017

# Unleashing the Potential of Consistency Learning for Detecting and Grounding Multi-Modal Media Manipulation

## Supplementary Material

## **001** A. Additional visualization

002As shown in Fig. 1, we give the additional visualizations on003 $DGM^4$  dataset, which demonstrate the effectiveness of our004proposed method in different scenarios.

## **B.** Future work

The future work could be carried out based on the following aspects. First, the usage of large vision-language model (LVLM) should be explored on DGM<sup>4</sup> [2] task, as LVLM contains rich real-world knowledge which contributes to forgery reasoning. Second, more advanced constructing and<br/>supervising functions should be introduced on DGM4 task,<br/>which enhances the ability to distinguish between genuine<br/>and forged information in extreme difficult cases. Third,<br/>due to the frequent issues of distortion and noise in image<br/>transmission on the internet, the robustness of detection is<br/>also worth exploring in the future.010<br/>012

## **C.** Failure cases

As shown in Fig. 2, we give some failure cases. It could be 018 observed that when the manipulated face is too small, the 019



Figure 1. Visualization of detecting and grounding results on  $DGM^4$  datasets. Here, red box and text indicate the prediction of manipulated faces and words, while green box and text represent the corresponding ground truth. The first line is the prediction for text forgery, the second line is the prediction for image forgery, and the third line is the prediction for multi-modal forgery.

032

037

045

052

058

059

060

061



Figure 2. Visualization of failure cases on  $DGM^4$  datasets. Here, red box and text indicate the prediction of manipulated faces and words, while green box and text represent the corresponding ground truth.

(1)

detecting or grounding tasks to image modality may fail.
What's more, when the expression of the sentence is relatively vague and there is no well-aligned in semantics between image-text pair, the detecting or grounding tasks to text modality may fail.

## 025 D. Additional loss details

1026 Here, we give additional details of losses used in different 1027 sub-tasks which are similar to [3].  $V_{cls}$  and  $T_{cls}$  represent 1028 the class embeddings of the outputs of multi-modal interac-1029 tion, while  $\tilde{V}_a$  and  $\tilde{T}_a$  represent the aggregated embeddings 1030 of the outputs of semantic consistency decoder. For binary 1031 classification, the loss  $L_{bcls}$  can be obtained by Eq. 1.

$$L_{bcls} = L_{ce}(C_b([V_{cls}, T_{cls}]), \hat{y}),$$

where  $L_{ce}$  is the cross entropy loss,  $C_b$  is the binary classifier and  $\hat{y}$  is the corresponding ground truth. [.,.] is the concatenation operator. For fine-grained manipulated type prediction, the loss  $L_{fcls}$  can be calculated by Eq. 2.

$$L_{fcls} = L_{ce}(C_i(\widetilde{V}_a), \hat{y}_i) + L_{ce}(C_t(\widetilde{T}_a), \hat{y}_t), \quad (2)$$

where  $C_i$  is the FS/FA classifier,  $C_t$  is the TS/TA classifier.  $\hat{y}_i$  and  $\hat{y}_t$  are the corresponding ground truth. For grounding face manipulation, the loss is composed of L1 and GIoU loss [1] as shown in Eq. 3.

042 
$$L_{img} = L_1(D_i(\widetilde{V}_a), \hat{b}) + L_{giou}(D_i(\widetilde{V}_a), \hat{b}), \quad (3)$$

where,  $D_i$  is the Bbox decoder, and  $\hat{b}$  is the ground truth of Bbox. The overall loss is shown in Eq. 4.

$$L = L_{bcls} + \alpha L_{fcls} + \beta L_{img} + \gamma (Lc + Ls), \quad (4)$$

046where Lc and Ls are the consistency loss of contextual and047semantic consistency matrices, respectively.  $\alpha$ ,  $\beta$  and  $\gamma$ 048are the hyper-parameters to balance different loss functions.049Based on the principle of achieving a uniform order of mag-050nitude for all losses to prevent preference bias, we set them051to 1, 0.1, and 10, respectively.

### References

- Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 658–666, 2019. 2
- [2] Rui Shao, Tianxing Wu, and Ziwei Liu. Detecting and grounding multi-modal media manipulation. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6904–6913, 2023. 1
- [3] Jiazhen Wang, Bin Liu, Changtao Miao, Zhiwei Zhao, Wanyi Zhuang, Qi Chu, and Nenghai Yu. Exploiting modality-specific features for multi-modal manipulation detection and grounding. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4935–4939. IEEE, 2024. 2