPEG compression quality factors and tampered region sizes remains a challenging task for forgery localization.

For future works, we will consider: first, exploring more advanced models to generate more realistic manipulated images to further challenge media forensics. Secondly, investigating approaches to improve the generalization performance of forgery detection and localization models to handle various types of image manipulations. Thirdly, conducting experiments to evaluate the performance of existing forgery detection and localization models on the proposed dataset under various testing scenarios and compare them with state-of-the-art approaches. Last but not least, exploring the potential of transfer learning approaches to enhance the performance of existing models on new datasets with limited training samples.

## 6. Impact Statement

The recent advancements in language-image models have led to highly realistic image generation from textual descriptions, which can pose a potential threat to media forensics. This paper provides a new image dataset called AutoSplice, created using the DALL-E2 language-image model to splice masked regions guided by a text prompt. The unsatisfactory generalization performance of existing forgery detection and localization models on the proposed dataset highlights the need for further investigation and improvement in this area. Future works proposed in this paper, such as exploring more advanced models and transfer learning approaches, aim to address these challenges and contribute to the advancement of media forensics research. Ultimately, this work will help detect image manipulations in various applications, including social media, journalism, and law enforcement, and contribute to ensuring the authenticity and reliability of digital media content.

# References

[1] Nist nimble 2016 datasets., 2016. Accessed: 2023-03-11. 2, 3

[2] Tiziano Bianchi and Alessandro Piva. Image forgery localization via block-grained analysis of jpeg artifacts. *IEEE Transactions on Information Forensics and Security*, 7(3):1003–1017, 2012. 3

[3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 1

[4] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 2

[5] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 2

[6] Tong Chen, Bin Li, and Jinhua Zeng. Learning traces by yourself: Blind image forgery localization via anomaly detection with vit-vae. *IEEE Signal Processing Letters*, 2023. 3, 5, 6, 7

[7] Xinru Chen, Chengbo Dong, Jiaqi Ji, Juan Cao, and Xirong Li. Image manipulation detection by multi-view multi-scale supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14185–14193, 2021. 5, 6, 7

[8] Xinru Chen, Chengbo Dong, Jiaqi Ji, Juan Cao, and Xirong Li. Image manipulation detection by multi-view multi-scale supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14185–14193, 2021. 6

[9] Riccardo Corvi, Davide Cozzolino, Giada Zingarini, Giovanni Poggi, Koki Nagano, and Luisa Verdoliva. On the detection of synthetic images generated by diffusion models. *arXiv preprint arXiv:2211.00680*, 2022. 1, 2

[10] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022. 2

[11] Davide Cozzolino and Luisa Verdoliva. Single-image splicing localization through autoencoder-based anomaly detection. In *2016 IEEE international workshop on information forensics and security (WIFS)*, pages 1–6. IEEE, 2016. 3

[12] Davide Cozzolino and Luisa Verdoliva. Noiseprint: A cnn-based camera model fingerprint. *IEEE Transactions on Information Forensics and Security*, 15:144–159, 2019. 3, 5, 6, 7

[13] Hao Dang, Feng Liu, Joel Stehouwer, Xiaoming Liu, and Anil K Jain. On the detection of digital face manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern recognition*, pages 5781–5790, 2020. 2

[14] Tal Daniel and Aviv Tamar. Soft-introvae: Analyzing and improving the introspective variational autoencoder. In *Pro-ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4391–4400, 2021. 1

[15] Sowmen Das, Md Islam, Md Amin, et al. Gca-net: utilizing gated context attention for improving image forgery localization and detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 81–90, 2022. 3

[16] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 1

[17] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021. 2

[18] Chengbo Dong, Xinru Chen, Ruohan Hu, Juan Cao, and Xirong Li. Mvss-net: Multi-view multi-scale supervised networks for image manipulation detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 6

[19] Jing Dong, Wei Wang, and Tieniu Tan. Casia image tampering detection evaluation database. In *2013 IEEE China summit and international conference on signal and information processing*, pages 422–426. IEEE, 2013. 2, 3

[20] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, pages 12873–12883, 2021. 2

[21] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. Leveraging frequency analysis for deep fake image recognition. In *International conference on machine learning*, pages 3247–3258. PMLR, 2020. 3

[22] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 89–106. Springer, 2022. 1, 2

[23] Diego Gragnaniello, Davide Cozzolino, Francesco Marra, Giovanni Poggi, and Luisa Verdoliva. Are gan generated images easy to detect? a critical analysis of the state-of-the-art. In *2021 IEEE international conference on multimedia and expo (ICME)*, pages 1–6. IEEE, 2021. 5, 6

[24] Yang He, Ning Yu, Margret Keuper, and Mario Fritz. Beyond the spectrum: Detecting deepfakes via re-synthesis. In *30th International Joint Conference on Artificial Intelligence (IJCAI)*, 2021. 3, 5, 6

[25] 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. 2

[26] Y.-F. Hsu and S.-F. Chang. Detecting image splicing using geometry invariants and camera characteristics consistency. In *International Conference on Multimedia and Expo*, 2006. 2, 3

[27] Mingzhen Huang, Shan Jia, Ming-Ching Chang, and Siwei Lyu. Text-image de-contextualization detection using vision-language models. In *ICASSP 2022-2022 IEEE Inter-*

*national Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8967–8971. IEEE, 2022. 4

[28] Yan Ju, Shan Jia, Jialing Cai, Haiying Guan, and Siwei Lyu. Glff: Global and local feature fusion for face forgery detection. *arXiv preprint arXiv:2211.08615*, 2022. 2, 3, 5, 6

[29] Yan Ju, Shan Jia, Lipeng Ke, Hongfei Xue, Koki Nagano, and Siwei Lyu. Fusing global and local features for generalized ai-synthesized image detection. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 3465–3469. IEEE, 2022. 1, 3, 5, 6

[30] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 1

[31] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *NeurIPS*, 34, 2021. 2

[32] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 1

[33] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 1, 2

[34] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. 2

[35] Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019. 1

[36] Chenqi Kong, Baoliang Chen, Haoliang Li, Shiqi Wang, Anderson Rocha, and Sam Kwong. Detect and locate: Exposing face manipulation by semantic-and noise-level telltales. *IEEE Transactions on Information Forensics and Security*, 17:1741–1756, 2022. 3

[37] Paweł Korus and Jiwu Huang. Evaluation of random field models in multi-modal unsupervised tampering localization. In *2016 IEEE international workshop on information forensics and security (WIFS)*, pages 1–6. IEEE, 2016. 2, 3

[38] Myung-Joon Kwon, Seung-Hun Nam, In-Jae Yu, Heung-Kyu Lee, and Changick Kim. Learning jpeg compression artifacts for image manipulation detection and localization. *International Journal of Computer Vision*, 130(8):1875–1895, Aug. 2022. 3, 8

[39] Myung-Joon Kwon, In-Jae Yu, Seung-Hun Nam, and Heung-Kyu Lee. Cat-net: Compression artifact tracing network for detection and localization of image splicing. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 375–384, 2021. 3, 5, 6, 7, 8

[40] 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

[41] Xun Lin, Shuai Wang, Jiahao Deng, Ying Fu, Xiao Bai, Xinlei Chen, Xiaolei Qu, and Wenzhong Tang. Image manipulation detection by multiple tampering traces and edge artifact enhancement. *Pattern Recognition*, 133:109026, 2023. 3

[42] Bo Liu and Chi-Man Pun. Deep fusion network for splicing forgery localization. In *proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018. 3

[43] Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. Visual news: Benchmark and challenges in news image captioning. *arXiv preprint arXiv:2010.03743*, 2020. 2, 4, 6

[44] Xiaohong Liu, Yaojie Liu, Jun Chen, and Xiaoming Liu. Pscc-net: Progressive spatio-channel correlation network for image manipulation detection and localization. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(11):7505–7517, 2022. 5, 6, 7, 8

[45] Hannes Mareen, Dante Vanden Bussche, Fabrizio Guillaro, Davide Cozzolino, Glenn Van Wallendael, Peter Lambert, and Luisa Verdoliva. Comprint: Image forgery detection and localization using compression fingerprints. *arXiv preprint arXiv:2210.02227*, 2022. 3

[46] O. Mayer and M. C. Stamm. Exposing fake images with forensic similarity graphs. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):1049–1064, 2020. 5, 6, 7

[47] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022. 2

[48] Gustav Müller-Franzes, Jan Moritz Niehues, Firas Khader, Soroosh Tayebi Arasteh, Christoph Haarburger, Christiane Kuhl, Tianci Wang, Tianyu Han, Sven Nebelung, Jakob Nikolas Kather, et al. Diffusion probabilistic models beat gans on medical images. *arXiv preprint arXiv:2212.07501*, 2022. 2

[49] Tian-Tsong Ng, Shih-Fu Chang, and Q Sun. A data set of authentic and spliced image blocks. *Columbia University, ADVENT Technical Report*, 4, 2004. 2, 3

[50] 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. 2

[51] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 2

[52] Adam Novozamsky, Babak Mahdian, and Stanislav Saic. Imd2020: A large-scale annotated dataset tailored for detecting manipulated images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision Workshops*, pages 71–80, 2020. 2, 3

[53] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1, 2, 3

[54] 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. 1, 2

[55] 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, 2

[56] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 1

[57] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 2, 4

[58] Rui Shao, Tianxing Wu, and Ziwei Liu. Detecting and grounding multi-modal media manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 4

[59] Abhishek Sinha, Jiaming Song, Chenlin Meng, and Stefano Ermon. D2c: Diffusion-decoding models for few-shot conditional generation. *Advances in Neural Information Processing Systems*, 34:12533–12548, 2021. 1

[60] Michał Stypułkowski, Konstantinos Vougioukas, Sen He, Maciej Zięba, Stavros Petridis, and Maja Pantic. Diffused heads: Diffusion models beat gans on talking-face generation. *arXiv preprint arXiv:2301.03396*, 2023. 2

[61] Chentian Sun, Xiaopeng Fan, and Debin Zhao. Lossless recompression of jpeg images using transform domain intra prediction. *IEEE Transactions on Image Processing*, 32:88–99, 2022. 5

[62] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021. 2

[63] Luisa Verdoliva. Media forensics and deepfakes: an overview. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):910–932, 2020. 1

[64] Junke Wang, Zuxuan Wu, Jingjing Chen, Xintong Han, Abhinav Shrivastava, Ser-Nam Lim, and Yu-Gang Jiang. Objectformer for image manipulation detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2364–2373, 2022. 1

[65] Menglu Wang, Xueyang Fu, Jiawei Liu, and Zheng-Jun Zha. Jpeg compression-aware image forgery localization. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5871–5879, 2022. 6

[66] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8695–8704, 2020. 1, 2, 3, 5, 6

[67] Xin Wang, Hui Guo, Shu Hu, Ming-Ching Chang, and Siwei Lyu. Gan-generated faces detection: A survey and new perspectives. *arXiv preprint arXiv:2202.07145*, 2022. 1

[68] Bihan Wen, Ye Zhu, Ramanathan Subramanian, Tian-Tsong Ng, Xuanjing Shen, and Stefan Winkler. Coverage—a novel database for copy-move forgery detection. In *2016 IEEE international conference on image processing (ICIP)*, pages 161–165. IEEE, 2016. 2, 3

[69] Yue Wu, Wael AbdAlmageed, and Premkumar Natarajan. Mantra-net: Manipulation tracing network for detection and localization of image forgeries with anomalous features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9543–9552, 2019. 3, 5, 6, 7

[70] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 2

[71] Peng Zhou, Xintong Han, Vlad I Morariu, and Larry S Davis. Learning rich features for image manipulation detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1053–1061, 2018. 3

[72] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 350–368. Springer, 2022. 4

[73] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Lafite: Towards language-free training for text-to-image generation. *arXiv preprint arXiv:2111.13792*, 2021. 1

[74] Wanyi Zhuang, Qi Chu, Zhentao Tan, Qiankun Liu, Haojie Yuan, Changtao Miao, Zixiang Luo, and Nenghai Yu. Uia-vit: Unsupervised inconsistency-aware method based on vision transformer for face forgery detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part V*, pages 391–407. Springer, 2022. 3