

# MONET: MULTI-SCALE OVERLAP NETWORK FOR DUPLICATION DETECTION IN BIOMEDICAL IMAGES

Ekraam Sabir\*, Soumyaroop Nandi\*, Wael AbdAlmageed\*†, Prem Natarajan\*

\*USC Information Sciences Institute, Marina del Rey, CA, USA

†Department of Electrical and Computer Engineering, USC, Los Angeles, USA

## ABSTRACT

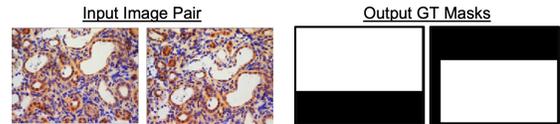
Manipulation of biomedical images to misrepresent experimental results has plagued the biomedical community for a while. Recent interest in the problem led to the curation of a dataset and associated tasks to promote the development of biomedical forensic methods. Of these, the largest manipulation detection task focuses on the detection of duplicated regions between images. Traditional computer-vision based forensic models trained on natural images are not designed to overcome the challenges presented by biomedical images. We propose a multi-scale overlap detection model to detect duplicated image regions. Our model is structured to find duplication hierarchically, so as to reduce the number of patch operations. It achieves state-of-the-art performance overall and on multiple biomedical image categories.

**Index Terms**— Biomedical forensics, image forensics, image manipulation, duplication detection

## 1. INTRODUCTION

Advancements in multimedia technology has enabled the proliferation of digitally manipulated misinformation. Prevalent manifestations of misinformation include fake news, digitally manipulated images, deepfake videos and more. Efforts towards the development of automated detection methods for fake news [1, 2, 3], natural-image forensics [4, 5, 6] and deepfakes [7, 8, 9] have gained traction. However, an important yet almost neglected field is that of biomedical image forensics. Misrepresentation of experimental results by manipulating biomedical images has been an issue of concern for a while in the biomedical community [10]. Unlike natural images, biomedical images often contain arbitrary patterns with no semantic context which allows for these manipulations to go undetected during the peer review process. Investigative or follow-up research can lead to the discovery of such manipulations which consequently leads to retractions [11]. However the entire discovery process is a loss of time and money [12].

Recently, a new biomedical image forensics dataset (BioFors) [13] was released to promote the development of automated detection methods. The dataset comprises images belonging to four categories collected from biomedical research documents. The paper also proposed three biomedical



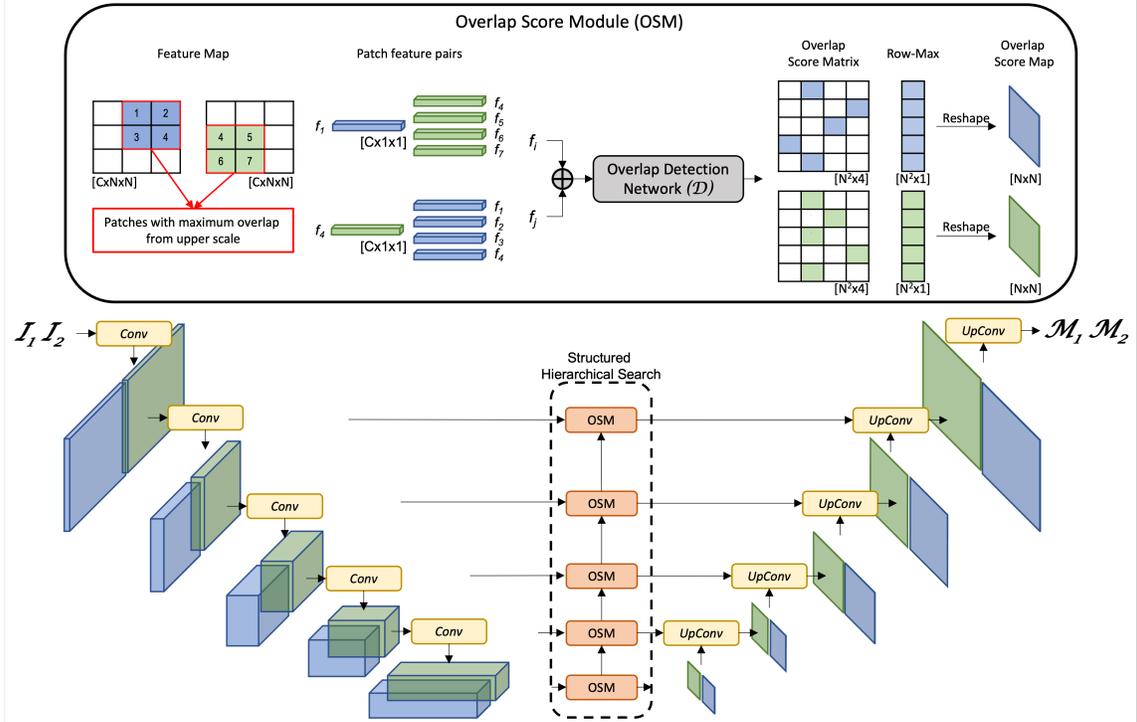
**Fig. 1:** An image pair sample from the EDD task in BioFors.

forensic tasks to overcome the lack of structured problem definitions in literature. There are three manipulation detection tasks described in [13] – external duplication detection (EDD), internal duplication detection (IDD) and cut/sharp transition detection (CSTD). These tasks together cover popular forms of semantic and digital manipulations found in biomedical literature. Of these problems, we focus on the largest task involving the detection of duplicated regions between images a.k.a the external duplication detection (EDD) task. The provenance of manipulations in this task is not known i.e. the images could be spliced, cropped with an overlap from a larger image or simply reused. Figure 1 shows a manipulated sample under the EDD task from [13].

Related research areas of image matching and splicing detection are matured. However, as shown in [13], traditional computer vision methods trained on natural images are not suitable for biomedical forensics. Difficulties in detecting keypoints from biomedical images limit keypoint-descriptor based methods [14, 15, 16]. Use of coarse feature maps limits the detail in deep-learning based splicing detection methods [17, 4] and dense matching of features is computationally expensive [18]. To overcome these challenges, we propose a multi-scale overlap detection network (MONet) that recursively finds overlap between patches to locate duplicated image regions. Recursive overlap detection is performed at multiple scales in an hierarchical manner from large to small image patches. Our model increases the matching detail from coarse to refined feature maps in a top-down approach, while simultaneously reducing the computational burden by making fewer patch-comparisons.

## 2. PROPOSED METHOD

The EDD task is structured to locate duplicated regions between image pairs. Given two input images  $I_1, I_2 \in \mathbb{R}^{H \times W \times 3}$  the objective of the EDD task is to predict two binary masks



**Fig. 2:** Illustration of MONet. The top shows details of overlap score module (OSM) and the bottom shows overall architecture.

$M_1, M_2 \in \mathbb{R}^{H \times W \times 1}$ , highlighting the duplicated image region as shown in Figure 1. Our model comprises convolutional encoding and decoding operations at multiple patch scales as shown in Figure 2. The overlap-score module (OSM) performs overlap detection at each scale. The search for overlapping image regions is structured hierarchically across scales by linked OSMs.

## 2.1. Architecture Overview

The general structure of our model resembles a U-Net [19] with a series of convolutional encoders at multiple scales  $s \in \{1, 2, 3, 4, 5\}$  that produce feature maps  $F \in \mathbb{R}^{N_h \times N_w \times C_s}$ . For notational convenience we consider that images are square with  $N \times N \times 3$  dimension. Consequently, the dimension of encoder feature maps at each scale is  $\frac{N}{2^s} \times \frac{N}{2^s}$ . The upsampling involves a series of convolutional decoders that produce feature maps at corresponding scales to that of the encoder. The final output of the decoder produces a pair of binary prediction masks. To find duplicated regions between images, we measure overlap between patches of two images at each scale within overlap-score modules (OSMs) in a top-down hierarchy. The maximum overlap score of a patch indicates the confidence with which all or a part of it is considered to have been repeated in the other image. A higher score indicates full or substantial repetition, while a low score represents negligible or no repetition. To minimize the number of patch comparisons, patch pairs in  $I_1$  and  $I_2$  with maximum overlap at the current scale  $s$  are used to guide the search among sub-patches at the lower scale  $s - 1$ .

Scale	Patch Dimension	Naive Comparisons	Ours
1	2x2	$\sim 268.43\text{M}$	$\sim 131\text{K}$
2	4x4	$\sim 16.77\text{M}$	32,768
3	8x8	$\sim 1.04\text{M}$	8,192
4	16x16	65,536	2,048
5	32x32	4,096	4,096

**Table 1:** Number of patch comparisons at each scale.

## 2.2. Overlap-Score Module (OSM)

The purpose of the OSM module is to predict two overlap score maps at each scale corresponding to the feature maps. Deconvolution layers upsample the overlap score maps sequentially to produce binary output masks. Score maps from previous and current scale are concatenated for upsampling. Overlap scores are produced by an overlap detection network  $\mathcal{D}$  which takes as input two patch feature vectors (one from each image). It is trained on patch feature triplets (anchor, overlapping and non-overlapping patches) generated from synthetic data at each scale. We consider a feature map  $F$  at scale  $s$ , to be composed of a grid of patch feature vectors  $f \in \mathbb{R}^{1 \times 1 \times C_s}$ , such that each feature vector represents a patch of dimension  $d_s \times d_s$  in the input image, where  $d_s = \frac{N}{2^s}$ . While the convolutional receptive field of a feature vector  $f$  is larger than the patch dimension  $d_s$  at any given scale, we implicitly limit the scope of each feature vector to its patch dimensions when measuring overlap. The overlap score map, is indexed similar to a feature map  $F$ . The score at each index represents the maximum overlap found for that patch feature vector when compared to vectors from the other image.

Method	Microscopy		Blot/Gel		Macroscopy		FACS		Combined	
	Image	Pixel								
SIFT [14]	0.180	0.146	0.113	0.148	0.130	0.194	0.11	0.073	0.142	0.132
ORB [15]	0.319	0.342	0.087	0.127	0.126	0.226	0.269	0.187	0.207	0.252
BRIEF [16]	0.275	0.277	0.058	0.102	0.135	0.169	0.244	0.188	0.180	0.202
DF - ZM [18]	<b>0.422</b>	0.425	0.161	0.192	0.285	0.256	<b>0.540</b>	<b>0.504</b>	0.278	0.324
DMVN [17]	0.242	0.342	0.261	0.430	0.185	0.238	0.164	0.282	0.244	0.310
<b>Ours - regular margin loss</b>	0.398	<b>0.435</b>	0.507	<b>0.520</b>	0.221	0.262	0.313	0.356	<b>0.410</b>	<b>0.438</b>
<b>Ours - flexible margin loss</b>	0.346	0.386	<b>0.520</b>	<b>0.520</b>	<b>0.309</b>	<b>0.281</b>	0.256	0.336	0.398	0.410

**Table 2:** MCC scores on external duplication detection (EDD) task in BioFors across image categories.

### 2.3. Structured Hierarchical Search

The OSMs are structurally linked from higher to lower scale such that patch comparisons can be made hierarchically. Sub-patches of a patch with maximum overlap at a higher scale, are candidate patches for overlap detection at a lower scale. Since, the spatial dimension of each feature map gets halved at each scale, a feature vector  $f$  at a higher scale overlaps with four feature vectors at the immediate lower scale. This observation is useful in limiting the number of patch comparisons to be made at a lower scale. For two patches with maximum overlap at a higher scale, each of their four sub-patches are compared only with each other. At the largest scale (lowest resolution feature map), with no prior scoring, overlap  $o$  is measured between all possible pairs to predict an overlap score map  $O^{N \times N \times 1}$ . Table 1 shows the reduction in patch comparisons at each scale for 256x256 image pairs.

### 2.4. Loss

We pretrain the endoder and overlap detection network jointly using the margin ranking loss function  $\mathcal{L}_o$ . The model is then trained end-to-end with mask output using binary cross-entropy loss. For two feature vectors  $x_1$  and  $x_2$  the regular margin ranking loss function is calculated as (1), where  $m$  is the margin hyper-parameter. In our experiments for an anchor, positive and negative patch triplet  $\langle a, a^+, a^- \rangle$ ,  $x_1$  and  $x_2$  represent the overlap scores between patch pairs  $\langle a, a^+ \rangle$  and  $\langle a, a^- \rangle$  respectively. Therefore the difference between  $x_1$  and  $x_2$  represents the difference in overlap between positive and negative patch pairs. As a result, we also experiment with a flexible margin that is measured as a function of overlap difference. Specifically, if the true overlap in pixels for  $\langle a, a^+ \rangle$  and  $\langle a, a^- \rangle$  is  $o^+$  and  $o^-$ , the flexible margin  $m_{flex}$  is shown in (2), where  $d$  is the patch dimension. Then the flexible margin ranking loss  $\mathcal{L}_{flex}$  is calculated as (3).

$$\mathcal{L}_o = \max(0, (x_2 - x_1) + m) \quad (1)$$

$$m_{flex} = \frac{o^+ - o^-}{d^2} \quad (2)$$

$$\mathcal{L}_{flex} = \max(0, (x_2 - x_1) + m_{flex}) \quad (3)$$

### 2.5. Implementation and Training Details

We resize all input images to  $256 \times 256 \times 3$  dimension. The largest scale has  $8 \times 8 \times 256$  dimension feature map for  $32 \times 32$  dimensional patches. The channel dimension is halved at each scale (256 at scale 5 and 16 at scale 1). The overlap detection layer is a two layer feed-forward network. We pretrain our model for 25 epochs with the margin ranking loss on overlapping and non-overlapping patch triplets generated from synthetic data. The model is trained end-to-end for 50 epochs after that with binary cross-entropy and margin ranking loss. We use the adam optimizer with a learning rate of  $1e-4$ .

## 3. EXPERIMENTS

Method	Image	Pixel
Ours w/o gating	0.340	0.398
Ours w/ dot product overlap	0.076	0.052
Ours - normalized margin	0.398	0.410
Ours - flexible margin	0.410	0.438

**Table 3:** Ablation of gates or using dot-product for overlap.

### 3.1. Dataset and Metrics

We use the BioFors dataset introduced in [13]. The EDD task has 1,547 manipulated images. Train and test splits have 30,536 and 17,269 images respectively, divided into four image categories – Microscopy, Blot/Gel, Macroscopy and FACS. Each category has a different origin or semantic meaning, which leads to diverse image properties and challenges. We evaluate at the image and pixel level, according to the protocol described in [13]. Image level evaluation assigns binary labels to images. Pixel level evaluation is performed on aggregated pixel statistics across images. We use matthews correlation coefficient (MCC) metric as reported in [13].

### 3.2. Synthetic Data Generation

BioFors dataset does not provide any manipulated samples for training. Hence, we train our model using synthetically generated samples similar to the process described in [13]. However, our model additionally requires joint pre-training

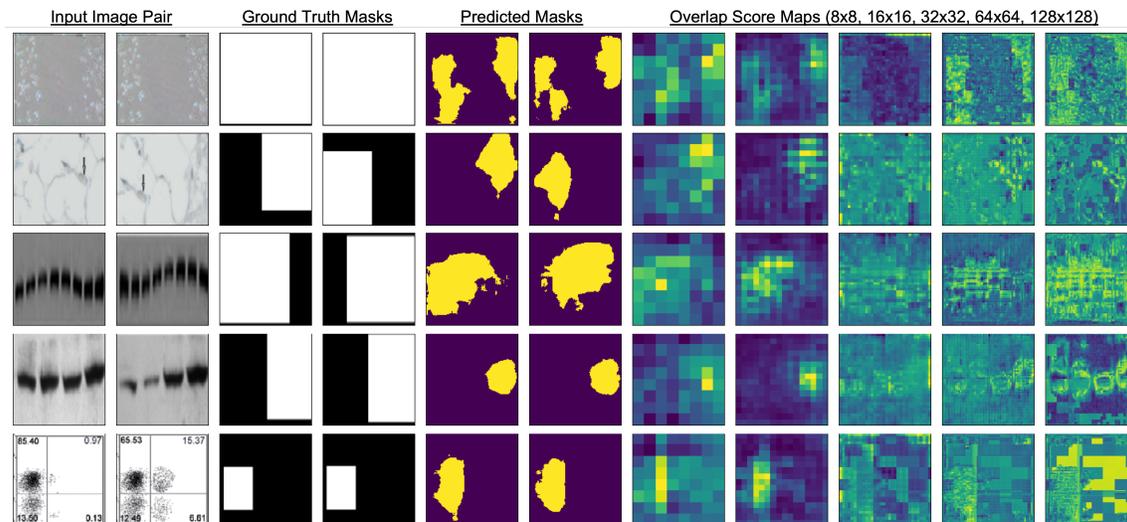


Fig. 3: Input images, ground truth masks, predicted masks and intermediate score maps from MONet.

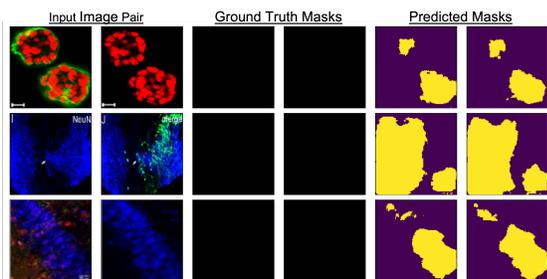


Fig. 4: False positive samples on Microscopy images.

of encoders and overlap scoring modules (OSMs). This requires extensive hierarchical annotation of patch overlap at each pixel i.e. patch pairs and their overlap scores at each scale. Generating such extensive annotation on the fly is computationally expensive. As a workaround, we generate predefined annotation templates, which can be used with random image-pairs on the fly to generate unique synthetic samples.

### 3.3. Results

Table 2 shows the performance of our model. Baseline results are presented as reported in [13]. Image and pixel columns denote corresponding evaluation protocol. We highlight two versions of our model – with regular margin ranking loss and with a flexible margin ranking loss. Our model achieves a new state-of-the art on blot/gel, microscopy, macroscopy image categories and also on the combined evaluation.

### 3.4. Analysis

As shown in Table 2, our model achieves state-of-the art result across multiple categories. However, the performance fluctuates across image categories. Additionally, a single model does not hold top-performance across categories. We believe that the unique characteristics of each category make it difficult to train a single outperforming model. Figure 3 shows

sample predictions from our model. Overlap score maps from each scale show the progression of patch overlap detection. Figure 4 also shows that our method generates false positives. As described in [13], these duplicated regions are not considered manipulations due to the semantics of experiments that produced them, such as image overlay or chemical staining. Overcoming these false positives requires either additional semantic information from source documents or the definition of manipulation needs rethinking.

**Ablation:** We perform an ablation analysis of our model in Table 3. The model performance degrades if we remove the gating operation when concatenating overlap score maps across scales. Additionally, performance drops drastically if we use feature dot products from literature [17, 4] instead of one hidden layer overlap detection network.

## 4. CONCLUSION

Duplication of images across biomedical experiments is a concerning issue. Our proposed model achieves state-of-art-performance on the EDD task for some image categories. Further investigation with dedicated model for each image category is required.

## 5. ACKNOWLEDGEMENT

This material is based on research sponsored by DARPA and Air Force Research Laboratory (AFRL) under agreement number FA8750-20-2-1004. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of DARPA and Air Force Research Laboratory (AFRL) or the U.S. Government.

## 6. REFERENCES

- [1] Ekraam Sabir, Wael AbdAlmageed, Yue Wu, and Prem Natarajan, “Deep Multimodal Image-Repurposing Detection,” in *Proceedings of the 26th ACM International Conference on Multimedia*, New York, NY, USA, 2018, MM ’18, pp. 1337–1345, ACM.
- [2] Eric Müller-Budack, Jonas Theiner, Sebastian Diering, Maximilian Idahl, and Ralph Ewerth, “Multi-modal analytics for real-world news using measures of cross-modal entity consistency,” in *Proceedings of the 2020 International Conference on Multimedia Retrieval*, 2020, pp. 16–25.
- [3] Ekraam Sabir, Ayush Jaiswal, Wael AbdAlmageed, and Prem Natarajan, “Meg: Multi-evidence gnn for multi-modal semantic forensics,” in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 9804–9811.
- [4] Yue Wu, Wael Abd-Almageed, and Prem Natarajan, “Busternet: Detecting copy-move image forgery with source/target localization,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [5] 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 (CVPR)*, June 2019.
- [6] 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*, 2022, pp. 2364–2373.
- [7] Ekraam Sabir, Jiaxin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan, “Recurrent convolutional strategies for face manipulation detection in videos,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2019, pp. 80–87.
- [8] Iacopo Masi, Aditya Killekar, Royston Marian Mascarenhas, Shenoy Pratik Gurudatt, and Wael AbdAlmageed, “Two-branch recurrent network for isolating deepfakes in videos,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [9] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Ting Zhang, Weiming Zhang, Nenghai Yu, Dong Chen, Fang Wen, and Baining Guo, “Protecting celebrities from deepfake with identity consistency transformer,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 9468–9478.
- [10] Elisabeth M Bik, Arturo Casadevall, and Ferric C Fang, “The prevalence of inappropriate image duplication in biomedical research publications,” *MBio*, vol. 7, no. 3, 2016.
- [11] Elisabeth M Bik, Ferric C Fang, Amy L Kullas, Roger J Davis, and Arturo Casadevall, “Analysis and correction of inappropriate image duplication: the molecular and cellular biology experience,” *Molecular and Cellular Biology*, vol. 38, no. 20, 2018.
- [12] Andrew M Stern, Arturo Casadevall, R Grant Steen, and Ferric C Fang, “Financial costs and personal consequences of research misconduct resulting in retracted publications,” *Elife*, vol. 3, pp. e02956, 2014.
- [13] Ekraam Sabir, Soumyaroop Nandi, Wael AbdAlmageed, and Prem Natarajan, “Biofors: A large biomedical image forensics dataset,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10963–10973.
- [14] David G Lowe, “Distinctive image features from scale-invariant keypoints,” *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [15] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski, “Orb: An efficient alternative to sift or surf,” in *Proceedings of the IEEE International Conference on Computer Vision*. Ieee, 2011, pp. 2564–2571.
- [16] Michael Calonder, Vincent Lepetit, Christoph Strecha, and Pascal Fua, “Brief: Binary robust independent elementary features,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2010, pp. 778–792.
- [17] Yue Wu, Wael Abd-Almageed, and Prem Natarajan, “Deep matching and validation network: An end-to-end solution to constrained image splicing localization and detection,” in *Proceedings of the 25th ACM international conference on Multimedia*, 2017, pp. 1480–1502.
- [18] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva, “Efficient dense-field copy-move forgery detection,” *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 11, pp. 2284–2297, 2015.
- [19] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.