



TextSleuth: Towards Explainable Tampered Text Detection

Chenfan Qu^{1,2}, Jian Liu², Haoxing Chen², Baihan Yu²,
Jingjing Liu², Weiqiang Wang², Lianwen Jin¹

¹South China University of Technology ²Ant Group
202221012612@mail.scut.edu.cn, rex.lj@antgroup.com, eelwjin@scut.edu.cn

Abstract

Recently, tampered text detection has attracted increasing attention due to its essential role in information security. Although existing methods can detect the tampered text region, the interpretation of such detection remains unclear, making the prediction unreliable. To address this problem, we propose to explain the basis of tampered text detection with natural language via large multimodal models. To fill the data gap for this task, we propose a large-scale, comprehensive dataset, ETDD, which contains both pixel-level annotations for tampered text region and natural language annotations describing the anomaly of the tampered text. Multiple methods are employed to improve the quality of the proposed data. For example, elaborate queries are introduced to generate high-quality anomaly descriptions with GPT4o. A fused mask prompt is proposed to reduce confusion when querying GPT4o to generate anomaly descriptions. To automatically filter out low-quality annotations, we also propose to prompt GPT4o to recognize tampered texts before describing the anomaly, and to filter out the responses with low OCR accuracy. To further improve explainable tampered text detection, we propose a **simple yet effective** model called TextSleuth, which achieves improved fine-grained perception and cross-domain generalization by focusing on the suspected region, with a two-stage analysis paradigm and an auxiliary grounding prompt. Extensive experiments on both the ETDD dataset and the public dataset have verified the effectiveness of the proposed methods. In-depth analysis is also provided to inspire further research. Our dataset and code will be open-source.

1 Introduction

Text image is one of the most important information carriers in today's society, containing a large amount of sensitive and private information [Chen *et al.*, 2024a]. With the rapid development of image processing technologies, sensitive text information can be more easily manipulated for malicious purposes, such as fraud, posing serious risks to information

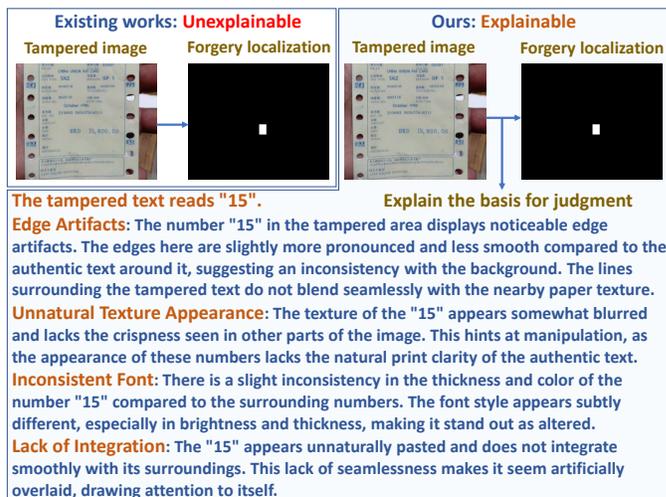


Figure 1: We propose to both detect the tampered text region and explain the basis for the detection in natural language, making the prediction more reliable. We construct the first dataset and propose a novel model for the explainable tampered text detection task.

security [Dong *et al.*, 2024]. Consequently, tampered text detection has become a major research topic in recent years [Qu *et al.*, 2024b]. It is crucial to develop effective and reliable methods for detecting tampered text images.

Existing works model tampered text detection as semantic segmentation [Shao *et al.*, 2023] or object detection [Qu *et al.*, 2024a], with the aim of interpreting the basis for image forgery classification by predicting tampered regions. Despite the progress made in recent years, such fine-grained predictions are still black-box and cannot provide a convincing basis to support the judgement, leading to unreliable results.

To provide more reliable predictions for tampered text detection, we propose to leverage multimodal large models to both detect tampered text regions and explain the basis for their detection in natural language. Given the absence of dataset for interpretable tampered text detection, we construct the Explainable Tampered Text Detection (ETDD) dataset. To ensure the comprehensiveness of the data, we collect multilingual card images, document images and scene text images from the Internet and the existing text-rich datasets such as ICDAR2017 [Nayef *et al.*, 2017] and LSVT [Sun

et al., 2019]. And we perform text tampering on the collected data with various methods, including traditional methods copy-move, splicing, and the deep generative method DiffUTE [Chen *et al.*, 2024b]. Poisson Blending [Pérez *et al.*, 2023] is employed to reduce the visual inconsistency around tampered region. Finally, we 12,000 tampered text images with accurate pixel-level annotations of the tampered region and 10,500 authentic text images. The large-scale of our data notably alleviates the data hunger of deep models. The images are split to three parts, ETDD-Train, ETDD-Test and ETDD-CD, the two test sets have the same and different distributions with the ETDD-Train respectively, allowing both in-domain and cross-domain evaluation.

With the obtained tampered text, we utilize GPT4o to generate the description of both visual and linguistic anomalies caused by text tampering, and to generate the text recognition result for specifying the target tampered text. To achieve this, we prompt the GPT4o with a novel elaborate query, the tampered image and its corresponding mask annotation indicating the tampered region. However, since text is mostly dense and has similar location and shape, directly inputting the binary mask, as existing work [Xu *et al.*, 2024] does will cause severe confusion to the GPT4o, making it unclear which is the actual tampered text. To solve this problem, we propose to fuse the binary mask into the original tampered image with pixel-wise weighting. With the proposed fused mask prompt, the GPT4o has a much better understanding of the location of the target region, which in turn significantly reduces the errors and obviously improves the annotation quality. In addition, the GPT4o’s output is not always correct and manual verification is costly. Inspired by the fact that incorrect detection of manipulated text leads to unclear perception and poor anomaly description, we further propose to address this issue by automatically filtering out the annotation based on the accuracy of the manipulated text OCR.

The tiny area and visual consistency of tampered text [Wang *et al.*, 2022] pose multiple challenges for explainable tampered text detection, making it difficult for existing methods to achieve good enough performance. For example, misidentification of tampered text leads to incorrect anomaly description, difficulty in finding tampered text weakens the analysis quality, and increases the risk of overfitting to unrelated background styles. To this end, we propose a novel **simple-yet-effective** model termed as TextSleuth. Specifically, an extra RCNN [Ren *et al.*, 2015] based text detection module initially scans the image and predicts the location of the tampered text with cascaded RoI heads. The initial prediction of tampered region is converted into a grounding prompt and fed into the large language model along with the image tokens and the original question to obtain the final prediction. The proposed two-stage analysis paradigm and auxiliary prompt in TextSleuth effectively minimizes errors, improves explanation quality and cross-domain generalization by drawing the model’s special attention to the anomaly region and helping it to learn more general features. In addition, since the reference grounding comprehension task is mostly involved in the pre-training stage of large models [Chen *et al.*, 2024c], the proposed auxiliary grounding prompt can reduce comprehension difficulty and alleviate forgetting.

Both our proposed ETDD dataset and TextSleuth model are the first efforts in the field of interpretable tampered text detection. Extensive experiments have confirmed that the proposed TextSleuth significantly improves upon the baseline model, outperforming existing methods by a large margin on both the proposed ETDD dataset and the public Tampered IC-13 [Wang *et al.*, 2022] dataset, demonstrating strong in-domain and cross-domain generalization capabilities. In-depth analysis is also provided to inspire further work in the field of interpretable tampered text detection.

In summary, the our main contribution is fourfold:

- We propose a novel **task**, explainable tampered text detection, which aims to provide reliable prediction by describing the anomalies of tampered text in natural language, serving as a pioneering effort in this field.
- We obtain the data annotation for this task by prompting GPT4o with elaborate queries. We propose effective methods to improve the quality of the annotations. For example, a fused mask prompt to reduce model confusion and a novel method to automatically filter out bad responses. Based on these, we construct the ETDD dataset, which is **the first**, large-scale and comprehensive **dataset** for explainable tampered text detection.
- We propose **the first** multimodal large **model** TextSleuth for interpretable tampered text detection, which achieves state-of-the-art performance with a two-stage analysis paradigm and a novel auxiliary prompt.
- Extensive experiments are conducted. Valuable conclusions and insights are provided through in-depth analysis, inspiring the further research in this field.

2 Related works

2.1 Tampered Text Detection

Early work on tampered text detection is achieved by printer classification [Lampert *et al.*, 2006] or template matching [Ahmed and Shafait, 2014], which is only applicable to scanned documents and does not work well for photographed documents [Dong *et al.*, 2024]. DocTamper [Qu *et al.*, 2023] is the first large-scale comprehensive dataset for tampered text detection in arbitrary style documents. It is created by a novel method Selective Tampering Generation [Qu *et al.*, 2023], which effectively reduces the cost of high-quality training data. DTD [Qu *et al.*, 2023] is proposed to detect visually consistent tampering in documents through examining the continuity of the block artifacts grids. The SACP [Alibaba Security, 2020] and RTM [Luo *et al.*, 2024] datasets provide large-scale manually forged documents. Tampered IC13 [Wang *et al.*, 2022] dataset is the first work on tampered scene text detection, containing about 300 scene text images tampered by SR-Net [Wu *et al.*, 2019]. The OSFT dataset [Qu *et al.*, 2024a] significantly improves Tampered IC-13 with realistic forgeries created by nine latest deep generative models. CAFTB-Net [Song *et al.*, 2024] benefits from noise-domain modeling and cross-attention mechanism. DTL [Shao *et al.*, 2025] improves model robustness with latent manifold adversarial training. Omni-IML [Qu *et al.*, 2024b] unifies image forensics across different text image types. Despite the

Table 1: A brief summary of the ETTD dataset statistics. "Forged Area" denotes the area ratio of tampered text.

Dataset	Image types	Image source	Languages	Tampering methods & Numbers	Authentic num.	Forged Area
ETTD-Train	Documents, ID cards,	Internet, ICDAR2013,	EN, CH	Total (10400): DiffUTE (800), CopyMove (4800), Splicing (4800)	9600	0.0268
ETTD-Test	scene texts, etc.	ICDAR2017, LSVT	EN, CH	Total (600): DiffUTE (200), CopyMove (200), Splicing (200)	400	0.0202
ETTD-CD	scene text	ICDAR2013	EN	Total (1000): Copymove (500), Splicing (500)	500	0.0608

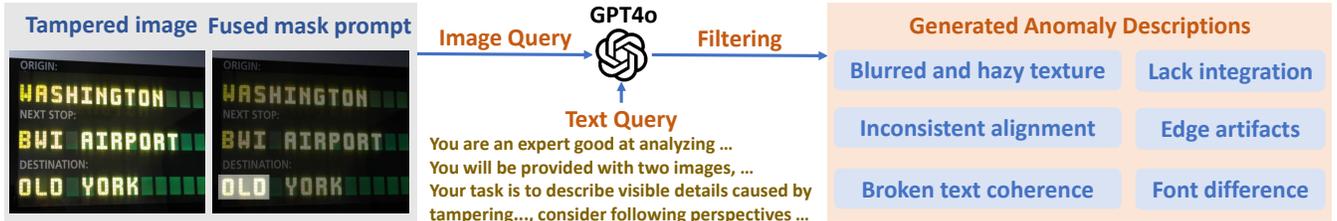


Figure 2: The pipeline for obtaining the textual anomaly description for the tampered text.

progress made in recent years, existing work on tampered text detection can still only localize the tampered region in an unreliable black-box manner, unable to explain the basis of its judgement in natural language.

2.2 Explainable Image Forgery Detection

Recently, some works achieve explainable image forensics through multimodal large language models. FFAA [Huang *et al.*, 2024a] utilizes GPT4o to generate detailed basis description for the identification judgment. MMTT [Lian *et al.*, 2024] further proposes to perform in-depth face forgery analysis by observing the facial organs one by one. FakeShield [Xu *et al.*, 2024] leverages GPT4o to create anomaly description for natural style image forgery and introduces a new model based on LISA [Lai *et al.*, 2024]. EditScout [Nguyen *et al.*, 2024] focuses on the detection of image forgeries that generated by diffusion models and also develops a new model based on LISA. ForgeryGPT [Li *et al.*, 2024a] improves interpretable natural image forensics with binary mask prompt. FakeBench [Li *et al.*, 2024b] fills the blank in explainable image forensics for forgeries generated by deep generative models (e.g. Diffusion models). ForgerySleuth [Sun *et al.*, 2024] obtains hierarchical forgery description annotation with the proposed Chain-of-Clues for the MIML [Qu *et al.*, 2024c], a high quality image forensics dataset with more than 123k manually forged images and their corresponding pixel-level annotations. It also improves LISA for natural image forensics with noise domain modeling. SIDA [Huang *et al.*, 2024b] achieves explainable forensics for natural image and AIGC image especially on social media scenario, an effective new model is introduced by improving LISA. Despite the progress made, none of the existing work achieves interpretable forensics on another important task, tampered text detection. Due to the tiny size and visual consistency of tampered text [Qu *et al.*, 2023], natural image forgery detection methods mostly do not work very well for tampered text detection [Luo *et al.*, 2024], leaving massive risks. It is crucial to develop explainable tampered text detection techniques for reliable text image forensics.

3 ETTD Dataset

To fill in the data gap for explainable tampered text detection dataset, we construct a large-scale comprehensive dataset called Explainable Tampered Text Detection (ETTD).

3.1 Text Tampering

To ensure the comprehensiveness of the proposed data, we collect multilingual document and card images from the Internet and scene text images from the existing datasets (e.g. ICDAR2013 [Karatzas *et al.*, 2013] and LSVT [Sun *et al.*, 2019]). We then forge some of the collected images with the widely-used methods, copy-move and splicing. Poisson Blending [Pérez *et al.*, 2023] is employed to reduce visual inconsistency. To further improve the data diversity, we manually edit the text with DiffUTE [Chen *et al.*, 2024b], a latest diffusion model for realistic tampered text generation.

3.2 Anomaly Description Generation

As shown in Figure 2, we leverage the GPT4o to generate the description of both visual and linguistic anomalies caused by text tampering. Given the different features between tampered text and tampered natural objects [Wang *et al.*, 2022], the textual queries in existing works [Xu *et al.*, 2024] can not work well for tampered text (e.g. "unnatural depth" is usually observed in tampered natural objects but not in tampered text). To address this issue, we propose an elaborate query that inspires the GPT4o to analyze anomalies for tampered text on six major perspectives, covering texture, integration, alignment, edge artifacts, text coherence, font, as shown in Figure 2. The detailed query is presented in the Appendix.

We further input this elaborate query along with the tampered image and its corresponding mask annotation into the GPT4o. However, due to the similarity in location and shape of the text instances in an image, directly inputting the binary mask as in existing work [Xu *et al.*, 2024] will cause considerable confusion to the GPT4o. As shown in Fig. 3, the annotator model usually struggles to identify the target text with the binary mask, often mistaking a nearby authentic text as a fake text. Analyzing anomaly on authentic text



Figure 3: The binary mask prompt as in existing work is confusing in text images. In contrast, our proposed fused mask prompt clearly indicates the content and the exact location of the tampered text.

undoubtedly produces incorrect anomaly descriptions. To address this issue, we propose the fused mask prompt, where the original image is fused with the binary mask by pixel-wise weighting. Specifically, given the input image $I \in \mathbb{R}^{H,W,3}$ and the binary mask annotation $M \in \mathbb{B}^{H,W}$, $\mathbb{B} \in \{0, 1\}$, the fused mask prompt $M^{fused} \in \mathbb{R}^{H,W,3}$ can be formulated as $M^{fused} = I * \lambda_1 + M * \lambda_2$. We set λ_1 and λ_2 to 0.5 in practical. With the proposed fused mask prompt, the annotator can clearly recognize the tampered text on the target region and better understand where the target region is by referring to the surrounding content. The proposed method significantly reduces hallucination and errors caused by frequent confusion.

Since the responses of GPT4o are not always correct, directly using the GPT4o responses as annotations leads to poor data quality, while manually verifying the annotation is costly. Our aim is to automatically filter out unsatisfactory responses. We empirically find that, the anomaly description from the GPT4o is also mostly accurate when the GPT4o can correctly recognize the tampered text. This means that the GPT4o is clear about the location of the tampered text and the visual details of it. Inspired by this, we propose to automatically filter out the bad responses with tampered text OCR accuracy [Zhang *et al.*, 2019] lower than 0.8. The OCR ground-truth is obtained from dataset annotation or OCR engine, and is used to replace the GPT4o OCR in the remaining samples to ensure accuracy. The proposed method effectively improves the quality of anomaly description for tampered text in an automatic manner. For authentic text images, the textual description is set to "There is no tampered text in this image."

3.3 Dataset Summary

As shown in Tab. 1, there are 5,500 text images tampered by copy-move, 5,500 text images tampered by splicing and 1,000 text images tampered by DiffUTE in our ETTD dataset. Another 10,500 images without text tampering serve as the authentic part. 20,000 images from the ETTD dataset are split as the training set (ETTD-Train), 1,000 images from the ETTD dataset are split as the test set (ETTD-Test) and another 1,500

images from the ETTD dataset are split as the cross-domain test set (ETTD-CD). The ETTD-CD consists of Copy-move forgeries, Splicing forgeries and authentic images from IC-DAR2013, which are not included in ETTD-Train. Therefore, the ETTD-CD has a different data distribution from ETTD-Train and can evaluate model performance on unknown scenarios. Accurate pixel-level annotations for tampered regions are provided to facilitate fine-grained analysis of the tampered text regions. The data hunger of large models can be effectively alleviated with our large-scale diverse data.

4 TextSleuth

The tampered text is mostly tiny in size and the the visual anomaly is mostly unobvious [Luo *et al.*, 2024]. Consequently, three major challenges are emerged for interpretable tampered text detection: 1. The multimodal large models suffer from more risks of misidentifying the tampered text, resulting in incorrect anomaly description. 2. A large part of the model parameters are used to find the tampered region, which weakens their ability to analyze and describe the tampered region. 3. The models are more likely to be disturbed by the irrelevant background style, which weakens their generalization on unseen tampering methods and scenarios. To this end, we propose a **simple-yet-effective** model termed as TextSleuth, which overcomes the above challenges through a two-stage analysis paradigm and a novel reference grounding auxiliary prompt.

As shown in Figure 4, given an input image, the suspected tampered text region is initially detected by a Swin-Transformer based detection model with cascaded RoI heads [Cai and Vasconcelos, 2018]. The predicted coordinates are then normalized to 0-1000 and are converted to the reference grounding auxiliary prompt "The suspected tampered text $\langle box \rangle [[x_{min}, y_{min}, x_{max}, y_{max}]] \langle /box \rangle$ ". Given that the reference grounding comprehension task is involved in the pre-training stage of most large models [Wang *et al.*, 2024], the large language model can effortlessly comprehend the target location in the proposed auxiliary prompt. In the auxiliary prompt, the large language model naturally pays special attention to the region represented by the coordinates, as it has learned in its pre-training stage. This differs from existing work [Li *et al.*, 2024a] that forces the model to look at the suspected region with binary mask embeddings, which is confusing in indicating tampered text, violates the pre-training paradigm and causes more forgetting. The auxiliary prompt is fed into the large language model along with the image tokens and the original question, to obtain the recognition and describe the anomaly for tampered text.

Despite its simplicity, the proposed method effectively addresses the three major challenges in explainable tampered text detection: 1. The initial prediction of the suspect region significantly reduces the risk of misidentifying the tampered region and reduces hallucination. 2. The detection prior alleviates the difficulty in detecting tampered text, models can save more parameters to obtain better anomaly analysis and description. 3. By focusing on the tampered region, the model gets rid of the interference from unrelated background styles, learns more general features during training,

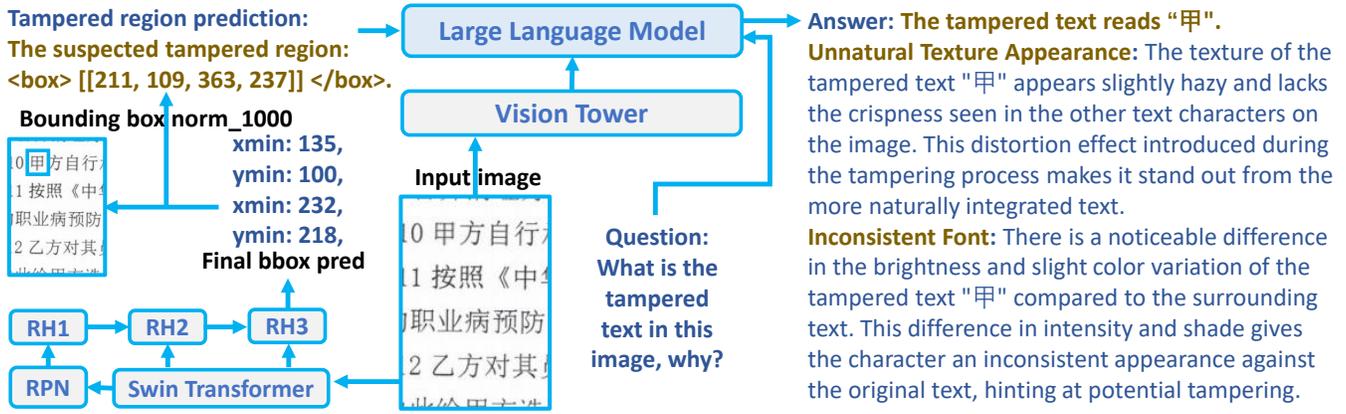


Figure 4: The overall pipeline of the proposed TextSleuth.

and thereby obtains improved generalisation to unseen tampering methods and scenarios.

5 Experiments

We conduct experiments on both the proposed ETTD dataset and the public Tampered IC-13 dataset [Wang *et al.*, 2022] with multiple advanced large multimodal models, including GPT4o [OpenAI, 2024], Yi-VL-6B [AI *et al.*, 2024], DeepSeekVL-7B [Lu *et al.*, 2024], MiniCPMV2.5 [Hu *et al.*, 2024], the 1B to 8B versions of Intern2VL [Chen *et al.*, 2024c] and the 2B, 7B versions of Qwen2VL [Wang *et al.*, 2024]. We fine-tune all models except GPT4o on the ETTD training set for 5 epochs with the same settings, which is sufficient for all models to achieve their best performance.

5.1 Evaluation Metric

To evaluate the similarity between the predicted anomaly description and the textual annotation, we calculate the OCR accuracy [Zhang *et al.*, 2019] for tampered text recognition and the paragraph cosine similarity for non-OCR parts. The weighted summary of OCR accuracy and paragraph similarity is used as the final similarity score. For misclassified samples, the paragraph cosine similarity is set directly to 0 as the gist is opposite. Specifically, we extract the content within the quotation marks from the first predicted sentence and use it to calculate the OCR accuracy Acc_{OCR} . We then remove stop-words and the content within the quotation marks from both prediction P_{pred} and ground-truth paragraphs P_{gt} for more accurate paragraph similarity calculation. The paragraph feature vectors V_{pred}, V_{gt} are obtained by averaging the word vectors in each paragraph, $V_{pred} = average([W2V(word) \text{ for } word \text{ in } P_{pred}])$, $V_{gt} = average([W2V(word) \text{ for } word \text{ in } P_{gt}])$, where $W2V$ is the pretrained word-to-vector function [Mikolov *et al.*, 2018]. Finally, the cosine similarity between the two paragraph feature vectors is used as the paragraph similarity score, $Sim_{para} = Cos(V_{pred}, V_{gt})$. We have manually verified that better predictions almost always lead to higher cosine similarity scores. The final similarity score Sim is calculated by $Sim = 0.3 * Acc_{OCR} + 0.7 * Sim_{para}$. The common accuracy metric [Guillaro *et al.*, 2023] is adopted for image forgery classification task.

5.2 Implement Details

The vision tower and projector of the large multimodal model are full-parameter fine-tuned and the large language model part is LoRA [Hu *et al.*, 2021] fine-tuned with rank 8 and alpha 16. We adopt AdamW [Loshchilov and Hutter, 2017] optimizer with a learning rate decaying linearly from $1e-4$ to 0. The batch-size is set to 16 for all models and the experiments are run on NVIDIA A100 80GB GPUs. We set the maximum input area to 1.2M pixels for the Qwen2VL models. In the proposed TextSleuth model, the Swin-Transformer based detection model is trained for 30 epochs on the ETTD training set, with a batch-size of 16 and a maximum input resolution of 1.2M pixels. The AdamW optimizer is employed with a learning rate that decays linearly from $6e-6$ to $3e-6$.

For all the fine-tuned models, the input text query is "What is the tampered text in this image, why?", which matches the training data. The image classification prediction is regarded as "tampered" if the edit distance between the model output and the string "There is no tampered text in this image." is greater than 3. For GPT4o and other pre-trained models, to output the most similar format with the annotation, the query is set to "Does this image have tampered text on it? Please start your answer with "Yes" or "No". If "Yes", then recognize the tampered text and describe the anomaly of the tampered region.". The image classification prediction is regarded as "tampered" if the output starts with "Yes".

5.3 Comparison Study

Anomaly Description. The comparison results of anomaly description on the ETTD dataset are shown in Table 2. Four conclusions can be drawn through analyzing the results:

(1) **High input resolution is essential for models to correctly recognize the tampered text and explain its anomaly, but it can also slightly weaken model generalization to unfamiliar scenarios.** On the ETTD-Test set, the Yi-VL-6B with the lowest input resolution 448×448 achieves the lowest final score of 68.5, which is notably lower than the fine-tuned models. Resizing the input image to such a small resolution makes the subtle visual anomaly hard to detect, thus Yi-VL-6B performs the worst. However, the Yi-VL-6B performs much better on ETTD-CD, where the tampered text is mostly larger. Most of the other models suf-

fer significant performance degradation due to the unfamiliar scenario. This indicates that high input resolution mostly weakens model generalization on unfamiliar scenarios.

(2) **The scaling law also applies to the explainable tampered text detection task.** Even within the same series (e.g. InternVL2 or Qwen2VL) where the vision tower is the same and the pre-training data is similar, models with larger LLMs mostly perform better. For example, Qwen2VL-7B achieves an average final score of 90.7, which is better than Qwen2VL-2B. This confirms that the scaling law behind our task.

(3) **Model performance on the explainable tampered text detection task is highly related to its pre-trained model performance on general task.** For example, on the OpenCompass, the pre-trained MiniCPMV2.5 model performs better than DeepSeekVL-7B, which is consistent with the model performance of their fine-tuned versions on the explainable forensic task. A better large multimodal model on general task is likely to be more intelligent and can better learn to describe the anomaly of tampered text.

(4) **The proposed TextSleuth performs the best on both ETDD-Test and ETDD-CD, significantly outperforming other models in both in-domain and cross-domain scenarios.** This verifies that the proposed auxiliary prompt improves model’s fine-grained perception and helps the model to produce high quality anomaly description by focusing its attention on the suspected region at start.

Image Forgery Classification. The comparison results of the image forgery classification are shown in Tab. 4. The public tampered IC-13 dataset used in evaluation consists of texts tampered by SR-Net [Wu *et al.*, 2019] and can also evaluate model’s cross domain performance on unknown tampering method. The proposed TextSleuth considerably outperforms existing methods on all the three benchmarks, and improves the Qwen2VL-7B baseline by +3.5 points, +7.9 points and +12.7 points on the three benchmarks respectively, demonstrating the effectiveness of the proposed method. Similar conclusions can be drawn as analyzed in the above paragraphs.

Auto-annotation. To verify the effectiveness of the proposed fused mask prompt, we manually obtain the OCR results of 100 random tampered texts from the collected data. We compare the tampered text OCR performance between the binary mask prompt as in existing work and the proposed fused mask prompt, the results are shown in Table 5. The annotator GPT4o has significantly higher OCR accuracy with the proposed prompt. The higher OCR accuracy indicates that the large model can better understand the actual location of the tampered text. Therefore, the anomaly description from GPT4o is also more accurate with the method.

5.4 Ablation Study

The ablation study of the proposed TextSleuth is shown in Table 3. We conduct experiments on three base multimodal models. For each base model (e.g. InternVL2-2B), there are four ablation settings. Setting (1) is the official pre-trained model performance. Setting (2) is the official model fine-tuned on the ETDD Train. Setting (3) is the TextSleuth fine-tuned with the proposed grounding auxiliary prompt. Setting (4) is the TextSleuth with the perfect tampered text detector. Four conclusions can be drawn through analysis:

(1) **The existing multimodal models do not have the ability to recognize tampered text and the anomaly.** All three base models perform poorly in setting (1), but much better in setting (2). This confirms that the official open-source models are mostly incapable of detecting tampered text. Training them on the ETDD data is essential for them to gain the ability for explainable tampered text detection.

(2) **The proposed auxiliary prompt can significantly improve model performance across different base models.** For each of the three base models, the model performance in setting (3) is significantly better than that in setting (2) (+10.8 points average final score for InternVL2-2B, +9.2 for Qwen2VL-2B and +6.5 for Qwen2VL-7B). These improvements are achieved by the proposed two-stage analysis paradigm and the auxiliary prompt in our TextSleuth. The proposed methods alleviate the difficulty in detecting tampered region and make the models better focused on analyzing the anomaly, resulting in an improved anomaly description quality. Additionally, by focusing on the tampered text with the proposed prompt, the models can learn more general features by reducing the interference from unrelated background styles. Consequently, the model’s cross-domain generalization is considerably improved. The huge improvements on different basic large multimodal models also demonstrate that our TextSleuth is widely applicable.

(3) **The performance of our TextSleuth can be further improved with better tampered region detectors.** For all of the three base models, model performance in setting (4) is better than those in setting (3). The improvement is achieved by eliminating the errors of the initial tampered text box prediction. Therefore, our TextSleuth can easily be improved in the future with an advanced tampered text region detector.

Robustness Evaluation. We evaluate the robustness of the TextSleuth under different JPEG compression quality factors and different resize factors on ETDD-Test and ETDD-CD. As shown in Table 6, the stable performance under common distortions has verified the robustness of our TextSleuth.

The **prediction visualization** is presented in the Appendix.

6 Conclusion

This paper is the first work that achieves explainable tampered text detection, by describing the anomaly of the tampered text image with natural language. Given the absence of dataset for this task, we construct a large-scale comprehensive dataset termed as ETDD, which consists of multilingual document and scene text images tampered by copy-move, splicing and AIGC-based text editing. Given the obtained tampered text, we obtain the anomaly description by prompting GPT4o with elaborate queries, the tampered image and tampered region annotation. However, due to the similar shape and position of the text, prompting GPT4o with a binary mask as in existing works mostly confuses the GPT4o and leads to incorrect responses. To address this issue, we propose to prompt GPT4o with a fused mask, which is obtained by weighting the image with the binary mask. Experiments verify that the proposed fused mask prompt significantly improves the annotation quality. Moreover, given that the incorrect recognition of tampered text means unclear perception and leads to

Table 2: Comparison study of the proposed method.

Methods	ETTD-Test (in-domain)			ETTD-CD (cross-domain)			Average
	OCR accuracy	Cosine similarity	Final score	OCR accuracy	Cosine similarity	Final score	Final score
GPT4o	48.3	66.1	60.7	74.6	78.0	77.0	68.9
Yi-VL-6B	49.9	76.5	68.5	64.3	81.4	76.2	72.4
DeepSeekVL-7B	66.6	86.9	80.8	37.9	64.7	56.7	68.8
MiniCPMV2.5	79.3	92.6	88.6	68.9	74.8	73.0	80.8
InternVL2-1B	77.8	89.1	85.7	79.2	84.0	82.5	84.1
InternVL2-2B	81.1	91.5	88.3	78.2	82.7	81.3	84.8
InternVL2-4B	75.8	82.4	80.4	91.4	94.0	93.1	86.8
InternVL2-8B	80.9	90.7	87.7	80.0	85.1	83.5	85.6
Qwen2VL-2B	84.8	93.7	91.0	82.1	85.0	84.1	87.6
Qwen2VL-7B	87.1	94.8	92.4	87.1	89.9	88.9	90.7
TextSleuth-7B (Ours)	92.6	98.3	96.5	97.7	98.1	97.9	97.2

Table 3: Ablation study of the proposed method. "SFT" denotes surprised fine-tuning. "TextSleuth" denotes equipping the model with the proposed TextSleuth method. "Perfect Detector" denotes using ground-truth tampered region boxes in the TextSleuth's auxiliary prompt.

Base Multi-modal Model	Ablation settings				ETTD-Test (in-domain)			ETTD-CD (cross-domain)			Average
	Num	SFT	Text Sleuth	Perfect Detector	OCR accuracy	Cosine similarity	Final score	OCR accuracy	Cosine similarity	Final score	Final score
InternVL2-2B	(1)	×	×	×	14.1	57.8	44.7	34.3	64.4	55.4	50.1
	(2)	✓	×	×	81.1	91.5	88.3	78.2	82.7	81.3	84.8
	(3)	✓	✓	×	83.1	98.3	93.6	96.5	98.1	97.5	95.6
	(4)	✓	✓	✓	83.7	99.4	94.6	97.7	99.4	98.8	96.7
Qwen2VL-2B	(1)	×	×	×	18.5	57.0	45.5	29.8	63.8	53.6	49.6
	(2)	✓	×	×	84.8	93.7	91.0	82.1	85.0	84.1	87.6
	(3)	✓	✓	×	90.4	98.2	95.8	97.2	98.0	97.7	96.8
	(4)	✓	✓	✓	91.3	99.3	96.8	98.5	99.3	99.0	97.9
Qwen2VL-7B	(1)	×	×	×	14.0	41.8	33.5	36.4	53.4	48.3	40.9
	(2)	✓	×	×	87.1	94.8	92.4	87.1	89.9	88.9	90.7
	(3)	✓	✓	×	92.6	98.3	96.5	97.7	98.1	97.9	97.2
	(4)	✓	✓	✓	93.6	99.4	97.6	99.0	99.4	99.2	98.4

Table 4: Accuracy performance of different large multimodal models on image forgery classification task.

Method	ETTD-Test (in-domain)	ETTD-CD (out-domain)	Tampered-IC13 (zero-shot)
GPT4o	67.3	79.3	82.8
Yi-VL-6B	76.9	81.9	45.9
DeepSeekVL-7B	87.4	66.7	76.4
MiniCPMV2.5	93.2	75.5	56.7
InternVL2-1B	89.7	84.6	59.2
InternVL2-2B	92.1	83.3	58.8
InternVL2-4B	82.8	94.5	36.1
InternVL2-8B	91.2	85.7	60.5
Qwen2VL-2B	94.3	85.7	73.8
Qwen2VL-7B	95.4	90.5	75.1
TextSleuth-7B	98.9	98.6	88.4

Table 5: Comparison study for the proposed fused mask prompt.

Method	OCR Accuracy	Perfect Match
Binary mask prompt	47.3	30.4
Fused mask prompt (Ours)	84.2	73.0

Table 6: Robustness evaluation under common distortions.

Average final score	JPEG Ori.	JPEG compress quality75	JPEG compress quality50	Image resize factor0.75	Image resize factor0.5
Qwen2VL	90.7	89.6	87.2	89.2	86.1
TextSleuth	97.2	96.3	94.4	95.8	93.0

bad anomaly description, we also propose to filter out the responses with low tampered text OCR accuracy, which can effectively improve annotation quality in an automatic manner. The proposed ETTD dataset has both in-domain and cross-domain test subsets, which allows a thorough evaluation of model generalization. In addition, a novel TextSleuth model is proposed to further improve explainable tampered text detection, which overcomes several major challenges in the field with a two-stage analysis paradigm and an auxiliary prompt. Experiments have confirmed that the proposed method considerably improves upon different baseline models, and that our TextSleuth notably outperforms existing methods in both in-domain and cross-domain evaluation on both the ETTD and public datasets. In-depth analysis is also provided to inspire further work. We believe that our **valuable** ETTD dataset and our **first-of-its-kind, simple-yet-effective** methods can shed light on the further research on interpretable tampered text detection.

References

- [Ahmed and Shafait, 2014] Amr Gamal Hamed Ahmed and Faisal Shafait. Forgery detection based on intrinsic document contents. In *2014 11th IAPR International Workshop on Document Analysis Systems*, pages 252–256. IEEE, 2014.
- [AI *et al.*, 2024] 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024.
- [Alibaba Security, 2020] Alibaba Security. Security ai challenger program. <https://tianchi.aliyun.com/competition/entrance/531812/introduction>, 2020.
- [Cai and Vasconcelos, 2018] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [Chen *et al.*, 2024a] Changsheng Chen, Liangwei Lin, Yongqi Chen, Bin Li, Jishen Zeng, and Jiwu Huang. Cma: A chromaticity map adapter for robust detection of screen-recapture document images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15577–15586, 2024.
- [Chen *et al.*, 2024b] Haoxing Chen, Zhuoer Xu, Zhangxuan Gu, Yaohui Li, Changhua Meng, Huijia Zhu, Weiqiang Wang, et al. Diffute: Universal text editing diffusion model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [Chen *et al.*, 2024c] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.
- [Dong *et al.*, 2024] Li Dong, Weipeng Liang, and Rangding Wang. Robust text image tampering localization via forgery traces enhancement and multiscale attention. *IEEE Transactions on Consumer Electronics*, 2024.
- [Guillaro *et al.*, 2023] Fabrizio Guillaro, Davide Cozzolino, Avneesh Sud, Nicholas Dufour, and Luisa Verdoliva. Trufor: Leveraging all-round clues for trustworthy image forgery detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20606–20615, 2023.
- [Hu *et al.*, 2021] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [Hu *et al.*, 2024] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zheng Leng Thai, Kaihuo Zhang, Chongyi Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of small language models with scalable training strategies, 2024.
- [Huang *et al.*, 2024a] Zhengchao Huang, Bin Xia, Zicheng Lin, Zhun Mou, Wenming Yang, and Jiaya Jia. Ffaa: Multimodal large language model based explainable open-world face forgery analysis assistant, 2024.
- [Huang *et al.*, 2024b] Zhenglin Huang, Jinwei Hu, Xiangtai Li, Yiwei He, Xingyu Zhao, Bei Peng, Baoyuan Wu, Xiaowei Huang, and Guangliang Cheng. Sida: Social media image deepfake detection, localization and explanation with large multimodal model, 2024.
- [Karatzas *et al.*, 2013] Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluís Gomez i Bigorda, Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazan Almazan, and Lluís Pere De Las Heras. Icdar 2013 robust reading competition. In *2013 12th international conference on document analysis and recognition*, pages 1484–1493. IEEE, 2013.
- [Lai *et al.*, 2024] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024.
- [Lampert *et al.*, 2006] Christoph H Lampert, Lin Mei, and Thomas M Breuel. Printing technique classification for document counterfeit detection. In *2006 International Conference on Computational Intelligence and Security*, volume 1, pages 639–644. IEEE, 2006.
- [Li *et al.*, 2024a] Jiawei Li, Fanrui Zhang, Jiaying Zhu, Esther Sun, Qiang Zhang, and Zheng-Jun Zha. Forgerygpt: Multimodal large language model for explainable image forgery detection and localization, 2024.
- [Li *et al.*, 2024b] Yixuan Li, Xuelin Liu, Xiaoyang Wang, Bu Sung Lee, Shiqi Wang, Anderson Rocha, and Weisi Lin. Fakebench: Probing explainable fake image detection via large multimodal models, 2024.
- [Lian *et al.*, 2024] Jingchun Lian, Lingyu Liu, Yaxiong Wang, Yujiao Wu, and Zhedong Zheng. A large-scale interpretable multi-modality benchmark for image forgery localization, 2024.
- [Loshchilov and Hutter, 2017] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [Lu *et al.*, 2024] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. Deepseek-vl: Towards real-world vision-language understanding, 2024.

- [Luo *et al.*, 2024] Dongliang Luo, Yuliang Liu, Rui Yang, Xianjin Liu, Jishen Zeng, Yu Zhou, and Xiang Bai. Toward real text manipulation detection: New dataset and new solution. *Pattern Recognition*, page 110828, 2024.
- [Mikolov *et al.*, 2018] Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhresch, and Armand Joulin. Advances in pre-training distributed word representations. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- [Nayef *et al.*, 2017] Nibal Nayef, Fei Yin, Imen Bizid, Hyunsoo Choi, Yuan Feng, Dimosthenis Karatzas, Zhenbo Luo, Umapada Pal, Christophe Rigaud, Joseph Chazalon, et al. Icdar2017 robust reading challenge on multi-lingual scene text detection and script identification-rrc-mlt. In *2017 14th IAPR international conference on document analysis and recognition (ICDAR)*, volume 1, pages 1454–1459. IEEE, 2017.
- [Nguyen *et al.*, 2024] Quang Nguyen, Truong Vu, Trong-Tung Nguyen, Yuxin Wen, Preston K Robinette, Taylor T Johnson, Tom Goldstein, Anh Tran, and Khoi Nguyen. Editscout: Locating forged regions from diffusion-based edited images with multimodal llm, 2024.
- [OpenAI, 2024] OpenAI. Gpt-4 technical report, 2024.
- [Pérez *et al.*, 2023] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 577–582. 2023.
- [Qu *et al.*, 2023] Chenfan Qu, Chongyu Liu, Yuliang Liu, Xinhong Chen, Dezhi Peng, Fengjun Guo, and Lianwen Jin. Towards robust tampered text detection in document image: new dataset and new solution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5937–5946, 2023.
- [Qu *et al.*, 2024a] Chenfan Qu, Yiwu Zhong, Fengjun Guo, and Lianwen Jin. Generalized tampered scene text detection in the era of generative ai, 2024.
- [Qu *et al.*, 2024b] Chenfan Qu, Yiwu Zhong, Fengjun Guo, and Lianwen Jin. Omni-impl: Towards unified image manipulation localization, 2024.
- [Qu *et al.*, 2024c] Chenfan Qu, Yiwu Zhong, Chongyu Liu, Guitao Xu, Dezhi Peng, Fengjun Guo, and Lianwen Jin. Towards modern image manipulation localization: A large-scale dataset and novel methods. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10781–10790, June 2024.
- [Ren *et al.*, 2015] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015.
- [Shao *et al.*, 2023] Huiru Shao, Kaizhu Huang, Wei Wang, Xiaowei Huang, and Qiufeng Wang. Progressive supervision for tampering localization in document images. In *International Conference on Neural Information Processing*, pages 140–151. Springer, 2023.
- [Shao *et al.*, 2025] Huiru Shao, Zhuang Qian, Kaizhu Huang, Wei Wang, Xiaowei Huang, and Qiufeng Wang. Delving into adversarial robustness on document tampering localization. In *European Conference on Computer Vision*, pages 290–306. Springer, 2025.
- [Song *et al.*, 2024] Yalin Song, Wenbin Jiang, Xiuli Chai, Zhihua Gan, Mengyuan Zhou, and Lei Chen. Cross-attention based two-branch networks for document image forgery localization in the metaverse. *ACM Transactions on Multimedia Computing, Communications and Applications*, 2024.
- [Sun *et al.*, 2019] Yipeng Sun, Zihan Ni, Chee-Kheng Chng, Yuliang Liu, Canjie Luo, Chun Chet Ng, Junyu Han, Errui Ding, Jingtuo Liu, Dimosthenis Karatzas, et al. Icdar 2019 competition on large-scale street view text with partial labeling-rrc-lsvt. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 1557–1562. IEEE, 2019.
- [Sun *et al.*, 2024] Zhihao Sun, Haoran Jiang, Haoran Chen, Yixin Cao, Xipeng Qiu, Zuxuan Wu, and Yu-Gang Jiang. Forgery sleuth: Empowering multimodal large language models for image manipulation detection, 2024.
- [Wang *et al.*, 2022] Yuxin Wang, Hongtao Xie, Mengting Xing, Jing Wang, Shenggao Zhu, and Yongdong Zhang. Detecting tampered scene text in the wild. In *European Conference on Computer Vision*, pages 215–232. Springer, 2022.
- [Wang *et al.*, 2024] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution, 2024.
- [Wu *et al.*, 2019] Liang Wu, Chengquan Zhang, Jiaming Liu, Junyu Han, Jingtuo Liu, Errui Ding, and Xiang Bai. Editing text in the wild. In *Proceedings of the 27th ACM International Conference on Multimedia*, MM ’19, page 1500–1508, New York, NY, USA, 2019. Association for Computing Machinery.
- [Xu *et al.*, 2024] Zhipei Xu, Xuanyu Zhang, Runyi Li, Zecheng Tang, Qing Huang, and Jian Zhang. Fakeshield: Explainable image forgery detection and localization via multi-modal large language models, 2024.
- [Zhang *et al.*, 2019] Rui Zhang, Yongsheng Zhou, Qianyi Jiang, Qi Song, Nan Li, Kai Zhou, Lei Wang, Dong Wang, Minghui Liao, Mingkun Yang, et al. Icdar 2019 robust reading challenge on reading chinese text on signboard. In *2019 international conference on document analysis and recognition (ICDAR)*, pages 1577–1581. IEEE, 2019.

7 The Proposed Textual Prompt

Due to the different characteristics of tampered text, existing textual prompts designed for natural objects or deepfakes cannot be directly used to generate high-quality anomaly descriptions for tampered text. To this end, we redesign the textual prompt by summarizing the possible anomalies caused by text tampering into six major perspectives and providing a detailed explanation for each of them.

The full version of our textual prompt is:

You are an expert good at analyzing tampered text images. You will be provided with two images, **the first is the tampered text image A and the second is the reference image B, with the tampered areas highlighted and the authentic areas darkened.**

Your task is to: **First, recognize the tampered text and output its OCR result. Second, Describe visible details in the image that have been tampered with.** Please consider the visible details caused by tampering from these perspectives.

- 1. Edge artifacts.** The background of the tampered text may be inconsistent with the authentic regions. Therefore, the edges around the tampered text region may be discontinuous and inconsistent with the background.
- 2. Unnatural texture appearance.** The texture appearance of the tampered text may be slightly blurred, hazy, jagged, have a distortion effect, or have an unnatural clarity.
- 3. Inconsistent font.** The font of the tampered text may be slightly different in color, size, thickness, brightness, or style from the surrounding authentic text.
- 4. Inconsistent alignment.** The tampered text may have inconsistent spacing with the surrounding text or a small offset to the text line.
- 5. Text incoherence.** Tampered text may break the coherence of the sentence.
- 6. Lack of integration.** The tampered text may appear unnaturally placed and not integrated with its surroundings, or it may not blend seamlessly with its surroundings, appearing artificially overlaid or unnaturally pasted. Don't mention the image B in your answer, always assume that you are only observing the input image A.

As shown in Figure 1, our proposed prompt can help GPT4o output a satisfactory anomaly description.

8 Detection Performance

We present the detection performance of the detector in the proposed TextSleuth in Table 7. The precision, recall and F1-score under the ICDAR2017 DetEval protocol [Nayef *et al.*, 2017] and an IoU threshold of 0.5 are used.

Table 7: The detection performance of the detector in the proposed TextSleuth. The IoU threshold is set to 0.5.

Dataset	Precision	Recall	F1-score
ETTD-Test	0.994	0.986	0.990
ETTD-CD	0.983	0.995	0.989

9 Visualization

The prediction visualization of GPT4o, Qwen2VL-7B and our TextSleuth is shown in Figures 6 and 7. Evidently, the proposed method can produce more accurate results. More samples of the ETTD dataset are shown in Figures 8 and 9.



User

You are an expert good at analyzing tampered text images. You will be provided with two images, **the first is the tampered text image A and the second is the reference image B, with the tampered areas highlighted and the authentic areas darkened.**

Your task is to: First, **recognize the tampered text and output its OCR result.** Second, **Describe visible details in the image that have been tampered with.** Please consider the visible details caused by tampering from these perspectives:

- 1. Edge artifacts.** The background of the tampered text may be inconsistent with the authentic regions. Therefore, the edges around the tampered text region may be discontinuous and inconsistent with the background.
- 2. Unnatural texture appearance.** The texture appearance of the tampered text may be slightly blurred, hazy, jagged, have a distortion effect, or have an unnatural clarity.
- 3. Inconsistent font.** The font of the tampered text may be slightly different in color, size, thickness, brightness, or style from the surrounding authentic text.
- 4. Inconsistent alignment.** The tampered text may have inconsistent spacing with the surrounding text or a small offset to the text line.
- 5. Text incoherence.** Tampered text may break the coherence of the sentence.
- 6. Lack of integration.** The tampered text may appear unnaturally placed and not integrated with its surroundings, or it may not blend seamlessly with its surroundings, appearing artificially overlaid or unnaturally pasted.

Don't mention the image B in your answer, always assume that you are only observing the input image A.



User



GPT4o

The tampered text reads "GRASS".

Edge Artifacts: The edges around the tampered text "GRASS" appear to be discontinuous with the surrounding background. There is a noticeable boundary around the letters that does not seamlessly blend with the text's immediate environment, suggesting a cut-and-paste effect commonly found in image tampering.

Unnatural Texture Appearance: The texture of the tampered text "GRASS" appears slightly blurred compared to the sharper appearance of the other texts. This distortion effect, where the edges of the letters do not have the same crispness or definition as the authentic text, makes it stand out unnaturally.

Inconsistent Alignment: The spacing between the tampered text "GRASS" and the preceding text "BREAK" appears irregular or has a slight offset when compared to the consistency seen in the other lines of text. This misalignment disrupts the uniformity typically expected in such signage.

Text Incoherence: The text "BREAK GRASS" does not create a coherent or meaningful message within the context of safety instructions, which typically involve actions like "BREAK GLASS." This incoherence adds to the suspicion of tampering, as the revised phrasing results in linguistic nonsense.

Lack of Integration: The tampered text "GRASS" appears as though it has been overlaid or pasted onto the image rather than being a naturally integrated part of the scene. This artificial placement makes the text appear as a foreign element within the image, lacking the integration seen in genuine, unaltered text regions.

Figure 5: Our proposed textual prompt are specially designed for tampered text can can guide GPT4o to generate high-quality anomaly descriptions for tampered text.

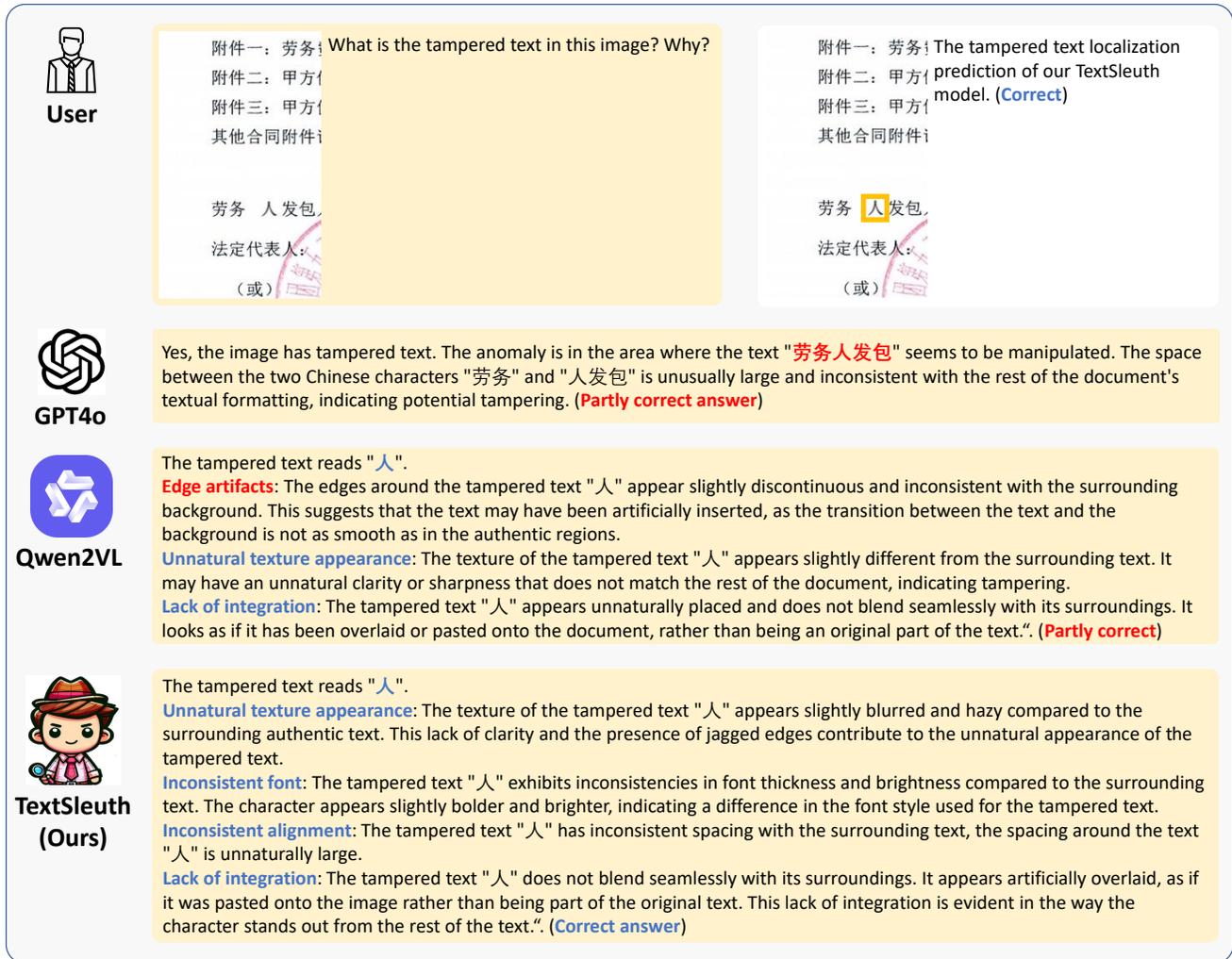


Figure 6: Our proposed TextSleuth can accurately detect the tampered text and better explain the anomalies of the tampered text.



What is the tampered text in this image? Why?

PARTY SUPPORT DEMOGRAPHICS														
		Age					Social Class				Region			
		18-24 %	25-34 %	35-49 %	50-64 %	65+ %	ABC1 %	C2DE %	F %	Other %	Dublin %	Midlands %	East of Eng %	West of Eng %
FINE GAEL	40%	+4	35	43	39	41	39	39	37	62	32	44	43	41
labour	20%	+1	15	22	20	17	25	19	23	6	29	18	20	6
SDP	11%	+1	19	7	33	10	5	8	14	7	8	11	11	12
FINIAN FAUL	10%	-7	10	9	7	8	21	10	9	21	5	11	12	16
green party	2%	=	2	4	2	2	2	3	2	-	5	-	1	1
Independent/Other	17%	+2	17	14	19	21	6	19	15	4	20	14	11	22

Our model's localization prediction of tampered text.

PARTY SUPPORT DEMOGRAPHICS														
		Age					Social Class				Region			
		18-24 %	25-34 %	35-49 %	50-64 %	65+ %	ABC1 %	C2DE %	F %	Other %	Dublin %	Midlands %	East of Eng %	West of Eng %
FINE GAEL	40%	+4	35	43	39	41	39	39	37	62	32	44	43	41
labour	20%	+1	15	22	20	17	25	19	23	6	29	18	20	6
SDP	11%	+1	19	7	33	10	5	8	14	7	8	11	11	12
FINIAN FAUL	10%	-7	10	9	7	8	21	10	9	21	5	11	12	16
green party	2%	=	2	4	2	2	2	3	2	-	5	-	1	1
Independent/Other	17%	+2	17	14	19	21	6	19	15	4	20	14	11	22



Yes, the text appears to have been tampered with. The manipulated area is the percentage figure next to "FINE GAEL" under the "65+ %" category in the Age section. The abnormality is in the appearance of the number "62" which seems larger, bolder, and slightly misaligned compared to the other numbers in the grid, indicating it might have been digitally altered or superimposed. Additionally, visually it looks out of place and indicates it has been tampered with. **(Incorrect answer)**



Qwen2VL

The tampered text reads "22".
 Edge artifacts: The edges around the tampered text "22" appear slightly discontinuous and inconsistent with the surrounding background. This suggests that the text may have been inserted or altered, as the transition between the text and the background is not seamless.
 Unnatural texture appearance: The texture of the tampered text "22" seems slightly different from the authentic text. It appears slightly blurred or hazy compared to the sharper, more defined text around it, indicating potential tampering.
 Lack of integration: The tampered text "22" appears unnaturally placed and does not blend seamlessly with its surroundings. It looks artificially overlaid, as if it was pasted onto the image rather than being part of the original design. **(Incorrect answer)**



TextSleuth (Ours)

The tampered text reads "33".
Edge artifacts: The edges around the tampered text "33" appear slightly inconsistent with the background. There are noticeable discontinuities or a change in texture in the immediate vicinity of the text, suggesting the insertion of new text.
Unnatural texture appearance: The texture of the number "33" seems to have a slightly different appearance compared to neighboring text. It indicates alterations.
Inconsistent font: The font style, thickness, or color of "33" differs subtly from the surrounding numbers. This inconsistency can serve as an indicator of tampering, suggesting that the font does not match the original style used in the rest of the document.
Lack of integration: The "33" does not blend seamlessly with its surrounding text and grid lines, making it appear artificially overlaid or pasted. It lacks the natural integration seen in untouched sections of the document, enhancing the perception of tampering. **(Correct answer)**

Figure 7: Our proposed TextSleuth can accurately detect the tampered text and better explain the anomalies of the tampered text.



The tampered text reads "PR".

Edge Artifacts: The tampered text "PR" shows visible discontinuities at the edges compared to the surrounding areas. The edges around the "PR" text show a break or discontinuity against the background, indicating it was not part of the original image.

Unnatural Texture Appearance: The texture of the tampered "PR" text appears somewhat hazy and less sharp than the surrounding authentic texts. This indicates manipulation, as the integrated texture should normally appear consistent.

Inconsistent Font: There is a noticeable inconsistency in size and color with the text just before it and the authentic "STOP" text below. The tampered text "PR" appears slightly more faded and not as bold as the surrounding text.

Inconsistent Alignment: The "PR" does not line up properly with the preceding text "DON," appearing offset both horizontally and vertically. This inconsistent alignment makes the text look less natural and more likely tampered with.

Text Incoherence: The inclusion of the letters "PR" within the context of the image does not form a coherent or expected phrase within the setting of a traffic sign, suggesting tampering in terms of intended message.

Lack of Integration: The tampered "PR" text does not smoothly integrate with the color and shading nuances of the stop sign's red background or the adjacent legitimate text, making it appear as an artificially placed element.



The tampered text reads "JAMES".

Edge artifacts: The edges around the tampered text "JAMES" appear to be inconsistent with the background, indicating discontinuity. This is evident where the boundary between the tampered region and the surrounding area is not as smooth or consistent as the original regions.

Unnatural texture appearance: The tampered text "JAMES" exhibits a slightly unnatural texture, where the clarity of the letters seems detached from the rest of the image. This could manifest as either an overly sharp appearance compared to the slightly blurred background or an inconsistent texture.

Inconsistent font: There is a noticeable difference in the font style or thickness of the tampered text "JAMES" compared to the surrounding authentic text. This can include variations in color, size, or brightness, which makes the tampered text stand out unnaturally from the rest.

Lack of integration: The tampered text "JAMES" appears unnaturally placed, as if overlaid on the image without seamlessly blending with its surroundings. This lack of integration points to an artificial addition which has a contrasting appearance compared to the authentic text.



The tampered text reads "7890".

Edge Artifacts: The edges around the tampered text "7890" appear to be slightly discontinuous and inconsistent with the background texture, which results in an unnatural demarcation between the tampered text and the surrounding authentic text. This discontinuity suggests tampering.

Unnatural Texture Appearance: The tampered text "7890" has a slightly different texture compared to the surrounding numbers. It appears somewhat blurred or artificially smoothed, which results in an unnatural clarity that is distinct from the rest of the text, indicating a potential tampering effect.

Lack of Integration: The tampered text "7890" does not blend seamlessly with the surrounding characters and the background. It appears somewhat as if it is placed on top of the card's design rather than integrated into it, giving an impression of being unnaturally overlaid or pasted into the image.

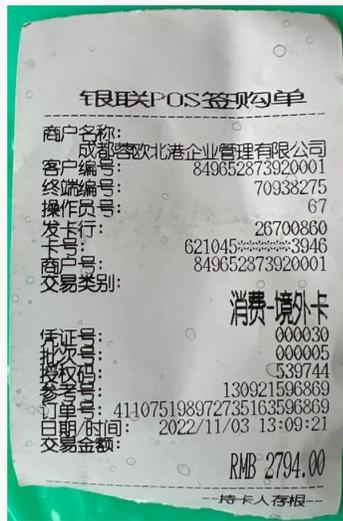
Figure 8: More data samples in the proposed ETTD dataset. The left image in each row is the original image and the tampered text region is highlighted in the right image.



The tampered text reads "B".

Unnatural Texture Appearance: The letter "B" displays an unnatural texture. This is evident through slight blurriness, jagged edges, or a distortion effect making it stand out against the surrounding authentic text.

Inconsistent Font: The character "B" differs in font size, thickness, brightness, and color compared to the surrounding text. This variance makes it appear distinct and separate from the rest of the text elements in the label.

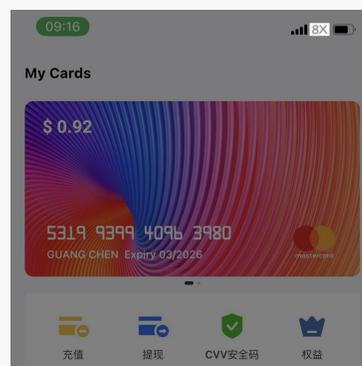
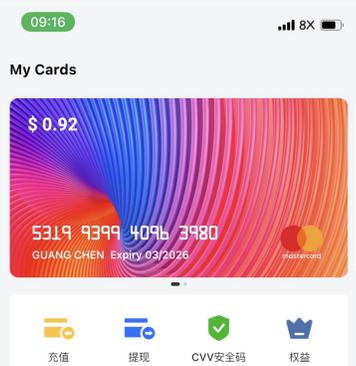


The tampered text reads "67".

Edge Artifacts: The number "67" has visible edge artifacts, indicating tampering. The boundaries around these characters appear more prominent or sharper than the surrounding text, suggesting they were added separately.

Inconsistent Font: The font used for "67" shows slight variances in thickness or brightness compared to the authentic text, indicating that it was not originally part of the document.

Unnatural Texture Appearance: The "67" appears slightly blurred or hazy compared to nearby numbers. This effect might result from editing that failed to seamlessly integrate the new text with the texture of the paper.



The tampered text reads "8X".

Edge Artifacts: The edges around the "8X" tampered text appear slightly inconsistent with the smooth and continuous background typically seen in authentic regions. This manifests in subtle edge discontinuities or abrupt transitions that do not blend seamlessly with the surrounding area.

Unnatural Texture Appearance: The "8X" text shows signs of unnatural texture compared to its original appearance. This includes potential blurriness, haziness, or jagged contours which are not naturally found on the finely rendered elements typical of this user interface.

Figure 9: More data samples in the proposed ETTD dataset. The left image in each row is the original image and the tampered text region is highlighted in the right image.